

**Improving the security and cyber
security of companies and
individuals using behavioural
sciences: a data-centric approach**

by

**Leonardo Mariano
Castro Gonzalez**

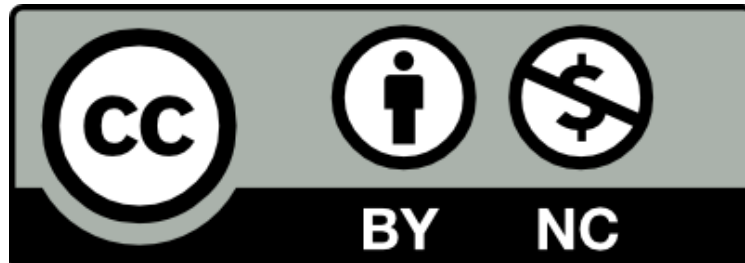
A thesis submitted to the University of Birmingham for the
degree of
Doctor of Philosophy



**UNIVERSITY OF
BIRMINGHAM**

Department of Economics
Birmingham Business School
College of Social Sciences
University of Birmingham
September 2022

University of Birmingham Research Archive e-theses repository



This unpublished thesis/dissertation is under a Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0) licence.

You are free to:

Share — copy and redistribute the material in any medium or format

Adapt — remix, transform, and build upon the material

The licensor cannot revoke these freedoms as long as you follow the license terms.

Under the following terms:



Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.



NonCommercial — You may not use the material for commercial purposes.

No additional restrictions — You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Notices:

You do not have to comply with the license for elements of the material in the public domain or where your use is permitted by an applicable exception or limitation.

No warranties are given. The license may not give you all of the permissions necessary for your intended use. For example, other rights such as publicity, privacy, or moral rights may limit how you use the material.

Unless otherwise stated, any material in this thesis/dissertation that is cited to a third-party source is not included in the terms of this licence. Please refer to the original source(s) for licencing conditions of any quotes, images or other material cited to a third party.

Abstract

While security and cyber security systems literature focus on how to detect threats at a logistics, software and hardware level, there is not enough work around how to improve the security by incorporating the understanding of the human behaviour for those individuals that form part of the system. The present dissertation focus in the latter problem and has it as main research question. To do so, we study three different security and cyber security problems. We study a problem of communication framing when training employees in cyber security by deploying a two-staged survey in a British financial institution to then analyse it with a behavioural segmentation model. We find that, depending on their risk-perception and risk-taking attitudes, employees can become better cyber security sensors when correctly framed. We also study a problem of illicit drugs distribution in England to understand the territorial logic of the operators. Using public data, we analyse the problem using Spatial Analysis models. We find that gangs avoid places with a high number of knife crime events and hospital admissions by misuse of drugs. Finally, we study the transition of companies to the “New Normal” when the pandemic started. Using a qualitative model to understand the cyber security culture within, we find that cyber security was not a priority of the narrative of big companies during the first months of 2020. The three essays contribute to the literature in behavioural sciences applied to security and cyber security by using modern tools and frameworks of statistical learning and Natural Language Processing. By incorporating these different resources, we show how to improve the efficiency of security and cyber security systems by analysing the behaviour data extracted from them.

Keywords: Behavioural sciences, Data science, Cyber Security, Security Systems,

Acknowledgements

Foremost, I would like to thank my supervisor Prof. Ganna Pogrebna for the opportunities and the support she granted me during these past four years. Without her, much of the professional opportunities I had during the PhD would have not happened at all. Thanks to my supervisors Prof. Oleksandr Talavera and Prof. Kimberley Scharf for their academic support. Thanks to Prof. Eiman Kanjo and Dr. Roshan Boojihawon for their time and agreeing on reading my thesis.

Thanks to CONACYT-SENER and to the University of Birmingham to financially support me during this journey. Quite literally, none of this could have been possible without their support.

This work would have not been the same if it were not for the support of different brilliant academics throughout the journey. I would like to thank Prof. Elsa Arcaute for giving me a warm academic support when I needed it the most. Thanks to Prof. Carsten Maple for allowing me to work with his great team at WMG. To Prof. Keri Pearlson and Dr. Daniel Gozman for the stimulating talks. Finally, thanks to Prof. Weisi Guo for the incredible help during the first half of my PhD.

Thanks to Prof. Dame Helen Margetts, Dr. Jonathan Bright and Dr. Omar Guerrero for their trust in me and my work during the last months of my PhD, and giving me the opportunity to work with them in my next professional steps at the Alan Turing Institute.

Not only the present dissertation, but my entire experience throughout these past 4 years would have not been the same without the presence of some people with who I had the pleasure to work with and with who I have the privilege to call friends. Thanks to my incredible Tortuga/Mexicovid team, Fabián Aguirre López, Blas Kólic, Rodrigo Leal Cervantes, Santiago Martínez Balvanera, Carlota Segura García and Diego Vidal Cruz-Prieto for the inspiration, the fun and the never-ending stimulation I receive with them. Thanks to Cécile de Bézenac and Giulia Occhini for the incredible discussions of what would later become my next professional steps. Thanks to Torty Sivill for allowing me to co-host with her the Turing Student Seminar, and then for the privilege of being her first guest at her amazing podcast.

Thanks also to all my friends that were always close to me: my London friends Sebastián “Bato” Balvanera Nadurille, Paula Blancarte Jaber, Sofía Bowen Silva, Magdalena Bravo Rojas, Victoria González Budilova, Zahra Jafari, Elena Luciano Suástegui, Maeva Mésanger Godemer, Jazmín Ruíz Díaz, Edgar René Ruíz López and Sergio Pérez Pérez; my virtual friends Andrea Calderón Stephens and Andrés

de los Ríos Sommer; and the friends I collected during this journey, Frida Xaman Ek Estrella García, Emanuela Carollo, Sofia Lykou, Ismael Rodríguez Fernández, Mariana Gamba Fadul, Alberto Nieto Tibaquirá and Emilia-bebé.

Thank you all for the love and support; thanks for all the fun and amazing time spent together; thanks for inspiring me, motivating me –and pressuring me– to always move towards a better place.

To my parents, Luz González Godínez and Roberto Castro Rodríguez, for their unconditional and endless love and support. None of this could have been without you. Thanks to my family also for always showing their support and love without borders. An eternal thanks and remembrance to those that could not see this achievement made true: Sergio González Corona, Zoila Rodríguez and Héctor Castro Rodríguez. I deeply miss you and would love for you to see this moment.

Last but definitely not least, thanks to Emilia Casillas Rodríguez for her invaluable and incommensurate support, her always accurate words, her warming hugs and her infinite love. Thanks for keep making, just as 6 years ago, an entire ocean to feel like a small pond.

Contents

Abstract	2
List of Figures	8
List of Tables	12
List of Equations	15
List of Acronyms	16
Introduction	18
1 The effects of a framed communication in cyber security risk perception and action: a data-centred case study using behavioural segmentation	28
1.1 Introduction	29
1.2 Literature review	31
1.3 Methodology	36
1.3.1 Data gathering	39
1.3.2 Methods	41
1.4 Results: case study	43
1.4.1 Phase 1	43
1.4.2 Phase 2	46
1.4.3 Work from home habits and cyber security higyene during the COVID-19 pandemic	52
1.5 Discussion	56
1.5.1 Framed communication strategies using a segmentation strategy	56
1.6 Conclusions	58
2 Understanding and nowcasting the illicit drug distribution in England: a data-centric approach to the County Lines Model	62
2.1 Introduction	63
2.2 Literature review	66

2.2.1	County Lines literature	66
2.2.2	Social Network Analysis	67
2.2.3	Criminology and other domains	69
2.2.4	Literature gap and hypotheses	70
2.3	Methodology	72
2.3.1	Retail model	74
2.3.2	Gravity model	75
2.3.3	Radiation model	75
2.3.4	Model selection process	76
2.3.5	Pipeline	77
2.3.6	Data	78
2.4	Results	79
2.4.1	Model selection	79
2.4.2	Model analysis and geographic distribution	81
2.5	Conclusions	84
3	The Impact of Corporate Values and Factors of Internal and External Culture on Formulating the Post-COVID “New Normal”: Implication for Cybersecurity and Information Systems	89
3.1	Introduction	90
3.2	Literature Review	92
3.3	Methodology	97
3.3.1	The CAMS-inspired New Normal Model extension	97
3.3.2	Factors of external culture	99
3.3.3	Factors of internal culture	101
3.3.4	Beliefs, attitudes, and values	101
3.3.5	COVID-19 business resilience	102
3.4	Processed data	102
3.5	Results	106
3.5.1	Business resilience with respect to beliefs, attitudes and values	106
3.5.2	Links between cultural factors and beliefs, attitudes and values	110
3.6	Discussion	114
3.7	Conclusions	116
4	Conclusions	120
	References	128
	Appendices	148

A	Chapter 1: Phase 1 survey and results	148
B	Chapter 1: Phase 2 survey and results	153
C	Chapter 1: Human-as-a-Cyber-Security-Sensor test	156
D	Chapter 2: Table of models	169
E	Chapter 2: Database	172
E.1	The database	172
E.1.1	Drug related hospital admissions data	172
E.1.2	Drug related deaths data	173
E.1.3	Number of hospital beds data	174
E.1.4	Police workforce data	174
E.1.5	Police numbers of drug seizures data	175
E.1.6	Knife crime related data	176
E.1.7	Disposable Income data	176
E.1.8	Demographic data	176
E.1.9	Geographic data	176
E.1.10	Territorial resolutions	177
E.1.11	England and its 9 regions	177
E.1.12	Local Police Forces	177
E.1.13	Merged local authorities	178
F	Chapter 3	179
F.1	Formation of sectors for both datasets	179
F.2	Topic Modelling	183
F.2.1	Perplexity plots to find optimal number of topics	183
F.2.2	Insight about resulting topics – Most important words in each of them	183
F.2.3	Frequency tables	185
F.3	LIWC	185
F.3.1	Dictionary of words for each dimension	185
F.3.2	Tables about averages of each dimension in both datasets . . .	187
F.4	Results	189

List of Figures

1.1	Design of the study done.	37
1.2	Schematic of the 4 segments (Relaxed, Opportunistic, Anxious and Ignorant), plus the fifth Well-calibrated area at the centre of the figure.	42
1.3	Average scores and standard deviation for each of the 13 activities of the CyberDoSpeRT.	43
1.4	Scatter plot of the risk taking and risk perception scores for the 605 respondents of the first survey. Colours represent the different segments (Anxious, Opportunistic, Relaxed, Ignorant) found performing a k -means++ algorithm.	44
1.5	Proportion of segments for Phase 1, Phase 2, and the different priming samples obtained in the latter phase.	49
1.6	(a) Proportion of correct answers by the participants when attempting the HaaCSS test in the second phase. (b) Distribution of incorrect answers divided by primings.	50
1.7	Results for each of the six questions of the HaaCSS test, showing the proportion of people answering if the presented screenshot was a potential start of a cyber attack, and the confidence of their answer. Results are shown for the total sample of Phase 2, as for each of the different treatments given to the respondents.	53
1.8	Heatmap matrix showing the different distributions of the satisfaction of life and work for the respondents of both surveys	54
1.9	Heatmap matrix showing the different distributions of the contracted and actual worked hours for the respondents of both surveys	54
1.10	Heatmap matrix showing the different distributions of the satisfaction of life and work for the respondents of both surveys	55
1.11	Levels of indicated confidence to detect cyber attacks by the respondents of both phases.	55
1.12	Answers to the frequency and type of cyber attacks detected by the respondents of Phase 1.	56
2.1	Schematic of models. Coloured boxes are those considered by the model.	73

2.2	Schematic of analysis pipeline	78
2.3	Results of BIC and the Sørensen-Dice index for the 68 different models tested. Zone 0 corresponds to the Gravity model. Zone 1 corresponds to the Radiation model. Finally, zones 2 and 3 correspond to the retail model with the Poissonian loss function and with the MSE loss function respectively. The annotated year corresponds to the data in which the model was trained on. Details of each model can be found in Appendix A.	80
2.4	MSE costs when comparing the trained models with the Metropolitan Police data.	81
2.5	Data points and modelled lines ordered by police force for 2019 (a) and 2020 (b). Note: for both years, there were 0 lines detected in Durham. As the plot is in log scale, this data point was not included.	82
2.6	Heatmaps for the Metropolitan Police data for 2019 (2.6a) and 2020 (2.6e), and the three different models tested: Retail (2.6b and 2.6f), Radiation (2.6c and 2.6g) and Gravity (2.6d and 2.6h).	83
2.7	Heatmaps showing the difference between the modelled distribution of lines with respect to the Metropolitan Police data for 2019 and 2020.	85
3.1	MIT CAMS Culture of Cybersecurity Model adapted [Huang and Pearlson, 2019a]	96
3.2	The CNNM structure, with details about data and processes	98
3.3	Profiles of different sectors and regions of the world by LIWC psychological dimensions. Left column depicts the US sample. The right column depicts the Global sample.	105
3.4	Sankey representation of the Tables obtained in Section 3.5. We highlight those statistically significant correlations ($p < 0.05$) with the LIWC dimension of Focus on the future as part of the organisational culture. The idea of these Sankey diagrams is to graphically represent the CNNM with the external factors and managerial mechanisms on the left, the values, beliefs and attitudes at the center of the figure and the behaviour on the right. Blue links refer to positive correlations while red links refer to negative correlations.	119
A.1	Distribution of the employees that took part of the Phase 1 of the study by (a) age and (b) gender.	149
A.2	Distribution of the employees that took part of the Phase 1 of the study by (a) life satisfaction and by (b) job satisfaction.	149

LIST OF FIGURES

A.3	Distribution of the employees that took part of the Phase 1 of the study by (a) confidence to detect cyber threats and (b) if they are in a cyber security role or not.	150
A.4	Distribution of the employees that took part of the Phase 1 of the study by (a) contracted hours per week (b) actual working hours per week and (c) days working from home.	150
A.5	Distribution of the employees that took part of the Phase 1 of the study by (a) years of experience in the company and (b) if they have a costumer-facing role within the company.	151
A.6	Distribution of employees by (a) relationship status and by (b) caring responsibilities. We also present the ages of those cared ones in (c). Some employees had more than one cared ones.	152
B.1	Distribution of the employees that took part of the Phase 2 of the study by (a) age and (b) gender.	153
B.2	Distribution of the employees that took part of the Phase 2 of the study by (a) life satisfaction and by (b) job satisfaction.	154
B.3	Distribution of the employees that took part of the Phase 2 of the study by (a) confidence to detect cyber threats and (b) if they have a costumer-facing role within the company.	155
B.4	Distribution of the employees that took part of the Phase 2 of the study by (a) contracted hours per week (b) actual working hours per week and (c) days working from home.	155
C.1	Screenshot presented as first question at the Human-as-a-Cyber-Security-Sensor test.	157
C.2	Results for the first question of the HaaCSS test.	158
C.3	Screenshot presented as second question at the Human-as-a-Cyber-Security-Sensor test.	159
C.4	Results for the second question of the HaaCSS test.	160
C.5	Screenshot presented as third question at the Human-as-a-Cyber-Security-Sensor test.	161
C.6	Results for the third question of the HaaCSS test.	162
C.7	Screenshot presented as fourth question at the Human-as-a-Cyber-Security-Sensor test.	163
C.8	Results for the fourth question of the HaaCSS test.	164
C.9	Screenshot presented as fifth question at the Human-as-a-Cyber-Security-Sensor test.	165
C.10	Results for the fifth question of the HaaCSS test.	166

C.11 Screenshot presented as sixth question at the Human-as-a-Cyber-Security-Sensor test. 167

C.12 Results for the sixth question of the HaaCSS test. 168

F.1 Perplexity plots for (a) US-based companies and for (b) Global companies. 184

F.2 Sankey diagram representing statistically significant correlations between internal factors (topics – on the left), the attitudes, values and beliefs (psychological dimensions – centre) by industrial sectors for the US sample. Blue links represent positive correlations, while red links represent negative correlations. We highlight those links presented in the main manuscript. 189

F.3 Sankey diagram representing statistically significant correlations between internal and external factors (topics, cultural value orientations, GCI – on the left), the attitudes, values and beliefs (psychological dimensions – centre) by industrial sectors for the Global sample. Blue links represent positive correlations, while red links represent negative correlations. We highlight those links presented in the main manuscript. 190

F.4 Sankey diagram representing statistically significant correlations between internal and external factors (topics, cultural value orientations, GCI – on the left), the attitudes, values and beliefs (psychological dimensions – centre) by region for the Global sample. Blue links represent positive correlations, while red links represent negative correlations. We highlight those links presented in the main manuscript. 201

List of Tables

1	Summary of studied systems and objectives for each Chapter.	26
1.1	Detail of the two surveys.	38
1.2	Determinants of the Anxious behavioural type. Results obtained performing a Probit regression.	45
1.3	Determinants of the Opportunistic behavioural type. Results obtained performing a Probit regression.	46
1.4	Determinants of the Relaxed behavioural type. Results obtained performing a Probit regression.	47
1.5	Determinants of the Ignorant behavioural type. Results obtained performing a Probit regression.	48
1.6	Results of the Poisson regressions the three different treatments applied in the second phase of the survey with respect to the behavioural segments.	51
2.1	Results for the best three models calibrated.	81
3.1	Subsets used for the CNN model from our data sets.	103
3.2	Different variables used in the CAMS-inspired New Normal Model.	104
3.3	Results for a Clustered OLS regression for the US sample, by industry sector. In this case the dependent variable is the business resilience to COVID-19 (profit change with respect to 2019).	107
3.4	Results for a Clustered OLS regression for the Global Sample, by region. In this case the dependent variable is the business resilience to COVID-19 (profit change with respect to 2019).	108
3.5	Results for a Clustered OLS regression for the Global sample, by industry sector. In this case the dependent variable is the business resilience to COVID-19 (profit change with respect to 2019).	109
3.6	Results for a clustered OLS regression for the US Sample, by industry sector. In this case the dependent variable are the psychological dimensions and the tested variables are the topics drawn from the Topic Modelling exercise.	111

LIST OF TABLES

3.7 Results for a clustered OLS regression for the Global sample (Asia subset). In this case the dependent variable are the psychological dimensions and the tested variables are the topics drawn from the Topic Modelling exercise, the Cultural Value Orientation and the Global Cybersecurity Index (GCI). 112

3.8 Results for a clustered OLS regression for the Global sample (sectors with future-oriented link). In this case the dependent variable are the psychological dimensions and the tested variables are the topics drawn from the Topic Modelling exercise, the Cultural Value Orientation and the Global Cybersecurity Index (GCI). 113

C.1 Topic modelling results. Reasoning for text example 1 157

C.2 Topic modelling results. Reasoning for text example 2 160

C.3 Topic modelling results. Reasoning for text example 3 162

C.4 Topic modelling results. Reasoning for text example 4 163

C.5 Topic modelling results. Reasoning for text example 5 165

C.6 Topic modelling results. Reasoning for text example 6 166

D.1 List of all trained models. 169

F.1 Industry sectors for US-based companies used with respect to the industry labels at the Forbes Fortune database. 179

F.1 Industry sectors for US-based companies used with respect to the industry labels at the Forbes Fortune database. 180

F.1 Industry sectors for US-based companies used with respect to the industry labels at the Forbes Fortune database. 181

F.2 Industry sectors for Global companies used with respect to the industry labels at the Forbes Global database. 181

F.2 Industry sectors for Global companies used with respect to the industry labels at the Forbes Global database. 182

F.3 Topics and 10 most important keywords for the US-based companies database. 184

F.4 Topics and 10 most important keywords for the Global companies database. 185

F.5 Sample of words from the LIWC dictionaries for each of the dimensions used. 186

F.5 Sample of words from the LIWC dictionaries for each of the dimensions used. 187

F.6 Proportion of dimensions in collection of texts for the US-based companies database per industrial sector. 187

LIST OF TABLES

F.6	Proportion of dimensions in collection of texts for the US-based companies database per industrial sector.	188
F.7	Proportion of dimensions in collection of texts for the global companies database per industrial sector.	188
F.8	Proportion of dimensions in collection of texts for the global companies database per region.	188
F.8	Proportion of dimensions in collection of texts for the global companies database per region.	189
F.9	Distribution of topics found in the companies of the US database. . .	190
F.9	Distribution of topics found in the companies of the US database. . .	191
F.9	Distribution of topics found in the companies of the US database. . .	192
F.9	Distribution of topics found in the companies of the US database. . .	193
F.9	Distribution of topics found in the companies of the US database. . .	194
F.9	Distribution of topics found in the companies of the US database. . .	195
F.9	Distribution of topics found in the companies of the US database. . .	196
F.9	Distribution of topics found in the companies of the US database. . .	197
F.9	Distribution of topics found in the companies of the US database. . .	198
F.9	Distribution of topics found in the companies of the US database. . .	199
F.9	Distribution of topics found in the companies of the US database. . .	200
F.10	Distribution of topics found in the companies of the Global database.	201
F.10	Distribution of topics found in the companies of the Global database.	202
F.10	Distribution of topics found in the companies of the Global database.	203
F.10	Distribution of topics found in the companies of the Global database.	204
F.10	Distribution of topics found in the companies of the Global database.	205

List of Equations

1.1	Probit model	42
1.2	Poisson model	43
2.1	Betweenness coefficient in a network	68
2.2	Retail model	74
2.3	Modified Retail model	74
2.4	Gravity model	75
2.7	Gravity model with outflow restriction	75
2.8	Radiation model	75
2.9	Modified Radiation model	76
2.10	Opportunity function	76
2.11	Sørensen-Dice index	76
2.12	Bayesian Information Criterion	76
2.13	Gaussian loss function	77
2.14	Poissonian loss functions	77
F.1	Perplexity equation	183

List of Acronyms

ABM Agent-Based Model.

BIC Bayesian Information Criterion.

BYOD Bring-your-own-Device.

CAMS Cybersecurity at MIT-Sloan.

CISO Chief Information Security Officer.

CLM County Lines Model.

CNNM CAMS-inspired New Normal Model.

COVID-19 Corona Virus Disease - 2019.

CyberDoSpERT Cyber-Domain-Specific-Risk-Taking.

GCI Global Cybersecurity Index.

GDHI Gross Disposable Household Income.

HaaCSS Human as a Cyber Security Sensor.

HaaSS Human as a Security Sensor.

ITU International Telecommunication Union.

LDA Latent Dirichlet Allocation.

LIWC Linguistic Inquiry and Word Count.

MSE Mean Squared Error.

NCA National Crime Agency.

NHS National Health System.

OLS Ordinary Least Squares.

ONS Office for National Statistics.

PDE Partial Differential Equation.

PPE Personal Protective Equipment.

SCVOC Schwartz Cultural Value Orientation Coefficients.

SNA Social Network Analysis.

SVS Schwartz Value Survey.

UBA User Behaviour Analysis.

UEBA User-Entity Behaviour Analysis.

UK United Kingdom.

US United States.

Introduction

We live in a particular moment of history where every year there is an increasing amount of data created, analysed and stored around the world [Buyya et al., 2016]. Society has moved beyond a *digitalisation*, transforming analogous information into digital information and entered into a state of *datification*, “*taking all aspects of life and turning them into data*” [Biltgen and Ryan, 2016]. From heartbeats [Xu and Liu, 2020] to global supply chains [Choi, 2018] and cities’ pollution [Zappi et al., 2012] and well-being [Johnson et al., 2020], almost every physical aspect and abstraction, such as emotions and sentiments [Pennebaker et al., 2007, Wang et al., 2004], can be turned into data [Biltgen and Ryan, 2016]. With this technological revolution [Skilton and Hovsepian, 2008], it is expectable that an important number of academic fields, companies and government seek to incorporate modern methods involving analysing a big number of data to solve problems in an innovative and more efficient way.

In particular, strategies to increase *security* have incorporated an increasing number of data-driven methods. Given the increasing use of technology in the field [Anderez et al., 2021], multiple types of data, such as WiFi logs, audio data and GPS positions have been introduced to problems of security, crime-prevention [Short et al., 2010, Caminha et al., 2017] and terrorism [Krebs, 2002, Lum et al., 2006], to name a few. In the same way, different methods coming from multiple disciplines (e.g. Physics or Statistics) have been included to the disciplines [D’Orsogna and Perc, 2015], like Agent-Based modelling [Hegemann et al., 2011] and Social Networks Analysis [Campana, 2016, Podolny and Page, 1998, Rostami and Mondani, 2015]. Most of the streams around criminology and terrorism are focused on prevention by mapping, *nowcasting* and forecasting future crime patterns.

Another particular aspect of security that has involved much of novel data-driven methods has been cyber security. The field moves fast as it has become one of the most important fields for governments’ and companies’ overall security systems [Goutam, 2015, Srinivas et al., 2019, Pogrebna and Skilton, 2019]. As new threats emerge, cyber security has included modern techniques in Machine Learning, Artificial Intelligence or Blockchain [Handa et al., 2019, Taylor et al., 2020]. Cyber security research puts special emphasis on hardware and software [Heartfield et al., 2016], with the idea on mind that humans are “the weakest link” in the overall system [Heartfield and Loukas, 2018].

While in the security and cyber security systems treated in this dissertation, humans take a fundamental role on how the system is correctly deployed and maintained, most of *the data-driven solutions* offered in the literature often forget to un-

derstand the behaviour of the different individuals as a group. Rather, particularly in cyber security, authors tend to take a “punitive” attitude towards humans when involving them as part of the cyber security prevention system, particularly when training employees in cyber security [Rege et al., 2020, Sabillon et al., 2019, Yeoh et al., 2021]. However, as companies and governments are exploring more data-centric solutions [MacArthur et al., 2022], there is also a need on not only involving humans as part of the final solution when talking of security systems, but also a need to understand them in order to make the whole system more efficient towards potential threats or attacks.

In this dissertation we present three different problems of security and cyber security, and their respective solutions that have as basis to understand the behaviour of humans as individuals or as a collective. By analysing the behaviour of the respective set of individuals, we then build solutions that are driven by the respective results of the analysis and framed with different methodologies from both qualitative or quantitative nature. In that sense, at the centre of the research is to understand *how* and *why* individuals or collectives are behaving as data shows to then find a more efficient solution within a (cyber)security system.

With respect to security, we study -at the best knowledge of the author- for the first time, the County Lines Model (CLM) [Crime Agency, 2019] from a quantitative perspective. The CLM is an illicit drugs distribution model in the UK that has brought important consequences on modern slavery and public health issues in different populations of Great Britain. In this case, we are interested in understanding the territorial logic of the gangs operating in this model.

With respect to cyber security, we study two different problems. First, we are interested in going beyond the homogeneous “punitive” approach [Sabillon et al., 2019] to train employees in cyber security. To do that, we test the capacity of employees in one of the largest financial companies in the UK by applying not only a “punitive” approach, but also a non-punitive one to then see how each approach works one on one. We are assuming that not all employees are the same, thus allowing some heterogeneity in our analysis.

Finally, we analyse how the transition for companies to the “New Normal” [Habersaat et al., 2020] was done when the COVID-19 pandemic started with respect to cyber security. We do this by analysing the official documents that companies released when the COVID-19 pandemic started using Topic Modelling [Blei et al., 2003], a Natural Language Processing (NLP) method that has been widely used in Social Sciences [Nikolenko et al., 2017]. In this case we are interested in expanding the Cybersecurity at MIT-Sloan (CAMS) model [Huang and Pearlson, 2019b] to understand how companies managed the cyber security transition from a managerial

point of view.

Given that each one of the three different problems studied in this dissertation are of different nature, every one of them has their own problematic, methodology and data used. Thus, in what is next we present each one of the problems. We will then set the general research questions and the structure of this dissertation.

Chapter 1

Literature often focuses on different methods and approaches to train employees/workers in a given institution when involving the human behaviour into a cyber security system, [Alshaikh and Adamson, 2021, Yeoh et al., 2021, Rege et al., 2020, Kwak et al., 2020, Sabillon et al., 2019]. However, the approach taken normally by this literature consists in two different elements: humans are the most vulnerable element of the security system and training has to have a punitive and threatening approach to the workers, highlighting the consequences of a successful cyber attack.

Although the first point has been tested before [Heartfield and Loukas, 2018], proving that humans can, with the appropriate training, become very efficient *sensors* that detect potential threats at a good accuracy rate, the second point has not been proved wrong. In Chapter 1, we test if workers do homogeneously respond to a “negative” or “punitive” training -highlighting consequences- or if they actually could also be trained with a “positive” training, highlighting the benefits of a good cyber security culture.

To do so, we partnered with the team on cyber security at the Warwick Manufacturing Group of the University of Warwick¹ to work with one of the most important financial institutions based in the UK, with more than 10 000 workers. The financial institution has had an active and constant cyber security awareness approach towards its employees, having different campaigns and trainings for workers. However, its innovation team of the cyber security branch is interested in knowing how to communicate in a more efficient way with the heterogeneous workforce of the bank.

We staged a 2-phase survey that was launched throughout the bank’s workforce. This allowed to obtain data from the employee’s demographic information and perspective and attitudes around different topics of cyber security. The survey included the Cyber-Domain-Specific-Risk-Taking scale [Kharlamov et al., 2018] and a Human-as-a-Security-Sensor test [Heartfield and Loukas, 2015], which allowed us

¹The project between the WMG and the financial institution had different streams of work, which included qualitative and quantitative work. The stream of work presented here was entirely analysed and written by the author of this dissertation alone, and managed by Prof. Ganna Pogrebna and Prof. Carsten Maple.

to understand how different types of employees reacted to different types of communication when asked to assess if a given situation could entail a potential cyber security threat. Our methodology allows us to understand the reaction at a disaggregated level, not supposing that all employees are homogeneous, thanks to a behavioural segmentation model from [Kharlamov et al., 2018] based on risk-taking and risk-perception scores.

In that sense, Chapter 1 intends to provide a methodology for institutions to (i) make better communication strategies to train their employees, to then have as a direct consequence an (ii) increased cyber security culture which would increase the accuracy of the employees to correctly detect potential cyber threats. The methodology has as a basis the behaviour of the employees with respect to different types of communication and taking into account their heterogeneity.

More generally, in this first chapter we want to study how to increase the cyber security of a company by increasing the accuracy of the employees to detect cyber threats. This is done by understanding the behaviour of the workforce towards their received training. Then, by implementing changes in the training, so it does a better work by making the employees better cyber security sensors. In that sense, we study and change the behaviour of the employees that make part of the whole cyber security system of the financial institution.

Chapter 2

The illicit drugs distribution model, County Lines Model (CLM), studied here contrasts with the traditional distribution/retail model seen in the UK, where a clear division between importers, medium-size sellers and small-size sellers is present [Coombes, 2018]. In the case of the CLM, medium- and small-size sellers are merged, thus making the operation more efficient [Black, 2020a]. As a result, the final consumer sees a drop in prices, a faster delivery and an increase in quality. [Rescue and Analysts, 2019, Rescue and Analysts, 2010, Coombes, 2018] On the other side, the CLM has provoked a major drawback in terms of intelligence for the police forces to tackle the problem [Crime Agency, 2019, Silver and Intelligence, 2021, Black, 2020b, Rescue and Analysts, 2019, Rescue and Analysts, 2010]. From a logistics point of view, the fact of merging different links of the supply chain of illicit drug made the whole operation more resilient. All of the links that make up the supply chain must be taken down to ensure that a CLM gang is no longer functional. However, as each of the three elements (consumer, delivery centre and operational centre) can be in different jurisdictions, an important effort of coordination between police forces is needed. Considering that in Great Britain (England, Scotland and Wales) there are 44 different local police forces, it is key to understand the territo-

rial logic of the gangs that operate under the CLM to build more efficient strategies when dismantling a CL operation.

The Metropolitan Police of London has declared in different reports that the CLM operators work in a basis of offer-demand [Rescue and Analysts, 2019, Rescue and Analysts, 2010]. However, evidence about the presence of gangs is not clear about this respect [Coombes, 2018, Andell and Pitts, 2018, Robinson, 2019, Stone, 2018, Madeley, 2018]. In Chapter 2, we study what is the *territorial logic* of the CL operators, thus understanding which are the incentives and the drawbacks when establishing a local distribution. In particular, we are interested in testing if, as the Metropolitan Police suggests, the CL operators are solely based on an offer-demand principle.

Given the territorial extension of the distribution network and the lack of available data from different gangs, we take a spatial analysis approach which would allow us to understand the *territorial logic* behind the CL operators. By using different spatial models, we are thus testing different approaches and understandings of how flows (of illicit drugs) between a place i and a place j would take place. We feed our models with public data from the Metropolitan Police of London, the Office for National Statistics, NHS digital and the House of Commons Library.

With the resultant work in Chapter 2, we provide a quantitative understanding of a problem that has brought public harm into the population of the UK, while also providing important information for authorities about how to tackle the CLM problem in a more intelligent and coordinated way. This is done by understanding the behaviour of the gangs via the data collected by different institutions of the British Government. Put it in another way, the work here intends to improve the illicit drug security system in the UK by implementing a study that understands the behaviour of the gangs at a national level, thus understanding the behaviour of the perpetrators to make the security system more efficient in decision-taking processes.

Chapter 3

When the COVID-19 pandemic hit in Western Europe during the first quarter of 2020 and lockdowns started to be imposed throughout the different countries and regions, working habits had to suddenly adapt to the new reality. Quickly called the “New Normal” [Habersaat et al., 2020], this new phase would involve a restricted mobility, restricted face-to-face contact, and particularly for this research, a Work-from-Home policy which -depending on the country and the job responsibilities- would make employees to remotely work from home until new notice.

However, as the “New Normal” was imposed, different industries adapted differently [Dwivedi et al., 2020]. Some literature was published around how cyber

security could be compromised by this new restrictions [Gerke et al., 2020], as the cyber security of companies and institutions would now largely depend on the cyber hygiene of the employees [Abukari and Bankas, 2020], thus leaving an important opportunity for cyber attackers to penetrate systems while the companies transitioned to the “New Normal”. In Chapter 3, we study what the top companies in the world were talking about when the pandemic started, and if cyber security was part of the priorities in that discussion. We frame our analysis in a validated model [Huang and Pearlson, 2019b, Marotta and Pearlson, 2019] of Cyber security culture for organisations to understand if the elements analysed actually resulted in a change of behaviour around cyber security.

We compiled a database of the official documents that these companies released when the COVID-19 pandemic globally started in the first half of 2020. Using Topic Modelling from Natural Language Processing [Blei et al., 2003, Nikolenko et al., 2017], we extract the different subjects that companies were talking about at that moment in time, not only by company but also by industrial sector and region of the world. That way, we make a disaggregated analysis that allow us to understand how each industrial sector of the different parts of the world reacted to the pandemic. As said, we base our framework in the Cybersecurity at MIT Sloan model (CAMS model) [Huang and Pearlson, 2019b], to extend it to the “New Normal” context.

With this study, we contribute by not only understanding how the companies were behaving when the pandemic hit, but we are also contributing by *a posteriori* knowing which of the industrial sectors studied financially persisted the most by thinking of cyber security when Work-from-Home policies started. That way, our research can benefit different companies of particular sectors to make a more resilient cyber security culture now that a hybrid working pattern has been adopted throughout the world [Sakurai and Chughtai, 2020, Mukherjee et al., 2020].

The work done in Chapter 3 thus intends to better understand the transition of different companies around the world during the COVID-19 pandemic to the “New Normal” with respect to cyber security by implementing the behaviour at the managerial level.

Behaviour, security and cyber security: Research gap, questions and contribution

In the present work we are studying three different problems in which, by understanding the behaviour of individuals or collectives at different stages, we improve the solutions around a particular aspect of a security system. In Chapter 1, we study the behaviour of employees towards cyber security risk to improve the training they receive so they can become more accurate sensors when facing a potential threat. In

Chapter 2, the focus is on understanding the territorial logic of gangs transporting illicit drugs from London to other parts of England so the different police bodies in England have better information on how to tackle them. Finally, in Chapter 3, we expand a proven cyber security culture model [Huang and Pearlson, 2019b] to the context of the COVID-19 “New Normal” in order to understand the behaviour of companies around cyber security when the pandemic started, from a managerial point of view.

When reviewing each one of the overall security systems we study, it is noticeable how literature published around does not take into account the heterogeneity of the different humans taking part in it. For Chapter 1, training of cyber security for employees takes the punitive approach without asking about the heterogeneity of workers that could arise in a big company [Heartfield and Loukas, 2015, Heartfield et al., 2016, Heartfield and Loukas, 2018, Yeoh et al., 2021, Rege et al., 2020, Sabilon et al., 2019]. For Chapter 2, the different police forces take an offer-demand perspective [Rescue and Analysts, 2019] towards the CLM, while not asking the different social and economic features that have allowed the model to prosper in such a successful way. It is only qualitative work from anthropology and social studies that has researched this [Coombes, 2018, Andell and Pitts, 2018, Stone, 2018], although their objectives do not seem to be found in governmental intelligence [Black, 2020a, Crime Agency, 2019, Silver and Intelligence, 2021]. Finally, for Chapter 3, we retake a qualitative behavioural model to increase the cyber security of a company [Huang and Pearlson, 2019b, Marotta and Pearlson, 2019] and extend it towards a data-driven model which could help to understand the cyber security culture of any set of companies. This extension comes as, although very successful qualitative models have been developed, as the one cited before, the behavioural element has not been seen taken in a more quantitative/data-wise framework.

Although each one of the three problems have their own research gap that will be developed respectively in each Chapter, we can detect a more general research gap from the literature and the problems described above in this introductory Chapter. As said, we detect a gap where little or no research has been done when studying security systems from a quantitative and data-driven point of view with a perspective involving human behaviour, particularly taking into account its heterogeneity.

The threading research question throughout this dissertation then comes to: how can we improve a security or cyber security system by taking into account the behaviour of the individuals or collectives within? This general research question is complemented with the three more specific questions that each one of the three problem has. In Chapter 1 we will focus on the question (i) how can we improve the training of the employees by understanding how they behave and react with respect to cyber security risk and the training they receive? By taking into account

the heterogeneities of the studied workforce using two surveys, we are introducing a behavioural analysis which allows to construct better communication strategies for a non-homogeneous set of workers in order to approach them in the best way to increase their efficiency to detect potential cyber threats. In this case we are using a behavioural segmentation model [Kharlamov et al., 2018] which we train and test using a statistical learning framework. In Chapter 2 we study the question around (ii) how can the English police forces better tackle the CLM problem by understanding the territorial logic of their operators? In this case we are interested in knowing the different social, geographical and demographic differences within the police forces territories in England that are taken into account by the criminals to choose one territory over other to establish a local market. By doing so, we are taking into account the mentioned variables to understand the behaviour of the criminal organisation so the law-enforcement bodies can have a better insight to tackle the illicit-drug problem in England, thus improving the security of the country. Finally, in Chapter 3 we focus on (iii) how can we help adapt a cyber security culture in companies that are transitioning to Work-from-Home by understanding their discussion when the pandemic started? In this case we are studying the response documents that the top global companies published when the pandemic started using Natural Language Processing. We link the results of the text analysis with a cyber security model [Huang and Pearlson, 2019b] and to the financial performance of the companies during the pandemic. The objective in this particular chapter is to understand if the topics, linguistic dimensions and values that the companies show in their documents are linked with how much they take into account cyber security.

As output of this dissertation we want to contribute to the theory (and literature) of Behavioural Data Sciences, Behavioural Sciences, Cyber Security and Criminology by providing different solutions to increase the efficiency of the security systems from a data-driven and human behaviour perspective. The main contribution is the implementation of the analysis of human behaviour as individuals or as collectives into the different security and cyber security systems studied here. As an extension of this contribution, our objective is to make these solutions practical enough so different parties can employ them, depending on their respective subject.

Two different characteristics are present as part of the solutions found for each one of the problems: (a) the making of policies for the stakeholder with respect to security and cyber security problem and (b) increasing the cyber security culture in companies. The first point is particularly present for the County Lines Model chapter (Chapter 2) and the “non-punitive” cyber security training (Chapter 1), while the second point is present in the latter project (Chapter 1) and the official COVID-19 documents chapter (Chapter 3).

While Chapter 1 refers to the employees of a company, Chapter 3 refers to the

managerial levels, and Chapter 2 studies perpetrators of an illicit activity. Thus, different areas of a given system where humans can be part of the solution to improve the security system are studied in the present dissertation. An emphasis is always made in obtaining research outputs that can become part of a policy/strategy decision [MacArthur et al., 2022]. This emphasis is particularly made in Chapter 2 where the outputs can help the Metropolitan Police to improve their efforts against the County Lines Model, and in Chapter 1 where our outputs helped the Financial Institution we worked with to improve their cyber security trainings. All of these ideas are summarised in Table 1.

Moreover, the work done in Chapter 1 mentioned has as output research policies that change the behaviour of the workforce in order to increase the cyber security culture. That is, changing their behaviour to make them a more efficient cyber threats sensor while also improving their *attitudes, values and beliefs* around cyber security [Huang and Pearlson, 2019b, Marotta and Pearlson, 2019]. The objective to also improve the cyber security culture is also present in Chapter 3.

Table 1: Summary of studied systems and objectives for each Chapter.

Chapter	Type of system	Level studied	Objective
1	Cyber security	Employees	Making strategies to increase the efficiency of threat detection systems; Increase the cyber security of the company.
2	Security	Threat perpetrator	Making strategies to increase the efficiency of threat detection systems.
3	Cyber security	Top Managerial employees	Increase the cyber security of the company.

Structure and summary

After this introduction, each one of the works described above are presented in Chapters 1-3. As stated before, in each one of them we expose their respective literature reviews, while also developing their own specific research gaps and questions, and also explain their particular methodologies and databases used. Each chapter has a similar structure in which we develop a literature review and a methodology, we present the results to later discuss them and finally conclude. We then proceed to the general Conclusions (Chapter 4). Different Appendices (A-F) are also presented to support the work of the main chapters.

As a summary of this introductory chapter, we motivated the need to research around data-driven solutions incorporating human behaviour to make security systems more efficient, as the existing solutions do not take into account the heterogeneity of behaviours/profiles, taking a one-lane strategy where technology now allows to incorporate a more diverse perspective. We set up the three different problems to develop, these being the efficient training of worker in cyber security using a behavioural segmentation model (Chapter 1), the territorial logic of the County Lines operators (Chapter 2), and the study of the cyber security culture in companies when the pandemic hit and the world transitioned to the “New Normal” (Chapter 3). We then set up the research gap, stating that data-driven solutions in security systems fail to incorporate human behaviour, to then set the research questions outlined above.

Chapter 1

The effects of a framed communication in cyber security risk perception and action: a data-centred case study using behavioural segmentation

Chapter Abstract

Cyber security threats are becoming more and more sophisticated. Although software and hardware improvements help overcome these threats, it is also important to efficiently train humans to properly detect possible attacks. We present a framework for developing disaggregated communication strategies to increase the cyber security in companies using data science methods. The framework is user-centred and is designed to minimise the risk of employees incorrectly assessing potential cyber threats. We worked with one of the most important financial companies in the UK) to develop a two-staged analysis to test our framework. For the first phase we obtained 605 respondents, while for the second 150. In phase 1, we trained a behavioural segmentation model based on the Cyber-Domain-Specific-Risk-Taking scale. The scale allows to create a 4-segments model based on the participant risk-taking and risk-perception scores. In phase 2, we validate our behavioural model and we assess a tailored Human-as-a-Security-Sensor test on each participants. The test is presented with one of three treatments at random. These are neutral (no treatment), positive (highlighting benefits of cyber security) and negative (highlighting consequences of a cyber attack). From the four segments, one is not affected by any treatment, a second is positively affected by both treatments, a third is negatively affected by both treatments, while the fourth is positively affected by the negative treatment and negatively affected by the positive treatment. This work uses novel techniques of data science to introduce a method to efficiently communicate and train a large volume of heterogeneous employees.

1.1 Introduction

When implementing cyber security strategies in large companies, the focus tends to lie on the implementation of software, firewalls and a generalised training for the employees [Alshaikh and Adamson, 2021, Yeoh et al., 2021, Rege et al., 2020, Kwak et al., 2020, Sabillon et al., 2019]. In addition to the proven benefits of software and hardware for a correct cyber security, these strategies usually adhere to the view of humans as “weak links” with a limited ability to detect and report cyber attacks to the company [Morgan et al., 2020, Stanton et al., 2005, D’Arcy et al., 2009, Sasse et al., 2001]. Although this view has been contested [Rege et al., 2020, Morgan et al., 2020, Heartfield et al., 2016, Heartfield and Loukas, 2018], to design and to implement an user-centred strategy that complements hardware and software can result challenging, particularly when it comes to companies with a large and heterogeneous workforce. In this article we present a comprehensive framework for implementing a disaggregated strategy that improves communication and engagement of cyber security threats in companies with a large number of employees. In short, the objective of this framework is to increase the capacity of employees to accurately assess potential cyber security threats. As said, the framework could be of particular interest for a company with a heterogeneous set of workers, as different kinds of employees could respond differently to a cyber security training, thus allowing for an unknown number of workers to have a deficient engagement with the company’s cyber security culture and standards.

We collaborate with one of the largest financial companies in the UK. The company is interested in increasing the cyber security culture and resilience. This aim is of particular interest during a hybrid time where companies do not only have to be aware of the cyber hygiene of employees at the company’s facilities, but also at home given the work-from-home policies around the globe by the COVID-19 pandemic [Caligiuri et al., 2020, Gerke et al., 2020, Dwivedi et al., 2020]. We accomplish this goal by (i) understanding the actual cyber security culture of the workforce, (ii) building a behavioural segmentation model which allows us to better map how the surveyed population perceives and reacts to cyber security risks, and (iii) design the right communication strategies using a tailored Human-as-a-Security-Sensor [Heartfield and Loukas, 2018] with different primings.

The behavioural model used is based on the one created in [Kharlamov et al., 2018], where the authors examine, using a 13-question probe related to different kinds of cyber security attacks, the perception associated to cyber security risk, and the level of engagement towards this risk. The authors named it the Cyber Domain Specific Risk Taking scale (CyberDoSpeRT), with which they map the results of the

risk perception and risk-taking scores to 4 different behavioural segments: Anxious (high risk perception, low risk taking), Opportunistic (high risk perception, high risk taking), Relaxed (low risk perception, high risk taking) and Relaxed (low risk perception, low risk taking).

The tailored Human-as-a-Security-Sensor test is based on the work of [Heartfield and Loukas, 2018], where the authors develop a framework with the same name with the objective of encouraging users to detect and report semantic cyber attacks. Using the taxonomy introduced in [Heartfield and Loukas, 2015], we choose 6 different cyber attacks (social media masquerading, WiFi speed masquerading, typosquatting and phishing, Microsoft Planner phishing, Phishing email and QRishing) to test the user’s ability to detect possible cyber threats and have an insight about their reasoning behind their decision. This is done by asking two follow-up questions: how confident is the user about the decision (if the presented situation is a cyber attack or not), and asking the user to briefly describe the reasons behind their decision. The second question is then analysed using Topic Modelling via LDA [Blei et al., 2003, Nikolenko et al., 2017], a tested method in Natural Language Processing. We call our tailored test the Human-as-a-cyber-Security-Sensor test (HaaCSS test).

When presented to the user, our HaaCSS test is primed with 3 different treatments. A neutral one (no treatment) which serves as our baseline. A “Positive” one highlighting the benefits of cyber security, and a “Negative” one highlighting the consequences of a possible cyber attack. The 3 different primings are allocated randomly to each of the respondents. By testing the user’s capacity to correctly assess cyber threats, we are putting the employees as the centre of the design of solutions to augment the cyber security culture and resilience of the company. This is in addition to other cyber security solutions that any company should implement.

Combining the elements described above, we design a two staged study to empirically test the effects in the company’s workforce. In the first phase of the study, a survey was launched to train our behavioural segmentation model using the CyberDoSperRT scale. We are also interested at this stage in studying the impact of the COVID-19 pandemic in cyber attacks. Healthcare organisations are a common target to cyber security attacks [Nifakos et al., 2021], particularly during the COVID-19 pandemic [De Cauwer and Somville, 2021]. We are interested in knowing if this is also the case for this particular financial company. The second phase of the study is designed to validate the behavioural segmentation model and to implement the HaaCSS test with different communication treatments. For the first phase of the study, we obtained 605 respondents, while 150 for the second one.

As said, this objective of this project is to present a framework to increase the capacity of employees to correctly assess cyber threats, taking into account their

heterogeneities present in a large workforce. To do so, we are using a *statistical learning* work pipeline which allows us to (i) train with the features obtained from the first survey our behavioural segmentation model to then (ii) validate it with our second survey to finally (iii) analyse the result of our HaaCSS test with respect to our trained and validated behavioural model.

Taking into account the above, the questions we are interested in answering are: which are the most important features in each of the behavioural segments in the financial company? Which are the effects of the different framed communications in the general workforce when doing the HaaCSS test? Can we see different effects in each of the segments depending on the primed presented to them in the test? Can we thus create a disaggregated communication strategy to make the company's workforce more engaged with cyber security?

We first introduce a literature review in Section 1.2, where we also explore the literature gap found and the hypotheses we work with. In Section 1.3.1 we present the CyberDoSpeRT scale and the HaaCSS test implemented with the three different primings (neutral, negative and positive). The general methods to analyse and interpret the collected data is presented in Section 1.3.2. The results of the two different stages are presented in Section 1.4.1 and Section 1.4.2 respectively, as those with respect to the COVID-19 pandemic (Section 1.4.3). We then discuss the results in Section 1.5, to then conclude in Section 1.6. This paper is accompanied with three appendices, where we present with more detail survey 1 (Appendix A), survey 2 (Appendix B), and the HaaCSS test (Appendix C).

1.2 Literature review

The Human-as-a-Security-Sensor is a paradigm defined by the authors in [Heartfield et al., 2016, Heartfield and Loukas, 2015] based on the need to construct security systems with a holistically approach, allowing the ability of the system's users to detect and report threats. The cited authors define it as "*The paradigm of leveraging the ability of human users to act as sensors that can detect and report information security threats.*" As the definition states, the used test in our study is based on the capacity of humans to discriminate different kinds of semantic attacks with respect to the technical and automatic tools, like security software. In that sense, we are interested in generating a test allowing to score the capacity of the company's employees to correctly assess different cyber security semantic attacks. The latter are defined as "*the manipulation of user-computer interfacing with the purpose to breach a computer system's information security through user deception*" [Heartfield

and Loukas, 2018]. The idea of humans as sensors to detect cyber security threats is in direct opposition with the idea of humans as the “weakest link” [Heartfield and Loukas, 2018]. However, the idea has been contested and humans have successfully been part of security detection systems for physical systems and cyber physical systems.

In [Zheng et al., 2014], the authors study data from a New York City app where citizens can complain about issues encountered in the city. The large and accurate number of complaints with respect to noise-pollution problems allowed the authors to generate an analysis that could sort the different sources with respect to category and importance, with a potential to make the city’s agencies to become more efficient when treating this kind of public problem. In [Jürrens et al., 2009], the authors propose a system for failure detection in water supply systems. The sensor distinguishes between two kinds of users: those reporting and those receiving the reports of failure, while supposing that any of the two kinds of users are trained to professionally detect water systems failures.

In recent years, the User-Entity Behaviour Analysis (UEBA) literature has emerged around cyber security systems [Revanth Filbert Raj and Babu, 2019, Salitin and Zolait, 2018]. As an extension of the User Behaviour Analysis (UBA), UEBA also incorporates the analysis of physical entities into the research. Thus, this literature intends to detect potential threats in users through their behaviour and the behaviour of the entities they used (IP logs, usernames, etc.). In the latter study, the authors perform a small literature review in which they outline the different objectives and methods from the recent UEBA literature. In [Salitin and Zolait, 2018], the authors perform a second literature review to understand which mechanisms of Machine Learning, such as Support Vector Machines and different methods involving Neural Networks could be implemented to real-life detection of cyber attacks using UEBA. However, the latter study lacks any application of the reviewed methods. In UEBA we find different sources of data as proxy of user’s activity, such as IP addresses logs [Raguvir and Babu, 2020]. This also allows to use different methods of Machine Learning and Data Science to analyse the users behaviour [Muliukha et al., 2020, Chen et al., 2017]. Within the known literature, none uses clustering analysis nor natural language processing methods as it is done in this paper.

An interesting article is the one of [Chowdhury et al., 2019] where the authors explore the importance of time pressure on individuals when a cyber threat is present. The authors present a framework with which this problem can be addressed.

Another interesting take when studying cyber security from a human behaviour perspective is the work in [Steingartner et al., 2021], where the authors study a cyber deception system approach rather than a cyber prevention system. Cyber deception deals with implementing different proxies, traps and decoys to *deceive* the

attacker, thus not allowing the cyber attack to ever happen.

Within the UK, study cases have been performed in local communities regarding the role of the security of the citizen with respect to the security of the technology when designing cyber security architectures. The authors in [Coles-Kemp et al., 2018] worked with a deprived local community in North Eastern England to develop a digital tool that could give communal support regarding public welfare. This was done also to study how individual and communal cyber security in digital platforms is understood. As part of their conclusions, the authors find as a structural absence to not knowing the target person to whose the cyber security is directed. Also, a study is performed in [Whitty, 2021] in order to develop a model of archetypal insiders after 99 case studies of British organisations and corporations receiving insider attacks. The author interviews figures of different managerial levels that were in contact with the attacker to then arrive to a list of archetypal behaviours of internal perpetrators of cyber attacks. The author distinguishes between 8 different archetypes.

Regarding our scope on cyber security culture and awareness within similar companies, only the work in [Marotta and Pearlson, 2019] was found. The authors base their research on the the Cybersecurity at MIT Sloan Model (CAMS) [Huang and Pearlson, 2019b] for cyber security culture at managerial level at the Banca Popolare di Sondrio and do not use any quantitative tool. The CAMS model objective is to increase the cyber security culture in a given company from an organisational perspective. In other words, to increase the awareness in cyber security in the different managerial levels taking into account internal culture factors, such as organisational mechanisms, and external culture factors, such as those given by the particular country where the company resides. The model was constructed parallel to a case study with Liberty Mutual’s managerial leaders to understand how the beliefs, values and attitudes inside the company -and particularly at managerial level- can positively change the behaviour with respect to cyber security. The CAMS model was then validated with an Italian Bank, thus also changing the external culture factors with respect to Liberty Mutual’s one (USA), in [Marotta and Pearlson, 2019]. The CAMS model, although comprehensive by taking into account different internal and external factors, was built as a qualitative model.

In recent years, studies have been published regarding the improvement of cyber security awareness within companies. However, all of the following studies are of qualitative nature. The authors in [Sabillon et al., 2019] present a model to efficiently train different kinds of employees within an organisation. Nevertheless, the focus the authors have is by assuming that humans are the “weakest link” in cyber security systems, which is the opposite to ours. The authors of the latter research validate their model in a Canadian company. In [Rege et al., 2020], the authors introduce

a social engineering awareness training for STEM students and professionals. In [Sillanpää and Hautamäki, 2020] the authors study the behaviour of the employees in a Finnish company towards social engineered cyber security threats by using several penetration methods. They show how, although employees can understand the different cyber attacks and know what to do when one is detected, there is a statistical significance with how employees behave when they face one. The study in [Alshaikh and Adamson, 2021] moves from the perspective of increasing awareness and presents a study case in Australia where the goal was to influence the workers' behaviour with respect to cyber security. This shift from increasing awareness to positively changing the workers' behaviour is more aligned with the goal of the present paper. This research is done by studying the effects of implementing different strategies from the psychological attachment theory.

Finally, an emphasis has been done in phishing email which, as the results in the present study shows, is the most recognised method of cyber threat by the surveyed participants. In [Kwak et al., 2020], the authors study the psychological factors that might affect the low records on reporting this kind of threat. Using survey data, they show how by increasing the cyber security awareness and by highlighting the negative outcomes of a cyber attack via a phishing email, the share of reported attacks increased. In [Yeoh et al., 2021], an awareness and training campaign around phishing emails is studied simulating a phishing attack. The authors arrive to the same conclusion that awareness and highlighting the costs of allowing a cyber attack changes the individual's behaviour to increase the awareness of this phishing emails vulnerabilities.

The CyberDoSpeRT scale that is used in this work [Kharlamov et al., 2018] is based on the DoSpeRT scale developed by [Weber et al., 2002] and then modified in 2006 [Blais and Weber, 2006]. The original scale was developed to test the difference between risk taking attitudes and risk perception with respect to very specific domains like ethics, health and social problems. The scale has since been adapted to different specific domains (like [Einav et al., 2012, Harris et al., 2006, Wilke et al., 2014]) and translated to different languages. In [Kharlamov et al., 2018], the authors adapt the original DoSpeRT to cyber security for the first time with 13 specific questions around the field of cyber security. The CyberDoSpeRT easily allows to apply quantitative methods as the surveyed population need to assess their risk perception and risk taking attitude on a scale from 1 to 7.

From the literature reviewed above, we can delimit a literature gap encountered and that we intend to fill in this work. Although there has been research done to create models that are tailored to increase the cyber security awareness and culture in companies [Rege et al., 2020, Marotta and Pearlson, 2019, Huang and Pearlson,

2019b, Yeoh et al., 2021] none of them has tried to implement a quantitative approach since as foundation which would allow to have a disaggregated study that could be scaled to large enough companies. Also, we found that in [Rege et al., 2020] and [Yeoh et al., 2021], the authors find that a “punitive” approach, highlighting the negative consequences of a cyber attack, increases the awareness and thus the efficiency of their frameworks. However, none of them try the opposite approach where highlighting the positive consequences of a cyber awareness and prevention has.

Different study cases were also found in our literature review, where the authors validate a given hypothesis or model in different companies around the world [Marotta and Pearlson, 2019, Sabillon et al., 2019, Rege et al., 2020, Sillanpää and Hautamäki, 2020, Alshaikh and Adamson, 2021]. The surveys or interviews are normally targeted to managerial posts, or *ad hoc* employees [Whitty, 2021]. However, none of the studies tries to understand the behaviour of a random sample of employees in the company.

Although different studies recognise and study the heterogeneity of the employees in a company, there has not been research works that actually correlate the difference within employees to cyber security behaviour.

Thus, the literature gap found is a study that presents a quantitative framework that is actually tested in a random sample of a given company. As intended, a study that also explores the effects of a “positive” communication (highlighting benefits of prevention), while still taking into account the different kinds of employees that can be found in a company, has not been found in the literature.

The CyberDoSpeRT [Kharlamov et al., 2018] and the HaaSS [Heartfield and Loukas, 2018] are a good fit to our study given the scalability that they have to a large volume of surveyed individuals. By combining the two, interesting results correlating the risk perception, the risk taking attitude and the ability to correctly detect a cyber threat could be found. Thus, from the literature reviewed, these two are the ones chosen for this study.

Following the literature on cyber attacks during COVID-19 [Nifakos et al., 2021, De Cauwer and Somville, 2021], we have as a first working hypothesis that

H1: The bank employees experienced an increased number of cyber attacks, particularly once they moved to a Work-from-Home routine.

Also, from the literature on semantic attacks [Heartfield and Loukas, 2015, Heartfield et al., 2016, Heartfield and Loukas, 2018, Zheng et al., 2014, Jürrens et al., 2009] where humans tend to be perceived as “weak links” when it comes to cyber security

attacks, thus we can expect that

H2: The Positive treatment will not have statistically significant worst results than the Negative treatment in our HaaCSS test.

That is, employees will not have a significant worst performance when receiving the Positive treatment than when receiving the Negative one. We can also extract a third hypothesis with respect to the Negative treatment by stating that:

H3: The Negative treatment will not act beneficially throughout the four behavioural segments of our model.

1.3 Methodology

The financial company we worked with has made an important effort in implementing a mandatory cyber security training for its employees. However, as cyber attacks always transform in form and content, the company is interested in keeping the cyber trainings and updates on employees as effective as possible. That is, by making them as accessible and understandable as possible.

However, as the company has a workforce of more than 10 000 employees, to design a unique training that is efficient for the entire employee population is virtually impossible given the heterogeneity of the employees.

Also, the financial company is interested in knowing the cyber hygiene habits of their employees during Work-from-Home times.

In order to design and implement a research that allow us to obtain an effective communications towards a disaggregated training for employees, we decide to undertake a two staged survey as part of a quantitative case study. During the first stage of the survey we also ask some questions related to the perceptions of cyber attacks during a Work-from-Home policy time.

For the case study we are supposing a myriad of affirmations: (1) the surveyed employees have a heterogeneous behaviour with respect to the cyber security prevention and detection. (2) Employees will react differently to each one of the treatments depending on the prior behaviour detected. (3) Employees are not “weak links” with respect to cyber security detection, but actually can be a good indicator of a cyber and/or semantic attack. (4) Employees will see more cyber threats since starting to work from home.

1.3. METHODOLOGY

Given the heterogeneity and the volume of the workforce, a survey is the most appropriate method to capture the biggest number of information from the employees.

The second survey involves a Human-as-a-Cyber-Security-Sensor (HaaCSS) test that has been adapted from [Heartfield and Loukas, 2015]. We chose to adapt this test to our case study as the authors follow the same understanding as the scope of this work about semantic attacks and the ability of humans to correctly detect cyber threats when correctly trained. Each one of the two steps is detailed in Section 1.3.1.

The assessment of the information obtained by the survey is analysed using the CyberDoSpeRT [Kharlamov et al., 2018]. To our knowledge, this is the only test that has been adapted for cyber security matters and that studies the risk perception and risk aversion founded in an already established theory (DoSpeRT tests [Weber et al., 2002]).

Finally, as our numerical data obtained from the participants consists of natural numbers (integrals), we must choose the adequate statistical tools for analysing the results.

The statistical method and CyberDoSpeRT model used for this work are described in Section 1.3.2.

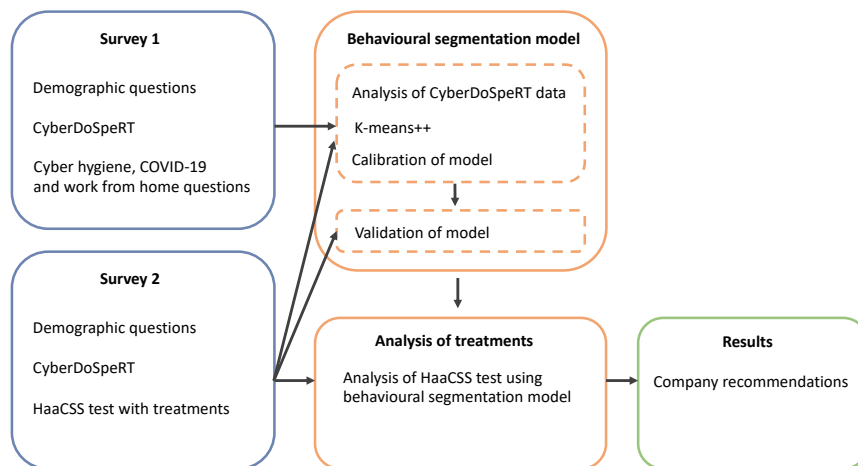


Figure 1.1: Design of the study done.

A diagram of the whole research design is shown in Figure 1.1. The Figure details each one of the surveys, stating which information was asked and which purpose it serves. Figure 1.1 is supported with Table 1.1.

Table 1.1: Detail of the two surveys.

Survey 1		
Question	Type	Use
Demographic data	Demographic and working habits questions (age, sex, number of working hours, etc.)	Used as features for our behavioural segmentation model using a Probit regression.
CyberDoSpeRT	13 different activities related to cyber security. Respondents have to rate from 1 to 7 their perception of risk and their likelihood to engage with the activity.	Used to obtain behavioural segment of respondent with respect to the risk taking and risk perception score.
Cyber hygiene, COVID-19 and work from home questions	Similar to demographic questions, but about the cyber hygiene at home while working from home during the COVID-19 pandemic.	Used for analysis of cyber security hygiene during the pandemic.
Survey 2		
Question	Type	Use
Demographic data	Demographic and working habits questions (age, sex, number of working hours, etc.)	Used as features to validate the behavioural segmentation model from Survey 1.
CyberDoSpeRT	13 different activities related to cyber security. Respondents have to rate from 1 to 7 their perception of risk and their likelihood to engage with the activity.	Used to obtain behavioural segment of respondent with respect to the risk taking and risk perception score and compare it with the one obtained in Survey 1.
Treated HaaCSS test	6 different scenarios with potential starts of cyber attacks. Respondents have to evaluate if (i) they think is a cyber attack or not, (ii) how confident they are about their answer and (iii) give in a sentence long their reasoning behind their answers. Tests are treated with three different primings	Used to obtain scores and then correlate results with behavioural segments using Poisson regression analysis.

1.3.1 Data gathering

Two different surveys were launched. The first one was launched from October 19 to November 6, 2020. We obtained results for 605 respondents. 26 different questions were asked with respect to demographics, working habits during the pandemic (March 23, 2020 onward) and cyber security. In particular, as part of the survey, the participants had to answer the cyber domain specific risk taking scale (CyberDoSpeRT) [Kharlamov et al., 2018]. In it, the respondents consider 13 different activities related to cyber security and scale their perceived risk to the activity and the likelihood to engage with it. For each activity, the allowed answers run from 1 (not at all risky/extremely unlikely to engage with) to 7 (extremely risky/extremely likely to engage with). The 13 activities listed in the CyberDoSpeRT are:

1. Not using anti-virus or anti-malware protection.
2. Enabling automatic uploading and/or automatic back-ups.
3. Linking multiple social media websites (e.g., linking Twitter, Facebook, and Instagram accounts, etc.)
4. Using a wearable device to collect your private data (e.g., FitBit, AppleWatch, etc.)
5. Letting web browser remember your passwords.
6. Not installing software updates as soon as they become available.
7. Not making hard drives unreadable before disposing of the old PC.
8. Providing private information (such as your email address) to obtain free WiFi in public places such as coffee shops, airports, train stations, etc.
9. Installing an Internet-connected security system in your home.
10. Not reading App permissions before uploading an App on your smartphone.
11. Not using tools which protect your browsing history (e.g., TorBrowser).
12. Keeping Location Services on your smartphone turned on.
13. Letting web browser remember your credit card information.

The second survey was launched from May 13 to June 3, 2021. We obtained information for 150 respondents. We also included different questions about demographics, working habits and the CyberDoSpeRT scale. The participants were also

presented with a HaaCSS test (Appendix C) [Heartfield and Loukas, 2018] consisting of six different screenshots replicating real potential starts of cyber security threats (social media masquerading, WiFi speed masquerading, typosquatting and phishing, Microsoft Planner phishing, Phishing email and QRishing). For each of the different cases, the respondents had to (i) answer if they believed or not it was the start of a potential cyber security attack, (ii) rate their confidence about their answers on a scale from 0 (not confident at all) to 10 (extremely confident), and (iii) describe in a few sentences why they chose their answer.

The data recorded for question (iii) is analysed using Latent Dirichlet Allocation for Topic Modelling [Blei et al., 2003, Nikolenko et al., 2017]. This analysis allows us to obtain a specific view of the different topics covered by the surveyed population with their respective keywords (most frequent and important words) and their frequency (how much these topics are mentioned).

For this HaaCSS test, participants were divided into three different groups. Each one of them received a different priming with respect to cyber security. These are:

Neutral treatment No priming

Positive treatment “Correct identification of threats by the employees (i) saves the COMPANY millions of pounds annually; it also (ii) helps to avoid wasting resources (both human and financial) on “false alarms” (i.e., on responding to potential cyber incidents, when harmless events are reported as potentially dangerous). Previous research found that people, on average, correctly detected over 70% of threats in similar tasks and often detected potential cyber threats better than technical (automated) security systems.”

Negative treatment “Wrong identification of threats by the employees (i) costs the COMPANY millions of pounds annually; it also (ii) results in wasted resources (both human and financial) on “false alarms” (i.e., on responding to potential cyber incidents, when harmless events are reported as potentially dangerous). Previous research found that people, on average, failed to detect approximately 30% of threats in similar tasks and sometimes detected potential cyber threats worse than technical (automated) security systems.”

As the answers for the specific questions not listed in the main body of this work, More details of each survey can be found in Appendices A, B and C.

1.3.2 Methods

Segmentation and clustering

As stated in Section 1.3.1, the participants of each of the two surveys had to answer a 13 activities CyberDoSpeRT test [Kharlamov et al., 2018]. For each of the activities, they had to scale from 1 to 7 their risk perception of the activity and their likelihood to engage with it. Adding the 13 scores, we obtain the final numbers for both risk perception and likelihood of engagement. Each of these scores is a natural number going from 13 (1x13) to 91 (7x13). With both scores, we can then map each of the participant’s behaviour around cyber security threats in a two-dimensional way.

Taking this two-dimensional mapping into account, we use the 4-types segmentation used in [Kharlamov et al., 2018]: anxious (high perception and low engagement scores), opportunistic (high perception and high engagement scores), relaxed (low perception and high engagement scores) and ignorant (low perception and low engagement scores). Each type has its own advantages and drawbacks. Even though risk averse types (i.e., Anxious and Ignorant) may be more preferable for companies as they generally avoid risky activities online, Anxious types may overestimate potential risks to such an extent that they will report too many false positives (cases where threat does not exist, but Anxious types believe that it is there) and Ignorant types might disengage with cyber domain to a degree that would prevent them from effectively fulfilling their work duties. In contrast, even though Opportunistic and Relaxed types are generally considered to take “too much” risk, these types might be desirable in some roles within organisations (e.g., money laundering investigators, technology developers within the company that may require engaging with cyber risk on a daily basis). The four segments can be schematically seen in Figure 1.2, where we plot risk perception and risk engagement scores on a 2D canvas. A fifth segmentation is included for those respondents who are close to have a balanced score. We call these respondents as Well-calibrated, as they are in the neighbourhood of 52 ± 5 .

Given that our sampled population has obtained previous training around cyber security threats and hygiene, we could expect a non-homogeneous dispersion of the sample over the different segments. Moreover, we are expecting a distribution that tends to be allocated in the Anxious segment, with high scores of risk perception and low risk taking scores. In order to take into account the former knowledge of the sample and have an adapted segmentation of it, we perform a k -means++ algorithm [Arthur and Vassilvitskii, 2007]. We are only interested in adapting the four original segments (anxious, opportunistic, relaxed and ignorant), as the well-calibrated scores are independent of the process.

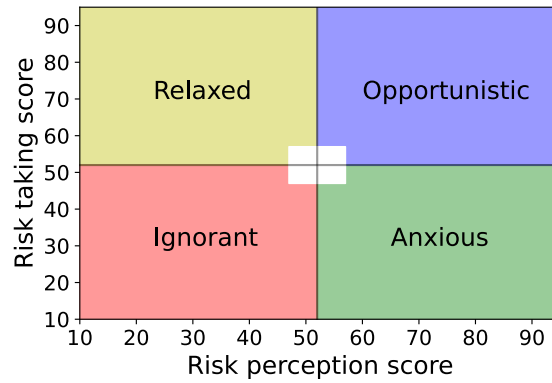


Figure 1.2: Schematic of the 4 segments (Relaxed, Opportunistic, Anxious and Ignorant), plus the fifth Well-calibrated area at the centre of the figure.

The purpose of the traditional k -means algorithm is to cluster a sample of data points with respect to k centroids. Each of these centroids minimises the distance between them and the data points within the cluster. One major limitation of the algorithm is that its results may heavily depend on the initial choice of the k starting centroids [Kanungo et al., 2004]. The k -means++ algorithm is an extension of the original k -means algorithm with an extra initial process to find optimal initial positions of the k centres [Arthur and Vassilvitskii, 2007]. This extension to the original k -means algorithm allows to eliminate the bias that arise from choosing arbitrary starting points [Shindler, 2008]. k -means++ has been theoretical reviewed [Nielsen and Nock, 2015, Bahmani et al., 2012]. A list of applications to geographical studies, financial analysis and general data science use can be found in [Lee et al., 2008, Li and Wang, 2022].

Regression analysis

We conduct a series of regressions using the probit model [Amemiya, 1978] in order to understand factors, which influence whether a particular individual is classified as Anxious, Ignorant, Opportunistic, or Relaxed. Specifically, the probit model takes the form of

$$\Pr(Y_k = 1|\mathbf{X}) = \Phi(\mathbf{X}^T\boldsymbol{\beta}). \quad (1.1)$$

Eq. (1.1) is translated as the probability of a person being part of a particular segment ($Y_k = 1$), given a number of explanatory variables \mathbf{X} coming from the survey (demographics, working habits, COVID-cyber security experience) is equal to the Normal distribution Cumulative Distribution Function (Φ) evaluated at $\mathbf{X}^T\boldsymbol{\beta}$, where $\boldsymbol{\beta}$ is a vector of free parameters.

We also perform Poisson regression in Section 1.4.2 to see the effect of the different behavioural segments in the scores obtained by the Humans-as-a-Cyber-Security-Sensor (HaaCSS) test. Given that the scores are natural numbers, we are using the Poisson model to model the impact factors in the counts. In this case,

$$\log E[y|\mathbf{X}] = \mathbf{X}^T \boldsymbol{\beta}, \quad (1.2)$$

which would mean that $E[y|\mathbf{X}] = \exp\{\mathbf{X}^T \boldsymbol{\beta}\}$.

1.4 Results: case study

1.4.1 Phase 1

In the first phase of the study we obtained the scores of 605 respondents. We segmented the population with respect to the 4 different types of behaviour plus the well-calibrated individuals. In Figure 1.3 we observe the mean answer for each of the 13 activities (listed in Section 1.3.2) with their respective standard deviations. We observe how there is a high concentration of answers in the bottom right quarter of the plane, which would correspond to a high risk perception and a low likelihood to engage with the activity (Anxious type). In order to not have an over representation of the Anxious profile, we perform a k -means++ clustering with $k = 4$. Results are shown in Figure 1.4.

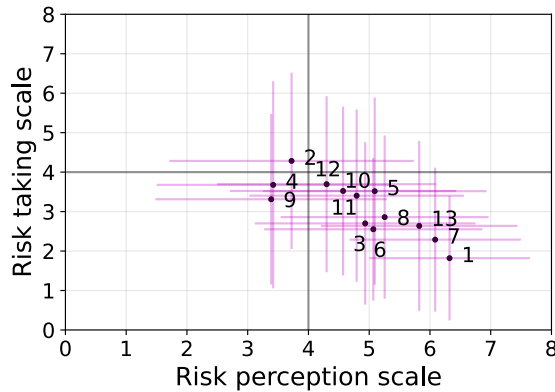


Figure 1.3: Average scores and standard deviation for each of the 13 activities of the CyberDoSpeRT.

In Figure 1.3 we can observe a negative relation between the risk perception and the risk engagement. Thus, generally speaking, the higher the risk perception is, the lower the risk engagement is. We observe that the three riskiest actions perceived are: 1 - not using anti-virus or anti-malware protection, 7 - not making hard drives

1.4. RESULTS: CASE STUDY

unreadable before disposing of an old PC and 13 - letting web browser remember your credit card information.

On the other side, the least risky actions obtained from the survey are: 9 - Installing an internet-connected security system in your home, 4 - Using a wearable device to collect private data (e.g. Fitbit, etc.) and 2 - Enabling automatic uploading and/or automatic back-ups.

With this first simple result, the company could detect a red flag with respect to any interested activities. In case that one of the 13 activities has an average behaviour that is not tolerated by the cyber security standards of the company, a more targeted approach can be taken.

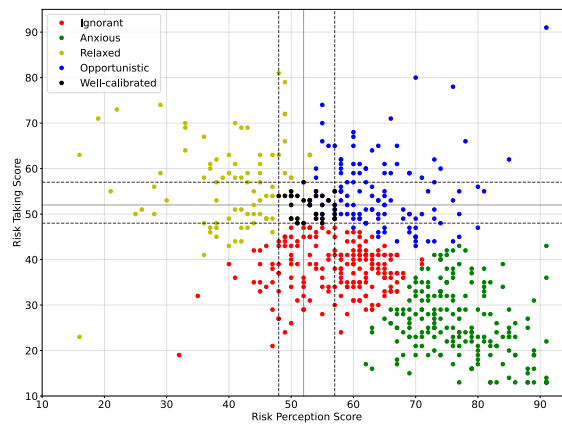


Figure 1.4: Scatter plot of the risk taking and risk perception scores for the 605 respondents of the first survey. Colours represent the different segments (Anxious, Opportunistic, Relaxed, Ignorant) found performing a k -means++ algorithm.

The next step is to have a deeper insight of the characteristics of each of the behavioural segments. In order to do so we present the different tables of results when performing a probit regression for each of the four behavioural types. As previously mentioned, a more in-depth analysis of each of the dependent variables presented in Tables 1.2-1.4 is detailed in Section 1.4.3 and Appendix A. The reader will note that the sample in which the different regressions are made is smaller than the total sample of respondents of Phase 1 (503 vs 605). This is because, given the internal policies of the company, it was not compulsory to answer to all of the questions of the survey. Thus, we are only taking into account those surveys that were fully responded.

In Table 1.2 we can observe the results for the Anxious behavioural type. Two interesting variables are found to be statistically significant: the confidence to detect cyber threats and the the age category, both with positive coefficients. This can be interpreted as a correlation between an anxious type (high score to perceive risk and low score to engage with risky activities), and older workers and employees with a

high confidence to detect cyber threats.

Results for the Opportunistic type are shown in Table 1.3. Here, we can observe opposite results to the Anxious type, as younger workers and people with less confidence in their detection abilities tend to be of the Opportunistic type (high perception score, high engagement score). Also, in this case we also observe other two important statistically significant variables, such as gender and an overtime work tendency. Self-identified women workers and employees frequently working overtime are more likely to belong to this behavioural type than to any other.

The Relaxed type (Table 1.4) is defined by a low score of risk perception and a high score of risk taking. This type correlated to the highest number of dependent variables. As with the Opportunistic type, we find a negative correlation with age and confidence in cyber threats detection. There are three extra variables worth mentioning: caring responsibilities, meaning that not having someone to take care of has a correlation with this particular profile. Also, having a low job satisfaction has an impact, as do not having a customer-facing role.

Finally, results for the Ignorant type are found in Table 1.5. Caring responsibilities in this case are positively correlated as it is the frequency of working from home. Gender is also statistically significant, finding male respondents more prone to be of the Ignorant type.

Table 1.2: Determinants of the Anxious behavioural type. Results obtained performing a Probit regression.

$N = 503$ $R^2 = 0.1230$	coef	std err	z	$P > z $	[0.025	0.975]
Frequency of working from home	0.0244	0.046	0.534	0.594	-0.065	0.114
Confidence in cyber threats detection	0.2785	0.047	5.906	0.000	0.186	0.371
Frequency of cyber attacks before pandemic	0.0197	0.081	0.242	0.809	-0.140	0.179
Frequency of cyber attacks during pandemic	-0.0263	0.085	-0.308	0.758	-0.193	0.141
Job satisfaction	0.0263	0.041	0.641	0.521	-0.054	0.107
Life satisfaction	-0.0364	0.041	-0.883	0.377	-0.117	0.044
Gender	0.0696	0.131	0.530	0.596	-0.188	0.327
Age	0.1951	0.038	5.071	0.000	0.120	0.271
Relationship status	0.0374	0.132	0.284	0.776	-0.221	0.296
Caring responsibilities	-0.0263	0.167	-0.158	0.875	-0.354	0.301
Age of cared ones	-0.0005	0.016	-0.032	0.974	-0.031	0.030
Cyber security role	0.0598	0.210	0.285	0.775	-0.351	0.470
Years of experience in company	-0.0366	0.033	-1.098	0.272	-0.102	0.029
Customer-facing role	0.2013	0.153	1.315	0.188	-0.099	0.501
Work over contracted hours	-0.0787	0.139	-0.568	0.570	-0.350	0.193
Constant	-3.3523	0.535	-6.265	0.000	-4.401	-2.304

1.4. RESULTS: CASE STUDY

Table 1.3: Determinants of the Opportunistic behavioural type. Results obtained performing a Probit regression.

$N = 503$ $R^2 = 0.09910$	coef	std err	z	$P > z $	[0.025	0.975]
Frequency of working from home	-0.0822	0.046	-1.779	0.075	-0.173	0.008
Confidence in cyber threats detection	-0.1643	0.043	-3.850	0.000	-0.248	-0.081
Frequency of cyber attacks before pandemic	0.0119	0.083	0.144	0.886	-0.150	0.174
Frequency of cyber attacks during pandemic	-0.0994	0.087	-1.137	0.255	-0.271	0.072
Job satisfaction	0.0470	0.045	1.054	0.292	-0.040	0.134
Life satisfaction	0.0496	0.046	1.071	0.284	-0.041	0.140
Gender	0.3274	0.142	2.306	0.021	0.049	0.606
Age	-0.1326	0.044	-3.015	0.003	-0.219	-0.046
Relationship status	-0.0462	0.137	-0.337	0.736	-0.315	0.223
Caring responsibilities	-0.0107	0.183	-0.058	0.953	-0.368	0.347
Age of cared ones	0.0185	0.017	1.102	0.270	-0.014	0.051
Cyber security role	-0.2188	0.246	-0.889	0.374	-0.701	0.264
Years of experience in company	0.0020	0.039	0.050	0.960	-0.075	0.079
Customer-facing role	0.1843	0.159	1.157	0.247	-0.128	0.496
Work over contracted hours	0.3816	0.145	2.625	0.009	0.097	0.667
Constant	0.1815	0.535	0.339	0.734	-0.867	1.230

1.4.2 Phase 2

In the second phase of the study we obtained 150 different respondents. Although significantly less answers than in Phase 1 (605 respondents), this second sample allow us to validate the behavioural segmentation done in the first phase. By validating it, we can also link the segmentation results with the obtained results for the Human-as-a-CyberSecurity-Sensor (HaaCSS) test.

Thus, as part of the second survey, respondents had to answer the Cyber-DoSpeRT scale. By using the same process as for Phase 1 (adding the recorded answers for each of the 13 activities for obtaining the risk perception and risk taking scores), we obtained the behavioural type of each of the 150 respondents. The segmentation was done using the same model from Phase 1, i.e. we do not perform a second k-means++ algorithm in order to obtain the respondent's type, but rather used the one done in Phase 1 to segment. The different proportions for Phase 1, Phase 2, and within the different priming samples can be seen in Figure 1.5.

By performing a Fisher's exact test, the difference between the Phase 1 and Phase 2 proportions of segments are not statistically significant, while a Kruskal-Wallis test allows us to say that the difference between the proportions for the three different treatments (Neutral, Positive and Negative) are not statistically significant either. These results allow us to (1) validate the proportion of segments obtained

1.4. RESULTS: CASE STUDY

Table 1.4: Determinants of the Relaxed behavioural type. Results obtained performing a Probit regression.

$N = 503$ $R^2 = 0.1018$	coef	std err	z	$P > z $	[0.025	0.975]
Frequency of working from home	-0.0689	0.049	-1.398	0.162	-0.166	0.028
Confidence in cyber threats detection	-0.1433	0.047	-3.080	0.002	-0.234	-0.052
Frequency of cyber attacks before pandemic	-0.0292	0.095	-0.308	0.758	-0.215	0.157
Frequency of cyber attacks during pandemic	0.0505	0.096	0.528	0.597	-0.137	0.238
Job satisfaction	-0.1139	0.044	-2.615	0.009	-0.199	-0.029
Life satisfaction	0.0794	0.047	1.687	0.092	-0.013	0.172
Gender	0.1045	0.150	0.695	0.487	-0.190	0.399
Age	-0.1223	0.046	-2.663	0.008	-0.212	-0.032
Relationship status	-0.0716	0.141	-0.508	0.612	-0.348	0.205
Caring responsibilities	-0.4851	0.203	-2.386	0.017	-0.883	-0.087
Age of cared ones	0.0213	0.019	1.114	0.265	-0.016	0.059
Cyber security role	-0.0058	0.242	-0.024	0.981	-0.480	0.468
Years of experience in company	0.0207	0.041	0.509	0.611	-0.059	0.100
Customer-facing role	-0.4671	0.187	-2.497	0.013	-0.834	-0.100
Work over contracted hours	0.0769	0.153	0.501	0.616	-0.224	0.378
Constant	1.2274	0.563	2.181	0.029	0.125	2.330

from our model in Phase 1 with an independent validation sample from Phase 2 and (2) validate the sample obtained in Phase 2 in order to generalise to the total sample (Phase 1 + Phase 2) the results obtained for the treatments with respect to the segmentation done.

Fisher's exact tests is chosen given the fact that it is used to test the statistical significance for categorical variables, such as the case to validate the proportions of segments between the two phases. The Fisher's exact test is a better fit to our endeavour given the lack of need for a large number of data points, in comparison to the Chi-square test where a large sample is needed [Ross, 2017]. The Kruskal-Wallis, on the other hand, is used to know if two samples come from the same distribution [Siegel and Castellan, 1988]. Both tests have been in the mainstream statistics literature since first published and their use are taken from statistical textbooks [Ross, 2017, Siegel and Castellan, 1988].

In Figure 1.6 we observe two different ways of visualising the scores obtained by the respondents when performing the HaaCSS test. In Figure 1.6a we observe the proportion of correct answers obtained in Phase 2 for the total of the respondents, as for each of the different primings (Neutral, Positive and Negative). A first important result is that the minimum score is 2 correct answers, meaning that all respondents are able to correctly assess at least two of the six situations presented during the

1.4. RESULTS: CASE STUDY

Table 1.5: Determinants of the Ignorant behavioural type. Results obtained performing a Probit regression.

$N = 503$ $R^2 = 0.04628$	coeff	std err	z	$P > z $	[0.025	0.975]
Frequency of working from home	0.0995	0.048	2.088	0.037	0.006	0.193
Confidence in cyber threats detection	0.0010	0.041	0.026	0.980	-0.079	0.081
Frequency of cyber attacks before pandemic	-0.0102	0.077	-0.132	0.895	-0.161	0.141
Frequency of cyber attacks during pandemic	0.0670	0.080	0.839	0.401	-0.089	0.223
Job satisfaction	0.0045	0.040	0.113	0.910	-0.074	0.083
Life satisfaction	-0.0599	0.040	-1.488	0.137	-0.139	0.019
Gender	-0.4137	0.126	-3.282	0.001	-0.661	-0.167
Age	-0.0077	0.037	-0.208	0.836	-0.080	0.065
Relationship status	0.0906	0.127	0.713	0.476	-0.158	0.339
Caring responsibilities	0.3376	0.159	2.119	0.034	0.025	0.650
Age of cared ones	-0.0208	0.015	-1.367	0.171	-0.051	0.009
Cyber security role	0.0515	0.199	0.258	0.796	-0.339	0.442
Years of experience in company	0.0330	0.033	1.004	0.315	-0.031	0.097
Customer-facing role	-0.0668	0.151	-0.444	0.657	-0.362	0.228
Work over contracted hours	-0.2456	0.134	-1.836	0.066	-0.508	0.017
Constant	-0.5791	0.503	-1.151	0.250	-1.565	0.407

test. That said, the average score for the whole of the Phase 2 is 4.12 ± 0.94 correct answers.

The authors in [Heartfield and Loukas, 2015] performed a HaaSS test over different populations around the world and found that the average result for the UK population is 74% of correct answers. Making a direct translation to our own HaaCSS test, a 74% of success would mean an average of 4.44 correct answers over a 6-questions test. Taking this as our benchmark, the average score in Phase 2 of our study is within the range of the found average results for the UK sample, although slightly lower than the benchmark. Given that the main point of this study is to increase the cyber security culture of the company, increasing the average score above the benchmark and having a larger number of respondents having at least 5 or 6 out of 6 correct answers is desirable.¹ In Figure 1.6a we picture the proportion of respondents passing the HaaCSS benchmark in blue, while those that do not in red.

As seen in Figure 1.6a, the smallest proportion of respondents with 5 or 6 correct answers is found for the Positive priming, with 15.2% of the respondents having 5 correct answers and 9.1% having a perfect score (24.3% having a score ≥ 5). This contrasts to the overall proportion of employees having a score higher than 5 for

¹Taking as the ideal case a 100% of respondents having a perfect score, we are interested in minimising the share of people not meeting the benchmark of 5 correct answers.

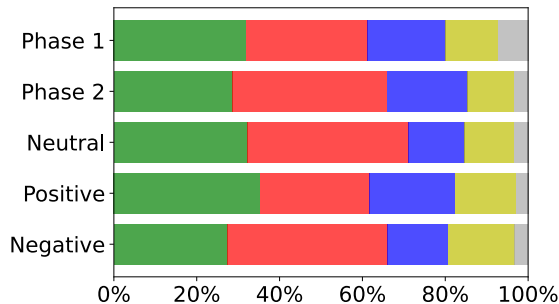


Figure 1.5: Proportion of segments for Phase 1, Phase 2, and the different priming samples obtained in the latter phase.

the Neutral priming (35.1%) and for the Negative priming (36.7%). However, as these two last framings considerably increase the proportion of employees meeting the benchmark, both non-neutral treatments increase the proportion of employees having a perfect score with respect to the neutral one. Indeed, for positive and negative primings we obtain a proportion of 9.1% and 11.7% respectively, while only 1.8% of employees having a perfect score with the neutral priming.

Although the positive treatment decreases the overall share of people having 5 correct answers or more, it increases the share of respondents with perfect score. This last effect is also seen for the Negative treatment, which also preserves the share of employees meeting the benchmark of 5 correct answers, but does not increase it considerably. In that sense, although we can say in a general way that framing the communication in a negative way could have the most beneficial effects for cyber protection of the three cases, the Negative priming does not get us close enough to a desirable state where the majority of employees correctly assess potential cyber threats.

This idea is also supported by the fact that, focusing now in the share of employees that only scores 2 or 3 correct answers, these vary only marginally along the three treatments: 3.5%, 3% and 1.7% for a score of 2, and 21.1%, 24.2% and 25% for a score of 3 (Neutral, Positive and Negative respectively).

Fortunately, the behavioural segmentation analysis allow us to have a deeper analysis in understanding the effect of each priming.

Another way of looking the success of a test is by looking at the fail rate. In this case, wanting to maximise the correct answers, thus maximising the ability of the employees to correctly assess possible cyber security threats, is equal to minimising the incorrect answers. This latter form of looking at the same results is shown in Figure 1.6b. In this case we are showing the number of respondents failing to correctly assess a possible cyber threat with respect to the priming that was

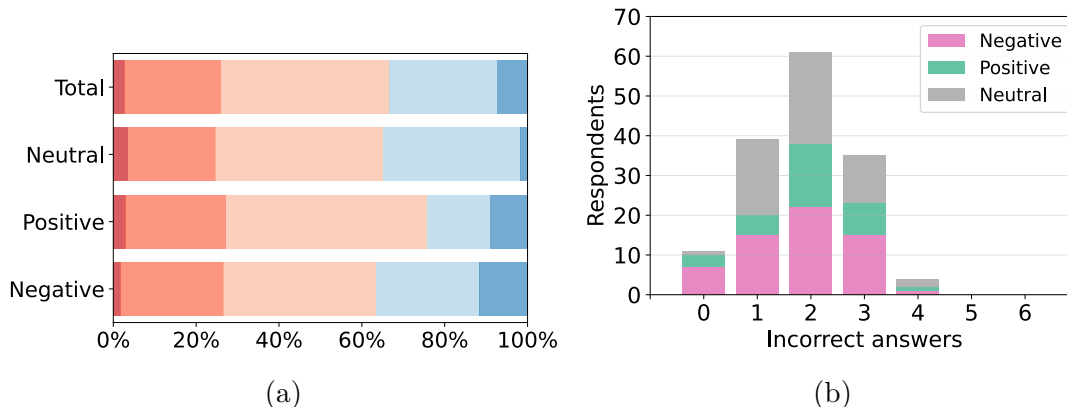


Figure 1.6: (a) Proportion of correct answers by the participants when attempting the HaaCSS test in the second phase. (b) Distribution of incorrect answers divided by primings.

presented to them.

Given the nature of the scores as data (integer numbers from 2 to 6), we are rather interested in knowing how the different primings affect the “likelihood” of failing. In order to obtain such information, we run three different Poisson regressions to the number of incorrect answers of participants with respect to their behavioural segment associated to them. Results are shown in Table 1.6.

As said, we cannot have a one to one correlation between the incorrect answers and the behavioural segmentation, i.e. we are not expecting a correlation such that, given a person belonging to segment x and receiving priming y would have n incorrect answers. What we are rather searching for is a likelihood that a treatment y have an impact (positive or negative) on a person’s incorrect answers, given that the employee belongs to segment x . In that sense, we first run a Poisson regression for the employees receiving the Neutral priming, using as dependent variable their number of incorrect answers. Results are observed in Table 1.6a. The four segments are statistically significant (at a 5% level), which is exactly what we are expecting to give robustness to the impact of the other two primings. At a 1% significance level, the Relaxed segment would not be considered as statistically significant. We take the four coefficients of Table 1.6a as our benchmarks with respect to the other two primings.

A first result from Tables 1.6b and 1.6c is that the Relaxed (low risk perception score and high risk taking score) segment is not statistically significant. Taking into account that this same segment could also be not taken as significant at a 1% level for the Neutral priming, then we can conclude that those employees belonging to the Relaxed segment are not affected by the priming framing presented in our study to correctly assess potential starts of cyber attacks.

1.4. RESULTS: CASE STUDY

Table 1.6: Results of the Poisson regressions the three different treatments applied in the second phase of the survey with respect to the behavioural segments.

Neutral priming						
$N = 57$, log-likelihood=-82.106	coef	std err	z	P> z 	[0.025	0.975]
Anxious	0.6190	0.196	3.156	0.002	0.235	1.003
Ignorant	0.6712	0.149	4.502	0.000	0.379	0.963
Opportunistic	0.6539	0.200	3.270	0.001	0.262	1.046
Relaxed	0.6190	0.277	2.232	0.026	0.075	1.163

(a)

Positive priming						
$N = 33$, log-likelihood=-48.759	coef	std err	z	P> z 	[0.025	0.975]
Anxious	0.8755	0.204	4.289	0.000	0.475	1.276
Ignorant	0.5261	0.213	2.468	0.014	0.108	0.944
Opportunistic	0.6931	0.267	2.594	0.009	0.169	1.217
Relaxed	0.5108	0.447	1.142	0.253	-0.366	1.387

(b)

Negative priming						
$N = 60$, log-likelihood=-88.520	coef	std err	z	P> z 	[0.025	0.975]
Anxious	0.6109	0.169	3.614	0.000	0.280	0.942
Ignorant	0.5664	0.164	3.445	0.001	0.244	0.889
Opportunistic	0.8473	0.218	3.883	0.000	0.420	1.275
Relaxed	0.3102	0.258	1.201	0.230	-0.196	0.816

(c)

However, for the other three segments we observe different behaviours from the results of our regressions. In a first place, the Anxious segment is correlated to make more mistakes when a Positive priming is presented to it (0.8744 ± 0.204 v. 0.6190 ± 0.196). On the other hand, a Negative priming would not create a significant impact on the Anxious segment, lightly lowering their impact on incorrectly assessing potential cyber threats (0.6109 ± 0.169 v. 0.6190 ± 0.196). A possible interpretation of this can be that Anxious segment (high risk perception score, low risk taking score) does not need an external factor to be in an alert state to detect threats. Nevertheless, by presenting a Positive framework, their detection skills decrease, thus being correlated to more mistakes.

Almost the opposite case is observed for the Opportunistic segment (high risk perception score and high risk taking score). By framing the HaaCSS test with a Negative priming, the likelihood to make more mistakes is increased (0.8473 ± 0.218

v. 0.6931 ± 0.267). A Positive priming also increases it, although marginally (0.6931 ± 0.267 v. 0.6539 ± 0.200). Employees belonging to the Opportunistic segment have a relative high score of risk perception, while a high score of risk taking. This can be seen as a person who is aware of the threats, but still engage with them. Both primings, but particularly the Negative, would have a negative impact on their ability to detect otherwise.

Finally, the Ignorant segment is positively affected by both primings. This can be seen by comparing both coefficients for Positive priming (0.5261 ± 0.213) and Negative priming (0.5664 ± 0.164) with respect to the Neutral priming (0.6712 ± 0.149). The Ignorant segment is characterised by a low risk perception score and a low risk taking score, thus being somewhat inattentive of cyber security threats. By simply framing the HaaCSS test with a non-neutral priming, the employees decreased their likelihood to incorrectly detect potential cyber threats.

General results for each of the 6 questions of the HaaCSS test are found in Figure 1.7. We remind that more in-depth results and observations are found in Appendix C.

1.4.3 Work from home habits and cyber security hygiene during the COVID-19 pandemic

During the two phases of the study we were interested in knowing the working habits during the COVID-19 pandemic. We set as starting date March 23, 2020, which is the day the UK government first announced a set of mobility restrictions due to the COVID-19 spread in the UK. Particularly, during the first phase of our study we asked different questions regarding working habits during the pandemic and potential cyber attacks detected before and after it. We present the respective answers to the questions regarding these two topics.

Most of the questions presented in what follows were asked in both phases of the survey. Given the results shown in Section 1.4.2, we are uniting both samples of Phase 1 and Phase 2. In case a particular question was only asked during Phase 1, it will be specifically noted. Most of the questions cited resulted to be statistically significant when included in the probit regression done to characterise the behavioural segments. Thus, in addition to have an insight about the working habits and cyber hygiene during the pandemic, the following results are also intended to give a more robust image to the results presented in Sections 1.4.1 and 1.4.2.

Given the accounted literature [Mills et al., 2021, Bhugra, 2021] around the effects of the COVID-19 measures on the wellbeing of the population in the UK, we asked the surveyed population to indicate their levels of general satisfaction with

1.4. RESULTS: CASE STUDY

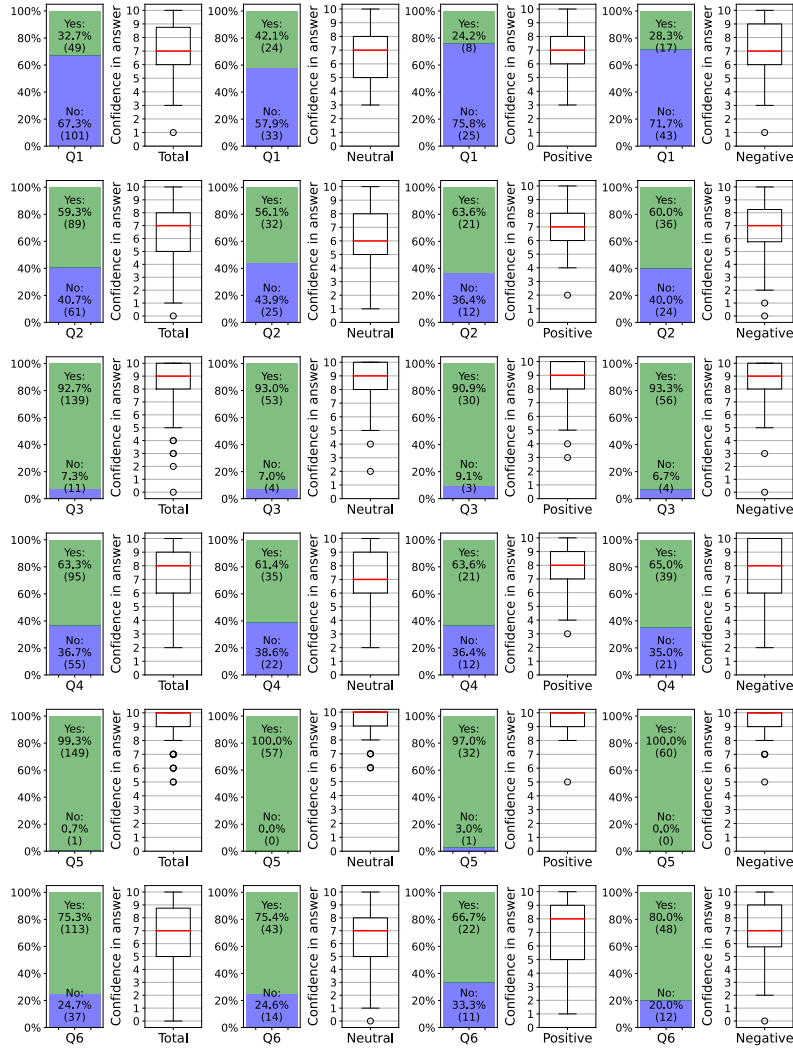


Figure 1.7: Results for each of the six questions of the HaaCSS test, showing the proportion of people answering if the presented screenshot was a potential start of a cyber attack, and the confidence of their answer. Results are shown for the total sample of Phase 2, as for each of the different treatments given to the respondents.

respect to their life and job. Results are shown in Figure 1.8. Results are almost identical, obtaining an average satisfaction of life of 7.51 ± 1.91 and an average job satisfaction of 7.72 ± 1.88 .

With respect to particular working habits during the pandemic, we asked three questions related to the number of days working from home per week, and the number of contracted hours against the number of actual working hours. Results are shown in Figure 1.9 and Figure 1.10. As expected, the great majority of the employees work 5 days a week from home. This majority is expected, as both surveys were done during periods of time where mobility restrictions were imposed in the UK. In particular, during both periods of time the UK government suggested offices and business to do not open and impose work-from-home policies. We also observe

1.4. RESULTS: CASE STUDY

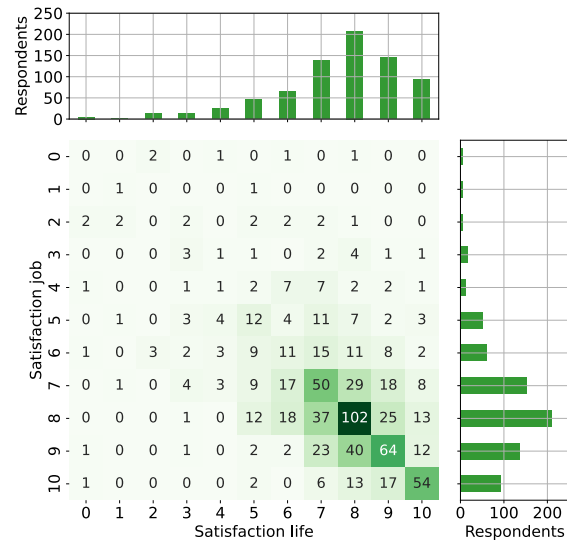


Figure 1.8: Heatmap matrix showing the different distributions of the satisfaction of life and work for the respondents of both surveys

an important number of employees who overwork with respect to their contracted hours.

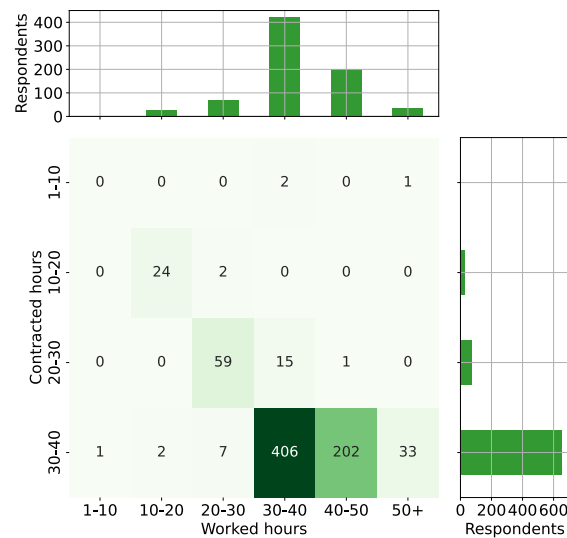


Figure 1.9: Heatmap matrix showing the different distributions of the contracted and actual worked hours for the respondents of both surveys

The previous results only refer to the working habits of the surveyed employees. In Figure 1.11 we observe the distribution of the confidence in detecting cyber attacks indicated by the employees.

During Phase 1, we asked two particular questions to the employees to indicate the frequency and the type of cyber attacks they have experiences before and during

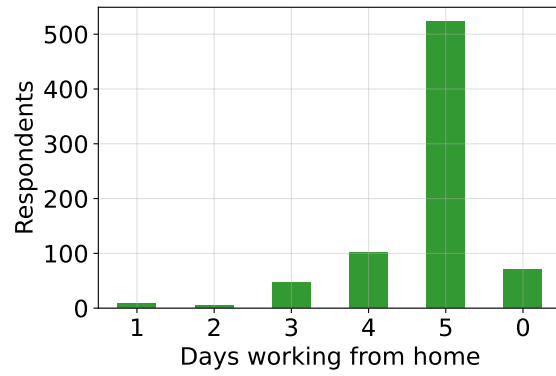


Figure 1.10: Heatmap matrix showing the different distributions of the satisfaction of life and work for the respondents of both surveys



Figure 1.11: Levels of indicated confidence to detect cyber attacks by the respondents of both phases.

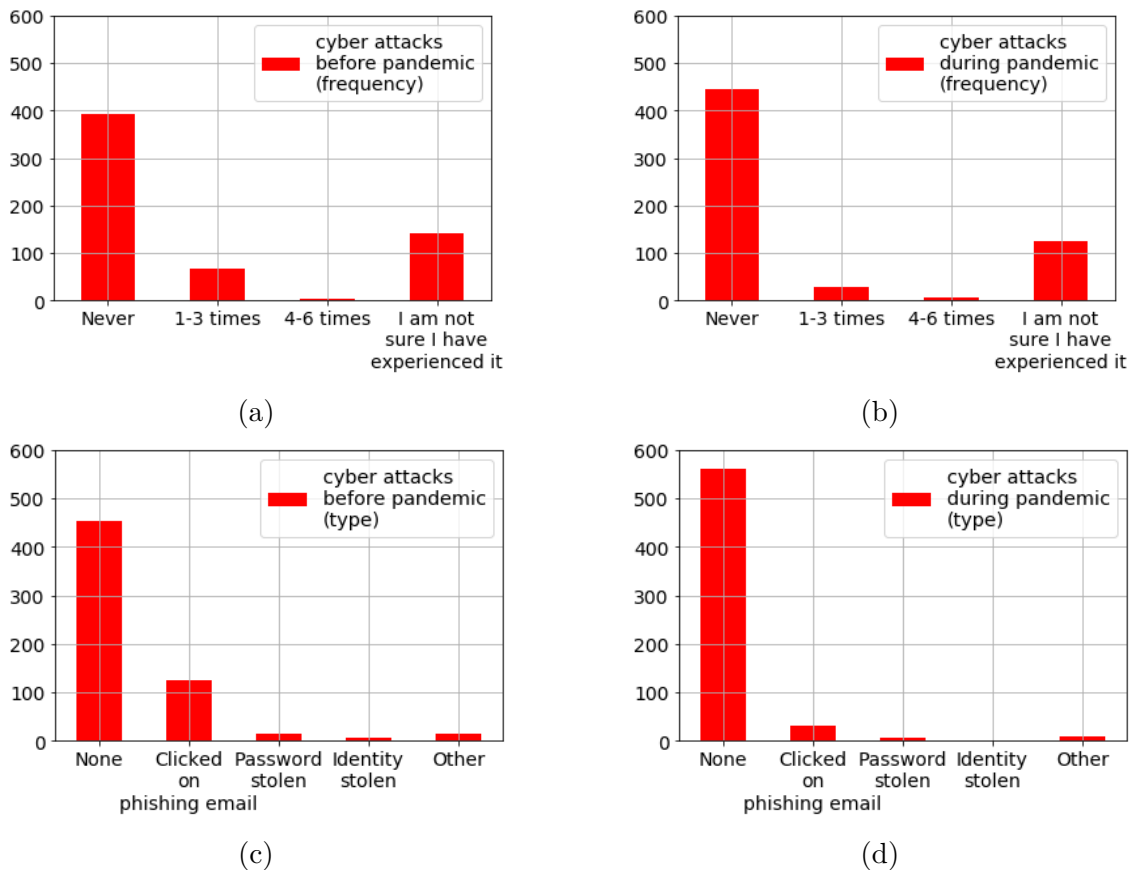


Figure 1.12: Answers to the frequency and type of cyber attacks detected by the respondents of Phase 1.

the pandemic. Results are shown in Figure 1.12.

1.5 Discussion

1.5.1 Framed communication strategies using a segmentation strategy

Taking the results of Phase 1 and Phase 2, along with the combined results from the behavioural segmentation exercise and the framed HaaCSS test, we have comprehensive grounds to create a disaggregated strategy to engage with employees in an efficient way. Indeed, a framed communication around cyber security can have different effects on a given heterogeneous workforce. While a negatively framed communication can motivate a high state of alert on its receivers [Rege et al., 2020, Yeoh et al., 2021], we tested if a positive framing can increase the confidence of employees to detect cyber attacks. However, as seen in the different results of Section 1.4, the effects of each one of the three treatments are not homogeneous within the surveyed

sample and rather have specific effects on each of the 4 behavioural segments. Given that we could obtain different demographic elements and working patterns for each of the 4 segments, we can thus devise different strategies in order to increase the likelihood of the surveyed employees to correctly assess potential cyber threats.

A Neutral (non-framed) communication is appropriate for Opportunistic types. A negative priming considerably increases the risk of making mistakes, while a positive priming does the same to a lesser degree. Considering that an Opportunistic employee is someone that has a high score of risk perception, thus being capable of potentially detecting cyber threats, and a high score of risk taking, we could assume that people belonging to this segment understand the risks they take when they engage with potential cyber threats. Reminding them of the potential consequences of this, by either a positive or a negative framing, could result in undesirable additional stress that increases their likelihood to falsely assess risky situations.

According to our results, the Opportunistic segment is characterised by young of age, a low self-perceived confidence for detecting cyber threats and its members tend to identify themselves as men. This segment is also characterised by an excess of working hours with respect to their contracted ones.

In order to decrease the likelihood of the Relaxed population to mistakenly assess potential cyber threats, further strategies should be explored. None of the framing was found statistically significant at 1% level. Like the Opportunistic segment, the Relaxed segment is characterised by a low self-perceived confidence for detecting cyber threats. In addition, employees belonging to this segment tend to identify themselves as men, have a low satisfaction with their current jobs, no caring responsibilities and are not in customer-facing roles.

A Negative communication is slightly beneficial for the Anxious segment, in contrast to positive priming which increases the likelihood of this segment to commit mistakes. Anxious employees are characterised by a high perception score and low risk taking score, which could be interpreted as having a state of high-alert towards cyber attacks. A positive priming could decrease their confidence or state of alert, while a negative framing could increase their awareness towards possible threats. Employees belonging to this segment are of older age and tend to have a high level of self-perceived confidence for detecting cyber threats.

Finally, any kind of non-neutral priming has beneficial effects on the Ignorant segment. This segment is defined by low level of risk perception and low level of risk taking. Thus, it makes sense to increase awareness of cyber attacks among employees in this segment. Whether a priming is positive or negative, its implementation decreases the likelihood of incorrectly assessing potential cyber attacks. “Ignorant”

employees are characterised by a high frequency of working from home, having caring responsibilities and tend to identify themselves as men.

Our results complement the existing literature that intends to create different strategies around how to better train employees with respect to cyber security in companies. Given that most of the employees recognised in the survey a phishing attack as the principal source of threat and that the results for the phishing questions in our HaaCSS tests were the most accurately answered, our results can also extend to the literature in this particular cyber threat. Our results show that, as stated in [Kwak et al., 2020] and [Yeoh et al., 2021], a reinforcement of the message indeed increases the awareness and, in our case, also the probability to correctly detect a cyber attack of this nature. However, our study shows how the fact that positive communication around cyber threats can *also* have positive results, while negative communication can in some cases have a non-desirable effect.

Our study also shows how, with respect to other training/awareness models [Rege et al., 2020, Sillanpää and Hautamäki, 2020, Marotta and Pearlson, 2019], the behaviour of our sample can be modified when applying our framework. This allows the interested company to easily detect changes in behaviour to then calibrate its communication strategies to better train the workforce.

1.6 Conclusions

For the present work we collaborate with a large financial company in the UK to, through a two phased survey, study a novel way to understand the culture and attitude around cyber security of their workforce. This study is done using a behavioural segmentation model based on the CyberDoSpeRT scale [Kharlamov et al., 2018] which allows to categorise the employees in four different profiles. Using the results of the latter segmentation, we are also able to study the effects of a Human-as-a-Cyber-Security-Sensor test [Heartfield and Loukas, 2018] on the employees with different communication framings. We are interested in testing if humans are the weakest link in cyber security sensors [Yeoh et al., 2021, Rege et al., 2020], or actually can work as very good threats detectors [Heartfield and Loukas, 2015, Heartfield et al., 2016, Heartfield and Loukas, 2018]. We are also trying to research if a positive communication highlighting the benefits of a cyber security prevention has any positive effects on security systems, which is opposed to a categorical vision where a punitive communication is preferred [Yeoh et al., 2021].

With the results of the first phase of the survey we are able to calibrate our behavioural model with four different behavioural segments. We obtain 605 different

answers. Each of the segments is obtained by computing a risk perception score and a risk taking score for each participant. A k -means++ [Arthur and Vassilvitskii, 2007] algorithm is then applied to the sample of scores, with $k = 4$. The different segments are: Anxious (high risk perception score and low risk taking score), Opportunistic (high risk perception score and high risk taking score), Relaxed (low risk perception score and high risk taking score) and Ignorant (low risk perception score and low risk taking score). Each of the four segments are studied using a probit regression analysis using demographic and work-related data from the 605 respondents of the survey.

The second survey (150 respondents) allows us to partially validate our behavioural model. This second sample does not have any statistically significant difference with respect to the proportions of the four different segments with respect to the first sample. It was during this second phase that each of the respondents underwent a HaaCSS test. For each respondent, the test was framed by one of three different treatments: neutral (non-framed), positive (emphasising on the benefits of cyber security) and negative (emphasising on the consequences of cyber security).

By merging both analyses, we move beyond generalising conclusions about the different primings, e.g. the negative priming is more beneficial to the workforce of the company. Instead, we are able to obtain a disaggregated strategy of communication that maximises the correct assessment of potential cyber threats according to different types of employees.

While a neutral approach works best for “Opportunistic” employees, any kind of framing has a positive effect on “Ignorant” employees. The former might be because of the conjunction of a high perception of risk and a high risk taking attitude (thus making any reminder of the risk taken to have a negative effect). The latter might have as a cause the low risk perception and the low risk engagement, thus allowing us to think of this segment as passive with respect to cyber security and positively reacting to any stimulus. The “Anxious” segment reacted negatively to the Positive framing, while marginally good to the Negative framing. Finally, we do not find any statistically significant effect from our framing strategy on the “Relaxed” segment.

The results prove two of our hypotheses as correct. The Negative treatment does not work homogeneously beneficial with all of the workforce, and for some might result counterproductive. On the other hand, the Positive treatment works also heterogeneously within the workforce.

By assessing the employee’s segment, or by using the demographic and work-related variables that are statistically significant to each segment (presented in Section 1.4), the company can then create different strategies to different audiences to tackle particular behaviours from their workforce.

Our study also does not find any particular change of treat levels or behaviour during the pandemic, then dismissing our first hypothesis (as of June 2020 when the first survey was launched).

This work researches the human side of cyber security, going farther than creating more awareness in a given workforce of a company, to actually change the behaviour of it to make them a better sensor of cyber threats. We are going away from the perspective of humans as the “weakest link” in a cyber security system, to actually empower the day-to-day employees to correctly assess possible threats, regardless of their position in the company.

We contribute to the literature of user-centred training for employees with respect to cyber security. The main contribution is to provide a model that can be scaled to companies with a large heterogeneous workforce, obtaining almost quick results in the behaviour of the sampled population, thus allowing the companies to quickly calibrate its communication strategy to obtain the best results as desired.

Also, our study proofs that the categorical affirmations that punitive communication is best when training employees [Rege et al., 2020, Yeoh et al., 2021] is not necessarily true. Our results show that, depending on how risk averse and how an employee can perceive risk, a given communication can work best than other when being trained to detect cyber threats. This in itself is useful, as it recognises that when designing a training in an important matter such as cyber security, the different behaviours that employees can have towards it must be taken into account.

The uniqueness of our work, in addition to the above mentioned, lies in the application of a *statistical learning pipeline* in order to analyse and test our hypotheses. As such, the fact of using a first survey to *train* our behavioural model to then *test* it with the data obtained from a second survey is indeed novel, and as far as the author knows, no other work has been done using this methodology.

Our work can have very practical implications for different companies whose workforce, as said, have different heterogeneous and volumes. By deploying a series of independent surveys and applying simple data science and statistical methods, the companies can try different communication strategies to their employees. This method could potentially be extended to other domains, although we keep our results to cyber security training.

An immediate implication is the understanding of security issues at individual and community levels. On one hand, our work reveals how there is not a unique and general way to better train individuals for security tasks. Thus opening a field of different ways to better train security at a personal level taking different methods

out of the established one. On the other hand, by incorporating methods from data science and statistical learning, we can then analyse those heterogeneities to a community level and thus create a more secure one. Looking at this same point from a different perspective, when talking about security at a community level, studies should now take into account the heterogeneity of it to implement different ways so individuals can be better trained into the needed tasks.

As further work, we are interested in implementing a potential third phase of the study, which is applying a HaaCSS test with the disaggregated strategy of communication. This would mean that a particular framed test would come to any participant, depending on the behavioural segment in which they would belong after computing their risk perception and risk taking scores. Also, novel strategies need to be designed in order to positively affect the Relaxed segment to have a more adequate response to potential cyber threats.

The present work belongs to the literature studying cyber security from the behavioural aspect of the users. We combine methods of data science, econometrics and behavioural science in order to help companies and organisations (in this case the Large Financial Company) to improve their cyber security putting an emphasis on their workforce.

Chapter 2

Understanding and nowcasting the illicit drug distribution in England: a data-centric approach to the County Lines Model

Chapter Abstract

The County Lines Model (CLM) is a relatively new illicit drugs distribution method found in Great Britain. The CLM has brought modern slavery and public health issues, while challenging the law-enforcement capacity to act, as coordination between different local police forces is necessary. In the present work we analyse the CLM using a data-driven approach for the first time. Our objective is to understand the territorial logic behind the line operators when establishing a connection between two places. We analyse three different spatial models (gravity, radiation and retail models), as each one of them understands flow from place i to j in a different way. Using public data from the Metropolitan Police of London, we train and cross-validate the models to understand which of the different physical and socio-demographic variables are considered when establishing a connection. We analyse hospital admissions by drugs, disposable household income, police presence and knife crime events, in addition to the population of a particular place and the distance and travel times between two different. Our results show that knife crime events and hospital admissions by misuse of drugs are the most important variables. Together, both elements could be interpreted as the operators avoiding conflict with other gangs and crowded markets by competing organisations. We also find that London operators distribute to the territory known as the “South” of England, as little presence of them is observed outside of it.

2.1 Introduction

During the last decade, a new illicit drugs distribution model has been developed in the UK. The model was baptised as the "County Lines Model" (CLM) by the UK government [Crime Agency, 2019] given its use of phone lines that are established between different counties. The problem has become increasingly worrying each year, becoming a top priority for security agencies given the limited ability to stop them, and the modern slavery and public health problems that the CLM brings to local communities [Crime Agency, 2019, Black, 2020a, Andell and Pitts, 2018, Robinson, 2019, Stone, 2018, Camber, 2020, Moyle et al., 2019].

The *modus operandi* can be described in the following way: a central hub is settled in big English cities like London, Birmingham, Manchester and Liverpool, from where drugs are sold and distributed [Coombes, 2018]. From these hubs, *lines* are settled to other parts of the country where a local market is established. So-called *settlers*, find a local accommodation (normally a flat belonging to consumers) in the destination market from which drugs can be distributed. Local *runners* are then hired to distribute the illicit merchandise to the consumers. Runners tend to be young people with knowledge of the local market whose tasks are to deliver merchandise and attract new clients. The distribution model increases the efficiency with respect to "the traditional model" [Coombes, 2018] where the "highstreet" illicit drug seller buys merchandise to a bigger distributor, to then sell it on the street. The improvement of the CLM is to merge both tiers (local and bigger seller) uniting both channels of distribution (hub-settler and runner-consumer).

Local consumers are given a phone number where they can place an order. The call is normally picked up in the central hub, from where they make the arrangements to distribute it to the consumer via the settler and the runner. The settler travels back and forth from the central hub and the local market bringing merchandise, while the runner distributes to the final consumer.

According to the National County Lines Coordination Centre [Silver and Intelligence, 2021], 3 cities account for more than 80% of the detected County lines in Great Britain in 2019 and 2020. These are, in respective importance, London, Birmingham and Liverpool. Public data is scarce, only having detailed records for London for those two years.

The implications of the proliferation of the CLM in the UK are multiple. Three are particularly highlighted in the literature [Andell and Pitts, 2018, Robinson, 2019, Black, 2020b]: (1) the rising of new illicit drug markets in small coastal towns and rural areas of England where illicit drugs problem were not found before. (2) Also, the involvement of young vulnerable people in the distribution scheme is of

worrisome for the UK Gov. This population is the most prone to be caught by law-enforcement bodies, while being involved in a modern slavery scheme making them hard to leave the CLM once they are involved. (3) Finally, a limited ability from the different police forces in England to dismantle any complete distribution channel between one place and another. Cooperation between different law-enforcement bodies is necessary, as every link in the distribution engine can work autonomously, making it hard for bodies to dismantle the whole distribution operation.

The fact that county line operators are found in small villages and coastal towns, far from local capitals and larger cities has risen different hypotheses about the logic behind establishing a *line*. Indeed, the size of the population of the target places seems not to be a primordial element, as large population centres (London, Manchester, etc.) do not attract a big number of lines according to public data shown in the strategic report from the National County Lines Council [Silver and Intelligence, 2021]. According to the same report, the logic behind the gangs operating county lines is a supply-demand balance.

The main objective of this work is to understand the territorial logic that county line operators follow to establish different distribution routes. To do so, we have to answer if the “traditional distribution model” has been broken as literature suggests [Stone, 2018], and if so, which are the new social, demographic and economic elements that are now taken into account to establish a new route. To answer both questions would help to obtain useful information for the Metropolitan Police to understand and tackle the county lines problem.

In order to answer these questions, we test three different spatial interactions models to compute flows from one place i to a second place j . We understand each of these models as different ways to understand the flow of persons/merchandise. Thus, by testing and comparing them we can extract information about which mechanisms could county lines operators follow. The models we use are the Gravity Model [Anderson, 2010], the Radiation Model [Simini et al., 2012] and the Retail Model [Wilson, 2008].

We use the classic gravity model as our benchmark, as it understands the flow from one place to another as proportional to the respective populations and inversely proportional to the distance between both places.

Radiation model understands flows as a process of sorting the available opportunities between i and j . To arrive to place j , the studied element (person/merchandise) should not be captured by the opportunities found in the way to it.

Finally, the retail model understands flow as a balance between the opportunities and the costs of going from one place to another, compared with all the other competing places in the given space. This latter model allows to test other kind of

dynamics involving different benefits and costs while considering competition too. We explore five different independent variables we expect to have some leverage for operators. These are knife crime events, number of police officers, gross disposable household income, hospital admissions by misuse and poisoning by drugs as possible costs.

The hospital admissions are taken as proxies for illicit drugs consumption, as no other data is available. Knife crime events are another high-priority incidents for UK Government [Bellis et al., 2019], which are reported to be related to gang rivalry. We are interested to see if the presence of this kind of event could be an element taken into account for operators as a disincentive for establishing a local market. In the same way, we are expecting police workforce to be a disincentive for gangs. Finally, we take the gross disposable household income as a measure of richness, as average income does not take into account regional disparities in rent prices, money transfers from the government and local taxes. We train and test the three models with public data from the Metropolitan Police of London [Rescue and Analysts, 2019] accounting the detected lines in other police force territories in Great Britain from London in 2019 and 2020.

This project is also found in the current context of need for better information for law-enforcement bodies in the UK, as there is an ongoing discussion about how Brexit and the COVID-19 pandemic will have a major effect on public spending, particularly in law enforcement bodies and the National Health Services (NHS), the public health body in the UK) [Roman-Urrestarazu et al., 2018]. In particular, reports state historical maximum numbers of drug-related deaths *per capita*, as a new generation of young consumers enters the market and an older generation requires more health care services [Black, 2020a]. Also, it has been discussed how Brexit would make more difficult for the United Kingdom to access and profit from European funding and infrastructure (like the European Monitoring Centre for Drugs and Drugs Addiction, EMCDDA) for better intel and tackling strategies for a better public health and general quality of life for its citizens [Coombes, 2018].

In the following, we present a literature review in Section 2.2. In this section we also present the literature gap found and the hypotheses we work with. The different models and data tested in Section 2.3. Results are presented in Section 2.4, to then discuss and conclude in Section 2.5. We also present two Appendices D and E. The former is a table to help the reader with the models tested, while the latter details the different sources and formats of the data used in this work. To the authors' knowledge, this is the first published work that studies the County Lines Model from a quantitative approach.

2.2 Literature review

A literature review is done to document the work done until now in the scope of this work. Firstly, we describe the literature that has been published around the County Lines Model. Then, we review some of the quantitative methods that might be applied to the kind of problem we are facing, such as network science, econometrics and others. Finally, we present literature around the topic of illicit drugs distribution and consumption in the UK to understand the context of this work. At the end of this section we discuss the literature gap found and how from it we can tackle the research questions of this project.

2.2.1 County Lines literature

In the case of the County Lines Model, only qualitative or official literature has been published. The Official literature includes documents and reports from different police agencies and the UK government. In particular, the NCA has published each year a statement regarding the views of the organisation about the County Lines Model [Crime Agency, 2019]. The document documents the findings from the NCA to understand the model and the different consequence it has had in the population.

In 2019, the UK Government's Home Office commissioned an up-to-date report to be done around the illicit drugs problem in the UK. The report was published in early 2020 [Black, 2020a, Black, 2020b] and reveals how the County Lines Model has evolved over the last decade. It also reports how the consumption of illicit drugs has changed in the population, stating that the UK faces an important challenge, as there currently are two peaks of consumers: one in their 20's and another in their 60's. Each one of those is of increasing worrisome, as the first one is the future workforce of the UK and the second represents an increasing pressure in the public services.

Two different police organisations have publicly published information about the County Lines Model information they have. These are the Metropolitan Police of London [Rescue and Analysts, 2019, Rescue and Analysts, 2010] and the West Midlands Police (Birmingham and metropolitan area) [Supt Mat, 2020]. Only the Metropolitan Police has published quantitative data about their detection of lines in other police territories.

In January 2018, a debate was held in the House of Commons (UK's lower parliamentary chamber) to discuss the exploitation and harms done by the County Lines Model in London [Pepin, 2018]. Different Members of the Parliament asked what has been done until that point to tackle the CLM problems in London, particularly

gang activity and exploitation.

Outside official documentation, academic literature about County Lines has mostly dedicated to report the child exploitation in different locations of England [Andell and Pitts, 2018, Robinson, 2019, Stone, 2018] and Scotland [Madeley, 2018]. In all of them we find a description of the model. An anthropological study can be found in [Coomber and Moyle, 2017], where the authors interview different consumers and victims of the CLM in South England.

Given the nature of the data by the Metropolitan Police that is used for this study, we are interested in knowing the different methods involving networks and methods to analyse flows/transportation from one location to another.

2.2.2 Social Network Analysis

Given the lack of information, we cannot perform a Social Network Analysis, simply because we do not have any information about the elements of the organisation. We rather have the amount of lines that have been found by the Metropolitan Police. Even so, we do a small literature review of what has been done in the past years.

Criminal networks analysis started to gain increasing popularity in criminology in the 90's [Campana, 2016]. Originally defined as a specific and fixed form of organisation [Podolny and Page, 1998], it slowly shifted to a more “instrumental” perspective [Campana, 2016] where it is seen as a fluid and complex structure allowing to understand a criminal organisation [Rostami and Mondani, 2015].

The more standard definition is to see a criminal network as a set of individuals (nodes) having a set of relations (links). In that sense, the criminal network can be defined by different kinds of participants, like offenders, consultants and *ad hoc* workers, cooperating together with different kinds of relationships. These relationships, implemented as weights, can cover different types of cooperation between two nodes, like advisory, subordination or equality of ranks [Campana, 2016]. However, different works have also studied more subjective attributes like trust between two individuals in the criminal network [Bright et al., 2019]. This has come to the establishment of analysing Crime Networks as Social Networks. Most of the literature found for the literature review uses Social Networks Analysis as its framework and tool set to understand these criminal organisations [Knoke, 2015, Campana, 2016, Bichler et al., 2017].

When characterising the edges and nodes, network analysis allows to know more about the structure of the criminal network, knowing its dimension and hierarchies (small gang, organised gang, organised crime structure) [Canter, 2004], its

efficiency and its robustness/security [Morselli, 2007]. To measure these three network's attributes, studies normally use the density of the network, and different centrality measures of the nodes [Bichler et al., 2017, Berlusconi, 2017, Bright et al., 2019]. In particular, authors use the betweenness centrality as a measure for security/efficiency. Having a relatively small number of nodes with high betweenness centrality would mean a very efficient process but also a very centralised one, where taking down the strategic nodes would mean a disaster for the network. On the other hand, a more decentralised network would mean a higher robustness, as a larger number of nodes would be necessary to take down in order to make the network inoperable. We remind that, given a network with nodes $i = 1, \dots, N$, then the betweenness centrality of node i is defined as

$$B_i = \sum_{j \neq k \neq i} \frac{\sigma_i(j, k)}{\sigma(j, k)}, \quad (2.1)$$

where $\sigma(j, k)$ is the total number of shortest paths going from node j to node k , and $\sigma_i(j, k)$ is the number of shortest paths going from node j to node k *passing through* i . In that sense, betweenness centrality is a measure of how important is a given node to pass information within the network.

It is here where we can also distinguish between two different broad topics in criminology literature: terrorist networks and criminal networks.

Terrorist networks became particularly relevant in academic literature after the 9/11 terrorists attacks [Krebs, 2002, Lum et al., 2006]. Criminologists study the different individuals involved in these attacks and their different dynamics. The biggest difference between these two networks is the driving motif of the organisation to stay formed. While in a non-terrorist criminal network the driving force is economic and social profit, in terrorist networks the force gravitates more around ideology [Krebs, 2002]. This elemental difference has important consequences in how the networks are structured. On the one hand, terrorist organisations care more about the robustness or security of their network. That is, if one or more elements fall down, the whole networks must not be exposed to fall down altogether. On the other hand, criminal organisations, by prioritising the economic profit of its network, their network would compromise robustness to increase their efficiency [Morselli, 2007]. However, this last point is still highly debated in literature, as some authors still present arguments and cases where criminal networks opt for security over efficiency [Berlusconi, 2017, Berlusconi et al., 2016].

The SNA found in the literature tends to do not incorporate the temporal dimen-

sion. Networks are studied as snapshots in time, working with the nodes and links at a certain date. Only a handful of papers [Bright et al., 2019, Berlusconi, 2017] have been found in which they incorporate different snapshots to the study. This kind of temporal study allows to better understand the key players and key relations that allow a criminal network to survive, while also understanding its fluidity and adaptation to internal changes [Rostami and Mondani, 2015].

While most of the criminal networks studies focus in the organisation, formation and evolution, the literature involving the different illicit activities into place is smaller. Instead of focusing in networks as organisations of different sizes with individuals as nodes, these studies focus on geographical networks where nodes are more organised gangs in different geographical locations [Berlusconi, 2017, Giommoni et al., 2017, Berlusconi et al., 2016, Dolliver et al., 2018]. The intention of the studies, located in an intersection between quantitative geography and economics, is to study the possible inflows and paths of illicit merchandise from one place to another tracking the price of the chosen merchandise. These can be MDMA in the US [Chandra et al., 2017] or heroin in Europe [Chandra, 2015].

2.2.3 Criminology and other domains

Criminology is the domain which has used a more diverse set of tools in order to analyse Crime Networks. In [D’Orsogna and Perc, 2015], an important review has been done compiling the different approaches taken in the last two decades. We can split these approaches into three different categories: Agent based models (ABMs) [Hegemann et al., 2011], analytical models [Short et al., 2010, Caminha et al., 2017] and network models [Bright et al., 2019, Berlusconi, 2017].

Analytical models try to understand the formation and dissipation of criminal hot spots and other phenomena using PDEs [Short et al., 2010]. ABMs and network models on the other hand have been developed to study the formation of criminal networks and the rehabilitation and recidivism of their members [Berlusconi, 2017].

While analytical methods do not incorporate real data from the criminal networks or their locations, ABMs and networks have been including more data from recorded cases in the US, Mexico and different European countries [Chandra et al., 2017, Chandra, 2015, Espinal-Enríquez and Larralde, 2015].

A different but important method to note that has been used in the last decades is the study of the water sewage systems in different locations of the world. In particular, the European Monitoring Centre for Drugs and Drugs Addiction publishes

every year the results obtained for different cities in Europe. The publication allows to monitor different drugs like cocaine, heroin and MDMA [European Monitoring Centre for Drugs and Drugs Addiction, 2020]. Normally, London appears in the top cities with cocaine and heroin consumption, although other British cities have appeared in the top rankings, such as Bristol [European Monitoring Centre for Drugs and Drugs Addiction, 2018].

A notable element in most of the papers read for this literature review is the lack of an exogenous/prior context or data which would explain the criminal networks. In that sense, all of the criminal networks, social or geographical, are supposed already formed and static. Only one article [Bright et al., 2019] was found in which the authors give a prior explanation about the context in which the analysed criminal social network was formed. However, this external information is only qualitative and does not have any data which can be used as input into the quantitative analysis, even though the necessity of it has been mentioned as a key element for analysing and predicting crime [Hardyns and Rummens, 2018].

2.2.4 Literature gap and hypotheses

The literature found for illicit drug networks in the UK studies the inflow of drugs and the organisation of the criminal networks from an international perspective, always studying the UK parallel to Europe, the US [Dolliver et al., 2018] or South Asia [Ruggiero and Khan, 2007].

In the case of the County Lines model, no quantitative analysis has been made. This is because of different elements: not enough data about gangs and the different routes used is publicly available; there is no data about the capacity and the volume of illicit drugs that these gangs are able to transport; there is no information about what is the share of sells in the illicit market attributed to the model.

Given the lack of information, we cannot perform a Social Network Analysis, simply because we do not have any information about the elements of the organisation. In that sense, we cannot approach the problem in a similar way as those used for the inference of inflow paths by price analysis [Berlusconi, 2017, Chandra et al., 2017, Chandra, 2015, Boivin, 2014], where specific information about the organisations is omitted, only supposing that the gangs/cartels exist and have enough capacity to move merchandise from one place to another and consumers which would buy it at the objective location. The problem we have with this latter approach is that we could not suppose which are the capacities for other important exporters

outside London.

The water sewage approach, although very interesting and conclusive, cannot be used in this framework given the specifications of the county lines model, tending to open local retailers in deprived and small towns, where the water sewage analysis has never been done. Indeed this particular analysis is only done in big cities in Europe, and does not have any data for more smaller and rural places in England.

The official documentation related to the County Lines Model reviewed in this Section shows different particular elements: (1) a shift from the traditional behaviour of the illicit drugs distribution organisations in England towards a more organised structure, (2) how much these new ways of distributing illicit drugs has had an important negative effect on local populations, and (3) the emergence and difficulty that Police forces in England have to tackle the problem that CLM poses. Point (1) and (2) have been supported by the qualitative studies found in the literature, showing modern slavery and decreased health conditions in the affected communities, while also understanding socially and logistically how CLM organisations work. However, these studies, given its framing and objectives, do not have any proposition that could directly benefit the law-enforcement bodies in a general way to tackle the problem in the UK.

By taking the official documentation and the qualitative studies as foundation to understand the needs of the police and the social structure of the CLM, and with help of the available data from the Metropolitan Police, we found that there is a literature gap which quantitatively studies the CLM to provide a general understanding of the problem and propose solutions or insights to the law-enforcement bodies.

A large number of network models can be adapted to our methodology. However, the work in [Piovani et al., 2018] has proven to be the most adaptable, given that the analysis can incorporate different variables in a non-linear way, and by being derived from different constraints with respect to demand and spending power. The different models presented there and that are introduced in a more detailed form in Section 2.3 has been used before for retail systems and for the London riots in 2012 [Davies et al., 2013], for human mobility studies [Yang et al., 2014, Simini et al., 2012], and more recently in exploring ancient trading routes in bronze-age Mesopotamia [Barjamovic et al., 2019], making them highly adaptable models.

From the literature cited above, we can also outline the different hypotheses we work with throughout this work. By the different official and academic literature around county lines [Black, 2020b, Coombes, 2018, Andell and Pitts, 2018, Robinson, 2019, Stone, 2018], we can suppose that

H1: county lines model does not take population and distance as a primordial

element to establish a connection

Thus making the traditional gravity model used here as the expected worse model to perform. On the other hand of the scale, given that the Retail Model allows to have as input different social and demographic models, we thus have as second hypothesis that:

H2: given the nature of the model, the Retail model will be the best performing of all.

Also, given the literature around county lines, we can expect that from the different variables used in the Retail model (disposable income, police workforce, knife crime events, hospital admissions by misuse and hospital admissions by poisoning of drugs).

H3: the gross dispensable household income is an incentive for operators while the other four are costs for the operators.

2.3 Methodology

The data obtained from the Metropolitan Police of London allows only to implement a handful of methods that were discussed in the Literature review in Section 2.2. Indeed, the territorial resolution used by the Metropolitan Police (police force territory) and the information extracted forces us to talk in terms of *flows*. We cannot speak of a flow of persons, but rather a number of detected *lines* (connections) established from a place i to another place j . In that sense, the data point is a natural number, T_{ij}^{data} , representing the detected connections.

The spatial resolution we work with is at police force territory, which in Great Britain account for 39 in England, 5 in Wales and 1 in Scotland. In our case, we work with the 39 territories in England only to train our models. We only train for England as not all features used in the models are available for the whole of Great Britain. We merge both territories in Greater London (Metropolitan Police + City of London Police) to work with London as a unique space.

To properly analyse the data obtained, answer the research questions we propose in Section 2.1 while also filling the literature gap detected in Section 2.2, and finally being able to translate our findings in policy proposals, we need to find appropriate methods and an analysis pipeline that incorporates the specific details of the available data.

2.3. METHODOLOGY

From the literature reviewed in Section 2.2, particularly in [Piovani et al., 2018], we find three different models that allow us to analyse the Metropolitan Police data in an appropriate way. These are the Gravity, the Radiation and the Retail model.

The three models allow to create an analysis pipeline in which they are trained using the available input and output data, while validating and assessing each one of the models. The pipeline is written and developed using the Python package `Pytorch`. In order to validate and assess each one of the different models, we need to also incorporate different metrics and methods to do so. As the dataset is limited in number of data points ($N = 74$) and we do not have a second dataset for validation, the best (and only) way to validate the results of our models is performing a cross-validation.

The assessments are done following the related literature [Altmann, 2020, Piovani et al., 2017, Piovani et al., 2018], using the Bayesian Information Criterion and the Sørensen-Dice index.

We present the three different general models used (Sections 2.3.1-2.3.3). We then present the model selection process (Section 2.3.4), to then present the data used throughout the work (section 2.3.6). We finish with a summary of the complete pipeline of analysis for this work in Section 2.3.5. A schematic of how each of the models work can be seen in Figure 2.1.

The retail model, although versatile with respect to the other two with given the amount of variables that can be taken into account and by considering the full topology of the space, it has not been tested as much as the retail model [Davies et al., 2013, Piovani et al., 2017, Wilson, 2008]. The Gravity model is by far the most tested of the three. However, as stated in [Simini et al., 2012, Noulas et al., 2012], the model has different limitations, like the limited number of variables allowed and the simplicity of the reasoning behind. Finally, the radiation model is created as a response to the Gravity model [Simini et al., 2012]. The model has as main limitation the fact that only relies on population and does not consider other geographical/social variables that the other two models do.

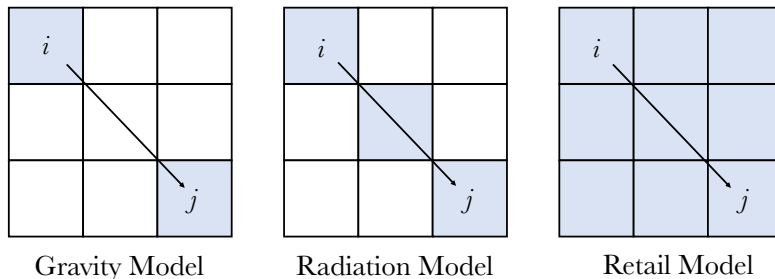


Figure 2.1: Schematic of models. Coloured boxes are those considered by the model.

2.3.1 Retail model

The Retail model was first presented in [Wilson, 2006] as an entropy-maximising model for the function T_{ij}^{retail} with three different conditions: (a) an outflow condition $\sum_j T_{ij}^{\text{retail}} = T_i$. (b) a Boltzmann-inspired energy conservation condition with respect to the travel time c_{ij} from i to j , $\sum_{i,j} T_{ij}^{\text{retail}} c_{ij} = C$, and (c) a similar conservation condition with respect to the total *benefit* found in the space, $\sum_{ij} T_{ij}^{\text{retail}} \log w_j = B$, where w_j is the benefit of place j to attract people.

Using the maximum entropy principle with the three constraints described above, we obtain the resulting function for T_{ij}^{retail} :

$$T_{ij}^{\text{retail}} = \frac{T_i \exp\{\alpha \log w_j - \beta c_{ij}\}}{\sum_k \exp\{\alpha \log w_k - \beta c_{ik}\}}. \quad (2.2)$$

Where α and β are two free parameters coming from the maximum entropy derivation. Notice how the exponent in the numerator represents the balance from the benefits at j and the cost to get to j from i , given by $\alpha \log w_j - \beta c_{ij}$. This latter balance competes with the other balances of going to the places k via the denominator of Eq. (2.2).

The retail system has been studied for different spatial dynamics in the past [Pivovani et al., 2018, Davies et al., 2013], allowing to include different types of data as benefit w_j . In this case, as we are interested in knowing if different social variables (police workforce, knife crime events, hospital admissions by misuse of or poisoning by drugs and drug-related deaths) might be relevant benefits or costs for the county lines operator, we thus replace condition (c) mentioned above by 5 analogous restrictions, one per variable, and use the different $w_j^{(n)}$ as the social/demographic variables. All of them (gross disposable household income, police workforce, knife crime events and hospital admissions) are normalised by the population of the police territory so they become per 100 000 inhabitants. We thus obtain as final solution

$$T_{ij}^{\text{retail}} = \frac{T_i \exp\{\sum_n \alpha_n \log w_j^{(n)} - \beta c_{ij}\}}{\sum_k \exp\{\sum_n \alpha_n \log w_k^{(n)} - \beta c_{ik}\}}. \quad (2.3)$$

By exploring the magnitude and sign of the different α_n , we can then have an insight about the elements that correlate to the detected lines from the Metropolitan Police, and if the variable is perceived as a benefit ($\alpha_n > 0$) or a cost ($\alpha_n < 0$).

2.3.2 Gravity model

The gravity model computes flows from i to j as proportional to the product of populations of i and j , and inversely proportional to the distance between them. The model has different expressions and different limitations [Simini et al., 2012, Noulas et al., 2012]. We take as basis for this work the following form [Anderson, 2010]

$$T_{ij}^{\text{gravity}} = G \frac{m_i^a m_j^b}{d_{ij}^c}. \quad (2.4)$$

We impose the outflow restriction $\sum_j T_{ij}^{\text{gravity}} = T_i$, which makes Eq. (2.4) become

$$\sum_j T_{ij}^{\text{gravity}} = G m_i^a \sum_j \left(\frac{m_j^b}{d_{ij}^c} \right) = T_i \quad (2.5)$$

$$\implies G m_i^a = T_i \left(\frac{m_j^b}{d_{ij}^c} \right)^{-1} \quad (2.6)$$

$$\implies T_{ij}^{\text{gravity}} = T_i \left(\sum_{k \neq i} \frac{m_k^b}{d_{ik}^c} \right)^{-1} \frac{m_j^b}{d_{ij}^c}. \quad (2.7)$$

2.3.3 Radiation model

The idea behind the radiation model originally comes from a particle transmission and absorption model in physics, where a particle is supposed to be emitted from place i and arriving to place j by sorting all *opportunities* in the way, i.e. not being absorbed in the way from one place to another. This idea has been applied for flow of persons in a given space, first used as a commuter model for job seeking in the US [Simini et al., 2012], to then being applied into different examples where commuters are modelled [Piovani et al., 2017, Masucci et al., 2013]. The original formulation of the radiation model is

$$T_{ij}^{\text{rad}} = T_i \frac{p_i p_j}{(p_j + p_{ij})(p_i + p_j + p_{ij})}, \quad (2.8)$$

where p_i and p_j are the populations of i and j , p_{ij} is the sum of populations between both places and T_i is given by the outflow constraint $T_i = \sum_{j \neq i} T_{ij}^{\text{rad}}$. In this particular project we work with a modified version from [Yang et al., 2014]:

$$T_{ij}^{\text{rad}} = T_i \frac{P(1|n_i, n_j, n_{ij})}{\sum_k P(1|n_i, n_k, n_{ij})}, \quad (2.9)$$

where n_i , n_j and n_{ij} are the opportunities in i , j , and between both places respectively. In this case, we simply suppose that $n_i = \rho p_i$, with $P(1|n_i, n_j, n_{ij})$ as the probability of the “particle” being absorbed in way from i to j given the opportunities n_i , n_j and n_{ij} .

$$P(1|n_i, n_j, n_{ij}) = \frac{[(n_i + n_j + n_{ij})^r - (n_i + n_{ij})^r](n_i^r + 1)}{[(n_i + n_{ij})^r + 1][(n_i + n_j + n_{ij})^r + 1]}. \quad (2.10)$$

2.3.4 Model selection process

The three models presented above represent different spatial interactions, interpreted in this context as different decision processes from the county lines operators to establish a connection between place i and j . To compare the different models and selecting the most appropriate one for our available data, we proceed using two different measures found in the literature: the Sørensen-Dice index S [Piovani et al., 2017], and the Bayesian information criterion (BIC) which is based in the maximum likelihood principle [Altmann, 2020].

The Sørensen-Dice S index measures the similarity between two different samples. Given a modelled number of detected lines T_{ij}^{model} after any of the models described above, and the observed data T_{ij}^{data} , we use the same formulation as in [Piovani et al., 2017]

$$S = \frac{2 \sum_{i,j} \min(T_{ij}^{\text{data}}, T_{ij}^{\text{model}})}{\sum_{i,j} T_{ij}^{\text{data}} + \sum_{i,j} T_{ij}^{\text{model}}}. \quad (2.11)$$

We perform a 2-fold cross-validation, splitting our database for 2019 and 2020. Thus, training with 2019 (2020) data to then validate with 2020 (2019) data. Over the results of the validation, we compute the Sørensen-Dice index, thus obtaining two measures of S for each model.

As an extra criterion to model selection, we also compute the BIC to the whole modelled sample by each of the models. BIC computes the log-likelihood and corrects with the size of the sample M for each model. In that sense,

$$BIC = 2 \log M - 2 \log \hat{L}. \quad (2.12)$$

M is the size of the sample and $\log \hat{L}$ represents the maximum value obtained for the log-likelihood when training the model. The log-likelihood is computed with

the parameters that minimise the loss functions used to calibrate the model. As discussed before, given the nature of the detected lines by the Metropolitan Police, we are interested in testing two different loss functions: the usual mean-square loss function derived from a Gaussian likelihood, shown in Eq. (2.13), and a loss function derived from a Poissonian likelihood, shown in Eq. (2.14). The choice of the Poissonian likelihood is given by the distribution of lines detected for both years, while the mean-square loss function is chosen to be a benchmark with respect to Eq. (2.14).

$$\begin{cases} \mathcal{L}_{\mathcal{G}}\left(\{T_{Lj}^{\text{model}}(\hat{\theta})\}_j \mid \hat{\theta}\right) = \frac{1}{2N} \sum_j \left(T_{Lj}^{\text{data}} - T_{Lj}^{\text{model}}\right)^2, & (2.13) \\ \mathcal{L}_{\mathcal{P}}\left(\{T_{Lj}^{\text{model}}(\hat{\theta})\}_j \mid \hat{\theta}\right) = \frac{1}{N} \sum_j T_{Lj}^{\text{model}} - T_{Lj}^{\text{data}} \log T_{Lj}^{\text{model}}, & (2.14) \end{cases}$$

where $\hat{\theta}$ is the vector of free parameters for each model. To each of both loss functions we are also adding an L2 regularisation term $\lambda \|\hat{\theta}\|^2$, with $\lambda = 1$. The subscript L in T_{Lj} represents London, thus showing the observation/model for London to any other police territory j .

2.3.5 Pipeline

The analysis pipeline is as follows: we perform a 2-fold cross-validation on each of the three types of models (gravity, radiation and retail). In total, we are training 1 gravity model, 1 radiation model and 32 retail models. The 32 retail models are a result of adding an offset to the 5 different free parameters $\{\alpha_n\}$ included in the Retail model of Eq. (2.3). Thus, the total number of models is $\sum_{i=0}^5 \binom{5}{i} = 32$. For all the 32 models we still take into account the β parameter which accounts for the travel times cost. A more detailed list of the models trained can be found in Appendix D. The models are trained using two different cost functions described in Section 2.3.4, and evaluated using the Sørensen-Dice index [Piovani et al., 2017] and the Bayesian Information Criterion (BIC) [Altmann, 2020].

The 2 folds correspond to the spatial data of England for the years 2019 and 2020. In that sense, we are performing a cross-validation on the temporal dimension, training the model using the whole of the England topology (space).

The main argument around why we perform a 2-fold cross-validation, and not an n -fold one with a higher n is that, in order to comply with an accurate comparison between the different models, the cross-validation must be performed in the same folds for each of the models. By including the radiation model in Eq. (2.9) which

works in slices of land rather than individual points, we would then have to correctly choose our different folds, so no information is lost when slicing. However, given the topology of England and the way the variable n_{ij} is constructed, we could only slice England in two different pieces, which by themselves are not well balanced (the south-east of England, and the rest of the country).

A schematic of the pipeline can be observed in Figure 2.2.

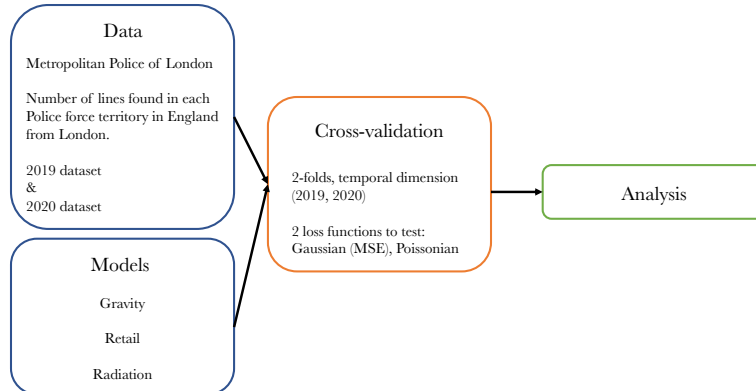


Figure 2.2: Schematic of analysis pipeline

2.3.6 Data

In this subsection we describe the different data that is implemented in the different tested models. In Appendix E we offer a more detailed description of the complete database used. The three models (gravity, radial and retail) have as one of the inputs the population of the police territories (directly or indirectly). These are public data from the Office of National Statistics (ONS), and by the time of submission, the last published update is of 2019¹.

The Gravity and the Retail model respectively use the distance and the travel time from one place to another. Given that the used resolution is at police territory level, we are using the distance/travel time from the most populous place in territory i to the most populous place in j . Data is drawn using the Google Maps©API.

The exponent of Eq. (2.3) allows to compute a balance between the different benefits and costs of going from i to j . The training and comparison process taken in this work allows to know if a given variable is a cost or a benefit, thus allowing to test between different variables.

¹In general, all data obtained from the digital platforms of the British Government (gov.uk) is used under the Open Digital Licence.

An important feature to test is the amount of potential costumers for the county lines operators. This accounts to current and potential consumers. We use two different measures as proxy to this consumption: finished hospital admissions [Lifestyles Team, 2019] by misuse of drugs and finished admissions by poisoning of drugs. Hospital admissions are normalised by population and by daytime hospital beds *per capita*.

Another feature we test is the police workforce in each territory. We use the number of average Full-time police officer over the British Fiscal year (May-April) which can be obtained from [Flatley, 2019].

To account for the disparities of richness in the different parts of England, we use the gross dispensable household income (GDHI). In comparison with the household income, the gdhi takes into account the amount of money that households have after local and national income taxes and benefits from the government. Data was obtained from the ONS [Fenton, 2021].

Finally, we are interested in testing the knife crime events *per capita* in each of the police territories. Knife crime events have been an increasingly worrying matter for the British Government, with numbers increasing 78% in England from 2014 to 2020 [Bellis et al., 2019].

2.4 Results

2.4.1 Model selection

Results for the BIC and the Sørensen-Dice index are found in Figure 2.3a and in Figure 2.3b respectively.

When comparing the Retail model calibrated with a Gaussian loss function (zone 3 in Figure 2.3) with respect to the other models, we observe how the former performs worst in all of its forms for both the BIC and the S index. We can thus proceed to discard these models.

With the models left, we perform a comparison by computing the MSE between the Metropolitan Police data for both years (2019 and 2020) and the predictions obtained from each model. The MSE is computed with the logarithms of the data points, so in this case $MSE = \frac{1}{N} \sum_{i,j} (\log T_{ij}^{\text{data}} - \log T_{ij}^{\text{model}})^2$ Results are seen in Figure 2.4.

The best performing model is the Retail model trained with the 2019 data and the Poissonian loss function. However, as it can be seen in the inset plot in Figure 2.4, the results can be differentiated in different levels. When examining each one of the four levels, we find that the hospital admissions by poisoning of drugs, the disposable income and the police presence variables do not have significant effect on the performance of the model. This can be seen in the upper level, as those

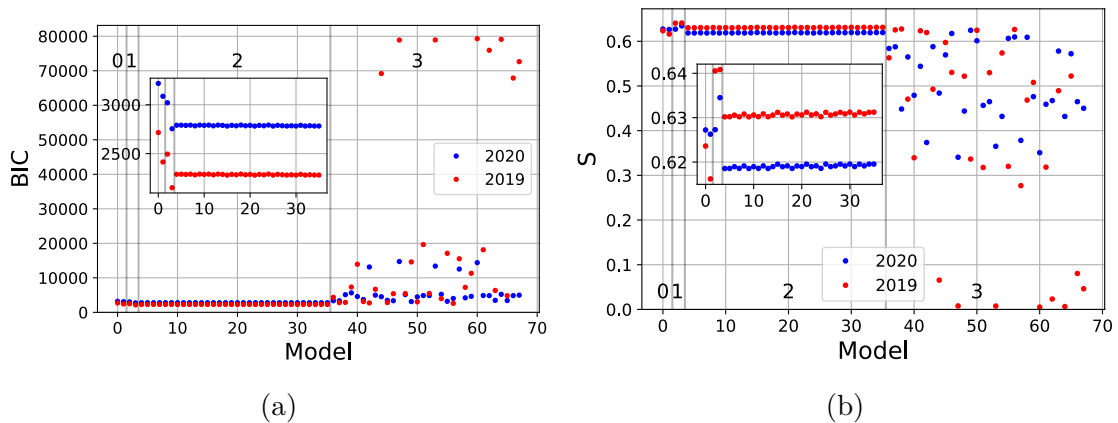


Figure 2.3: Results of BIC and the Sørensen-Dice index for the 68 different models tested. Zone 0 corresponds to the Gravity model. Zone 1 corresponds to the Radiation model. Finally, zones 2 and 3 correspond to the retail model with the Poissonian loss function and with the MSE loss function respectively. The annotated year corresponds to the data in which the model was trained on. Details of each model can be found in Appendix A.

combinations not containing the knife crimes and hospital admissions by misuse of drugs variables are those present there (all the different models are in Appendix D). The fact that the combination without any of the social variables and only the travel times is there allow us to interpret that any of the three mentioned variables before do not have any particular effect on the performance of the model. The hospital admissions by misuse of drugs seem to have an impact on the cost, although not as important as the knife crime variable. When combining both variables we obtain the most important effect on the MSE cost and the best performing models.

The Radiation model follows as best performing when trained with the 2019 data and Poisson loss function. Finally, we obtained the Gravity model trained in the same way.

In Table 2.1 we detail all the selected models. To keep the selected models as simple as possible, we filter out all the different Retail models and keep only those with the minimum number of variables. That is, one with both the misuse and the knife crime variables in addition to the travel times, one with only the knife crime variable and travel times, one with only the misuse variable and travel times, and finally one with only travel times.

From the exponent in Eq. (2.3), α_2 to the hospital admissions by misuse of drugs, and α_4 to the knife crime events. All variables are normalised by population.

We also select the best performing Radiation and Gravity models as we are interested in comparing them with respect to the Retail model.

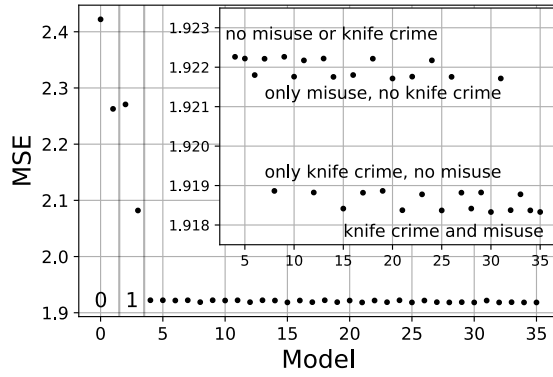


Figure 2.4: MSE costs when comparing the trained models with the Metropolitan Police data.

As it can be seen from Table 2.1, the five exponents α_n are negative, which is interpreted as all of the variables to represent a cost to county lines operators. This will be discussed in Section 2.5.

Table 2.1: Results for the best three models calibrated.

Ranking	Model	Loss function	Training year	Parameters
1	Retail	Poisson	2019	$\alpha_2 = -0.774e-2$, $\alpha_4 = -0.013$, $\beta = 0.014$
2	Retail	Poisson	2019	$\alpha_4 = -0.013$, $\beta = 0.014$
3	Retail	Poisson	2019	$\alpha_2 = -0.777e-2$, $\beta = 0.014$
4	Retail	Poisson	2019	$\beta = 0.014$
5	Radiation	Poisson	2019	$\rho = 2.085$, $n = 1.038$
6	Gravity	Poisson	2019	$b = 0.697$, $c = 0.368$

2.4.2 Model analysis and geographic distribution

Once we have obtained the best performing models, we proceed to compare and analyse the simulated distribution of modelled lines. In Figure 2.5 we present the different models compared with the Metropolitan Police Data for 2019 and 2020.

We observe outside this work that the four best performing models (Retail models) act almost identically, so we only depict models 1, 5 and 6 from Table 2.1.

The three models tend to overestimate the detected connections to places with less than 70 lines, while tending to underestimate them in police territories with more than 100 lines detected.

Each one of the models have different ways of understanding the dispersion of flow in a given space. On the one hand, the calibrated Radiation model sees the

2.4. RESULTS

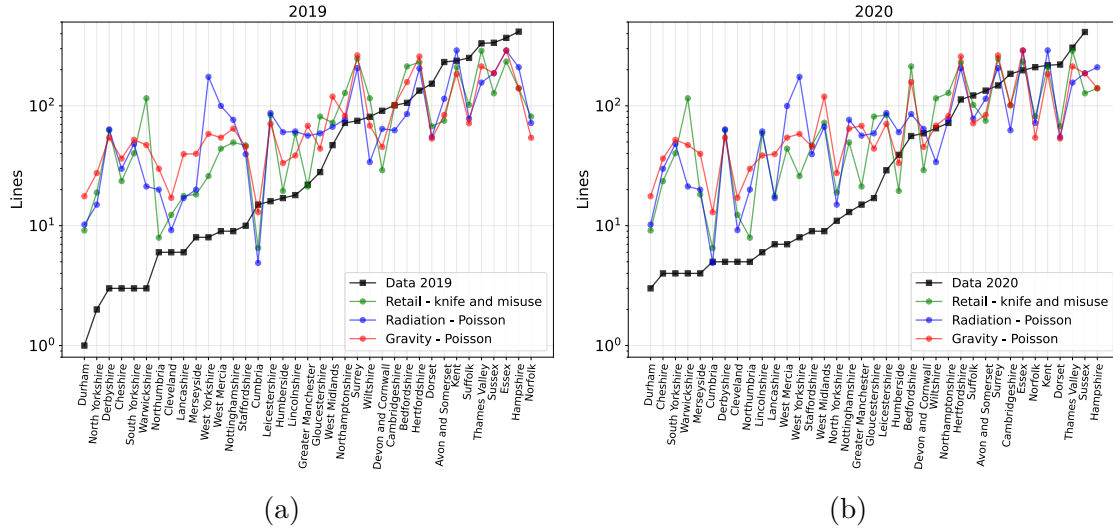


Figure 2.5: Data points and modelled lines ordered by police force for 2019 (a) and 2020 (b). Note: for both years, there were 0 lines detected in Durham. As the plot is in log scale, this data point was not included.

flow from London to another given police territory as a process of sorting opportunities presented on the way. After our calibration, opportunities here are seen as proportional to the population by the value of ρ given in Table 2.1. Thus, we are actually exploring how the population is distributed in England.

On the other hand, the Retail model understands flow as a balance of with respect to travel times and the other social variables using an exponential distribution. This means that flow from London to another police territory is given by how much time is spent commuting with respect to the other police territories and how much the other benefits/cost relate to it. Thus, a closer place from London would be favoured with respect to a farther one. However, given that this consideration is given by an exponential distribution, we can expect a slow decrease of lines when increasing travel times (light tail distribution).

Finally, the Gravity model explores the flow with respect to the distance between two places and the population of the target place. In that sense, closer and more populous locations would take most of the outflow, while distant and less populated locations would be disfavoured by the model.

To further understand the different results shown in Table 2.1 and in Figure 2.5, we map the different models and compare them with the Metropolitan Police data. This is shown in Figure 2.6. While Figure 2.6a and 2.6e present the Metropolitan Police data for 2019 and 2020, the rest present the modelled spatial distribution of lines. We also present the differences between the Metropolitan Police data and the models in Figure 2.7. Red zones correspond to territories overestimated by the

2.4. RESULTS



Figure 2.6: Heatmaps for the Metropolitan Police data for 2019 (2.6a) and 2020 (2.6e), and the three different models tested: Retail (2.6b and 2.6f), Radiation (2.6c and 2.6g) and Gravity (2.6d and 2.6h).

model, while blue zones correspond to territories underestimated by the model.

We start by analysing Figures 2.6a and 2.6e corresponding to the Metropolitan Police data. The first thing to notice is the decrease of detected lines in 2020 with respect to 2019. This effect can be given by mainly two factors taking into account the COVID-19 situation throughout 2020: the police had a smaller capacity to detect, or indeed the reduced mobility in the country reduced the number of connections. However, the decrease is not generalised and we can observe an increase in some police territories from 2019 to 2020, as in Hampshire (South of England) where we find the maximum number for 2020.

An important second element to note from the ground truth data is a very high share of the total lines (94.02% for 2019 and 93.77% for 2020) concentrated in 16 out of the 37 police territories considered. This set of 16 police territories, in addition to London, is considered to be the “South” of England, a social region with no administrative recognition which encloses the most developed parts of England and which opposes the “North” of England, where more industrial cities like Manchester and Liverpool are found (for a study using percolation theory please refer to [Arcaute et al., 2016]).

The “North-South” division is an element which none of the models captured.

However, we can still see different ways of simulating the problem in Figure 2.6. As discussed before, the Retail model distributes the lines in what appears a concentric fashion with respect to London, leaning towards the centre of England. This can be seen more clearly in Figures 2.7a and 2.7d, where we observe an overestimation in the Midlands and an underestimation of the coastal territories of the “South”. Note how the far South West of England (Cornwall and Devon), which is farther away in travel times than the centre of England from London, is underrepresented. This fact accounts for an argument in which the operators in London would not have as primordial element for establishing connections the travel times to the different territories. This argument is supported by the opposite fact, where we observe an overestimation by the retail-gravity model in more connected places from London, like the West Midlands (Birmingham) and Warwickshire (south of Birmingham).

The radiation model understands the flow in a different fashion, as seen in Figures 2.6c and 2.6g. In a similar way as the Retail model, the ring surrounding London is still catching an important number of lines. However, we can also observe a number of relatively large hotspots, particularly in West Yorkshire (North of England) and in West Mercia (border with Wales). While the former territory includes important cities and urban centres such as Leeds and Bradford, West Mercia is a diverse territory with dense suburban counties belonging to the Birmingham metropolitan area and more rural areas towards Wales, like Shropshire. In Figures 2.7b and 2.7e we observe also how the territories between West Yorkshire and London were filled with lines by the Radiation model. It is also important to note how the big metropolitan areas in England such as Birmingham do not appear as hotspots in Figures 2.6c and 2.6g.

Both models described above tend to distribute the number of lines in the centre of England, while avoiding the big cities. This is in contrast with the Gravity model (Figures 2.6d and 2.6h) where we observe the appearance of Birmingham and Manchester (2nd and 3rd most populous cities in the UK) as county lines hotspots.

The three models do not detect the territories where the maximum number of lines are detected, like Norfolk in 2019 and Hampshire in 2020. On one hand this is a sign of no overfitting from both models, but on the other hand makes very difficult for the models to detect future hotspots in the South of England.

2.5 Conclusions

In the present work we study the County Lines Model (CLM) distribution method of illicit drugs in England. Our aim is to shed some light around the territorial logic behind the data accounting the detected connections (lines) by the Metropolitan



Figure 2.7: Heatmaps showing the difference between the modelled distribution of lines with respect to the Metropolitan Police data for 2019 and 2020.

Police of London in other police territories [[Rescue and Analysts, 2019](#), [Rescue and Analysts, 2010](#), [Black, 2020a](#), [Black, 2020b](#)].

We understand the number of detected lines as a flow of people/merchandise that starts in London and finishes in a given police territory. In that sense, by modelling the flow from one place to another and compare it with the available data, we obtain information about which elements are present when establishing a local market.

Three different models are studied and compared. Each one of them follow, by construction, different logics about how they understand a flow from one place to another. The first one, the Gravity model [[Noulas et al., 2012](#), [Anderson, 2010](#)], sees the flow as proportional to the population of both places, while inversely proportional to the distance between them. We take this model as our benchmark as represents the classic idea that populous places would draw more attention than others given the same distance. The second, the Radiation model [[Simini et al., 2012](#)], understands flow to a given place as a process of sorting the opportunities before arriving to the final destination. With this second model we tested if the distribution of the population (in comparison to a single population spot) in England was involved in the decision making. Finally, the Retail model [[Wilson, 2006](#), [Wilson,](#)

[2008] takes into account the balance between the benefits and costs of establishing a flow towards a particular place. This final model allows to include as potential benefits/costs different social variables that we tested, like police workforce, knife crime events, hospital admissions for drug poisoning and for drug misuse.

We train the models using the Metropolitan Police data, and compare them using the Bayesian Information Criterion [Altmann, 2020], and the Sørensen-Dice index [Piovani et al., 2018] over a cross-validation. We also test two different loss functions, the classic Mean-square error and the one derived from a Poissonian likelihood.

The best performing model is the Retail model with different combinations of social variables and trained with the 2019 data and the Poisson loss function. We find that for certain combinations of social variables, the Retail model would give better results than others. Indeed, the hospital admissions by misuse of drugs *per capita* and knife crime events *per capita* are the two most influential variables to obtain better results. In particular, knife crimes shows to be more important to hospital admissions when compared one to one.

The Radiation and Gravity model also perform correctly when trained with the 2019 data and the Poisson loss function. However, when comparing the geographical distribution of both model to the Metropolitan Police data, these two models predict hotspots in populous regions of England where no important number of lines was detected by the Met. Police.

According to our ground truth, the distribution of the great majority of lines (93%) is over 16 of the 37 police territories in England, which form the union of the South West, the South East and the East of England. This territory is known as the “South” of England.

While the Gravity and Radiation model overestimated different territories outside the South of England with large populations, the Retail model did it in a more diffused way. This is due to the exponential form of the model.

None of the three models could capture the hard border that the data shows between the South of England and the rest of the country. This raises the question about the characteristics in the 16 police territories that represent the ‘South’ of England that make them so attractive to CL operators. A first hypothesis is that the CLM, although not reported in literature, actually acts within a more organised structure which can restrict itself to distribute in a given territory, as seen for other criminal organisations. In other words, even though not mentioned in the public information by the UK government, different CLM gangs operating from London could restrict themselves to these 16 police territories as a measure to do not enter in an open conflict with other gangs from other CLM hubs. This hypothesis could be studied by having data from the detected lines from the other important CLM hubs, like Birmingham, Liverpool, Manchester and West Yorkshire. In that sense,

we could expect a localised distribution in the 'North' of England, obtaining then a polycentric structure within the territory.

However, if data from other CLM hubs would not comply with the segregation and rather concentrate in a subset of the 16 police territories considered as the "South", then we would have a particularity of the consumers in those areas. This would also be of interest, as the population in this subset would have to have a distinction with respect to the other big metropolitan and rural areas the 21 police territories left. This distinction, although might be related to a particularity of the consumers, would also address the findings already obtained before in quantitative studies [Arcaute et al., 2016], where a clear distinction between the urban network between the South and the rest of England was found using percolation analysis.

The hypothesis about a polycentric structure could be supported by our findings on how knife crime events and hospital admissions by misuse of drugs are a cost to line operators. The fact that knife crime events appear as a cost might point to an avoidance from the operators to certain gangs so conflict is spared. Hospital admissions, on the other hand, are used as a proxy to illicit drugs consumers given the lack of public information about it. In that sense, the fact that the hospital admissions variable is one of the two most influential variables, combined with the knife crime variable, could be interpreted as county line operators avoiding places where there already is enough competition for them to handle. This competition can be regarded as possible origin of conflicts (knife crimes) and responsible of having a greater share of the illicit drug consumption market in a given territory.

The implications of our findings are then quite straightforward. We demonstrate that the logic behind the county line operators is not as simple as an offer-demand one [Rescue and Analysts, 2019, Rescue and Analysts, 2010, Black, 2020a, Black, 2020b, Supt Mat, 2020, Crime Agency, 2019, Silver and Intelligence, 2021], but actually might follow a social structure of the country given the lack of significance that distance, population distribution and the dispensable income have in our results. The result also shows a logic in which conflict with other gangs and other markets already filled with competition are avoided. This by itself can be of great help for law-enforcement bodies, as it gives a good lead on where to look for the presence of county lines from London: places within the 16 police territories where there is not an important number of knife crimes *per capita*. There is no mention of these factors in the reviewed literature.

This work also allows to implement a better coordination between local police forces, as the Metropolitan Police of London would only need to coordinate with 43% of the English police forces to tackle the 93% of the lines detected.

More in general, we have shown how with a combination of particular needs from

a public body we can show different solutions and insights by using a combination of public data, mathematical modelling, statistical learning and basing our pipeline in social qualitative work which frames our work.

The main limitation of this work is the lack of data. Having a larger dataset both in the temporal dimension and the territorial origin could make us obtain a more complete analysis of the county lines model not only in England but in Great Britain as a whole (no county lines has been identified in Northern Ireland). By obtaining data from other police territories like Liverpool (Merseyside police), Manchester, Birmingham (West Midlands Police) or West Yorkshire we could, as possible extension of this work, to analyse the hypothesis described above. However, we understand that given the different jurisdictions, each police force has the capacity to decide if that data is published or not, particularly given its sensitive nature. This fact limits the ability of this study to conduct a more comprehensive analysis.

This work also contributes to the literature of spatial analysis and quantitative security by implementing a comparative analysis between different spatial models to understand the territorial logic behind an illicit drugs distribution model in England that has brought important social problems to certain parts of the population. The present work also implements modern statistical learning techniques such as cross-validation, the introduction of multiple metrics to assess the trained models and the comparison between different cost function to determine the best one possible [Altmann, 2020, Piovani et al., 2018]. The results from our analysis can be extended in order to understand the distribution problem in England from a social data-driven approach, while also being able to deliver public policy suggestions for the law-enforcement bodies. Also, the present work delivers and fills the literature gap with respect to analysing the County Lines Model from a quantitative perspective, using novel data-driven models and tools to nowcast the Metropolitan Police data and understand the territorial logic of the county lines operators.

Chapter 3

The Impact of Corporate Values and Factors of Internal and External Culture on Formulating the Post-COVID “New Normal”: Implication for Cybersecurity and Information Systems

Chapter Abstract

This paper investigates whether and to what extent corporate values as well as factors of external culture impact on the companies’ ability and commitment to formulate effective and realistic “new normal” post-COVID strategies with particular focus on cybersecurity and information systems priorities. Using COVID-19 response documents from top 100 companies featured on the Forbes Fortune 500 global as well as US lists, we employ topic modelling to map top priority themes in the COVID responses mentioned by the companies and explore whether and how these priority themes, together with factors of external culture (Schwartz cultural value orientation, Global Cybersecurity Index) influence business financial success and resilience at times of uncertainty. We find that while cybersecurity and network security are rarely a subject of corporate focus, reaching a successful new normal requires businesses to concentrate on management of risks, risk and uncertainty aversion, as well as on tackling (digital) fraud. The originality of our research is focused on understanding and calibrating our model from a cyber security culture perspective by studying what was CEOs and high managerial posts were thinking by the time the pandemic started using statistical analysis and natural language processing techniques.

3.1 Introduction

For more than a year, our planet has been living in conditions of a global pandemic. This pandemic, caused by the 2019 coronavirus disease (i.e., COVID-19), affected all facets of human life, including the day-to-day functioning of businesses around the globe. While it is hard to fully appreciate the long-term economic impact of COVID-19, it is evident that most economies have already felt the consequences, showing high spikes in unemployment rates. The purpose of this paper is to formulate and test a new approach to understanding and mapping the future, so-called “New Normal”. We also look for strategic priorities for businesses around the globe, with particular emphasis on cyber security and information systems using corporate values as well as factors of internal and external business culture.

Under pandemic circumstances, many businesses require or recommend that their employees (dependent on the local area pandemic situation) work from home. Naturally, many companies are concerned about cyber security of remote work [Caligiuri et al., 2020]. The main problem with cyber security while working remotely is that business systems still rely on personal cyber hygiene of employees [Dwivedi et al., 2020], many businesses still not having a clear plan of what happens in case of a cyber security breach. As a result, employees do not know whom to contact to report cyber incidents, especially during a cyber emergency, and this is exasperated in a remote work environment. Who-does-what and who-reacts-to-what is not clearly identified [Gerke et al., 2020].

As companies get comfortable with Work-from-Home arrangements, their boundaries have now extended to their employee’s home and the personal technologies in their homes. This is a vulnerability companies now must manage [Abukari and Bankas, 2020].

Employees face different situations that can raise important issues to their companies’ cyber security, such as with an appropriate security of their home system [Forte and Power, 2007], e.g. Wi-Fi systems [Jang-Jaccard and Nepal, 2014], or an appropriate use of mixed devices in a “bring-your-own-device” policy [Yong Wang et al., 2014].

Nevertheless, businesses are already starting to prepare for the end of the pandemic and the “New Normal” by creating a set of processes and routines that will persist beyond the pandemic and help them be better prepared for the future [Sakurai and Chughtai, 2020, Hacker et al., 2020]. Much of the success of the “New Normal” business design and implementation heavily depends on the acceleration of the Industrial Revolution 4.0 technological advances such as AI, data-driven analytics and processes, as well as intelligent automation supported by the next-generation information systems (e.g. [Skilton and Hovsepian, 2008]). At the same time, in

their formulation of what the “New Normal” after the pandemic could look like, businesses can learn a lot from the social experiences in dealing with natural disasters, where information systems may be used to foster resilience against any source of crisis [Sakurai and Chughtai, 2020]. Businesses can also achieve greater resilience through nurturing the new virtual sense of togetherness through the use of the web conferencing systems [Hacker et al., 2020].

This paper uses a combination of qualitative and quantitative techniques to better understand corporate preferences, priorities, and culture, which impact on the formulation of the “new normal” and foster anti-crisis resilience. We are particularly interested on a cyber security and information systems emphasis. The qualitative basis of our approach is inspired by the Massachusetts Institute of Technology (MIT) Cyber security at MIT Sloan (CAMS) model [Huang and Pearlson, 2019a] of cyber security culture (henceforth, MIT CAMS model). The goal of the model (depicted on Figure 3.1) is to link cyber security behaviours with managerial influences. The model suggests that organisational culture influences behaviours for cyber security. Organisational culture can be described as the beliefs, values and attitudes held by leaders, groups, and individuals who make up the organisation. These beliefs, values and attitudes are shaped by external influences such as country norms, industry norms, regulations and other construct, and by internal managerial mechanisms (such as training, awareness programs, performance reviews, rewards, consequences, and corporate communications) that are directly under the control of organisational leaders.

We adapt the MIT CAMS model to the “New Normal” and test it with different data obtained from different sources. We call the adaptation of the model CNNM (CAMS-inspired New Normal Model), and we test it for two different examples: a data set from Fortune 500 for US companies only, and in a data set from Fortune Global 500, thus including companies whose HQ are not in the US. External cultural values are implemented via two different data sets: the Cultural Value Orientations coefficients [Schwartz, 2009], and the Global Cybersecurity Index (GCI). Internal cultural values are drawn from the corporate COVID-19 response documents using Natural Language Processing algorithms, while the beliefs, values and attitudes are obtained using the Linguistic Inquiry and Word Count (LIWC) framework. Business resilience is applied using the Forbes 500 database.

This paper tries to answer the following research questions: What are the top priority themes in the COVID-19 response documents of major global businesses in 2020, which are likely to determine their post-pandemic response and formulation of a new normal? Are (any of) these priorities related to cyber security or, more generally, to information systems? What is the correlation between company’s financial success and resilience and priority themes in their COVID-19 responses? How do

COVID-19 responses and strategies depend on external and internal cultural values? We divide our results and answers for different clusters of industries, thus gaining important insights for different sectors.

This paper contributes to the emerging literature on the global pandemic business strategy, information systems and the new normal. This emerging literature covers a wide range of COVID-19 impacts from environmental and urban to economic.

Yet, cyber security and information systems aspects of the “New Normal” remain under-researched. The main contribution of this paper is to fill this gap by developing a valid theoretical methodology supported by feasible and realistic empirical test. The remainder of this paper is organised as follows: In Section 3.2 we discuss the literature relevant to this work and the literature gap found. It is also in this section where we develop the hypotheses we work with throughout the work. Section 3.3 provides our methodology with an insight about our adaptation of the MIT CAMS model. A first visualisation of the data set is shown in Section 3.4, while the full results are provided in Section 3.5. Finally, Section 3.6 discusses the results to then conclude in Section 3.7. We are also including an Appendix F where all technical procedures are detailed.

3.2 Literature Review

In this section we discuss the relevant literature for this work. First we review the literature published until now around the Post-COVID 19 “New Normal”. We then review the literature published around cyber security when working remotely, to finally present the Cybersecurity at MIT Sloan model (CAMS model).

Since the pandemic was declared on March 11, 2020 by the World Health Organization, an increasing number of countries around the world started to see the number of infected people rise, leading most of them to impose lockdowns and curfews in different ways. Quite early, the term “New Normal” was coined to describe what would refer the life with the COVID-19 virus (SARS-CoV-2) around. While part of this new way of life referred to adapting to being at home for an unexpected number of days, much of the academic and industrial world, such as the different government bodies around the world. In [Habersaat et al., 2020], the authors outline 10 different considerations for companies, governments and individuals to transition to the “New Normal”.

By the end of 2020, peer-reviewed literature around the “New Normal” was mostly centred in the topic of tele-medicine and health care systems [Lanham et al.,

2020, Jiang et al., 2020]. As the pressure in Health Care systems around the world increased with the number of patients needing hospitalisation due to COVID-19 [Balakrishnan et al., 2020, Retzlaff, 2020, Eardley, 2020], different fields in the medicine world published articles and guidelines about how to control the need of surgery in their specialisation while dealing with COVID-19 patients. To name a few, papers were published around orthopaedic [Anoushiravani et al., 2020], dental [Tandale et al., 2020], arthroplasty [Zeegen et al., 2020], cardiovascular [Tamagnini et al., 2020, Zoghbi et al., 2020], radiology [Siegal et al., 2020] and gastrointestinal [Holtmann et al., 2020, Sethi et al., 2020], in addition to midwifery during the pandemic [Walton, 2020], application of local anaesthesia [Lie et al., 2020], dermatology treatments [Ng et al., 2020] and general surgery [Cobianchi et al., 2020].

From this particular literature, we expand the work of [Lanham et al., 2020], where the authors highlight that the success of tele-medicine in the future would depend on the *satisfaction of the patient*, making an emphasis on an user-centred perspective around the new technologies in medicine. Also, in [Tandale et al., 2020], the authors highlight how doctors and health care personnel must be careful, as any patient that arrives into any of the health care systems should be treated as a potential threat for infection of COVID-19. In [Tamagnini et al., 2020], the authors make the observation how for different medical specialisations such as cardiovascular surgery, there is an important overlap between the demographic sector that would require a medical intervention and the one who is more vulnerable to important consequences due to COVID-19. Finally, in [Balakrishnan et al., 2020], the authors talk about how the lack of PPE and medical equipment in general could lead to potential threats in the health care system as a whole.

Outside the medicine literature, other authors talked about the “New Normal” during 2020. In [Hesse and Rafferty, 2020] the authors try to understand how the pandemic could shape important urban hubs such as Luxembourg and Dublin from an urban perspective. In [Lappan et al., 2020] states some guidelines for academics working with primates, while in [Bloomquist, 2020] the author reflects about the “New Normal” from a theological point of view. In [Doolittle, 2020] the authors present a reflection about how COVID-19 could affect biological evolution in our world.

In [Yang, 2020], the author proposes a view about how the Chinese government responds to the pandemic and how it implements some first steps to the New Normal. In [Reuter et al., 2020], the authors analyse the results in the general health of the population in South Africa resulting of a alcohol prohibition in the same country. The authors finalise by suggesting an implementation of the same prohibition in other African countries given the benefits of the policy in the general population.

In [Mukherjee et al., 2020], the authors research around how the pandemic and the “New Normal” could help countries and international organisations to prepare for future disasters in our world, while in [Sakurai and Chughtai, 2020] the authors make a comparison between the 2011 Japan earthquake and the pandemic to draw some similarities to increase the resilience of society to face the 2020 pandemic.

Finally, there is an emphasis in the literature about how the pandemic and the “New Normal” also implies the change in different aspects of our daily habits. In [Hacker et al., 2020], the authors talk about the effects in our daily life of implementing an online conferencing system when working remotely, while [Dwivedi et al., 2020] shows a general perspective about how the “New Normal” transformed work, life and education patterns. In this last topic, [Greenhow and Chapman, 2020] and [Triyason et al., 2020] expand on the shift to an online education and a hybrid classroom for children and students. In [Larcher and Brierley, 2020], the authors make an analysis about the role that the children play during the pandemic in this change of life patterns.

Respecting the “New Normal” in businesses and cyber security, [Harwood, 2020] talks about the importance of maintaining a cash flow throughout the lockdowns, as different unexpected consequences could come if the flow is blocked throughout the world. In [Caligiuri et al., 2020], the authors research about how the international business must adapt to the “New Normal”. The authors of the latter research document show the decisions seen by international companies at managerial level at the moment of publication (June 2020) were mostly related to managing the imposed distancing and “rethinking boundaries”. However, more decisions should be thought and taken to deal with a new breed of remote workers throughout the pandemic and possible post-COVID times.

From the very beginning of the pandemic, authors started researching and publishing about the different challenges that the situation would represent for cyber security. In [Gerke et al., 2020], the authors discuss privacy laws of the USA and Europe around the topic of home monitoring techniques applied to decrease the interpersonal contact. In [Abukari and Bankas, 2020], the authors explore and outline protocols for cyber hygiene for remote worker during the pandemic. They outline the importance of educating and training workers, plus the importance of implementing policies within companies. Both works are important taking into account the relevance of Work-from-Home policies implemented since the pandemic started, and the different issues found in privacy and security before 2020 with respect to *Bring-your-own-Device* (BYOD) policies [Yong Wang et al., 2014].

The model we base the present work, the Cybersecurity at MIT-Sloan model

(CAMS model) first presented in [Huang and Pearlson, 2019a] explores how to create a cyber security culture in an organisation thinking in all of the possible managerial stages possible. The authors define cyber security culture as “*the beliefs, values, and attitudes that drive employee behaviors to protect and defend the organization from cyber attacks*” in the cited article. They centre their work on this “unwritten rules”, represented as the beliefs, values and attitudes of a company to then explore how to obtain a positive behaviour with respect to cyber security culture. The resultant behaviour, explain the authors, can be of *in-role* nature (where it is part of the job description) or *extra-role* nature (where it is not part of the job description). A schematic of the complete model is shown in Figure 3.1.

The authors in [Huang and Pearlson, 2019a] base their model in different works around organisational culture literature and studies about what this *culture* would mean in an information security context. On a first place, they base their definition of organisational culture from [Schein and Schein, 1985] in which the authors define it as “a pattern of shared basic assumptions that a group learns as it solves its problems of external adaptation and internal integration, that has worked well enough to be considered valid and, therefore, to be taught to new members as the correct way to perceive, think and feel in relation to those problems.” This definition is complemented with three different components that are found in any organisational culture: (i) the belief system forming the basis for collective action; (ii) the values representing people’s principles and (iii) the artefacts and creations related to any sensible behaviour as well as any myth, legend and common language used [Schein and Schein, 1985].

When adapting Schein’s definition of organisational culture to the cyber security domain, the work of [Da Veiga and Eloff, 2010] is taken as a point of reference, as the authors refer to the beliefs and attitudes that employees and stakeholders use to interact with the information systems at any given point. They put a particular emphasis on the fact that this culture can change in time.

[Huang and Pearlson, 2019b] construct the CAMS model with the behaviour (in-role and extra-role ones) from employees, leaders and stakeholders with respect to the cyber security of the organisation and link it with the organisational culture (values, beliefs and attitudes) of the company.

The two last components are external influences that have an effect on the organisation cyber security culture. This could be of legislative and/or regulatory nature, of peer institution-pressure nature or even societal nature, where different elements of the societal culture influence the internal cyber security culture. Finally, the authors include Managerial Mechanisms as the last element that influences and can be influenced by the cyber security culture of an organisation (beliefs, values and attitudes). This, as the name says, refers to specific mechanisms as training given,

3.2. LITERATURE REVIEW

leadership decisions, communication channels established and rewards/punishments put in place. Notice how in Figure 3.1, the External influences influence the organisational culture, but not the other way around. This is different to the Managerial mechanisms, where they influence and are influenced by the organisational culture. The authors explain that culture is created with leadership actions. However, as stated above, organisational culture has the capacity to change over time, thus also influencing future managerial mechanisms from leadership positions.

The CAMS model was constructed with help of senior executives in international companies, such as the Chief Information Security Officer (CISO) of Liberty Mutual, and then validated in Banca Popolare di Sondrio, another Financial Institution, in [Marotta and Pearlson, 2019].

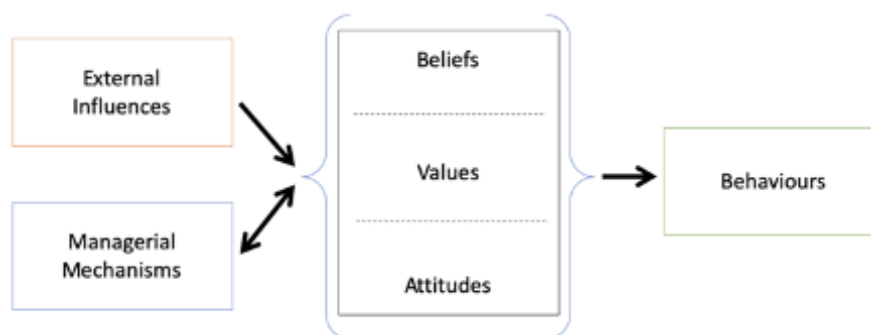


Figure 3.1: MIT CAMS Culture of Cybersecurity Model adapted [Huang and Pearlson, 2019a]

From the literature cited above we find that, although there is research done around the transition to the “New Normal” during the COVID-19 pandemic, there is no much related to cyber security, particularly referring to how the companies need to adapt to a Work-from-Home environment. We find an interesting literature gap where there is no research around how companies are tackling the transition. More specifically, we are interested in knowing not how they managed/are managing the transition, but if they thought about the dimension of the transition once the pandemic started, particularly referring to cyber security. Relevant literature written before the pandemic around cyber security like [Jang-Jaccard and Nepal, 2014, Yong Wang et al., 2014, Huang and Pearlson, 2019a] allows us to find the tools to fill it.

Given the above literature and the different works relating crisis control, COVID-19 and the “New Normal”, we can observe how having a quick adaptation and problem-solving plan when externals shock come have a direct impact on a given system’s resilience [Sakurai and Chughtai, 2020]. Also, as seen in the organisational culture literature [Schein and Schein, 1985, Da Veiga and Eloff, 2010, Huang and Pearlson, 2019b, Marotta and Pearlson, 2019], the culture itself is subject to change

over time, having a direct effect on the sensible behaviours as with different managerial mechanisms. As such:

H1: Those companies that make reference to the transition to a new way of working proceeding will be positively correlated to a better financial resilience.

However, as described above, there is a disproportionate mention to the “New Normal” within the industrial sectors, with a great majority of literature from medicine and health care systems. Given the lack of literature referring to the COVID-19 pandemic and the transition given during the first months of the pandemic. This leads us to our second hypothesis:

H2: Some industry sectors are more prone to think of a transition than others by the time the documents were released.

This second hypothesis is also supported by the fact that, as mentioned very early during the pandemic [Habersaat et al., 2020], when prioritising the different actions related to tackle the effects of the lockdowns, different industries will do it in a different way. Particularly, those relying on healthcare workers will have more emergency in transitioning, in the same way that those that rely on high-risk persons. As the authors in [Habersaat et al., 2020] propose, a phased transition plan allows margin to industries to take as much action as they want in a more hastened way as needed.

3.3 Methodology

In the following section we detail the methodology followed for this work. We adapt the MIT CAMS to the context of the “New Normal” in such a way that would allow us to make a quantitative study of the cyber security behaviour of companies.

First, we explain the adaptation we do to the MIT CAMS model in Section 3.3.1. A representation of the final adapted model used is shown in Figure 3.2. In Section 3.3.2-3.3.5 we detail the different resources used to analyse our new model.

3.3.1 The CAMS-inspired New Normal Model extension

Our methodological approach combines qualitative and quantitative methodology. The MIT CAMS model was successfully applied to many contexts and case studies as well as achieved real-world impact in public and private sectors, becoming a hit amongst practitioners (see e.g., [Huang and Madnick, 2019]; [Marotta and Pearson, 2019]; [Macedo and Menting, 2019]). The main advantage of the MIT CAMS model

3.3. METHODOLOGY

is that it allows to capture the factors of corporate culture and establish the causal links between these factors and observed behaviours in highly uncertain conditions with many unknowns.

In order to adapt the MIT CAMS model rationale such that it would provide an insight into the post-COVID-19 “new normal” business priorities, which, in turn, would foster corporate business resilience, we extend the MIT CAMS model to the CAMS-inspired New Normal Model (CNNM) and propose a feasible and logical way of testing it. The CNNM is summarised on Figure 3.2. The CNNM links factors of corporate culture with corporate success in the time of uncertainty through corporate beliefs, value, and attitudes. Observable factors of external corporate culture (such as cultural value orientations of the country, where the company has its headquarters; cyber security culture of that country, etc.) together with factors of internal culture (such as managerial mechanisms of addressing pandemic challenges) form the latent corporate beliefs, values and attitudes, which, in turn, influence observable corporate success (resilience) during the pandemic.

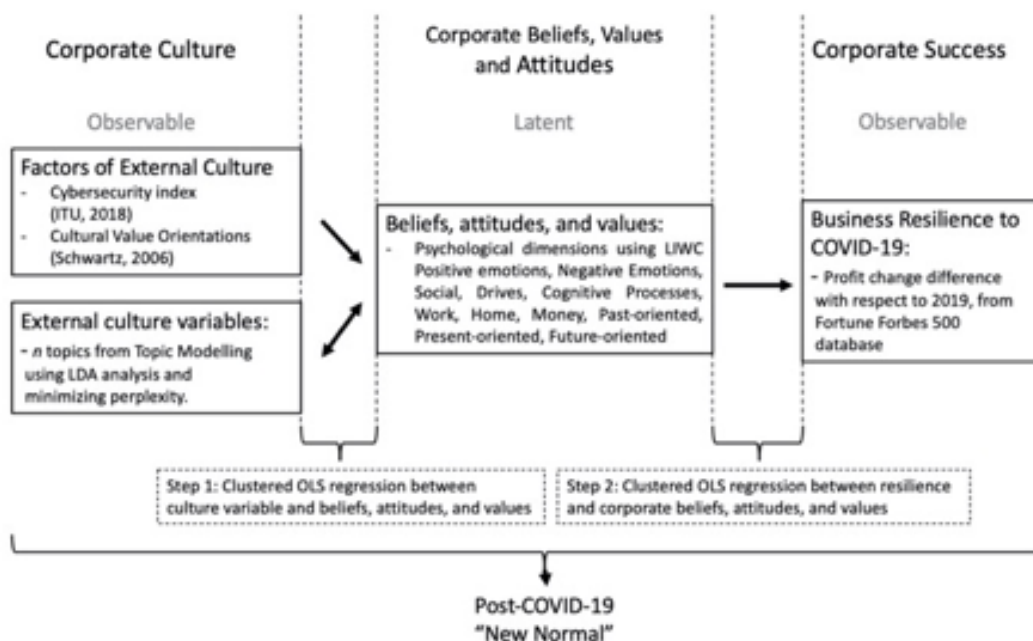


Figure 3.2: The CNNM structure, with details about data and processes

As proxies for external culture factors (Section 3.3.2), we include the Schwartz Cultural Value Orientation Coefficients (SCVOC) and the Global Cybersecurity Index (GCI).

Internal culture factors (Section 3.3.3) are obtained from a broad Topic Modelling exercise done for the COVID-19 response documents from the chosen companies. From the exercise we can then obtain the broad topics and the most influential words that the companies wrote in their documents.

Beliefs, attitudes, and values are also obtained from a second text analysis using the LIWC framework (Section 3.3.4). In this case we are looking for different psychological dimensions that the company expresses and reveals in its response document.

Finally, the Business resilience to COVID-19 (Section 3.3.5) is included by the difference of profit share of a company in 2020 with respect to 2019.

The testing of our model is divided in two steps. First, we analyse the links between the internal and external culture factors with corporate beliefs, attitudes and values using a clustered regression analysis. Second, we test the resulting significant beliefs, attitudes and values with the business success resilience observable, performing a second clustered regression.

We test the model for two different data sets. First, we focus our work in 217 US-based companies using the Forbes Fortune 500 database (focused on the USA). This allows to test the link between the internal culture values and the beliefs, attitudes, and values without considering the external cultural values. Given the nature of our data, we assume that the external culture is homogeneous for all US-based companies.

Once the link between internal culture values and beliefs, attitudes and values is tested, we used the Forbes Global Fortune 500 database to obtain a subset of the first 109 ranked international companies. In this second case we are including the external culture values (SCVOC and GCI).

For each data sets we present the results for different clusters of industry sectors. These clusters are formed by the similarity between the industries in our dataset, whilst ensuring a minimum number of companies so our analysis is done with enough data entries. The details of the different sectors can be consulted in Appendix F.

3.3.2 Factors of external culture

Two factors - Cultural Value Orientations as well as Exogenous Cybersecurity Culture – are used to capture external culture. While Cultural Value Orientations measures general culture in the company’s country of registration (headquarters’ country), Exogenous Cybersecurity Culture reflects the cyber security climate in that country.

Cultural value orientations coefficients for different countries [Schwartz, 2008] are used as proxies of Cultural Value Orientations factor. The Schwartz Value Survey (often referred to as SVS) asks respondents to rate 57 (general human) values according to their importance as a “*guiding principle*” of a respondent’s life on a scale from -1 to 7, where the answer “*opposed to my values*” scores -1; “*not important*” scores 0; “*important*” has a score of 3; and “*of supreme importance*” yields a

3.3. METHODOLOGY

score of 7. SVS split all obtained values into 6 broad value clusters: Embeddedness, Autonomy, Harmony, Mastery, Egalitarianism, and Hierarchy, each representing a dimension of the cultural value orientation. These value orientations are polarized and form pairs of antipodes [Schwartz, 2008]: Embeddedness (people consider themselves to be part of the collective) versus Autonomy (people view themselves as autonomous individuals); Harmony (desire to blend in with the nature) versus Mastery (desire to conquer the nature); and Egalitarianism (belief that all people are moral equals) versus Hierarchy (beliefs that hierarchy is necessary. Interestingly, according to Schwartz (2006) Autonomy can be Affective (concentration of own utility) and Intellectual (concentration in increasing own educational capability, following own ideas and creativity). These scores have been proved to be robust over time [Schwartz, 2008, Schwartz, 2013, Lee et al., 2010], and Professor Schwartz provided us the scores for 74 different countries.

Exogenous Cyber Security Culture is captured by the Global Cybersecurity Index (GCI) [International Telecommunication Union, 2017], calculated by the International Communications Union (ITU) based on 25 indicators forming 5 pillars (Legal, Technical, Organisational, Capacity and Coordination) from where 5 respective coefficients are computed. The GCI coefficient for a given country is then obtained by adding the 5 coefficients. According to ITU, the GCI allows to measure “cyber security commitment” in different countries. The index allows to understand the relative strength of commitment to cyber security governance and regulation in different parts of the world from hundreds of countries. Generally, the higher the index, the more committed a nation is to regulating and governing cyber security.

According to the ITU [International Telecommunication Union, 2017], the Legal and Technical pillars measurements are based “. . . on the existence of legal [technical] institutions and frameworks dealing with cyber security and cyber crime.” The Organisational pillar is based “. . . on the existence of policy coordination institutions and strategies for cyber security at national level.” The Capacity Building pillar contains “. . . measures based on the existence of research and development, education and training programs; certified professionals and public sector agencies fostering capacity building.” Finally, the Cooperation pillar is based on “. . . measures based on the existence of partnerships, cooperative frameworks and information sharing networks.”

This research used the latest version of the GCI index available in the public domain since it was released at the end of 2018 by the International Communications Union.

3.3.3 Factors of internal culture

In order to capture managerial mechanisms in the COVID-19 conditions, we use corporate COVID-19 response documentation, which summarises main corporate priorities and actions of the management, employees and customers at the time of the pandemic. To that end, the website of each company in our 317-companies' list was searched for the COVID-19 response documentation and the text of the main response document was copied and stored. This text was used to map managerial priorities during COVID-19 using the topic modelling exercise (e.g., [Hacker et al., 2020]).

We use a Latent Dirichlet Analysis (LDA) [Blei et al., 2003, Nikolenko et al., 2017] for topic modelling. To obtain the appropriate number of topics, we realised a Monte Carlo exercise (N=20) to know which number of topics minimises the perplexity. In other words, what is the minimum number of topics allowing to describe the complete set of documents we have. The technical details to find the number of topics and the details of the topics by themselves can be consulted in the Appendix F, as the frequency tables to see how these topics are distributed over the different industries.

For the US sample, the optimal number of topics is 6, these being: *Help and support to families, and donations to healthcare organisations; Office protocols for workers; Office protocols for clients and suppliers; Enabling Work from Home protocols; Financial statement for investors and markets; Economic waivers and insurance coverage.*

For the global sample, the optimal number of topics is 4, these being: *Ensure production chain and equipment at work; cost reduction, focus on income and production; Work from Home policy; Financial statement for investors and markets.*

3.3.4 Beliefs, attitudes, and values

To obtain the beliefs, attitudes and values, we use the Linguistic Inquiry and Word Count (LIWC) framework and software [Tausczik and Pennebaker, 2010, Pennebaker et al., 2007]. The LIWC allows to obtain the different psychological dimensions of a text, categorising the present words into different psychological emotions through an inbuilt dictionary. LIWC has already been successfully used in other COVID-19 studies, particularly studying emotions over time in Twitter [Dyer and Kolic, 2020], or in predicting results in German elections [Tumasjan et al., 2011].

As defined in [Huang and Pearlson, 2019a] in the original MIT CAMS paper, “beliefs, attitudes and values comprise the unwritten rules and therefore the culture of the organisation, ...”. In this case, the psychological dimensions present in the

COVID-19 responses can represent the unwritten attitude towards the pandemic.

LIWC allows to have a myriad of dimensions, ranging from “anxiety” to “body” and “power” dimensions. An advantage of LIWC is that a given score in a particular psychological dimension does not mean that the input text talks about *that* particular dimension. Instead, given the number of dimensions given, we can create a profile of that text by analysing a set of dimensions.

The dimensions chosen for this analysis are: *positive emotions; negative emotions; cognitive processes; drives; social; work; home; past-oriented; present-oriented* and *future-oriented*. The latter dimension is particularly interesting for our endeavour, as it is directly related to the preparation of a post-COVID-19 New Normal. Assuming most of these dimensions’ names are self-explanatory, we just expand on a couple of them. Cognitive processes are related to an argumentative communication, spanning from insights to disagreements. Drives, on the other side, suggest a communication using the narrator’s motivations, like power, ambition or hope. A set of words related to each of these dimensions can be consulted in the Supplementary Material, as they give a deeper insight about which information they reveal.

3.3.5 COVID-19 business resilience

The Forbes Fortune 500 databases contain different features from companies to create its ranking. These features include total assets, market value, number of employees, total revenues and profits, and these last two percent changes. To have an insight about a company’s resilience during the 2020 global pandemic year, we take as dependent variable the profit percent change between 2020 and 2019.

3.4 Processed data

In order to make the analysis more accessible to understand, we present in this Section a summary and visualisations of the different data formats used as input for our final Linear Regressions presented in Section 3.5. We are following the structure presented in Figure 3.2.

The whole analysis pipeline is done to the US-based companies data set and the global companies data set. Both databases are split by industry sector, while the global one is also divided in different geographical sub-regions. The detailed list is presented in Appendix F, although a summary table is also shown in Table 3.1.

In Table 3.2 we present the different variables used in the two steps of our analysis and that are shown in Figure 3.2.

Also, in In Figure 3.3 we visualise the process data for the LIWC dimensions. We observe the profiles for the different sectors at the US sample and the global sample.

3.4. PROCESSED DATA

Table 3.1: Subsets used for the CNN model from our data sets.

Data set	Subset variable	Subsets	N
US-based companies	Industry	Finance	39
		Food & Wholesale	16
		Health Care & Pharma	26
		Energy	23
		Chemicals	12
		Heavy Industry	49
		Services	13
		Retail	23
		Miscellaneous	16
Global companies	Industry	Finance	24
		Health Care & Pharma	14
		Energy	17
		Automotive	13
		Computing	14
		Conglomerate	15
		Construction	12
	Geographical region	USA	12
		Western Europe	12
		China	12
		Asia (without China)	12

The score plotted is how much do a sector reveals a psychological dimension with respect to the average in our database. That is, e.g., how much does the Retail sector in the US talks about positive emotions with respect to the other US based corporations. A score close to 0 does then mean that a particular sector reveals about the same as the average of a particular dimension.

3.4. PROCESSED DATA

Table 3.2: Different variables used in the CAMS-inspired New Normal Model.

Block in CNNM	Variable	Description	Use limitations
Factors of external culture	Cyber security index [International Telecommunication Union, 2017]	A real number between 0 and 1. A weighted average between the 5 pillars (numbers between 0 and 1): legal, technical, organizational, cooperation and capital building.	Only used for global data set.
	Culture Value Orientation coefficients [Schwartz, 2008]	7 real numbers between 0 and 10. One number for each of the dimensions (Harmony, embeddedness, hierarchy, mastery, affective and intellectual autonomy, and egalitarianism).	Only used for global data set.
Factors of internal culture	Topics from documents [Blei et al., 2003]	n different topics obtained from Topic Modelling. 6 Different for the US-based companies data set, 4 for the global companies one.	None.
Corporate Beliefs, Values and Attitudes	LIWC dimensions [Pennebaker et al., 2007]	10 real numbers between -1 and 1. One for each of the LIWC dimensions used: positive emotions, negative emotions, cognitive processes, drives, social, work, home, past-oriented, present-oriented, future-oriented.	None.
Corporate Success	Forbes companies profit change between 2019 and 2020	A percentage representing the change in a company's profit.	None.

3.4. PROCESSED DATA

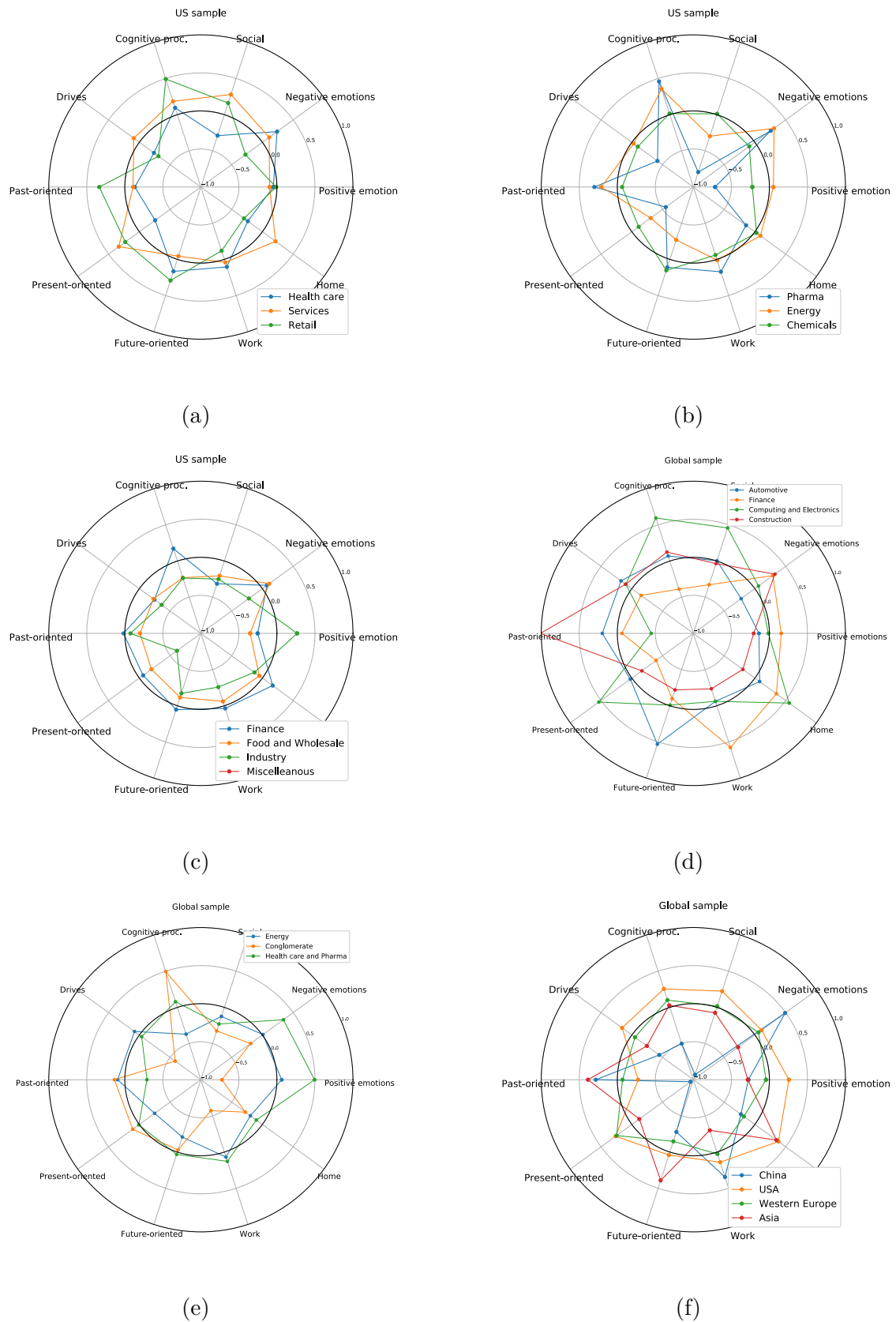


Figure 3.3: Profiles of different sectors and regions of the world by LIWC psychological dimensions. Left column depicts the US sample. The right column depicts the Global sample.

3.5 Results

Results are presented in an inverted way. Rather than examining the Step 1 shown in Figure 3.2, we first present results for the Step 2 of our process. This is to keep the thread of ideas easier to follow.

3.5.1 Business resilience with respect to beliefs, attitudes and values

We perform three different Clustered OLS regressions, one for the US sample and two for the global sample, respectively shown in Table 3.3, 3.4 and 3.5.

The heavy industry sector in the US has no statistically significance correlation with any of the psychological dimensions (Table 3.3). This is also true for Computing and Electronics and the Automotive sector in the global sample (Table 3.5). This is a nice reality check as both global sectors share an important number of companies with the US sector, and as these sectors largely depend on global supply chains, which could be interrupted by exogenous factors. In that sense, it is expected that none or a small number of psychological dimensions are significant for these sectors.

Positive and negative emotions are correlated on the business resilience according to our analysis. In the case of negative emotions, the Financial, Healthcare and Pharmaceutical and Miscellaneous sectors have a positive correlation when revealing negative emotions in their statements. This is complemented by the fact that the Finance sector is negatively affected when exposing positive emotions in the global sample. On the opposite side, Food and Wholesale, Energy and Chemicals are positively correlated when revealing positive emotions.

When focusing on the future-oriented attitude, we can only find a handful of companies in which this attitude is statistically significant. These are Chemicals, Services and Retail for the US sample, and Finance and Health Care and Pharma in the Global Sample. While Chemicals and Health Care and Pharma sectors have a positive impact with a future oriented vision, the other listed sectors have a negative impact. This is because of the different audiences towards their COVID-19 response documents are targeted. While Chemicals and Health Care and Pharmaceuticals are focused in both suppliers and clients, Finance, Services and Retail are targeted at clients.

Finally, at region level we could only find statistically significant coefficients for Asia (India, Russia, Saudi Arabia, Japan, Taiwan and South Korea companies). In that sense, no other region showed a statistically significant behaviour between the psychological dimensions and their business resilience.

Table 3.3: Results for a Clustered OLS regression for the US sample, by industry sector. In this case the dependent variable is the business resilience to COVID-19 (profit change with respect to 2019).

	Finance $R^2 = 0.16$	Food and Wholesale $R^2 = 0.86$	Health Care $R^2 = 0.49$	Energy $R^2 = 0.51$	Chemicals $R^2 = 0.97$	Industry $R^2 = 0.22$	Services $R^2 = 0.98$	Retail $R^2 = 0.49$	Misc. $R^2 = 0.37$
Positive emotions	0.48 (0.42)	1.88 (0.00)	-0.02 (0.92)	0.20 (0.00)	1.02 (0.00)	-0.39 (0.32)	13.15 (0.00)	-0.09 (0.63)	0.28 (0.49)
Negative emotions	0.22 (0.00)	0.20 (0.72)	0.50 (0.00)	-0.10 (0.34)	-0.09 (0.42)	0.68 (0.35)	-19.25 (0.00)	0.27 (0.39)	0.96 (0.00)
Social	0.23 (0.42)	0.39 (0.10)	0.08 (0.79)	0.07 (0.55)	4.20 (0.04)	-0.52 (0.43)	-64.41 (0.00)	0.28 (0.01)	-0.54 (0.79)
Cognitive processes	-0.17 (0.61)	1.20 (0.00)	-0.30 (0.24)	-0.12 (0.00)	-2.66 (0.04)	0.70 (0.36)	24.11 (0.00)	-0.04 (0.74)	-0.17 (0.78)
Drives	-0.42 (0.65)	-0.51 (0.01)	-0.44 (0.00)	-0.36 (0.13)	-1.94 (0.05)	0.23 (0.75)	36.77 (0.00)	0.15 (0.59)	.69 (0.66)
Past-oriented	0.10 (0.72)	0.46 (0.30)	0.10 (0.72)	0.06 (0.37)	-1.64 (0.05)	0.14 (0.78)	-14.28 (0.02)	-0.07 (0.53)	-1.15 (0.11)
Present-oriented	-0.22 (0.12)	-1.24 (0.00)	0.45 (0.60)	0.16 (0.00)	0.71 (0.08)	1.09 (0.11)	3.24 (0.62)	-0.08 (0.74)	0.44 (0.50)
Furutre-oriented	0.31 (0.39)	-0.34 (0.25)	0.30 (0.43)	0.15 (0.41)	1.53 (0.02)	-0.24 (0.50)	-12.12 (0.03)	-0.39 (0.00)	-0.17 (0.79)
Work	-0.13 (0.02)	-1.03 (0.00)	0.76 (0.01)	0.05 (0.71)	1.79 (0.04)	-0.48 (0.53)	-19.79 (0.00)	-0.36 (0.00)	-0.64 (0.58)
Home	-0.03 (0.81)	-0.80 (0.02)	-1.32 (0.17)	-0.07 (0.69)	-1.33 (0.01)	-0.52 (0.36)	25.13 (0.00)	-0.15 (0.11)	-0.21 (0.67)
Money	-0.02 (0.94)	-1.41 (0.00)	-0.29 (0.00)	0.20 (0.45)	-2.14 (0.10)	1.32 (0.22)	25.36 (0.00)	0.08 (0.56)	0.31 (0.80)

3.5. RESULTS

Table 3.4: Results for a Clustered OLS regression for the Global Sample, by region. In this case the dependent variable is the business resilience to COVID-19 (profit change with respect to 2019).

	China $R^2 = 0.47$	USA $R^2 = 0.22$	Western Europe $R^2 = 0.31$	Asia $R^2 = 0.98$
Positive emotions	0.041 (0.50)	0.943 (0.32)	-0.007 (0.49)	0.634 (0.00)
Negative emotions	-0.025 (0.31)	2.085 (0.25)	0.060 (0.29)	-0.258 (0.06)
Social	-0.089 (0.36)	-4.425 (0.13)	-0.147 (0.31)	2.123 (0.00)
Cognitive Processes	0.011 (0.86)	1.077 (0.37)	0.063 (0.06)	0.631 (0.00)
Drive	-0.020 (0.75)	0.682 (0.56)	0.118 (0.59)	-2.085 (0.00)
Past-oriented	-0.006 (0.91)	-1.234 (0.30)	0.196 (0.80)	-0.117 (0.01)
Present-oriented	-0.018 (0.86)	2.672 (0.46)	-0.071 (0.45)	-1.340 (0.00)
Future-oriented	-0.039 (0.49)	-1.094 (0.48)	-0.038 (0.87)	-0.449 (0.00)
Work	-0.064 (0.17)	-0.786 (0.81)	-0.072 (0.83)	-0.002 (0.89)
Home	0.122 (0.08)	0.525 (0.76)	0.164 (0.62)	-0.498 (0.00)
Money	0.024 (0.65)	0.286 (0.68)	-0.013 (0.70)	1.089 (0.00)

Table 3.5: Results for a Clustered OLS regression for the Global sample, by industry sector. In this case the dependent variable is the business resilience to COVID-19 (profit change with respect to 2019).

	Automotive $R^2 = 0.86$	Finance $R^2 = 0.56$	Energy $R^2 = 0.56$	Computing &Electronics $R^2 = 0.77$	Construction $R^2 = 0.98$	Conglomerate $R^2 = 0.73$	Health Care & Pharma $R^2 = 0.92$
Positive emotions	0.687 (0.44)	-0.177 (0.00)	0.059 (0.65)	-0.525 (0.27)	0.647 (0.22)	0.739 (0.79)	2.712 (0.39)
Negative emotions	-0.180 (0.76)	0.176 (0.00)	-0.073 (0.66)	0.811 (0.35)	-0.992 (0.33)	0.981 (0.42)	10.423 (0.00)
Social	2.121 (0.42)	0.092 (0.01)	-0.508 (0.00)	-0.508 (0.60)	1.568 (0.07)	1.798 (0.23)	13.068 (0.14)
Cognitive Processes	-0.833 (0.49)	0.208 (0.00)	-0.065 (0.02)	0.054 (0.93)	-2.318 (0.02)	1.888 (0.33)	-3.274 (0.42)
Drive	-1.690 (0.36)	0.000 (0.99)	0.095 (0.56)	0.314 (0.43)	2.478 (0.14)	-4.291 (0.00)	-2.399 (0.80)
Past-oriented	1.164 (0.59)	-0.354 (0.00)	-0.117 (0.16)	0.360 (0.16)	0.841 (0.16)	-0.768 (0.67)	-6.424 (0.00)
Present-oriented	0.224 (0.77)	-0.214 (0.00)	0.076 (0.65)	0.730 (0.50)	-0.387 (0.68)	1.431 (0.53)	-9.790 (0.08)
Future-oriented	-0.319 (0.45)	-0.468 (0.00)	0.313 (0.66)	0.047 (0.90)	1.388 (0.16)	-1.703 (0.05)	2.941 (0.00)
Work	0.896 (0.55)	0.053 (0.00)	-0.140 (0.27)	0.092 (0.73)	0.092 (0.27)	-1.120 (0.05)	-4.796 (0.56)
Home	-0.958 (0.26)	-0.019 (0.31)	0.255 (0.02)	0.504 (0.36)	1.631 (0.39)	-1.755 (0.00)	-16.138 (0.00)
Money	0.400 (0.44)	-0.123 (0.00)	0.065 (0.04)	0.309 (0.47)	-1.333 (0.07)	6.801 (0.02)	-8.095 (0.01)

3.5.2 Links between cultural factors and beliefs, attitudes and values

Given the research questions and the endeavour of this research, on this subsection we only focus our attention in those sectors and regions where a future-oriented attitude is statistically significant in Section 3.5.1. Therefore, we focus in: Chemicals, Services and Retail for the US sample, and Finance and Health Care and Pharma in the Global Sample. We also focus in the Asian-based companies (without including China).

In Table 3.6 we present the results for the US sample. In this case we are testing how influential are the 6 topics obtained from the COVID-19 response documents using a Topic Modelling analysis, using LDA.

In Table 3.7 and 3.8 we present the results for the Global sample. In this case we are also testing the external cultural factors, such as the Cultural Value Orientation coefficients and the Global Cybersecurity Index. For the Global sample, the number of topics is 4 and not 6.

Table 3.6: Results for a clustered OLS regression for the US Sample, by industry sector. In this case the dependent variable are the psychological dimensions and the tested variables are the topics drawn from the Topic Modelling exercise.

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home	Money
Help	1.314	-0.53	0.240	0.075	0.477	-0.165	-0.184	-0.527	0.291	-0.05	-0.45
Families	(0.00)	(0.00)	(0.38)	(0.84)	(0.29)	(0.51)	(0.41)	(0.11)	(0.34)	(0.86)	(0.13)
Protocol	-0.118	-0.71	-0.2	-0.59	-0.12	-0.115	-0.062	-0.352	-0.44	0.074	1.075
Office-Worker	(0.52)	(0.00)	(0.23)	(0.01)	(0.63)	(0.13)	(0.70)	(0.11)	(0.00)	(0.66)	(0.00)
Protocol	-0.412	-0.57	0.840	0.049	0.880	1.370	1.93	0.117	0.306	0.285	-0.36
Office-Client	(0.28)	(0.16)	(0.20)	(0.97)	(0.34)	(0.09)	(0.00)	(0.87)	(0.41)	(0.80)	(0.46)
Work from	-3.555	0.042	-1.6	-1.57	-0.83	-0.088	0.061	1.122	0.103	-0.20	-0.89
Home	(0.00)	(0.98)	(0.09)	(0.11)	(0.59)	(0.92)	(0.92)	(0.50)	(0.91)	(0.87)	(0.09)
Financial	2.281	1.671	-0.3	1.455	0.730	-0.770	-4.252	0.263	0.826	-1.12	3.688
Statement	(0.07)	(0.27)	(0.78)	(0.27)	(0.61)	(0.52)	(0.03)	(0.86)	(0.64)	(0.45)	(0.00)
Economic	-5.382	1.655	-0.2	5.625	-3.84	-0.223	2.876	3.428	-2.03	0.944	1.079
Waiver	(0.00)	(0.39)	(0.88)	(0.14)	(0.02)	(0.87)	(0.31)	(0.00)	(0.22)	(0.63)	(0.47)

Table 3.7: Results for a clustered OLS regression for the Global sample (Asia subset). In this case the dependent variable are the psychological dimensions and the tested variables are the topics drawn from the Topic Modelling exercise, the Cultural Value Orientation and the Global Cybersecurity Index (GCI).

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home	Money
Work from Home	1.244 (0.00)	-0.32 (0.00)	1.527 (0.00)	0.372 (0.33)	1.131 (0.00)	-0.68 (0.00)	0.861 (0.00)	-0.388 (0.43)	-0.95 (0.00)	0.618 (0.00)	-0.91 (0.00)
Ensure production	1.724 (0.00)	0.395 (0.00)	0.710 (0.00)	-0.75 (0.40)	2.140 (0.00)	0.989 (0.00)	-0.571 (0.00)	-0.726 (0.06)	-0.52 (0.36)	-0.20 (0.28)	-1.13 (0.00)
Financial Statement	1.667 (0.00)	-0.63 (0.07)	-3.95 (0.05)	-6.29 (0.00)	-1.52 (0.18)	-0.74 (0.87)	-4.237 (0.23)	-1.111 (0.21)	-0.47 (0.85)	-3.16 (0.00)	2.382 (0.00)
Cut costs	-6.62 (0.00)	1.897 (0.00)	-2.33 (0.27)	6.942 (0.06)	-9.70 (0.00)	12.733 (0.00)	1.078 (0.76)	11.763 (0.00)	2.573 (0.40)	-1.89 (0.16)	-5.77 (0.00)
Harmony	-0.27 (0.55)	-2.52 (0.00)	-4.89 (0.00)	-5.96 (0.00)	-1.61 (0.04)	-6.39 (0.04)	-4.497 (0.06)	-3.391 (0.00)	-0.06 (0.97)	1.143 (0.09)	3.523 (0.00)
Embedded	-0.11 (0.65)	1.049 (0.00)	2.292 (0.23)	3.943 (0.03)	0.389 (0.72)	5.619 (0.21)	2.397 (0.48)	-0.907 (0.27)	0.324 (0.89)	0.857 (0.30)	-3.87 (0.00)
Hierarchy	0.691 (0.01)	-1.15 (0.00)	-2.33 (0.09)	-4.45 (0.00)	0.478 (0.53)	-4.79 (0.13)	-2.995 (0.22)	-2.286 (0.00)	0.083 (0.96)	0.476 (0.43)	3.480 (0.00)
Mastery	0.718 (0.00)	-1.19 (0.00)	-2.49 (0.18)	-5.00 (0.01)	0.524 (0.61)	-5.89 (0.17)	-3.220 (0.33)	-1.480 (0.06)	-0.02 (0.99)	0.025 (0.98)	4.364 (0.00)
Affective Autonomy	-0.21 (0.00)	-0.06 (0.11)	-0.16 (0.32)	-0.03 (0.82)	-0.29 (0.00)	-0.38 (0.31)	0.004 (0.99)	0.650 (0.00)	-0.08 (0.73)	-0.26 (0.00)	0.210 (0.00)
Intelligent Autonomy	-0.76 (0.07)	1.516 (0.00)	2.899 (0.00)	4.801 (0.00)	-0.29 (0.60)	4.347 (0.04)	3.498 (0.03)	4.195 (0.00)	-0.22 (0.86)	-1.35 (0.00)	-2.87 (0.00)
Egalitarianism	-0.10 (0.69)	1.415 (0.00)	2.818 (0.02)	4.022 (0.00)	0.612 (0.37)	4.614 (0.10)	2.818 (0.19)	1.503 (0.01)	0.106 (0.95)	-0.31 (0.58)	-2.84 (0.00)
GCI	-0.31 (0.30)	1.535 (0.00)	3.047 (0.02)	4.657 (0.00)	0.325 (0.66)	5.108 (0.10)	3.269 (0.17)	2.251 (0.00)	0.018 (0.99)	-0.55 (0.36)	-3.30 (0.00)

Table 3.8: Results for a clustered OLS regression for the Global sample (sectors with future-oriented link). In this case the dependent variable are the psychological dimensions and the tested variables are the topics drawn from the Topic Modelling exercise, the Cultural Value Orientation and the Global Cybersecurity Index (GCI).

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home	Money
Work from Home	26.318 (0.79)	18.633 (0.86)	-135 (0.04)	49.941 (0.00)	-223.7 (0.00)	-79.22 (0.00)	-81.320 (0.00)	26.923 (0.65)	-135 (0.04)	-95 (0.10)	-226.5 (0.00)
Ensure production	27.602 (0.77)	16.199 (0.88)	-136 (0.04)	50.507 (0.00)	-224.7 (0.00)	-79.33 (0.00)	-82.004 (0.00)	28.583 (0.63)	-134 (0.03)	-97 (0.09)	-225.4 (0.00)
Financial Statement	25.150 (0.79)	18.074 (0.86)	-137 (0.04)	49.739 (0.00)	-225.1 (0.00)	-78.85 (0.00)	-84.048 (0.00)	27.359 (0.64)	-135 (0.03)	-96 (0.09)	-224.3 (0.00)
Cut costs	26.187 (0.79)	18.594 (0.86)	-136 (0.04)	48.798 (0.00)	-224.7 (0.00)	-78.89 (0.00)	-82.867 (0.00)	26.644 (0.65)	-135 (0.03)	-957 (0.10)	-224.0 (0.00)
Harmony	-5.25 (0.62)	1.285 (0.71)	2.486 (0.75)	6.768 (0.00)	2.890 (0.67)	3.391 (0.06)	2.229 (0.00)	8.270 (0.01)	-3.40 (0.32)	-8.2 (0.07)	-2.630 (0.55)
Embedded	-0.56 (0.94)	24.089 (0.00)	9.005 (0.10)	1.816 (0.04)	19.607 (0.00)	7.502 (0.00)	12.019 (0.00)	6.361 (0.02)	-8.15 (0.00)	-7.2 (0.09)	-13.02 (0.00)
Hierarchy	-1.45 (0.71)	-4.68 (0.01)	-0.47 (0.87)	2.745 (0.00)	-3.428 (0.20)	1.555 (0.06)	-1.458 (0.03)	1.337 (0.40)	1.021 (0.56)	-3.1 (0.07)	2.422 (0.12)
Mastery	-0.72 (0.85)	-15.1 (0.22)	11.202 (0.00)	-12.6 (0.00)	18.780 (0.00)	1.886 (0.44)	3.103 (0.33)	-11.52 (0.15)	23.446 (0.01)	23.917 (0.00)	37.445 (0.00)
Affective Autonomy	-1.16 (0.75)	8.704 (0.00)	3.936 (0.13)	2.717 (0.00)	7.676 (0.00)	2.625 (0.00)	5.329 (0.00)	4.350 (0.00)	-4.39 (0.00)	-5.1 (0.01)	-8.359 (0.00)
Intelligent Autonomy	2.391 (0.73)	-2.01 (0.83)	3.646 (0.51)	-10.6 (0.00)	7.711 (0.27)	1.121 (0.52)	0.569 (0.78)	-10.60 (0.07)	9.497 (0.15)	13.658 (0.00)	13.813 (0.00)
Egalitarianism	0.110 (0.98)	-6.88 (0.00)	2.919 (0.43)	0.039 (0.97)	2.953 (0.23)	1.884 (0.07)	1.014 (0.26)	-0.909 (0.69)	7.855 (0.00)	2.553 (0.41)	13.471 (0.00)
GCI	-5.29 (0.68)	-34.0 (0.00)	5.918 (0.52)	-1.07 (0.13)	-2.160 (0.76)	4.876 (0.04)	-6.992 (0.00)	-3.376 (0.37)	20.766 (0.00)	11.465 (0.06)	38.576 (0.00)

3.6 Discussion

The LIWC analysis allows to compare the different companies sectors through their profiles. We chose this method to capture the different sentiments, emotions and semantics from the *organisation culture* of each industry, given by the beliefs values and attitudes present in [Huang and Pearlson, 2019b].

Depending on the sector and the target reader of the COVID-19 response document, companies can reveal an inclination towards negative emotions (like Finance, Pharmaceuticals and Energy sectors in the US), or towards positive emotions (like the Heavy industry sector). This comparison can also be made by countries and regions. For example, from our results we can compare the profiles of Chinese-based companies with US-based companies. While in the US subset we find a heterogeneous distribution around the different psychological dimensions, in the China subset we find a more profiled narrative, centred between negative emotions, work and past-oriented sentences. This result must be carefully interpreted, as it does not mean that Chinese companies' discourses only talk about these dimensions. It rather means that, within our subset of companies, Chinese companies' response documents use more words or sentences related to these psychological dimensions *with respect to* the rest of our global sample of companies. For example, if in our global data set the word "fear", which is associated to the *negative emotion* dimension, is mentioned 5 times on average, then on the China subset it would be mentioned 7 times on average.

On the other hand, taking into account that US-based companies conform the most numerous subset of companies in our global set (38 out of 109), a more heterogeneous profile is expected as it captures a more diverse set of companies.

Looking at the complete analysis made, we can then have different insights about which are the top priority themes that are related to the post-pandemic thinking and the business resilience.

Only a handful of sectors have a statistically significant link between their resilience and a future-oriented narrative in their cultural aspects. These are the Services, Chemicals, Retail sectors in the US sample, and Finance and Healthcare and Pharma sectors in the global sample. For these sectors, the most important topics present in their COVID-19 responses related to future-oriented thinking are: *Ensuring economic waivers to their employees or clients* in the US sample, and *cutting costs in the company* for the global sample. The two latter are what we would see as managerial mechanisms from the CAMS model [Huang and Pearlson, 2019b]. Literature around Healthcare and Pharma during the pandemic discuss how these sectors had to adapt to the new conditions, be it as a whole system [Lanham et al., 2020], with respect to their workers [Tandale et al., 2020] or in the actual medi-

cal practice ([Cobianchi et al., 2020] more literature in Section 3.2). Although not specifically about Finance, [Harwood, 2020] talks about the different financial characteristics that businesses around the world should look out during the pandemic.

Although any of these two topics are directly connected to cyber security, we tested how the external culture around cyber security affected those global companies with a future-thinking narrative in the global sample. The correlation is statistically significant when talking about the Asian (without China) subset. This means that in those countries (Japan, Taiwan, South Korea, Russia, Saudi Arabia and India), the higher the Global Cybersecurity Index, the higher the probability there is that you have a future-oriented narrative in your COVID-19 response document. This is a very interesting result, as the companies in this subset heavily rely in technology, as it is with Samsung (SK), Toyota or Honda (Japan).

The fact of not finding direct mentions to cyber security in our documents is of no surprise taking into account how the literature about it and the “New Normal” was absent during the first months of the pandemic. Only a couple of works around privacy laws in the USA [Gerke et al., 2020] and cyber hygiene [Abukari and Bankas, 2020] were found.

Also, the above fact of not finding enough mentions about cyber security in our US and Global data sets make impossible to discuss, interpret and validate our CAMS-inspired New Normal Model. In order to do so, we would have had to observe more mention of cyber security in the documents that companies wrote at the beginning of the pandemic.

Focusing on the external cultural values -implemented via the Schwartz Cultural Orientation Coefficients-, we observe how a future-oriented narrative is present in those countries which have a more autonomic and a more egalitarian orientation. On the other side, a future-oriented narrative is less present for those countries with a more hierarchical and harmonical orientation.

Returning to our original idea of the CNNM, the above discussion can be synthesised with Figure 3.4. In it we can observe how the different elements of the CNNM can be seen: The managerial mechanisms by the topic modelling, the external factors by the GCI and the Schwartz CVO; the beliefs, values and attitudes by the LIWC dimensions in the centre and the behaviour represented by the financial profit of each sector/region on the right.

Topic Modelling has proven to be an efficient tool to extract information, as the different topics extracted to represent the managerial mechanisms have been used throughout the pandemic [Habersaat et al., 2020, Lanham et al., 2020, Harwood, 2020]. In that sense, the future-prone culture is only positively correlated

to economic waivers from American Energy companies having an economic waiver program as a managerial mechanism.

On the other hand, while for the Asian subset there is a negative correlation between the financial performance and the future prone culture, when looking at the global sample by sector we observe the global Healthcare and Pharma sector being positively correlated with the mentioned culture factor. This correlation is also positively correlated with to harmonical, autonomically affective, egalitarian societies where companies have a managerial mechanism of handling financial statements.

3.7 Conclusions

In the present research we adapt the MIT CAMS model [Huang and Pearlson, 2019a], which studies cyber security culture in companies, to the present COVID-19 pandemic. In particular, we are interested to test how companies are prepared to the “New Normal” [Habersaat et al., 2020] studying their COVID-19 responses. We call this adapted model the CAMS-inspired New Normal Model (CNNM). We intend to create a link between cyber security culture, preparation to the New Normal, and the companies’ resilience to the pandemic.

While the MIT CAMS model has been tested qualitatively, in this case we are taking a quantitative approach, using Natural Language Processing framework to process the COVID-19 response documents. We extract two different sources of information: Internal culture factors (Managerial processes, communication channels), and beliefs, attitudes and values (unwritten rules). The former information is extracted using Topic Modelling using Latent Dirichlet Allocation [Blei et al., 2003]. The latter is extracted using the LIWC framework [Tausczik and Pennebaker, 2010], which looks for different psychological dimensions in a given text using a predefined dictionary.

The first quarter of 2020 could be described by a single word: uncertainty. In that context, and with the available information, companies around the world had to write official response documents stating what could be the road map for the near future. The presented research used those texts to understand and analyse if and how those companies were preparing for the “New Normal”, and if they included a cyber security approach into it.

More than a year after the start of the pandemic, the business resilience was strong in those future-thinking companies that stated economic waivers for their employees and clients in the US or stated cost cuts throughout the company globally. This however is only true for a handful of sectors, like Chemicals, Retail and Services

3.7. CONCLUSIONS

in the US sample, and for Financial, Healthcare and Pharmaceuticals in the global sample.

No explicit mention for cyber security is found. This is particularly interesting, as the most cyber security-related measure of the pandemic, the transition to home office, was not an important topic in our analysis, even though it appeared in our topic modelling analysis. However, cyber security was found to be important as an external culture factor in Asia (without China), where an important cyber security culture has an impact in the future-thinking narrative. Given this absence of cyber security in most of our data set, we could not fully validate the CNNM and the different factors in it. By only obtaining one subset (Asia without China), not much can be said about the validity of our model.

Our research puts an emphasis on how companies around the world were reacting on a written form when the pandemic started, analysing who was the target, the topic to be communicated and the tone in which it was communicated with respect to different linguistic dimensions. We also focus in studying if any of those elements found in our text analysis is correlated to the resilience of the companies later into the future.

However, it is noticeable how cyber security is lacking in the general narrative, particularly in a context where the New Normal is defying the *status quo* of the worker-office dynamic, people commuting every day to a cyber-controlled space.

Our research contributes to the literature on quantitative methods to understand international businesses behaviour with respect to cyber security. Using both statistical analysis and Natural Language Processing, we can extend our analysis to different topics that could be intended to explore by any research team beyond cyber security. It also contributes by exploring official companies' documentation, which in this particular case is the official documents that companies released when the pandemic started. In that sense, our research explores how the content of these documents are correlated with future behaviour of the company by bridging Natural Language Processing, Economic data, econometric methods and Management literature and models around cyber security.

In general, we contribute with a methodology to understand the cyber security culture in companies around the world that is expandable to no matter the industrial sector or country where the Head Quarters reside.

As outlook of our work, it would be interesting to expand the database to a higher number of companies around the world. We could also find different sources of document in order to have a larger diversity of documents which would allow to make our framework more rich in analysis. This is in contrast to the main limitation of our work which was the lack of diversity of documents from where extract information.

3.7. CONCLUSIONS

Having, e.g., internal manuals of conducts or particular cyber security guidelines would have made our model much more accurate. We could also further expand our research by adapting a more appropriate “behaviour of cyber security resilience” other than the economic profit of the studied companies. This work could also be improved by using new techniques of Natural Language Processing, such as Weak Supervision [Ratner et al., 2018, Mekala and Shang, 2020, Fries et al., 2021], which allows to have labelled data without sacrificing time annotating.

Also, although this work is particularly targeted towards COVID-19 response documents, the same analysis could be done for other kind of corporate document, like purpose and vision statements.

3.7. CONCLUSIONS

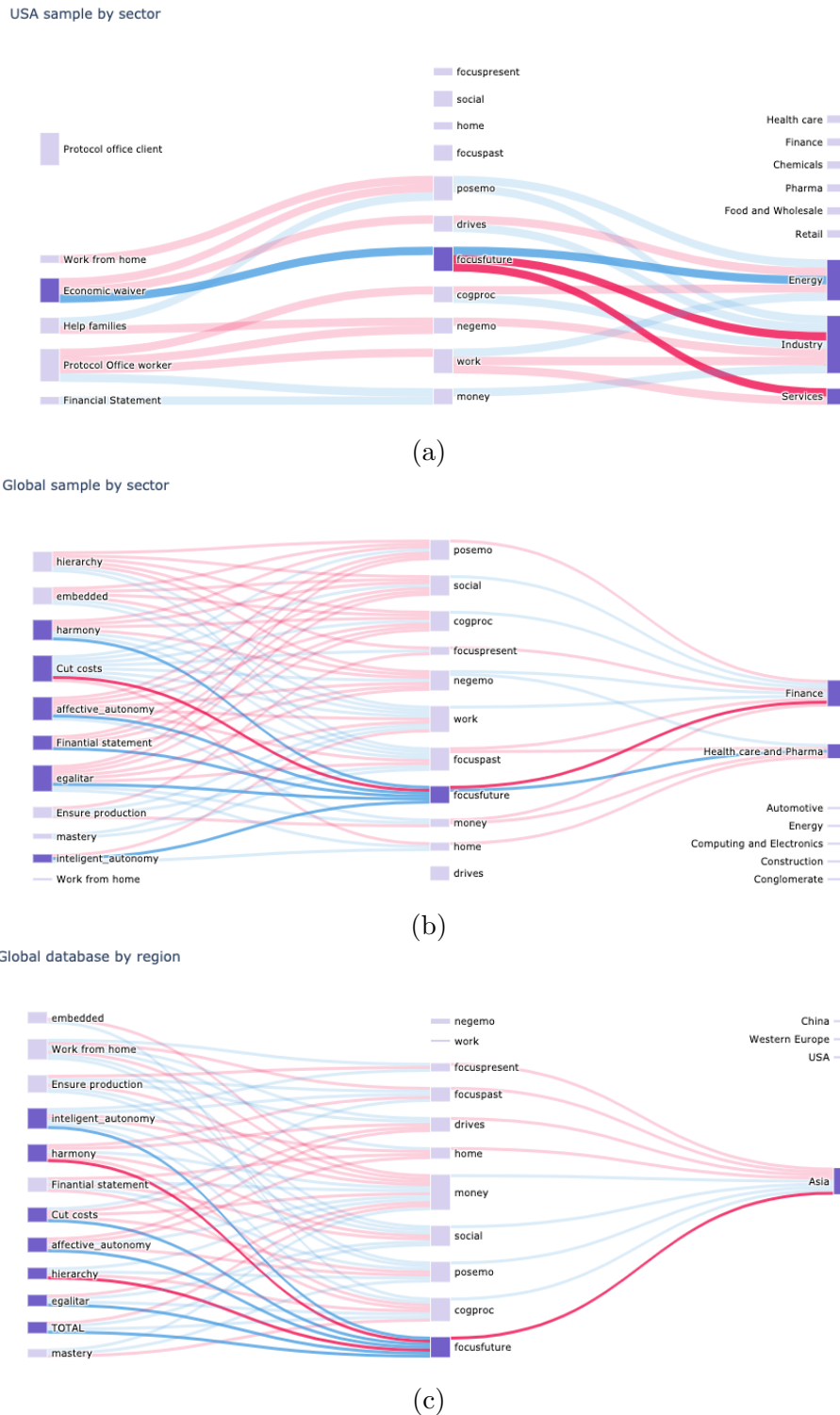


Figure 3.4: Sankey representation of the Tables obtained in Section 3.5. We highlight those statistically significant correlations ($p < 0.05$) with the LIWC dimension of Focus on the future as part of the organisational culture. The idea of these Sankey diagrams is to graphically represent the CNNM with the external factors and managerial mechanisms on the left, the values, beliefs and attitudes at the center of the figure and the behaviour on the right. Blue links refer to positive correlations while red links refer to negative correlations.

Chapter 4

Conclusions

The general objective throughout the present dissertation is to incorporate the analysis of human behaviour, be it of individuals or of collectives, into different security or cyber security systems. That way, it is expected, the system itself would be more efficient in finding and detecting potential threats. To prove it, we asked ourselves the question of *how can we improve such a system by taking into account the behaviour of particular individuals or collectives involved in it*, and analysed three different problems regarding security or cyber security.

With each of these three different problems and the resulting dissertation -which hopefully is greater than the sum of its parts-, we contribute to the literature in Behavioural Data Sciences, Behavioural Sciences, Cyber Security and Criminology literature by showing how, by implementing the behavioural analysis, an identifiable, quantifiable and measurable improvement can be obtained in each of the different studied cases.

In this concluding chapter, we will review each one of the three problems, summarising the findings, the limitations and the implications of the work done with respect to the particular research questions of the problem and the general ones of this dissertation. Once we go through every one of the chapters, a concluding remark stating the contributions and the limitations of this work is presented.

Chapter 1

In Chapter 1, the focus is put into understanding how different ways of communication when training employees to accurately detect cyber threats can change their behaviour, thus making them more efficient cyber security sensors. We work along the idea that humans are not “the weakest link” [Heartfield and Loukas, 2018] and can actually, when properly trained, become fundamental detectors of cyber threats. However, to “properly train” employees in such a subject can be extremely difficult and confusing, particularly when talking of a company with a heterogeneous workforce which accounts to more than 1000 employees. Although research has been done around *how* to train a workforce in cyber security [Rege et al., 2020, Yeoh et al., 2021, Sabillon et al., 2019], the approach taken is aligned with the “weakest link” idea, thus always suggesting a punitive approach towards employees highlighting the consequences than a successful attack could have.

However, is this punitive approach true? Does it work with an entire workforce, even though it might be incredibly diverse in terms of demographics and job skills and objectives? or can we actually find other ways of communicating that are non-punitive that also work to improve the workers' training? These are –paraphrasing them– the research questions we proposed ourselves to answer in Chapter 1.

We worked with one of the largest financial companies in the UK. We deployed a two staged study that included the Cyber-Domain-Specific-Risk-Taking (Cyber-DoSpeRT) scale test [Kharlamov et al., 2018] and a Human-as-a-Security-Sensor [Heartfield and Loukas, 2015] which was framed with three different communication treatments. These were a Positive one highlighting the benefits of a cyber security culture, a Negative one highlighting the consequences of a cyber attack, and a Neutral one, which is just the absence of a communication framing and that we took as control. From the CyberDoSpeRT scale we can apply a behavioural segmentation model based on the risk perception of cyber threats and the engagement (risk taking) with such potential threats. These 4 segments are called Anxious (high risk perception, low risk taking), Opportunistic (high risk perception, high risk taking), Relaxed (low risk perception, low risk taking) and Ignorant (low risk perception, high risk taking). The model was trained in the first survey and validated in the second. It is also in this second survey where we introduce to the surveyed employees a tailored Human-as-a-Security-Sensor test (HaaSS) with different communication framings.

Our results show that each of the 4 behavioural segments that are modelled into the company's workforce reacts differently to each one of the primings. While Anxious workers react better to a negative framing, Opportunistic employees react negatively to any kind of framing, thus preferring a neutral one. Ignorant workers, on the other hand, will react positively to any kind of non-neutral priming. Finally, we could not have any statistically significant result for Relaxed workers.

The work in Chapter 1 thus exposes how, in a heterogeneous workforce, complying with only one type of cyber security training approach does not have an homogeneous effect on the target audience. Moreover, the punitive approach does not only work homogeneously, but can actually be counter productive in a cyber security system that requires humans to accurately detect cyber threats.

Apart from the theoretical contribution described above, the work in Chapter 1 has immediate practical applications. Companies and institutions can test different ways of approaching their workforce when training it for cyber security in an ordered and easy way, ultimately calibrating the best way to communicate with their specific set of employees to efficiently train them to become accurate cyber security sensors.

As main limitation of this work, there was no the time to make a third survey which could allow us to validate the pipeline of the past two surveys. However,

the results obtained by the work done throughout the Chapter shows how being trained about cyber security with a behavioural analysis perspective, employees become better cyber security sensors, increasing their accuracy of correctly detecting potential cyber threats, thus having an overall positive impact on the cyber security system and making the interested company more secure.

Chapter 2

In Chapter 2, we study the illicit drug distribution problem of the County Lines Model (CLM) in England. The CLM has brought an important number of public health and modern slavery in deprived and vulnerable places around the UK [Black, 2020a, Black, 2020b, Silver and Intelligence, 2021]. Although the Government has done multiple efforts to tackle the problem from a national perspective [Bellis et al., 2019, Silver and Intelligence, 2021, Crime Agency, 2019], the logistic structure of the distribution network requires for the different law-enforcement bodies to have a good coordination, thus also needing of a good understanding of the problem. However, there has only been reports from the Metropolitan Police of London stating that County Lines operators guide themselves by an offer-demand principle [Rescue and Analysts, 2019, Rescue and Analysts, 2010]. Sociological and Anthropological literature around the public health and modern slavery issues that the CLM has brought to communities [Coombes, 2018, Andell and Pitts, 2018, Stone, 2018] speak differently, highlighting the fact that CL distribution points appear in small and remote towns in GB.

As such, in Chapter 2 we focus in going beyond the offer-demand principle and ask ourselves the question of what is the territorial logic of the CL operators in London. To reply to the latter question we need to ask also which are the social, geographic and demographic elements that can act as incentives or as costs for the operators.

Using data from the Metropolitan Police of London about the police territories where distribution hubs have been found throughout Britain, we train and cross-validate three different models of spatial analysis. The first one, the Gravity model [Noulas et al., 2012, Anderson, 2010] acts as a benchmark model given its understanding of flow of merchandise as proportional to the population and inversely proportional to the distance between two places. The second one, the Radiation model [Simini et al., 2012], understands the flow as a process of sorting “opportunities”, here represented as population, from one place to another. With this second model we want to test if the *distribution* of the population in England has something to do with how operators allocate their connections. Finally, the Retail model [Wilson, 2008] understands the flow from a place to another as process of selecting the

destination based on the incentives and the costs of going there with respect to the other competing places where to go. With the latter model, we include 6 different variables of different origin as input: the travel times from one place to another, the hospital admissions by misuse and by poisoning of drugs *per capita* (as proxy for number of illicit drug consumers), the knife crimes events *per capita*, the number of police officers working full-time *per capita*, and the gross disposable household income *per capita*.

Our results show that the Retail model is the best performing with respect to the data of the Metropolitan Police. Within the best performing model, the travel times, the hospital admissions by misuse and the knife crime events are the most significant variables as *costs* to the CL operators. This uncovers a territorial logic in which CL operators avoid places where violence between gangs exists, while also avoiding places where an excess of offer is also present. However, our best performing model could not replicate the phenomenon seen in the data from the Metropolitan Police. More the 90% of the lines detected in 2019 and 2020 were found in 16 out the 37 police territories. These 16 are what historically has been considered the “South of England”. The fact of finding this also uncovers an important element to understand the distribution network, as it might also indicate the presence of a poly-centric structure of the different gangs of CL throughout the UK, or a particularity of the consumers of these 16 territories for the CLM to operate.

To clarify the latter point, we would have needed an important volume of additional data. This was the main limitation of the work in Chapter 2. With more data from different police forces, in addition with the one obtained by the Met. Police, at a higher resolution and throughout a larger time frame, we could actually be able to respond to the points outlined above.

However, the limitations do not stop this work to have an important contribution. On the one hand, theoretically this is the first quantitative study around the County Lines Model known until now. We also contribute to the Spatial Analysis literature with an application of an important problem that affects the quality of life of the most vulnerable population in the UK. Practically, our contribution allows the different police bodies to have a better understanding of the problem, thus not only allowing them to build better strategies to bring down large CLM operations, but most importantly we have outlined a guideline of which data is necessary for them to better understand the problem. This by itself is an important practical contribution.

More generally, by understanding the territorial logic of the county lines operators from a behavioural perspective, we add more clarity to the whole security system in order to bring down entire operations that affect the quality of the population in Britain.

Chapter 3

In Chapter 3, we focus in the transition that companies had around cyber security when the COVID-19 pandemic started and we entered into a “New Normal” stage, characterised by a restricted mobility and remote working habits [Habersaat et al., 2020]. During this transition phase, the cyber security of companies and institutions could have been compromised, as the security system would depend on the personal cyber hygiene of the employees at home [Abukari and Bankas, 2020]. Although not many data around cyber security during that particular time frame exists, we are still interested in knowing if and how were the companies thinking about the transition to the “New Normal”, and if cyber security was one of their priorities.

Analysing the official documents that companies released when the pandemic started using Natural Language Processing, we ask ourselves the question which were the most important topics that companies were talking about at that point in time, and if cyber security was one of them. If so, can we detect a change in behaviour with respect to those that did not in their yearly financial outcomes?

We compiled and studied a database of the official COVID-19 response documents from different companies of the Top1000 US Forbes fortune ranking and Global Forbes fortune ranking. We frame the analysis using the Cybersecurity at MIT Sloan model (CAMS) [Huang and Pearlson, 2019b], which considers different external and internal elements as important factors to understand the cyber security culture and behaviour around it within a given company. As external cultural factors (different elements that are embedded in the culture of the place where the HQs are located) we implemented into our model the Cultural Value Orientation coefficient [Schwartz, 2009] and the Global Cybersecurity Index [International Telecommunication Union, 2017]. To proxy the Internal cultural factors, we do a Topic Modelling exercise with each one of the databases using Latent Dirichlet Allocation [Blei et al., 2003]. Finally, we test our model using a Linguistic Inquiry Word Count (LIWC) [Pennebaker et al., 2007] as internal managerial factors and the net change in profits as the emergent behaviour.

Our results show that only a handful of sectors had a future-oriented language present in their official documents. For the US-based sample, these were the Services, Chemicals and Retail sectors, and the Financial, Healthcare and Pharma sectors in the Global sample. Cyber security was not explicitly found in any of the documents, and only the Global Cybersecurity Index was significantly correlated to the future-oriented language in the Asian (excluding China) companies’ subset.

The limitations of this work were the lack of diversity of data that was available. A more appropriate data which would allow us to study the resultant behaviour in cyber security when a change in cyber security culture/internal or external factor

happens skewed the interpretation of the results obtained. An extension of this work should do in order to incorporate this element.

However, the work done has valuable contributions. From a theoretical point of view, we are contributing to the literature of understanding how the companies reacted when the pandemic started [Harwood, 2020, Caligiuri et al., 2020] from a top managerial point of view, with a disaggregated analysis by industrial sector and by region of the world. This is done by expanding the CAMS model [Huang and Pearlson, 2019b] to the context of the “New Normal” by quantifying each of its parts. Thus, the theoretical contribution of this work is to expand the CAMS model and present a “New Normal” extension, called the CNNM. From a practical point of view, the work in Chapter 3 can help the companies to understand *a posteriori* the outcomes of the pandemic from a cyber security point of view. The CNNM can be tested for different sets of data, thus making it transferable and scalable to different geographical and sectorial sets of industries.

Contribution, limitations, further work and concluding remarks

We have responded to our general research questions by looking at three different security and cyber security problems, improving the detection of potential threats by incorporating behavioural analysis from a modern data science perspective. We have used methods from Natural Language Processing, Statistical Learning, Spatial Analysis, Econometrics and qualitative models to uncover and explore the behaviour of gangs operating in Great Britain, in addition to employees and top managers in companies. By exploring and understanding how these subjects of analysis work within their respective parts of the security or cyber security system, we incorporate their behaviour to increase the accuracy of detection of threats.

In that sense, we manage to contribute to the literature by giving three solid examples of how to incorporate behavioural analysis in security and cyber security problems. Each one of them with theoretical and practical implications, which in particular for the case of Chapter 1 were already applied to the training of the financial institution’s training on cyber security.

With respect to the current published literature, our research compares in the following ways: The work in Chapter 1 contributes to the stream of literature that goes against the “weakest link” idea where humans should be taken of faulty and inaccurate nature, thus needing a punitive approach when training them, almost without repair [Heartfield and Loukas, 2018, Sabillon et al., 2019, Rege et al., 2020, Yeoh et al., 2021]. In that sense, our contribution is evidence of how not only humans can be accurate sensors of cyber threats, but also how the traditional approach to

train them is not necessarily the most adequate. This is something that has never been tested in cyber security until now. In Chapter 2 we apply different methods from Spatial Analysis and Statistical Learning to crime detection that have not been applied before [Dolliver et al., 2018] which help to understand the territorial logic of the operators with a small number of data available, with the potential to forecast or predict as more data is added to the problem. It is also the first quantitative work that has been done around County Lines Model, contributing to the qualitative research done until now [Stone, 2018, Coombes, 2018, Andell and Pitts, 2018]. The work in Chapter 2 also contributes by challenging the official public information around the CLM around the understanding of the operators, challenging the traditional 'offer-demand' statement to understand any underlying social factors the distribution network [Rescue and Analysts, 2019, Black, 2020a, Silver and Intelligence, 2021]. Finally, in Chapter 3 we contribute by understanding the transition towards the "New Normal" from a cyber security perspective. This has been moderately studied [Abukari and Bankas, 2020, Caligiuri et al., 2020, Gerke et al., 2020], as the "New Normal" was mostly studied from a Medicine perspective [Lanham et al., 2020, Jiang et al., 2020] and an online-education perspective [Dwivedi et al., 2020, Hacker et al., 2020]. We also contribute by extending the CAMS model [Huang and Pearson, 2019b] in a quantitative place, and by analysing the database of official COVID-19 documents using NLP.

The main limitation of this work is that, for the three chapters, one or more than one steps were missing to fully validate that the solutions obtained give the expected results. In Chapter 1, a third survey is missing to deploy the framed trainings and see their results. In Chapter 2, the police could implement the strategies to focus resources more intelligently. The results of this however could only be seen in a time scale of years. Finally, in Chapter 3, a more diverse number of data is missing in order to detect changes in cyber security culture within the different companies. Apart from this main limitation, there is also a second important limitation which is a lack of data available. This is particularly true for Chapter 2 and 3. However the limitations just described are also the source of the further work envisaged for this dissertation. For each one of the projects we understand which would be the next steps to round the results and fully have a deployment of the behavioural tools here produced to increase the security and cyber security of the respective institutions.

In [MacArthur et al., 2022], the authors comment how new technologies such as Artificial Intelligence and Data Sciences should not solely focus on automating what humans already do well, but rather do what humans cannot do. Going beyond this idea, new technologies can also support humans in doing what could poorly do to achieve it in a more correct and efficient way. Put it in another way, technology can

help humans in changing their behaviour towards being more healthy [Xu and Liu, 2020] or detect their emotions [Wang et al., 2004]. In the case of this dissertation we are centring ourselves in implementing these same technologies and methods to help humans detect different security threats. We take the perspective of a given institution or company to increase the security of a given system. However, trying to amplifying the scope of the research done, we are also helping humans to better detect cyber security threats, which could have beneficial consequences in their personal lives. We are helping British government instances on identifying underlying social dynamics which could be correlated with other social problems that could not see before the analysis done at that spatial resolution. Always looking at the ethics surrounding human data, a topic that was not discussed here given the anonymity of the data, but always taken into account, behavioural sciences have proven to be a game-shifting element in what is known as the 4th revolution [Skilton and Hovsepian, 2008]. The findings of this research, although modest, try to account for that.

References

- [Abukari and Bankas, 2020] Abukari, A. M. and Bankas, E. K. (2020). Some Cyber Security Hygienic Protocols For Teleworkers In Covid-19 Pandemic Period And Beyond. *International Journal of Scientific & Engineering Research*, 11(4):7. 23, 90, 94, 115, 124, 126
- [Alshaikh and Adamson, 2021] Alshaikh, M. and Adamson, B. (2021). From awareness to influence: toward a model for improving employees' security behaviour. *Personal and Ubiquitous Computing*, 25(5):829–841. 20, 29, 34, 35
- [Altmann, 2020] Altmann, E. G. (2020). Spatial interactions in urban scaling laws. *PLOS ONE*, 15(12):1–12. 73, 76, 77, 86, 88
- [Amemiya, 1978] Amemiya, T. (1978). The estimation of a simultaneous equation generalized probit model. *Econometrica*, 46(5):1193–1205. 42
- [Andell and Pitts, 2018] Andell, P. and Pitts, J. (2018). The end of the line? the impact of county lines drug distribution on youth crime in a target destination. *Youth & Policy*. 22, 24, 63, 67, 71, 122, 126
- [Anderez et al., 2021] Anderez, D. O., Kanjo, E., Amnwar, A., Johnson, S., and Lucy, D. (2021). The Rise of Technology in Crime Prevention: Opportunities, Challenges and Practitioners Perspectives. arXiv:2102.04204 [cs]. 18
- [Anderson, 2010] Anderson, J. E. (2010). The gravity model. Working Paper 16576, National Bureau of Economic Research. 64, 75, 85, 122
- [Anoushiravani et al., 2020] Anoushiravani, A. A., Barnes, C. L., Bosco, J. A., Bozic, K. J., Huddleston, J. I., Kang, J. D., Ready, J. E., Tornetta, P., and Iorio, R. (2020). Reemergence of Multispecialty Inpatient Elective Orthopaedic Surgery During the COVID-19 Pandemic: Guidelines for a New Normal. *The Journal of Bone and Joint Surgery. American Volume*, 102(14):e79. 93
- [Arcaute et al., 2016] Arcaute, E., Molinero, C., Hatna, E., Murcio, R., Vargas-Ruiz, C., Masucci, A. P., and Batty, M. (2016). Cities and regions in britain through hierarchical percolation. *Royal Society Open Science*, 3(4):150691. 83, 87
- [Arthur and Vassilvitskii, 2007] Arthur, D. and Vassilvitskii, S. (2007). k-means++: the advantages of careful seeding. In *SODA '07: Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035, New Orleans, Louisiana. Society for Industrial and Applied Mathematics. 41, 42, 59

- [Bahmani et al., 2012] Bahmani, B., Moseley, B., Vattani, A., Kumar, R., and Vasilvitskii, S. (2012). Scalable k-means++. *Proceedings of the VLDB Endowment*, 5(7):622–633. 42
- [Balakrishnan et al., 2020] Balakrishnan, A., Lesurtel, M., Siriwardena, A. K., Heinrich, S., Serrablo, A., Besselink, M. G. H., Erkan, M., Andersson, B., Polak, W. G., Laurenzi, A., Olde Damink, S. W. M., Berrevoet, F., Frigerio, I., Ramia, J. M., Gallagher, T. K., Warner, S., Shrikhande, S. V., Adam, R., Smith, M. D., Conlon, K. C., and E-AHPBA Scientific and Research Committee (2020). Delivery of hepato-pancreato-biliary surgery during the COVID-19 pandemic: an European-African Hepato-Pancreato-Biliary Association (E-AHPBA) cross-sectional survey. *HPB: the official journal of the International Hepato Pancreato Biliary Association*, 22(8):1128–1134. 93
- [Barjamovic et al., 2019] Barjamovic, G., Chaney, T., Coşar, K., and Hortaçsu, A. (2019). Trade, Merchants, and the Lost Cities of the Bronze Age*. *The Quarterly Journal of Economics*, 134(3):1455–1503. 71
- [Bellis et al., 2019] Bellis, A., Allen, G., and Audickas, L. (2019). Knife crime statistics. Technical report, House of commons Library. 65, 79, 122
- [Berlusconi, 2017] Berlusconi, G. (2017). The determinants of heroin flows in europe: A latent space approach. *Social Networks*, 51:104–118. 68, 69, 70
- [Berlusconi et al., 2016] Berlusconi, G., Calderoni, F., Parolini, N., Verani, M., and Piccardi, C. (2016). Link prediction in criminal networks: A tool for criminal intelligence analysis. *PloS one*, 11(4):e0154244. 68, 69
- [Bhugra, 2021] Bhugra, D. (2021). COVID-19 pandemic, mental health care, and the UK. *Industrial Psychiatry Journal*, 30(Suppl 1):S5–S9. 52
- [Bichler et al., 2017] Bichler, G., Malm, A., and Cooper, T. (2017). Drug supply networks: a systematic review of the organizational structure of illicit drug trade. *Crime Science*, 6(1):1–23. 67, 68
- [Biltgen and Ryan, 2016] Biltgen, P. and Ryan, S. (2016). *Activity-based intelligence : principles and applications*. Artech House electronic warfare library. Boston, [Massachusetts] : Artech House. 18
- [Black, 2020a] Black, C. (2020a). *Review of Drugs: evidence relating to drug use, supply and effects, including current trends and future risks*. Home Office. 21, 24, 63, 65, 66, 85, 87, 122, 126

REFERENCES

- [Black, 2020b] Black, C. (2020b). *Review of Drugs. Executive Summary*. Home Office. 21, 63, 66, 71, 85, 87, 122
- [Blais and Weber, 2006] Blais, A.-R. and Weber, E. U. (2006). A Domain-Specific Risk-Taking (DOSPERT) scale for adult populations. *Judgment and Decision Making*, 1(1):33–47. Place: US Publisher: Society for Judgment and Decision Making. 34
- [Blei et al., 2003] Blei, D., Ng, A., and Jordan, M. (2003). Latent Dirichlet Allocation. *The Journal of Machine Learning Research*, 3:601–608. 19, 23, 30, 40, 101, 104, 116, 124, 156
- [Bloomquist, 2020] Bloomquist, K. L. (2020). What is being revealed in 2020? *Dialog*, 59(3):184–187. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/dial.12585>. 93
- [Boivin, 2014] Boivin, R. (2014). Risks, prices, and positions: A social network analysis of illegal drug trafficking in the world-economy. 70
- [Bright et al., 2019] Bright, D., Koskinen, J., and Malm, A. (2019). Illicit network dynamics: The formation and evolution of a drug trafficking network. *Journal of quantitative criminology*, 35(2):237–258. 67, 68, 69, 70
- [Buyya et al., 2016] Buyya, R., Calheiros, R. N., and Dastjerdi, A. V., editors (2016). *Big data : principles and paradigms*. Cambridge, MA : Morgan Kaufmann. 18
- [Caligiuri et al., 2020] Caligiuri, P., De Cieri, H., Minbaeva, D., Verbeke, A., and Zimmermann, A. (2020). International HRM insights for navigating the COVID-19 pandemic: Implications for future research and practice. *Journal of international business studies*, 51(5):697–713. 29, 90, 94, 125, 126
- [Camber, 2020] Camber, R. (2020). 1,500 more county lines drugs gangs in just a year. *Daily Mail*, pages 40–41. 63
- [Caminha et al., 2017] Caminha, C., Furtado, V., Pequeno, T. H. C., Ponte, C., Melo, H. P. M., Oliveira, E. A., and Andrade, Jr, J. S. (2017). Human mobility in large cities as a proxy for crime. *PloS one*, 12(2):e0171609. 18, 69
- [Campana, 2016] Campana, P. (2016). Explaining criminal networks: Strategies and potential pitfalls. *Methodological Innovations*, 9:205979911562274. 18, 67
- [Canter, 2004] Canter, D. (2004). A partial order scalogram analysis of criminal network structures. *Behaviormetrika*, 31:131–152. 67

REFERENCES

- [Chandra, 2015] Chandra, S. (2015). Transnational cocaine and heroin flow networks in Western Europe: A comparison. 69, 70
- [Chandra et al., 2017] Chandra, S., Yu, Y.-L., and Bihani, V. (2017). How MDMA flows across the USA: evidence from price data. *Global Crime*, 18(2):122–139. 69, 70
- [Chen et al., 2017] Chen, X., Zeng, X., Wang, W., and Shao, G. (2017). Big data analytics for network security and intelligence. *Gongcheng Kexue Yu Jishu/Advanced Engineering Science*, 49(3):1–12. 32
- [Choi, 2018] Choi, T.-M. (2018). A System of Systems Approach for Global Supply Chain Management in the Big Data Era. *IEEE Engineering Management Review*, 46(1):91–97. Conference Name: IEEE Engineering Management Review. 18
- [Chowdhury et al., 2019] Chowdhury, N., Adam, M., and Skinner, G. (2019). The impact of time pressure on human cybersecurity behavior: An integrative framework. In *26th International Conference on Systems Engineering, ICSEng 2018 - Proceedings*, pages 1–10, Sydney, Australia. Institute of Electrical and Electronics Engineers Inc. 32
- [Cobianchi et al., 2020] Cobianchi, L., Pugliese, L., Peloso, A., Dal Mas, F., and Angelos, P. (2020). To a New Normal. *Annals of Surgery*, 272(2):e49–e51. 93, 115
- [Coles-Kemp et al., 2018] Coles-Kemp, L., Ashenden, D., and O’Hara, K. (2018). Why should I? Cybersecurity, the security of the state and the insecurity of the citizen. *Politics and Governance*, 6(2):41–48. 33
- [Coomber and Moyle, 2017] Coomber, R. and Moyle, L. (2017). The Changing Shape of Street-Level Heroin and Crack Supply in England: Commuting, Holidaying and Cuckooing Drug Dealers Across ‘County Lines’. *The British Journal of Criminology*, 58(6):1323–1342. 67
- [Coombes, 2018] Coombes, R. (2018). Brexit: the clock is ticking. *The BMJ*, 362:k4057. 21, 22, 24, 63, 65, 71, 122, 126
- [Crime Agency, 2019] Crime Agency, N. (2019). *Intelligence Assessment: County Lines Drug Supply, Vulnerability and Harm 2018*. 19, 21, 24, 63, 66, 87, 122
- [Da Veiga and Eloff, 2010] Da Veiga, A. and Eloff, J. H. P. (2010). A framework and assessment instrument for information security culture. *Computers & Security*, 29(2):196–207. 95, 96

REFERENCES

- [D’Arcy et al., 2009] D’Arcy, J., Hovav, A., and Galletta, D. (2009). User awareness of security countermeasures and its impact on information systems misuse: A deterrence approach. *Information Systems Research*, 20(1):79–98. 29
- [Davies et al., 2013] Davies, T. P., Fry, H. M., Wilson, A. G., and Bishop, S. R. (2013). A mathematical model of the London riots and their policing. *Scientific Reports*, 3:2045–2322. 71, 73, 74
- [De Cauwer and Somville, 2021] De Cauwer, H. and Somville, F. (2021). Health care organizations: Soft target during COVID-19 pandemic. *Prehospital and Disaster Medicine*, 36(3):344–347. 30, 35
- [Dolliver et al., 2018] Dolliver, D. S., Ericson, S. P., and Love, K. L. (2018). A geographic analysis of drug trafficking patterns on the TOR network. *Geographical review*, 108(1):45–68. 69, 70, 126
- [Doolittle, 2020] Doolittle, W. F. (2020). Could this pandemic usher in evolution’s next major transition? *Current biology: CB*, 30(15):R846–R848. 93
- [D’Orsogna and Perc, 2015] D’Orsogna, M. R. and Perc, M. (2015). Statistical physics of crime: A review. *Physics of Life Reviews*, 12:1 – 21. 18, 69
- [Dwivedi et al., 2020] Dwivedi, Y. K., Hughes, D. L., Coombs, C., Constantiou, I., Duan, Y., Edwards, J. S., Gupta, B., Lal, B., Misra, S., Prashant, P., Raman, R., Rana, N. P., Sharma, S. K., and Upadhyay, N. (2020). Impact of covid-19 pandemic on information management research and practice: transforming education, work and life. *International Journal of Information Management*, 55:102211. 22, 29, 90, 94, 126
- [Dyer and Kolic, 2020] Dyer, J. and Kolic, B. (2020). Public risk perception and emotion on Twitter during the Covid-19 pandemic. *Applied Network Science*, 5(1):1–32. Number: 1 Publisher: SpringerOpen. 101
- [Eardley, 2020] Eardley, I. (2020). A New Normal?: The COVID-19 pandemic has heralded different ways of working, triage of workload, collaborative research and cold-site surgery. *BJU international*, 126(2):215–217. 93
- [Einav et al., 2012] Einav, L., Finkelstein, A., Pascu, I., and Cullen, M. R. (2012). How General Are Risk Preferences? Choices under Uncertainty in Different Domains. *American Economic Review*, 102(6):2606–2638. 34
- [Espinal-Enríquez and Larralde, 2015] Espinal-Enríquez, J. and Larralde, H. (2015). Analysis of México’s narco-war network (2007–2011). *PLOS ONE*, 10(5):1–15. 69

REFERENCES

- [European Monitoring Centre for Drugs and Drugs Addiction, 2018] European Monitoring Centre for Drugs and Drugs Addiction (2018). *Perspectives on Drugs - Wastewater analysis and drugs: a European multi-city study*. European Monitoring Centre for Drugs and Drugs Addiction. 70
- [European Monitoring Centre for Drugs and Drugs Addiction, 2020] European Monitoring Centre for Drugs and Drugs Addiction (2020). *Perspectives on Drugs - Wastewater analysis and drugs: a European multi-city study*. European Monitoring Centre for Drugs and Drugs Addiction. 70
- [Fenton, 2021] Fenton, T. (2021). Regional gross disposable household income, uk: 1997 to 2019. Technical report, Office for National Statistics. 79
- [Flatley, 2019] Flatley, J. (2019). Police workforce england and wales statistics. Technical report, Home Office. 79
- [Forte and Power, 2007] Forte, D. and Power, R. (2007). The ultimate cybersecurity checklist for your workforce. *Computer Fraud & Security*, 2007(9):14–19. 90
- [Fries et al., 2021] Fries, J. A., Steinberg, E., Khattar, S., Fleming, S. L., Posada, J., Callahan, A., and Shah, N. H. (2021). Ontology-driven weak supervision for clinical entity classification in electronic health records. *Nature Communications*, 12(1):1–11. Number: 1 Publisher: Nature Publishing Group. 118
- [Gerke et al., 2020] Gerke, S., Shachar, C., Chai, P. R., and Cohen, I. G. (2020). Regulatory, safety, and privacy concerns of home monitoring technologies during COVID-19. *Nature Medicine*, 26(8):1176–1182. 23, 29, 90, 94, 115, 126
- [Giommoni et al., 2017] Giommoni, L., Aziani, A., and Berlusconi, G. (2017). How do illicit drugs move across countries? a network analysis of the heroin supply to europe. *Journal of Drug Issues*, 47(2):217–240. 69
- [Goutam, 2015] Goutam, R. K. (2015). Importance of Cyber Security. 18
- [Greenhow and Chapman, 2020] Greenhow, C. and Chapman, A. (2020). Social distancing meet social media: digital tools for connecting students, teachers, and citizens in an emergency. *Information and Learning Sciences*, 121(5/6):341–352. Publisher: Emerald Publishing Limited. 94
- [Habersaat et al., 2020] Habersaat, K. B., Betsch, C., Danchin, M., Sunstein, C. R., Böhm, R., Falk, A., Brewer, N. T., Omer, S. B., Scherzer, M., Sah, S., Fischer, E. F., Scheel, A. E., Fancourt, D., Kitayama, S., Dubé, E., Leask, J., Dutta, M., MacDonald, N. E., Temkina, A., Lieberoth, A., Jackson, M., Lewandowsky, S.,

REFERENCES

- Seale, H., Fietje, N., Schmid, P., Gelfand, M., Korn, L., Eitze, S., Felgendreff, L., Sprengholz, P., Salvi, C., and Butler, R. (2020). Ten considerations for effectively managing the COVID-19 transition. *Nature Human Behaviour*, 4(7):677–687. 19, 22, 92, 97, 115, 116, 124
- [Hacker et al., 2020] Hacker, J., vom Brocke, J., Handali, J., Otto, M., and Schneider, J. (2020). Virtually in this together – how web-conferencing systems enabled a new virtual togetherness during the COVID-19 crisis. *European Journal of Information Systems*, 29(5):563–584. Publisher: Taylor & Francis _eprint: <https://doi.org/10.1080/0960085X.2020.1814680>. 90, 91, 94, 101, 126
- [Handa et al., 2019] Handa, A., Sharma, A., and Shukla, S. K. (2019). Machine learning in cybersecurity: A review. *WIREs Data Mining and Knowledge Discovery*, 9(4):e1306. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/widm.1306>. 18
- [Hardyns and Rummens, 2018] Hardyns, W. and Rummens, A. (2018). Predictive policing as a new tool for law enforcement? recent developments and challenges. *European Journal on Criminal Policy and Research*, 24(3):201–218. 70
- [Harris et al., 2006] Harris, C. R., Jenkins, M., and Glaser, D. (2006). Gender differences in risk assessment: Why do women take fewer risks than men? *Judgment and Decision Making*, 1(1):48–63. Place: US Publisher: Society for Judgment and Decision Making. 34
- [Harwood, 2020] Harwood, M. (2020). Cash flow, forecasting and future action plans post Covid-19. *In Practice*, 42(6):357–360. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1136/inp.m2513>. 94, 115, 125
- [Heartfield and Loukas, 2015] Heartfield, R. and Loukas, G. (2015). A taxonomy of attacks and a survey of defence mechanisms for semantic social engineering attacks. *ACM Computing Surveys*, 48(3):1–39. 20, 24, 30, 31, 35, 37, 48, 58, 121
- [Heartfield and Loukas, 2018] Heartfield, R. and Loukas, G. (2018). Detecting semantic social engineering attacks with the weakest link: Implementation and empirical evaluation of a human-as-a-security-sensor framework. *Computers & Security*, 76:101–127. 18, 20, 24, 29, 30, 32, 35, 40, 58, 120, 125, 156
- [Heartfield et al., 2016] Heartfield, R., Loukas, G., and Gan, D. (2016). You are probably not the weakest link: Towards practical prediction of susceptibility to semantic social engineering attacks. *IEEE Access*, 4:6910–6928. 18, 24, 29, 31, 35, 58

REFERENCES

- [Hegemann et al., 2011] Hegemann, R. A., Smith, L. M., Barbaro, A. B., Bertozzi, A. L., Reid, S. E., and Tita, G. E. (2011). Geographical influences of an emerging network of gang rivalries. *Physica A: Statistical Mechanics and its Applications*, 390(21):3894 – 3914. 18, 69
- [Hesse and Rafferty, 2020] Hesse, M. and Rafferty, M. (2020). Relational Cities Disrupted: Reflections on the Particular Geographies of COVID-19 For Small But Global Urbanisation in Dublin, Ireland, and Luxembourg City, Luxembourg. *Tijdschrift voor Economische en Sociale Geografie*, 111(3):451–464. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/tesg.12432>. 93
- [Holtmann et al., 2020] Holtmann, G., Quigley, E. M., Shah, A., Camilleri, M., Tan, V. P., Gwee, K. A., Sugano, K., Sollano, J. D., Fock, K. M., Ghoshal, U. C., Chen, M., Dignass, A., and Cohen, H. (2020). "It ain't over ... till it's over!" Risk-mitigation strategies for patients with gastrointestinal diseases in the aftermath of the COVID-19 pandemic. *Journal of Gastroenterology and Hepatology*, 35(7):1117–1123. 93
- [Huang and Madnick, 2019] Huang, K. and Madnick, S. E. (2019). Does High Cybersecurity Capability Lead to Openness in Digital Trade? The Mediation Effect of E-Government Maturity within Cross-border Digital Innovation. *SSRN Electronic Journal*. 97
- [Huang and Pearlson, 2019a] Huang, K. and Pearlson, K. (2019a). For What Technology Can't Fix: Building a Model of Organizational Cybersecurity Culture. In *Proceedings of the 52nd Hawaii International conference on System Sciences*. 9, 91, 95, 96, 101, 116
- [Huang and Pearlson, 2019b] Huang, K. and Pearlson, K. (2019b). For What Technology Can't Fix: Building a Model of Organizational Cybersecurity Culture. In *Proceedings of the 52nd Hawaii International Conference on System Sciences*. 19, 23, 24, 25, 26, 33, 35, 95, 96, 114, 124, 125, 126
- [International Telecommunication Union, 2017] International Telecommunication Union (2017). *Global Cybersecurity Index 2017*. 100, 104, 124
- [Jang-Jaccard and Nepal, 2014] Jang-Jaccard, J. and Nepal, S. (2014). A survey of emerging threats in cybersecurity. *Journal of Computer and System Sciences*, 80(5):973–993. 90, 96
- [Jiang et al., 2020] Jiang, D. M., Berlin, A., Moody, L., Kumar, R., Hannon, B., Krzyzanowska, M. K., Dhani, N., Cole, H., Elliott, M., and Sridhar, S. S. (2020).

REFERENCES

- Transitioning to a New Normal in the Post-COVID Era. *Current Oncology Reports*, 22(7):73. 93, 126
- [John, 2019] John, E. (2019). Drug-related deaths by local authority, england and wales. Technical report, ONS. 173
- [Johnson et al., 2020] Johnson, T., Kanjo, E., and Woodward, K. (2020). Sensor Data and the City: Urban Visualisation and Aggregation of Well-Being Data. arXiv:2007.02674 [cs]. 18
- [Jürrens et al., 2009] Jürrens, E., Broering, A., and Jirka, S. (2009). A human sensor web for water availability monitoring. In *OneSpace 2009 : 2nd international workshop on blending physical and digital spaces on the internet*, Berlin, Germany. 32, 35
- [Kanungo et al., 2004] Kanungo, T., Mount, D. M., Netanyahu, N. S., Piatko, C. D., Silverman, R., and Wu, A. Y. (2004). A local search approximation algorithm for k-means clustering. *Computational Geometry*, 28(2):89–112. 42
- [Kharlamov et al., 2018] Kharlamov, A., Jaiswal, A., Parry, G., and Pogrebna, G. (2018). A cyber domain-specific risk attitudes scale to address security issues in the digital space. Working paper, The Alan Turing Institute, London. 20, 21, 25, 29, 34, 35, 37, 39, 41, 58, 121
- [Knoke, 2015] Knoke, D. (2015). *Emerging Trends in Social Network Analysis of Terrorism and Counterterrorism*, pages 1–15. American Cancer Society. 67
- [Krebs, 2002] Krebs, V. (2002). Mapping networks of terrorist cells. *Collections INSA*, 24. 18, 68
- [Kwak et al., 2020] Kwak, Y., Lee, S., Damiano, A., and Vishwanath, A. (2020). Why do users not report spear phishing emails? *Telematics and Informatics*, 48. 20, 29, 34, 58
- [Lanham et al., 2020] Lanham, N. S., Bockelman, K. J., and McCriskin, B. J. (2020). Telemedicine and Orthopaedic Surgery: The COVID-19 Pandemic and Our New Normal. *JBJS reviews*, 8(7):e2000083. 93, 114, 115, 126
- [Lappan et al., 2020] Lappan, S., Malaivijitnond, S., Radhakrishna, S., Riley, E. P., and Ruppert, N. (2020). The human–primate interface in the New Normal: Challenges and opportunities for primatologists in the COVID-19 era and beyond. *American Journal of Primatology*, 82(8):e23176. 93

REFERENCES

- [Larcher and Brierley, 2020] Larcher, V. and Brierley, J. (2020). Children of COVID-19: pawns, pathfinders or partners? *Journal of Medical Ethics*, 46(8):508–509. Publisher: Institute of Medical Ethics Section: Current controversy. 94
- [Lee et al., 2010] Lee, C.-T., Beckert, T. E., and Goodrich, T. R. (2010). The relationship between individualistic, collectivistic, and transitional cultural value orientations and adolescents’ autonomy and identity status. *Journal of Youth and Adolescence*, 39(8):882–893. 100
- [Lee et al., 2008] Lee, S., Won, D., and McLeod, D. (2008). Discovering Relationships among Tags and Geotags. *Proceedings of the International AAAI Conference on Web and Social Media*, 2(1):202–203. Number: 1. 42
- [Li and Wang, 2022] Li, H. and Wang, J. (2022). Collaborative annealing power k-means++ clustering. *Knowledge-Based Systems*, 255:109593. 42
- [Lie et al., 2020] Lie, S. A., Wong, S. W., Wong, L. T., Wong, T. G. L., and Chong, S. Y. (2020). Practical considerations for performing regional anesthesia: lessons learned from the COVID-19 pandemic. *Canadian Journal of Anaesthesia*, 67(7):885–892. 93
- [Lifestyles Team, 2019] Lifestyles Team, N. D. (2019). Statistics on drug misuse. Technical report, NHS. 79
- [Lum et al., 2006] Lum, C., Kennedy, L. W., and Sherley, A. (2006). Are counter-terrorism strategies effective? the results of the campbell systematic review on counter-terrorism evaluation research. *Journal of Experimental Criminology*, 2(4):489–516. 18, 68
- [MacArthur et al., 2022] MacArthur, B. D., Dorobantu, C. L., and Margetts, H. Z. (2022). Resilient government requires data science reform. *Nature Human Behaviour*, 6(8):1035–1037. 19, 26, 126
- [Macedo and Menting, 2019] Macedo, C. and Menting, J. (2019). Building a cybersecurity culture in the industrial control system environment. *International Journal of Information Security and Cybercrime*, 8:39–44. 97
- [Madeley, 2018] Madeley, G. (2018). County lines gangs using children to shift drugs in scotland. *Daily Mail*, pages 18–20. 22, 67

- [Marotta and Pearlson, 2019] Marotta, A. and Pearlson, K. E. (2019). A culture of cybersecurity at banca popolare di sondrio. In *Proceedings of the Twenty-fifth Americas Conference on Information Systems*, Cancun, Mexico. Americas Conference on Information Systems. 23, 24, 26, 33, 35, 58, 96, 97
- [Masucci et al., 2013] Masucci, A. P., Serras, J., Johansson, A., and Batty, M. (2013). Gravity versus radiation models: On the importance of scale and heterogeneity in commuting flows. *Phys. Rev. E*, 88:022812. 75
- [Mekala and Shang, 2020] Mekala, D. and Shang, J. (2020). Contextualized Weak Supervision for Text Classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 323–333, Online. Association for Computational Linguistics. 118
- [Mills et al., 2021] Mills, S., Albani, V., and Brown, H. (2021). Estimating the direct and indirect risks of becoming food insecure during the COVID-19 pandemic: a cross-sectional analysis using UK cohort data. *The Lancet*, 398(Special Issue):S67. 52
- [Morgan et al., 2020] Morgan, P., Asquith, P., Bishop, L., Raywood-Burke, G., Wedgbury, A., and Jones, K. (2020). A new hope: Human-centric cybersecurity research embedded within organizations. In A., M., editor, *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, pages 206–216, Copenhagen, Denmark. Springer Cham. 29
- [Morselli, 2007] Morselli, C. (2007). The efficiency/security trade-off in criminal networks. *Social Networks*, 29(1):143–154. 68
- [Moyle et al., 2019] Moyle, L., Childs, A., Coomber, R., and Barratt, M. J. (2019). Drugsforsale: An exploration of the use of social media and encrypted messaging apps to supply and access drugs. *International Journal of Drug Policy*, 63:101–110. 63
- [Mukherjee et al., 2020] Mukherjee, M., Chatterjee, R., Khanna, B. K., Dhillon, P. P. S., Kumar, A., Bajwa, S., Prakash, A., and Shaw, R. (2020). Ecosystem-centric business continuity planning (eco-centric BCP): A post COVID19 new normal. *Progress in Disaster Science*, 7:100117. 23, 94
- [Muliukha et al., 2020] Muliukha, V., Lukashin, A., Utkin, L., Popov, M., and Meldo, A. (2020). Anomaly detection approach in cyber security for user and entity behavior analytics system. In *ESANN 2020 - Proceedings, 28th European Symposium on Artificial Neural Networks, Computational Intelligence and*

REFERENCES

- Machine Learning*, pages 251–256, Bruges, Belgium. European Symposium on Artificial Neural Networks. 32
- [Ng et al., 2020] Ng, J. N., Cembrano, K. A. G., Wanitphakdeedecha, R., and Manuskiatti, W. (2020). The aftermath of COVID-19 in dermatology practice: What’s next? *Journal of Cosmetic Dermatology*, 19(8):1826–1827. 93
- [Nielsen and Nock, 2015] Nielsen, F. and Nock, R. (2015). Total Jensen divergences: Definition, Properties and k-Means++ Clustering. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2016–2020. arXiv:1309.7109 [cs, math]. 42
- [Nifakos et al., 2021] Nifakos, S., Chandramouli, K., Nikolaou, C., Papachristou, P., Koch, S., Panaousis, E., and Bonacina, S. (2021). Influence of human factors on cyber security within healthcare organisations: A systematic review. *Sensors*, 21(15). 30, 35
- [Nikolenko et al., 2017] Nikolenko, S. I., Koltcov, S., and Koltsova, O. (2017). Topic modelling for qualitative studies. *Journal of Information Science*, 43(1):88–102. 19, 23, 30, 40, 101, 156
- [Noulas et al., 2012] Noulas, A., Scellato, S., Lambiotte, R., Pontil, M., and Mascolo, C. (2012). A tale of many cities: Universal patterns in human urban mobility. *PLOS ONE*, 7(5):1–10. 73, 75, 85, 122
- [Pennebaker et al., 2007] Pennebaker, J., Chung, C., Ireland, M., Gonzales, A., and Booth, R. (2007). *The Development and Psychometric Properties of LIWC2007*. 18, 101, 104, 124
- [Pepin, 2018] Pepin, S. (2018). County lines exploitation in London. Technical report, House of Commons Library. 66
- [Piovani et al., 2018] Piovani, D., Arcaute, E., Uchoa, G., Wilson, A., and Batty, M. (2018). Measuring accessibility using gravity and radiation models. *Royal Society Open Science*, 5(9):171668. 71, 73, 74, 86, 88
- [Piovani et al., 2017] Piovani, D., Molinero, C., and Wilson, A. (2017). Urban retail location: Insights from percolation theory and spatial interaction modeling. *PLOS ONE*, 12(10):1–13. 73, 75, 76, 77
- [Podolny and Page, 1998] Podolny, J. M. and Page, K. L. (1998). Network forms of organization. *Annual Review of Sociology*, 24(1):57–76. 18, 67

REFERENCES

- [Pogrebna and Skilton, 2019] Pogrebna, G. and Skilton, M. (2019). *Navigating New Cyber Risks: How Businesses Can Plan, Build and Manage Safe Spaces in the Digital Age*. Palgrave MacMillan. 18
- [Raguvir and Babu, 2020] Raguvir, S. and Babu, S. (2020). Detecting anomalies in users - An UEBA approach. In *Proceedings of the International Conference on Industrial Engineering and Operations Management*, pages 863–876, Dubai, UAE. Industrial Engineering and Operations Management Society. 32
- [Ratner et al., 2018] Ratner, A., Hancock, B., Dunnmon, J., Sala, F., Pandey, S., and Ré, C. (2018). Training Complex Models with Multi-Task Weak Supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence. AAAI Conference on Artificial Intelligence*. arXiv. arXiv:1810.02840 [cs, stat]. 118
- [Rege et al., 2020] Rege, A., Nguyen, T., and Bleiman, R. (2020). A social engineering awareness and training workshop for stem students and practitioners. In *2020 IEEE Integrated STEM Education Conference (ISEC)*, pages 1–6, Princeton, NJ, USA. Institute of Electrical and Electronics Engineers. 19, 20, 24, 29, 33, 35, 56, 58, 60, 120, 125
- [Rescue and Analysts, 2010] Rescue and Analysts, R. P. (2010). Rescue and response county lines project. supporting young Londoners affected by county lines exploitation. 21, 22, 66, 85, 87, 122
- [Rescue and Analysts, 2019] Rescue and Analysts, R. P. (2019). Rescue and response county lines project. supporting young Londoners affected by county lines exploitation. 21, 22, 24, 65, 66, 85, 87, 122, 126
- [Retzlaff, 2020] Retzlaff, K. J. (2020). Lessons Learned From COVID-19 and the New Normal. *AORN Journal*, 112(3):212–215. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/aorn.13147>. 93
- [Reuter et al., 2020] Reuter, H., Jenkins, L. S., De Jong, M., Reid, S., and Vonk, M. (2020). Prohibiting alcohol sales during the coronavirus disease 2019 pandemic has positive effects on health services in South Africa. *African Journal of Primary Health Care & Family Medicine*, 12(1):e1–e4. 93
- [Revanth Filbert Raj and Babu, 2019] Revanth Filbert Raj, I. and Babu, S. (2019). User-entity behavior analytics (UEBA) - A systematic review of literatures. In *Proceedings of the International Conference on Industrial Engineering and Operations Management*, pages 3620–3621, Bangkok, Thailand. Industrial Engineering and Operations Management Society. 32

REFERENCES

- [Robinson, 2019] Robinson, G. (2019). Working county lines: Child criminal exploitation and illicit drug dealing in Glasgow and Merseyside. *International Journal of Offender Therapy & Comparative Criminology*, 63(5):694–712. 22, 63, 67, 71
- [Roman-Urrestarazu et al., 2018] Roman-Urrestarazu, A., Robertson, R., Yang, J., McCallum, A., Gray, C., McKee, M., and Middleton, J. (2018). European monitoring centre for drugs and drug addiction has a vital role in the UK’s ability to respond to illicit drugs and organised crime. *The BMJ*, 362. 65
- [Ross, 2017] Ross, S. M. (2017). *Introductory Statistics (4th ed.)*. Elsevier Inc. 47
- [Rostami and Mondani, 2015] Rostami, A. and Mondani, H. (2015). The complexity of crime network data: A case study of its consequences for crime control and the study of networks. *PLOS ONE*, 10(3):1–20. 18, 67, 69
- [Ruggiero and Khan, 2007] Ruggiero, V. and Khan, K. (2007). The organisation of drug supply: South asian criminal enterprise in the uk. *Asian Journal of Criminology*, 2(2):163–177. 70
- [Sabillon et al., 2019] Sabillon, R., Serra-Ruiz, J., Cavaller, V., and Cano, J. (2019). An effective cybersecurity training model to support an organizational awareness program: The Cybersecurity Awareness Training Model (CATRAM). A case study in Canada. *Journal of Cases on Information Technology*, 21(3):26–39. 19, 20, 24, 29, 33, 35, 120, 125
- [Sakurai and Chughtai, 2020] Sakurai, M. and Chughtai, H. (2020). Resilience against crises: COVID-19 and lessons from natural disasters. *European Journal of Information Systems*, 29(5):585–594. Publisher: Taylor & Francis eprint: <https://doi.org/10.1080/0960085X.2020.1814171>. 23, 90, 91, 94, 96
- [Salitin and Zolait, 2018] Salitin, M. and Zolait, A. (2018). The role of user entity behavior analytics to detect network attacks in real time. In *2018 International Conference on Innovation and Intelligence for Informatics, Computing, and Technologies (3ICT)*, pages 1–5, Sakhier, Bahrain. Institute of Electrical and Electronics Engineers Inc. 32
- [Sasse et al., 2001] Sasse, M., Brostoff, S., and Weirich, D. (2001). Transforming the ‘weakest link’ - a human/computer interaction approach to usable and effective security. *BT Technology Journal*, 19(3):122–131. 29
- [Schein and Schein, 1985] Schein, E. H. and Schein, P. (1985). *Organizational culture and leadership*. 95, 96

- [Schwartz, 2008] Schwartz, S. (2008). Cultural Value Orientations: Nature & Implications of National Differences : Psychology. *Journal of Higher School of Economics. Psychology. Journal of Higher School of Economics*, 5(2):37–67. 99, 100, 104
- [Schwartz, 2009] Schwartz, S. (2009). A theory of cultural value orientations: Explication and applications. *Comparative Sociology*, 109:173–219. 91, 124
- [Schwartz, 2013] Schwartz, S. (2013). National Culture as Value Orientations: Consequences of Value Differences and Cultural Distance. In *Handbook of the Economics of Art and Culture*, volume 2, pages 547–586. Elsevier. Journal Abbreviation: *Handbook of the Economics of Art and Culture*. 100
- [Sethi et al., 2020] Sethi, A., Swaminath, A., Latorre, M., Behin, D. S., Jodorovsky, D., Calo, D., Aroniadis, O., Mone, A., Mendelsohn, R. B., Sharaiha, R. Z., Gonda, T. A., Khanna, L. G., Bucobo, J. C., Nagula, S., Ho, S., Carr-Locke, D. L., Robbins, D. H., and New York Society for Gastrointestinal Endoscopy (2020). Donning a New Approach to the Practice of Gastroenterology: Perspectives From the COVID-19 Pandemic Epicenter. *Clinical Gastroenterology and Hepatology: The Official Clinical Practice Journal of the American Gastroenterological Association*, 18(8):1673–1681. 93
- [Shindler, 2008] Shindler, M. (2008). Approximation algorithms for the metric k-median problem. 42
- [Short et al., 2010] Short, M. B., Brantingham, P. J., Bertozzi, A. L., and Tita, G. E. (2010). Dissipation and displacement of hotspots in reaction-diffusion models of crime. *Proceedings of the National Academy of Sciences*, 107(9):3961–3965. 18, 69
- [Siegal et al., 2020] Siegal, D. S., Wessman, B., Zadorozny, J., Palazzolo, J., Montana, A., Rawson, J. V., Norbash, A., and Brown, M. L. (2020). Operational Radiology Recovery in Academic Radiology Departments After the COVID-19 Pandemic: Moving Toward Normalcy. *Journal of the American College of Radiology: JACR*, 17(9):1101–1107. 93
- [Siegel and Castellan, 1988] Siegel, S. and Castellan, N. (1988). *Nonparametric statistics for the behavioral sciences (2nd ed.)*. New York: McGraw-Hill. 47
- [Sillanpää and Hautamäki, 2020] Sillanpää, M. and Hautamäki, J. (2020). Social engineering intrusion: A case study. In *Proceedings of the 11th International Conference on Advances in Information Technology*, Bangkok, Thailand. Association for Computing Machinery. 34, 35, 58

REFERENCES

- [Silver and Intelligence, 2021] Silver, N. and Intelligence, B. (2021). Nclcc county lines strategic assessment 2020/2021. Technical report, National County Lines Coordination Centre. 21, 24, 63, 64, 87, 122, 126
- [Simini et al., 2012] Simini, F., González, M. C., Maritan, A., and Barabási, A.-L. (2012). A universal model for mobility and migration patterns. *Nature*, 484:96–100. 64, 71, 73, 75, 85, 122
- [Skilton and Hovsepian, 2008] Skilton, M. and Hovsepian, F. (2008). *The 4th Industrial Revolution*. Palgrave Macmillan. 18, 90, 127
- [Srinivas et al., 2019] Srinivas, J., Das, A. K., and Kumar, N. (2019). Government regulations in cyber security: Framework, standards and recommendations. *Future Generation Computer Systems*, 92:178–188. 18
- [Stanton et al., 2005] Stanton, J., Stam, K., Mastrangelo, P., and Jolton, J. (2005). Analysis of end user security behaviors. *Computers and Security*, 24(2):124–133. 29
- [Steingartner et al., 2021] Steingartner, W., Galinec, D., and Kozina, A. (2021). Threat defense: Cyber deception approach and education for resilience in hybrid threats model. *Symmetry*, 13(4). 32
- [Stone, 2018] Stone, N. (2018). Child criminal exploitation: ‘county lines’, trafficking and cuckooing. *Youth Justice*, 18(3):285–293. 22, 24, 63, 64, 67, 71, 122, 126
- [Supt Mat, 2020] Supt Mat, S. (2020). County Lines. 66, 87
- [Tamagnini et al., 2020] Tamagnini, G., Biondi, R., Ricciardi, G., Rutigliano, R., Trias-Llimós, S., Meuris, B., Lamelas, J., and Del Giglio, M. (2020). Cardiac surgery in the time of the novel coronavirus: Why we should think to a new normal. *Journal of Cardiac Surgery*, 35(8):1761–1764. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/jocs.14741>. 93
- [Tandale et al., 2020] Tandale, A. B., Khade, S. S., and Krishnakumar, K. (2020). Dental clinical practice changes needed during the COVID-19 pandemic: The ‘new normal’. *Journal of the Indian Medical Association*, pages 29–35. 93, 114
- [Tausczik and Pennebaker, 2010] Tausczik, Y. R. and Pennebaker, J. W. (2010). The Psychological Meaning of Words: LIWC and Computerized Text Analysis Methods. *Journal of Language and Social Psychology*, 29(1):24–54. Publisher: SAGE Publications Inc. 101, 116

REFERENCES

- [Taylor et al., 2020] Taylor, P. J., Dargahi, T., Dehghantanha, A., Parizi, R. M., and Choo, K.-K. R. (2020). A systematic literature review of blockchain cyber security. *Digital Communications and Networks*, 6(2):147–156. 18
- [Triyason et al., 2020] Triyason, T., Tassanaviboon, A., and Kanthamanon, P. (2020). Hybrid Classroom: Designing for the New Normal after COVID-19 Pandemic. In *Proceedings of the 11th International Conference on Advances in Information Technology*, IAIT2020, pages 1–8, New York, NY, USA. Association for Computing Machinery. 94
- [Tumasjan et al., 2011] Tumasjan, A., Sprenger, T. O., Sandner, P. G., and Welpe, I. M. (2011). Election Forecasts With Twitter: How 140 Characters Reflect the Political Landscape. *Social Science Computer Review*, 29(4):402–418. Publisher: SAGE Publications Inc. 101
- [Walton, 2020] Walton, G. (2020). COVID-19. The new normal for midwives, women and families. *Midwifery*, 87:102736. 93
- [Wang et al., 2004] Wang, Y., Ai, H., Wu, B., and Huang, C. (2004). Real time facial expression recognition with AdaBoost. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 3, pages 926–929 Vol.3. ISSN: 1051-4651. 18, 127
- [Weber et al., 2002] Weber, E., Blais, A.-R., and Betz, N. (2002). A domain-specific risk-attitude scale: Measuring risk perceptions and risk behaviors. *Journal of Behavioral Decision Making*, 15. 34, 37
- [Whitty, 2021] Whitty, M. (2021). Developing a conceptual model for insider threat. *Journal of Management and Organization*, 27(5):911–929. 33, 35
- [Wilke et al., 2014] Wilke, A., Sherman, A., Curdt, B., Mondal, S., Fitzgerald, C., and Kruger, D. J. (2014). An evolutionary domain-specific risk scale. *Evolutionary Behavioral Sciences*, 8(3):123–141. Place: US Publisher: Educational Publishing Foundation. 34
- [Wilson, 2008] Wilson, A. (2008). Boltzmann, Lotka and Volterra and spatial structural evolution: an integrated methodology for some dynamical systems. *Journal of The Royal Society Interface*, 5(25):865–871. 64, 73, 86, 122
- [Wilson, 2006] Wilson, A. G. (2006). Ecological and urban systems models: Some explorations of similarities in the context of complexity theory. *Environment and Planning A: Economy and Space*, 38(4):633–646. 74, 86

REFERENCES

- [Xu and Liu, 2020] Xu, X. and Liu, H. (2020). ECG Heartbeat Classification Using Convolutional Neural Networks. *IEEE Access*, 8:8614–8619. Conference Name: IEEE Access. 18, 127
- [Yang, 2020] Yang, K. (2020). Unprecedented Challenges, Familiar Paradoxes: COVID-19 and Governance in a New Normal State of Risks. *Public Administration Review*, 80(4):657–664. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/puar.13248>. 93
- [Yang et al., 2014] Yang, Y., Herrera, C., Eagle, N., and González, M. C. (2014). Limits of predictability in commuting flows in the absence of data for calibration. *Scientific Reports*, 4:5662. 71, 75
- [Yeoh et al., 2021] Yeoh, W., Huang, H., Lee, W.-S., Al Jafari, F., and Mansson, R. (2021). Simulated phishing attack and embedded training campaign. *Journal of Computer Information Systems*, pages 1–20. 19, 20, 24, 29, 34, 35, 56, 58, 60, 120, 125
- [Yong Wang et al., 2014] Yong Wang, Jinpeng Wei, and Vangury, K. (2014). Bring your own device security issues and challenges. In *2014 IEEE 11th Consumer Communications and Networking Conference (CCNC)*, pages 80–85, Las Vegas, NV. IEEE. 90, 94, 96
- [Zappi et al., 2012] Zappi, P., Bales, E., Park, J., Griswold, W., and Rosing, T. (2012). The CitiSense Air Quality Monitoring Mobile Sensor Node. *Proceedings of the IPSN'12 Workshop on Mobile Sensing*. 18
- [Zeegen et al., 2020] Zeegen, E. N., Yates, A. J., and Jevsevar, D. S. (2020). After the COVID-19 Pandemic: Returning to Normalcy or Returning to a New Normal? *The Journal of Arthroplasty*, 35(7S):S37–S41. 93
- [Zheng et al., 2014] Zheng, Y., Liu, T., Wang, Y., Zhu, Y., Liu, Y., and Chang, E. (2014). Diagnosing new york city’s noises with ubiquitous data. In *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, UbiComp '14, page 715–725, New York, NY, USA. Association for Computing Machinery. 32, 35
- [Zoghbi et al., 2020] Zoghbi, W. A., DiCarli, M. F., Blankstein, R., Choi, A. D., Dilsizian, V., Flachskampf, F. A., Geske, J. B., Grayburn, P. A., Jaffer, F. A., Kwong, R. Y., Leipsic, J. A., Marwick, T. H., Nagel, E., Nieman, K., Raman, S. V., Salerno, M., Sengupta, P. P., Shaw, L. J., Chandrashekhara, Y. S., and ACC Imaging Council (2020). Multimodality Cardiovascular Imaging in the Midst

REFERENCES

of the COVID-19 Pandemic: Ramping Up Safely to a New Normal. *JACC. Cardiovascular imaging*, 13(7):1615–1626. 93

Appendices

Appendix A

Chapter 1: Phase 1 survey and results

In this section we present the structure and results obtained from the Phase 1 of our study. This section is intended to complement the results shown in Section 1.4 to give a complete idea of the survey done for Phase 1. We remind that the survey was launched from October 19, 2020 to November 6, 2020. We obtained 605 employees, from which 503 individuals responded the complete survey¹.

Given the internal policies of the company we worked with, for some of the questions of the survey we cannot present the results obtained by the respondents. Thus, we only present those questions that are relevant to our study and are allowed to be published. Particularly, we present the results for the questions used for the regression analysis done in Section 1.4.1 (Tables 1.2-1.5).

In this case we are only presenting the answers for Phase 1, in contrast to some results in Section 1.4.3 where we present results for both survey. In the present Appendix A we are only omitting the results of both questions regarding the frequency and type of cyber attacks before and after the pandemic, as they have already been presented in Figure 1.12.

¹Given internal policies of the company, not all questions were compulsory to answer.

Age and Gender

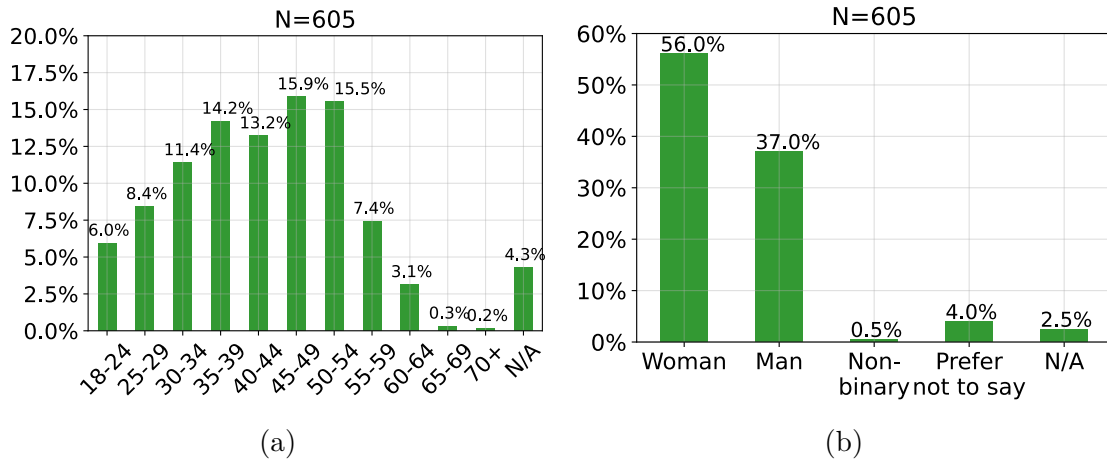


Figure A.1: Distribution of the employees that took part of the Phase 1 of the study by (a) age and (b) gender.

Life and Job Satisfaction

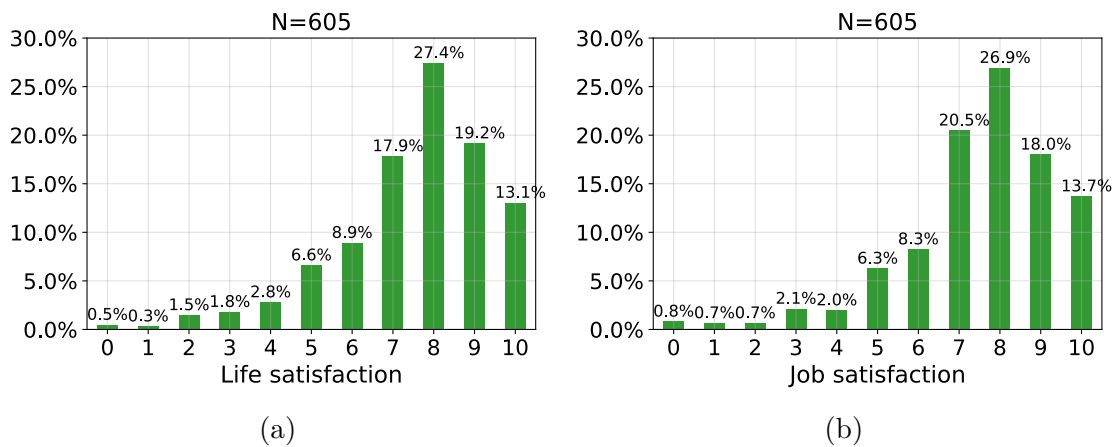


Figure A.2: Distribution of the employees that took part of the Phase 1 of the study by (a) life satisfaction and by (b) job satisfaction.

Confidence in cyber threats detection and cyber security role

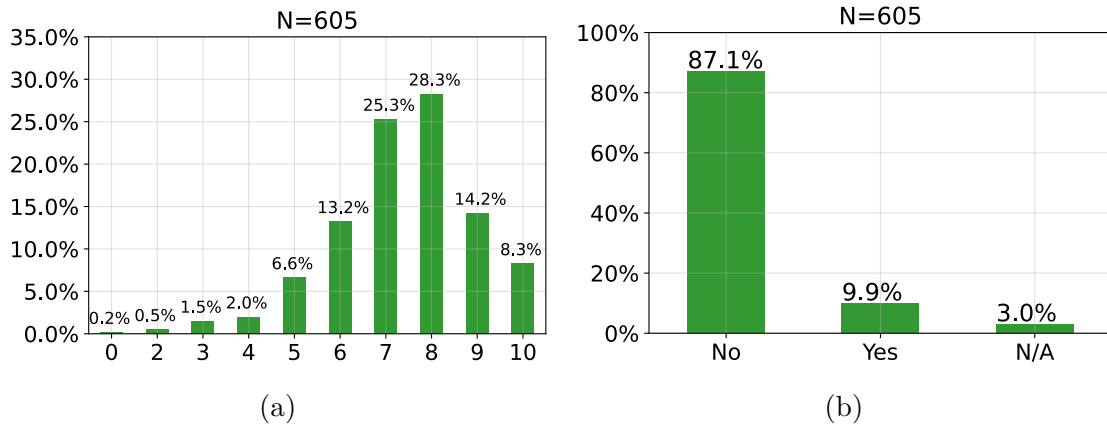


Figure A.3: Distribution of the employees that took part of the Phase 1 of the study by (a) confidence to detect cyber threats and (b) if they are in a cyber security role or not.

Frequency of working from home and over-work

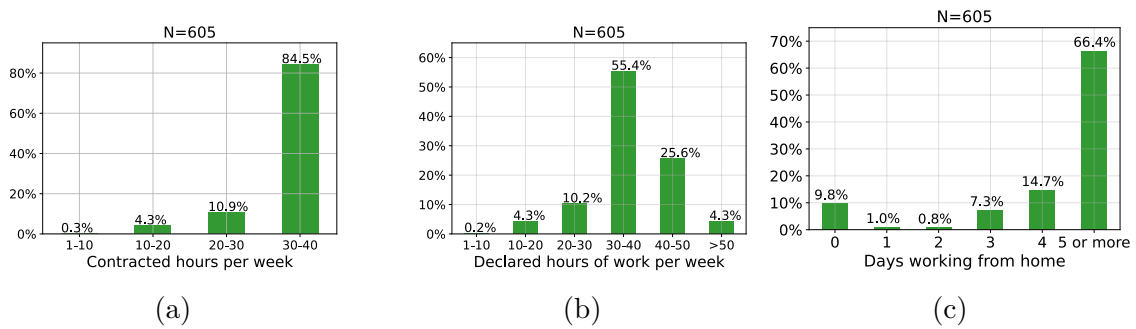


Figure A.4: Distribution of the employees that took part of the Phase 1 of the study by (a) contracted hours per week (b) actual working hours per week and (c) days working from home.

Years of experience in company and costumer-facing role

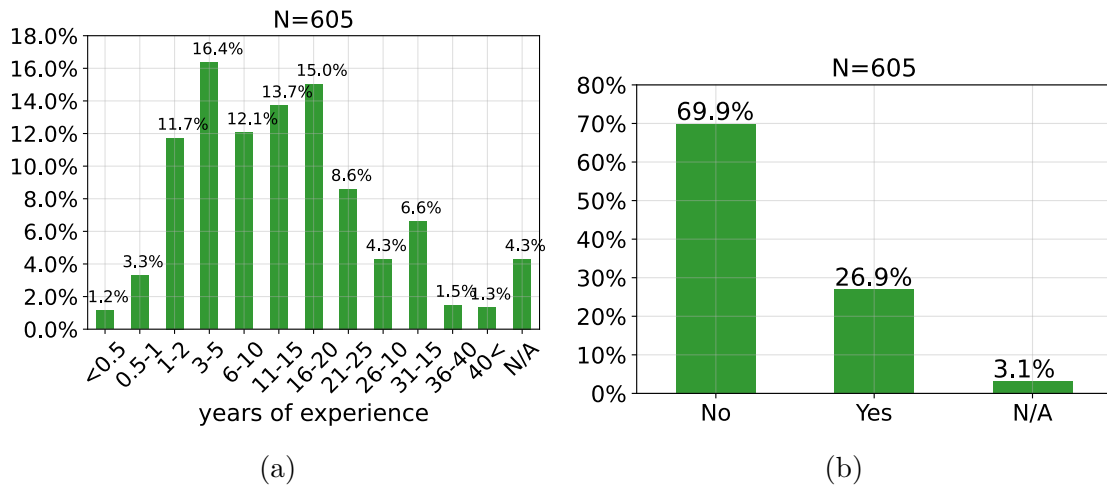


Figure A.5: Distribution of the employees that took part of the Phase 1 of the study by (a) years of experience in the company and (b) if they have a customer-facing role within the company.

Relationship status and caring responsibilities

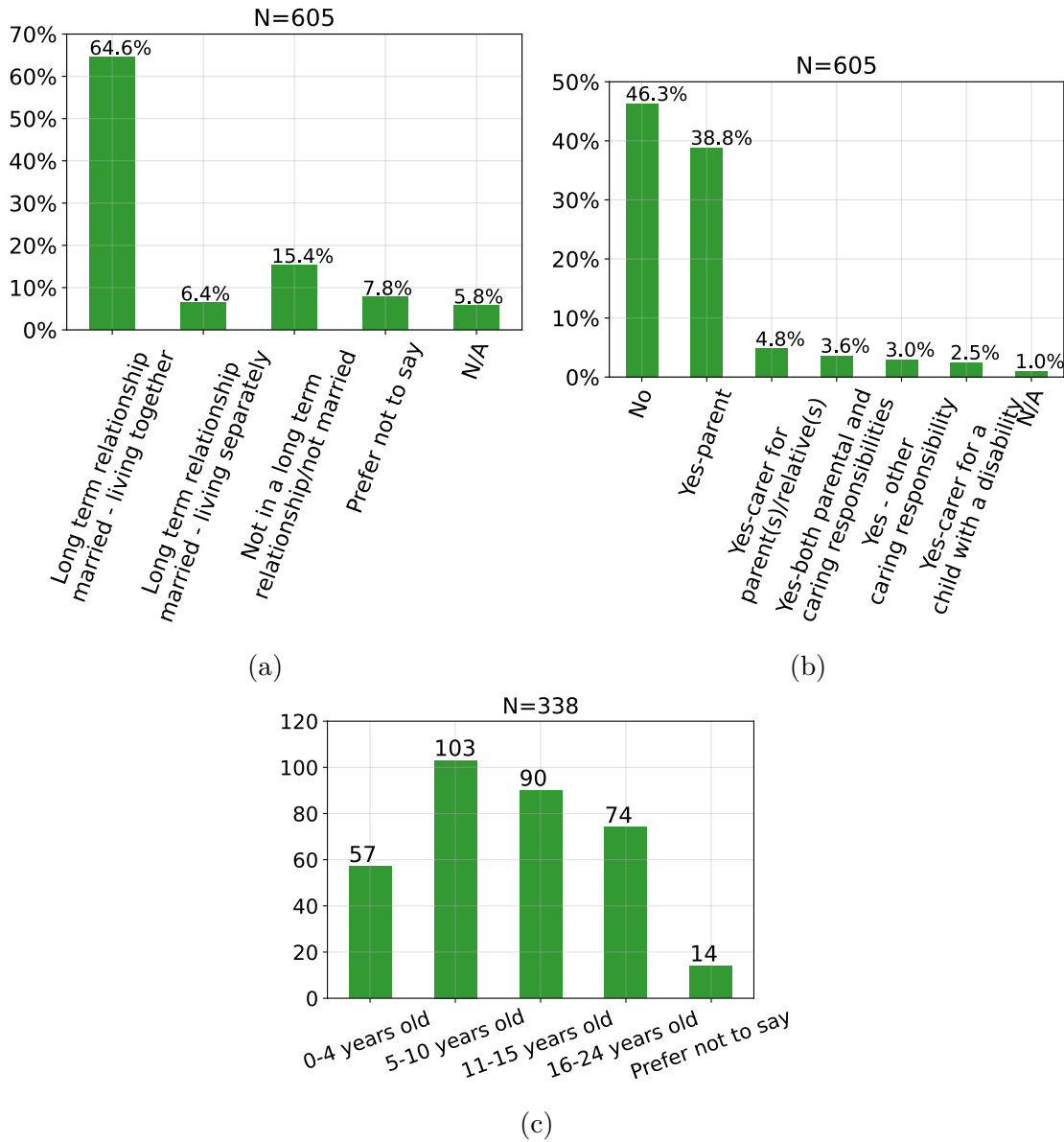


Figure A.6: Distribution of employees by (a) relationship status and by (b) caring responsibilities. We also present the ages of those cared ones in (c). Some employees had more than one cared ones.

Appendix B

Chapter 1: Phase 2 survey and results

In this section we present part of the results obtained from the Phase 2 of our study. This section is intended to complement the results shown in Section 1.4 to give a complete idea of the survey done for Phase 2. We remind that the second survey was launched from May 13, 2021 to June 3, 2021. We obtained 150 employees responses. In this case we are only focusing in the demographic and work-related questions found in the second survey. For a detailed analysis of the the Human-as-a-Cyber-Security-Sensor test, please refer to Appendix C.

Given the internal policies of the company we worked with, for some of the questions of the survey we cannot present the results obtained by the respondents. Thus, we only present those questions that are relevant to our study and are allowed to be published. Particularly, we present the results for the questions used for the regression analysis done in Section 1.4.1 (Tables 1.2-1.5).

In this case we are only presenting the answers for Phase 2, in contrast to some results in Section 1.4.3 where we present results for both survey.

Age and gender

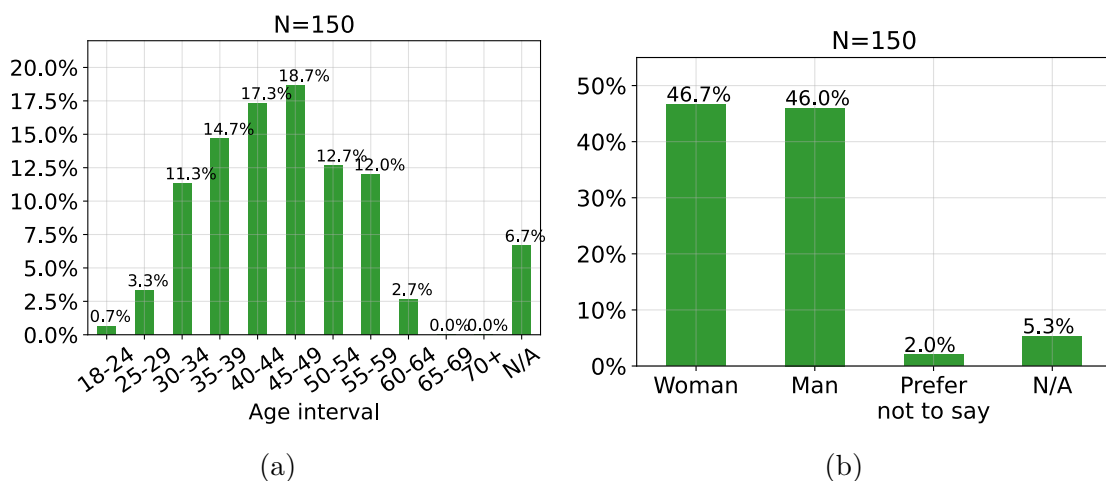


Figure B.1: Distribution of the employees that took part of the Phase 2 of the study by (a) age and (b) gender.

Life and Job Satisfaction

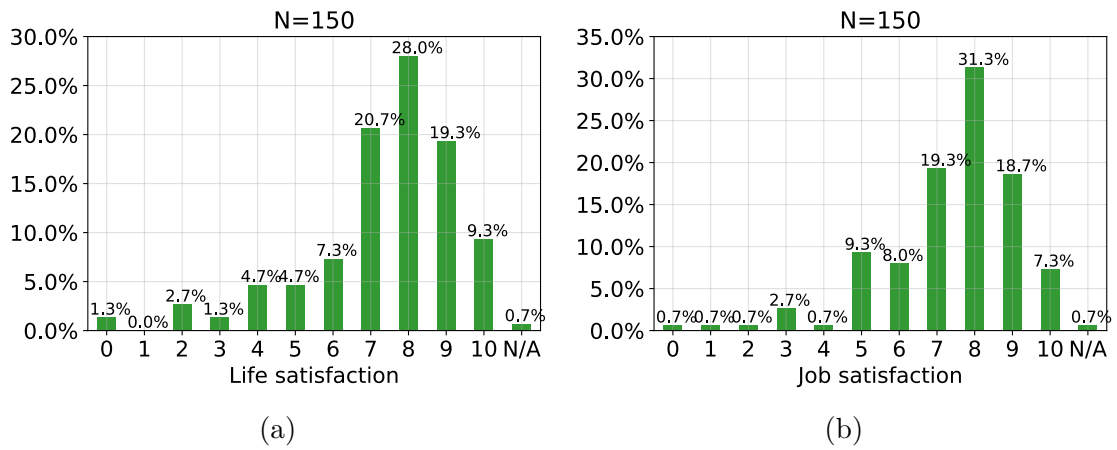


Figure B.2: Distribution of the employees that took part of the Phase 2 of the study by (a) life satisfaction and by (b) job satisfaction.

Confidence in detecting cyber threats and costumer-facing role

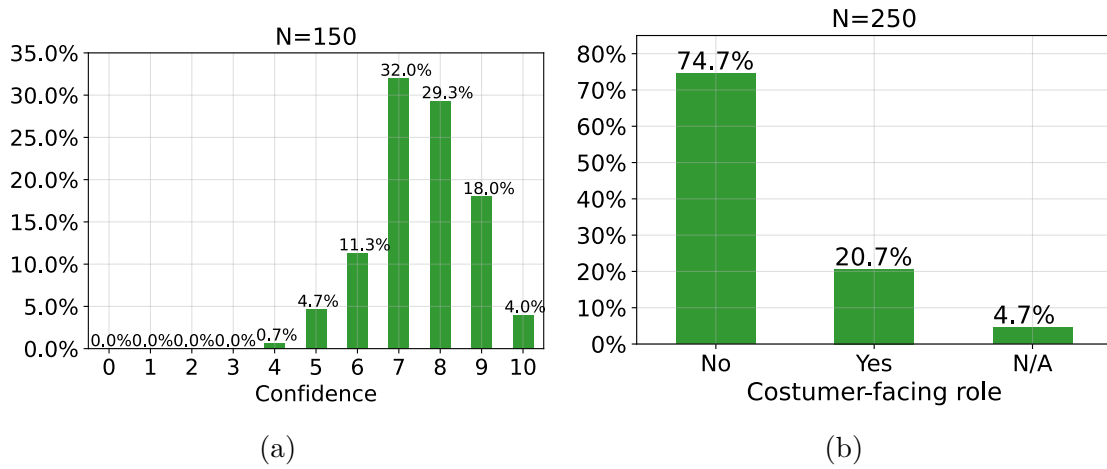


Figure B.3: Distribution of the employees that took part of the Phase 2 of the study by (a) confidence to detect cyber threats and (b) if they have a costumer-facing role within the company.

Frequency of working from home and over-work

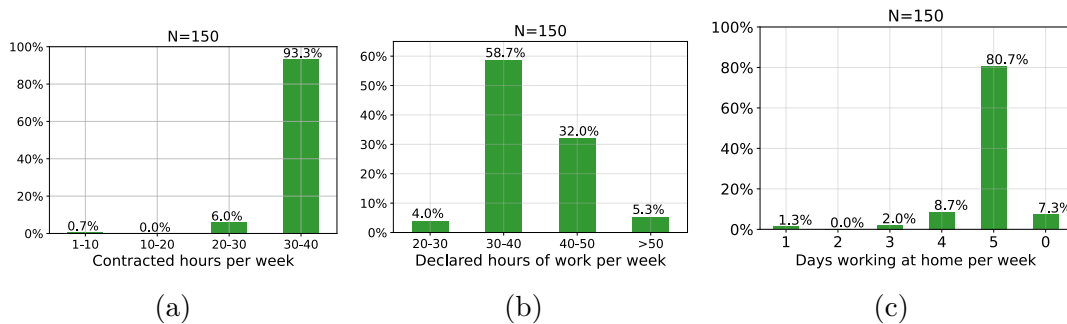


Figure B.4: Distribution of the employees that took part of the Phase 2 of the study by (a) contracted hours per week (b) actual working hours per week and (c) days working from home.

Appendix C

Chapter 1: Human-as-a-Cyber-Security-Sensor test

As part of the second phase of the study, the participants had to answer a Human-as-a-Cyber-Security-Sensor test [Heartfield and Loukas, 2018]. The HaaCSS test is composed of 6 different questions related to 6 different kinds of cyber security threats. For each one of them, a screenshot is presented to the participants, followed by 3 different questions:

1. Is this a start of a cyber attack? (Yes/no)
2. In a scale from 0 to 10, being 0 not confident at all and 10 being completely confident, how confident are you of your previous answer?
3. Briefly explain how you formed your opinion.

In the following we present each of the 6 questions with their respective results. In the case of questions 1 and 2, we present the results for each of the communication treatments. In the case of question 3, we present the results for the overall surveyed population.

The data recorded for question 3 is analysed using Latent Dirichlet Allocation for Topic Modelling [Blei et al., 2003, Nikolenko et al., 2017]. This analysis allow us to obtain a specific view of the different topics covered by the surveyed population with their respective keywords (most frequent and important words) and their frequency (how much these topics are talked about).

Exercise 1

The first example is a screenshot of a Facebook wall from a proxy profile showing a video of the World Health Organization with information regarding the COVID-19 pandemic. This screenshot does not represent the start of a cyber-attack, yet, many social media links (especially those when the user is nudged to click on a video) may represent social media masquerading. The screenshot presented to the respondents is observed in Figure C.1, while the results for the first two questions are presented in Figure C.2. Results for the topic modelling are presented in Table C.1.

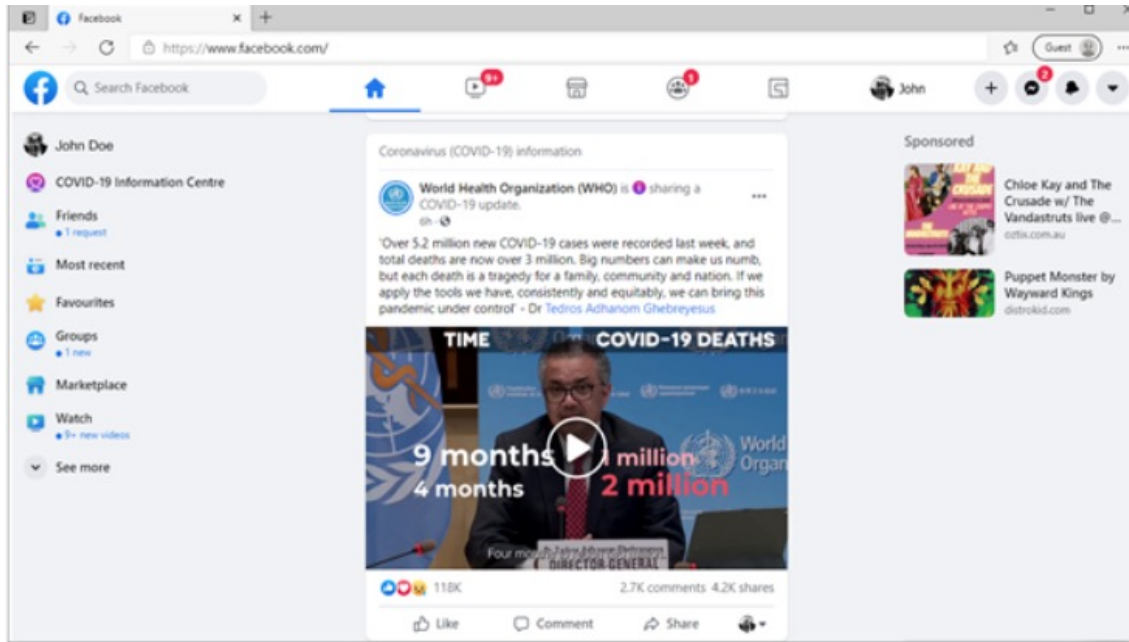


Figure C.1: Screenshot presented as first question at the Human-as-a-Cyber-Security-Sensor test.

Does represent the start of a cyber attack? No

Type of cyber attack: Social media masquerading

Invariantly, the number of respondents who declared that the screenshot was not a cyber-threat and their mean confidence increase for both Positive and Negative treatments. However, there is a mixed result for the mean confidence of those respondents declaring that the screenshot is a start of an attack. The confidence is increasing in the Positive treatment and decreasing in the Negative treatment.

Table C.1: Topic modelling results. Reasoning for text example 1

Answer	Topics obtained from LDA	Number	Frequency
No	Facebook website looks normal and straightforward.	27	0.27
	Website elements look normal (url, padlock)	44	0.44
	WHO is a legitimate source of information	30	0.30
Yes	The WHO video looks suspicious.	21	0.43
	The Facebook website looks suspicious.	28	0.57

“No” reasons include:

- Facebook website looks normal and straightforward
Example feedback: “*Look like a normal Facebook post.*”

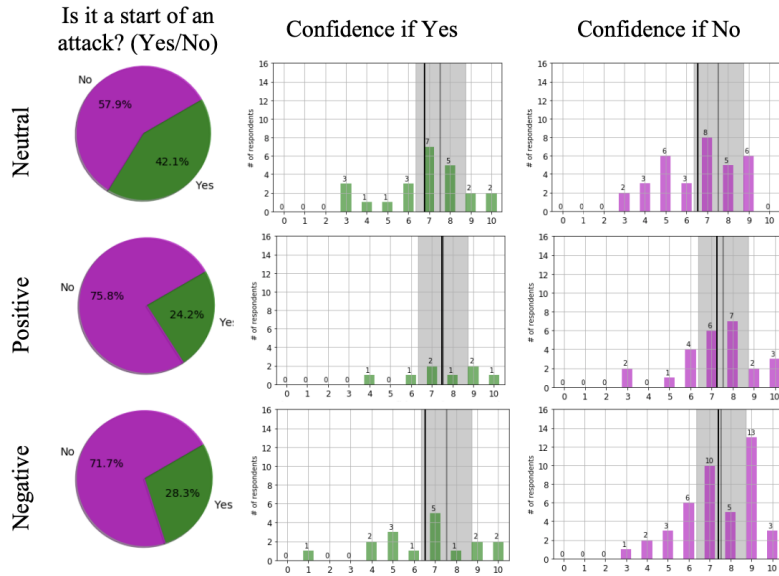


Figure C.2: Results for the first question of the HaaCSS test.

- Website elements look normal (url, padlock).
Example feedback: “*The page has the secure padlock symbol followed by https. There doesn't seem to be any push for clicking on links etc other than to watch the video.*”
- WHO is a legitimate source of information.
Example feedback: “*The WHO are a recognised organisation who have communicated via various mediums during the pandemic - including facebook. Lack of confidence comes as I know there were scams where fraudster posed as the WHO during the pandemic. There is a new friend request but that isn't an issue unless you add someone you don't know. Plus this is a secure link (Https)*”

“Yes” reasons include:

- The WHO video looks suspicious.
Example feedback: “*It looks like it has come from WHO. But the words used... Big Numbers and “Numb” - I don't think would be used by WHO.*”
- The Facebook website looks suspicious.
Example feedback: “*Facebook account holders name is John Doe (standard name for fake accounts), lack of Facebook friends I guess would also be a sign.*”

Exercise 2

The second test example is a screenshot of a legitimate wifi speed test. This screenshot does not represent a potential start of a cyber-attack. Fast.com is a service launched by Netflix. The screenshot presented to the respondents is observed in Figure C.3, while the results for the first two questions are presented in Figure C.4. Results for the topic modelling are presented in Table C.2.

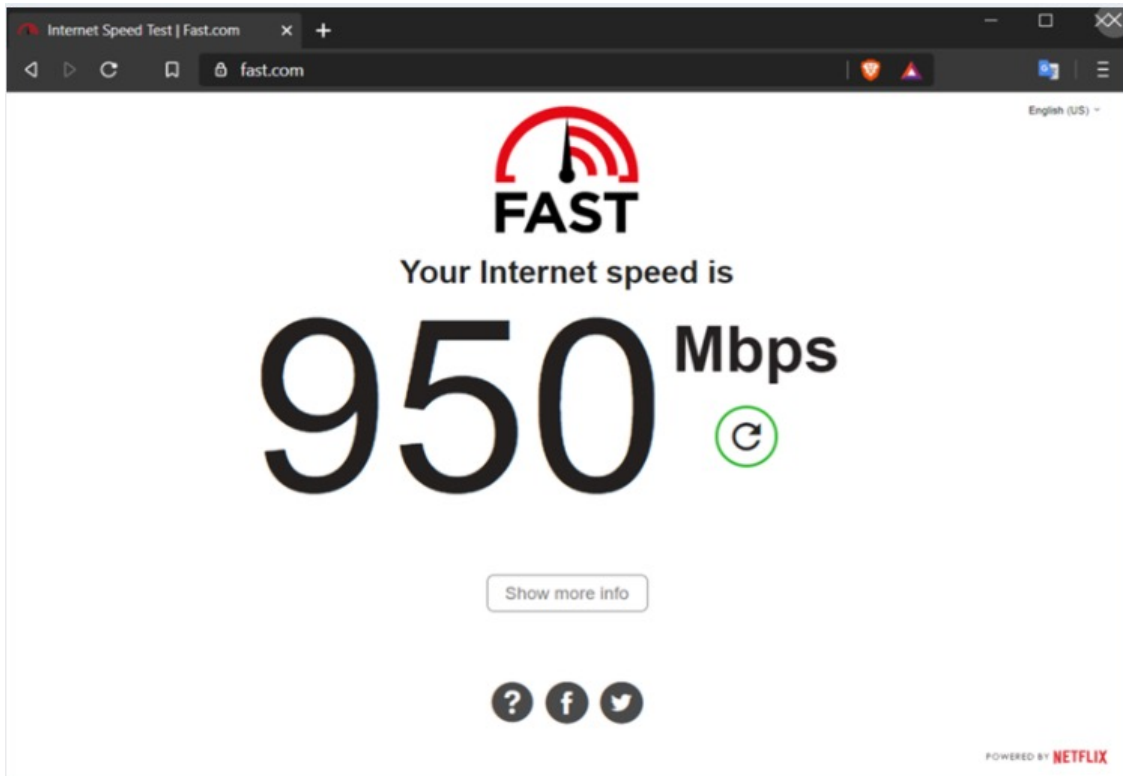


Figure C.3: Screenshot presented as second question at the Human-as-a-Cyber-Security-Sensor test.

Does represent the start of a cyber attack? No

Type of cyber attack: WiFi speed masquerading

“No” reasons include:

- It is powered by Netflix.
Example feedback: “*I have heard of this site as a main one to test broadband speed, though 950 mbs looks unusual.*”
- Website elements look normal (url, padlock).
Example feedback: “*The domain looks secure (padlock next to domain name).*”

Table C.2: Topic modelling results. Reasoning for text example 2

Answer	Topics obtained from LDA	Number	Frequency
No	It is powered by Netflix.	20	0.22
	Website elements look normal (url, padlock).	69	0.78
Yes	Netflix name looks like a scam.	19	0.31
	Website information is misleading.	18	0.29
	Website elements look suspicious.	24	0.40

“Yes” reasons include:

- Netflix name looks like a scam.
Example feedback: “*It says powered by Netflix but the website is fast.com, there is also a link to click for further information.*”
- Website information is misleading.
Example feedback: “*Not HTTPS and never seen speeds like that.*”
- Website elements look suspicious.
Example feedback: “*there is no https link, it only shows fast.com and also no privacy note on the right where it has English US - but am not 100% sure so would to click a nything open.*”

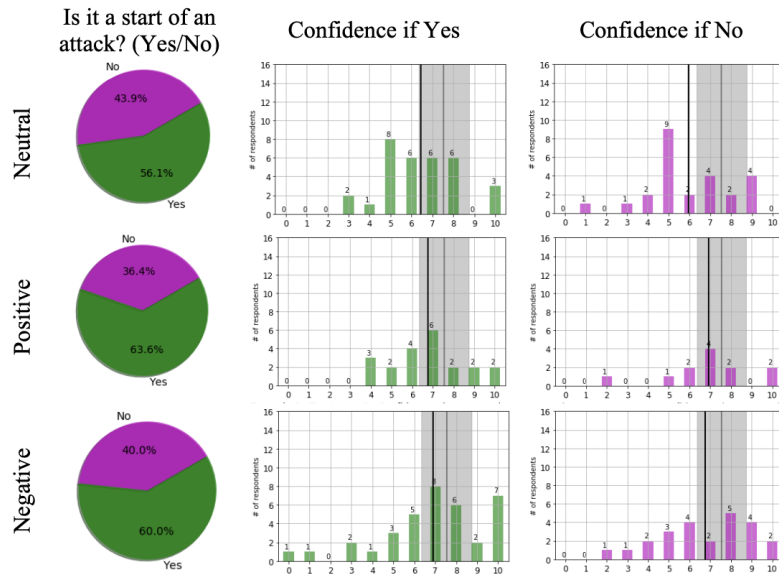


Figure C.4: Results for the second question of the HaaCSS test.

Exercise 3

The third example is a screenshot of a supposedly native Apple website offering to scan the user's computer to look for viruses. The website states that the user's system has been infected. It is an example of “typosquatting” and “phishing” and does represent the start of a cyber-attack. The screenshot presented to the respondents is observed in Figure C.5, while the results for the first two questions are presented in Figure C.6. Results for the topic modelling are presented in Table C.3.

In this case, the overwhelming majority of the surveyed employees have correctly detected the start of a cyber-attack.

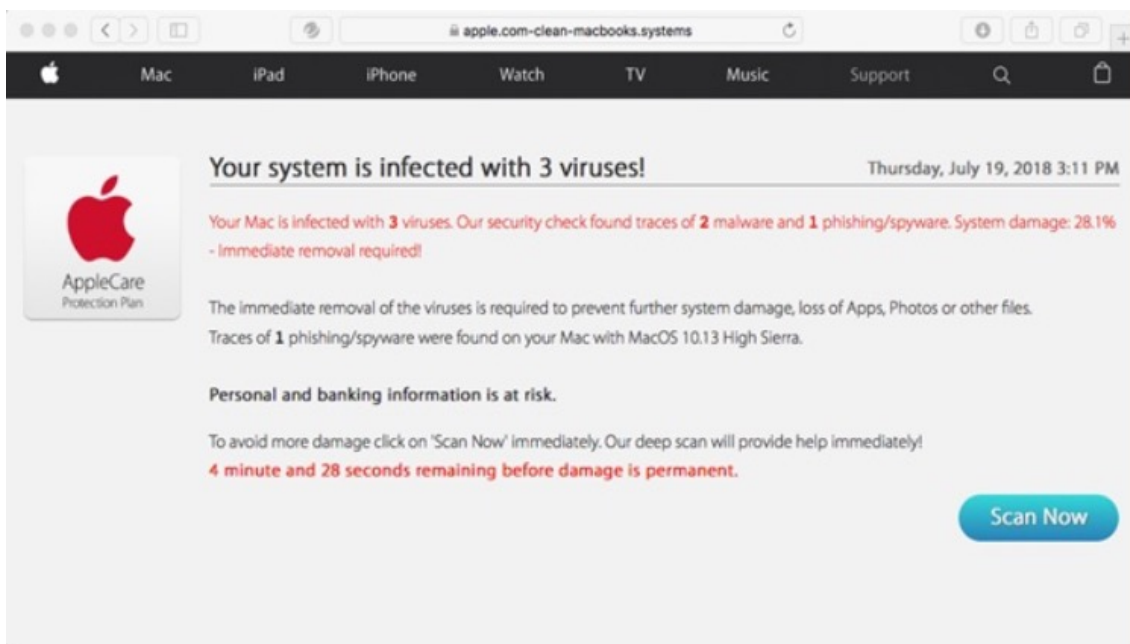


Figure C.5: Screenshot presented as third question at the Human-as-a-Cyber-Security-Sensor test.

Does represent the start of a cyber attack? Yes

Type of cyber attack: typosquatting and phishing

“No” reasons include:

- Website elements look normal (url, padlock), but unsure about content. Example feedback: “*again the padlock but unsure you would see a time limit so feel like I need to sit on the fence with this one!!*”
- Not familiar with Apple, but suspicious. Example feedback: “*The address is apple & linked to the Apple protection plan therefore does not suggest scam, unless Apple site has been hacked.*”

Table C.3: Topic modelling results. Reasoning for text example 3

Answer	Topics obtained from LDA	Number	Frequency
No	Website elements look normal (url, padlock), but unsure about content.	8	0.72
	Not familiar with Apple, but suspicious. (url, padlock).	3	0.28
Yes	Content and phrasing of message (urgency, wording).	78	0.56
	Pressure to click a button.	39	0.28
	Website elements look suspicious.	22	0.16

“Yes” reasons include:

- Content and phrasing of message (urgency, wording).
Example feedback: “*Wording on the screenshot is intended to cause panic (“damage is permanent” “loss of banking details” etc) and obviously trying to scare you into clicking on a phishing link. Layout of message looks pretty standard for a phishing email.*”
- Pressure to click a button.
Example feedback: “*Being asks to click to scan now is making me suspicious that this could be an attack.*”
- Website elements look suspicious.
Example feedback: “*The padlock is greyed out.*”

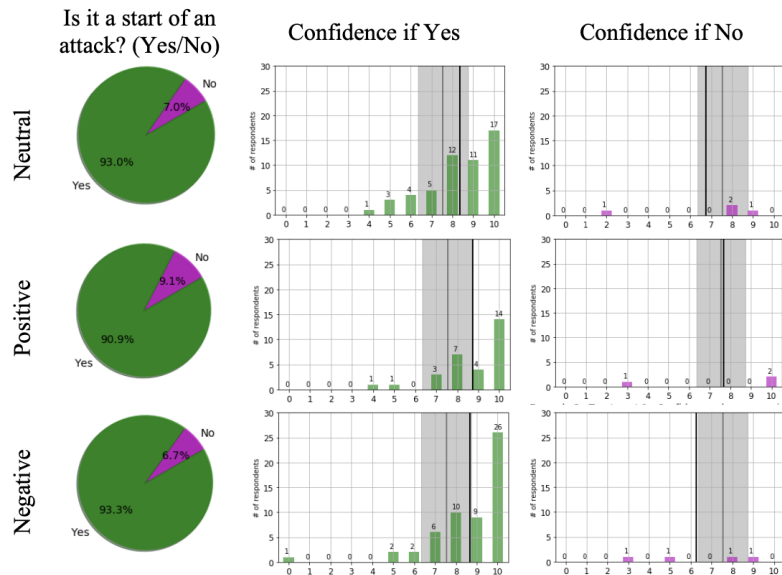


Figure C.6: Results for the third question of the HaaCSS test.

Exercise 4

The fourth example shows a screenshot from an email received from Microsoft Planner, stating a new task has been assigned to the user. This screenshot does not represent the start of a cyber-attack. The screenshot presented to the respondents is observed in Figure C.7, while the results for the first two questions are presented in Figure C.8. Results for the topic modelling are presented in Table C.4.

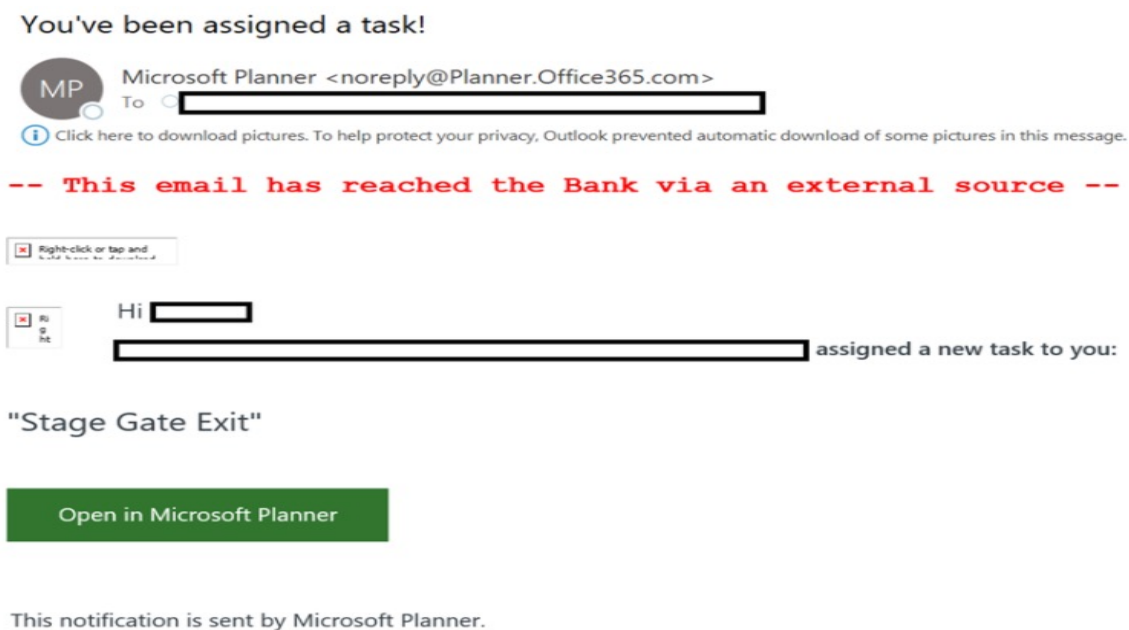


Figure C.7: Screenshot presented as fourth question at the Human-as-a-Cyber-Security-Sensor test.

Does represent the start of a cyber attack? No

Type of cyber attack: Microsoft Planner phishing

Table C.4: Topic modelling results. Reasoning for text example 4

Answer	Topics obtained from LDA	Number	Frequency
No	I am familiar with sender/software.	55	1
Yes	I am not familiar with the sender/External email.	33	0.35
	Action required inside the email.	62	0.65

“No” reasons include:

- I am familiar with sender/software.
Example feedback: “*This is a standard O365 email. However, it has an external banner so I would do more checks to make sure it’s legitimate.*”

“Yes” reasons include:

- I am not familiar with the sender/External email.
Example feedback: “*Email from an external address (red warning at top).*”
- Action required inside the email.
Example feedback: “*Taking action upon receipt of an email sending “do not reply” may be problem.*”

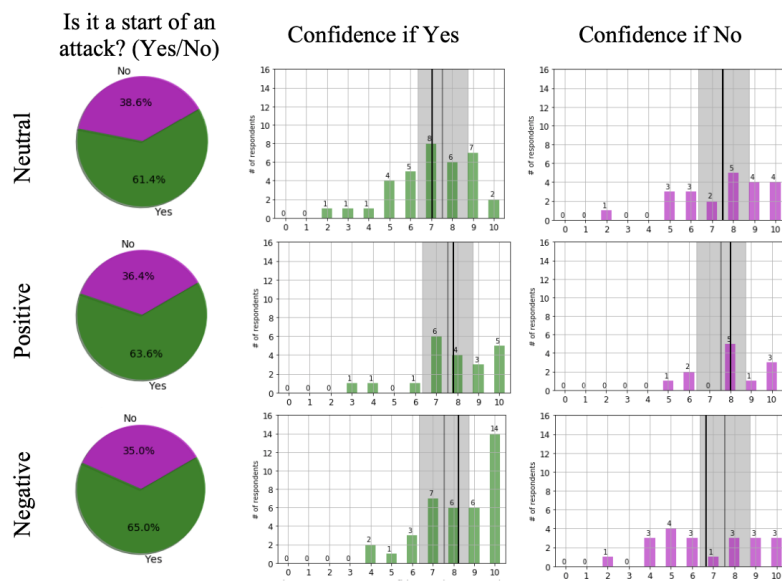


Figure C.8: Results for the fourth question of the HaaCSS test.

Exercise 5

The fifth screenshot shows a typical phishing email. The sender claims to be HMRC and requires urgent action from the user. This does represent the start of a cyber-attack. The screenshot presented to the respondents is observed in Figure C.9, while the results for the first two questions are presented in Figure C.10. Results for the topic modelling are presented in Table C.5.

In this case, the overwhelming majority of the surveyed employees have correctly detected the start of a cyber-attack.

Does represent the start of a cyber attack? Yes

Type of cyber attack: Phishing email

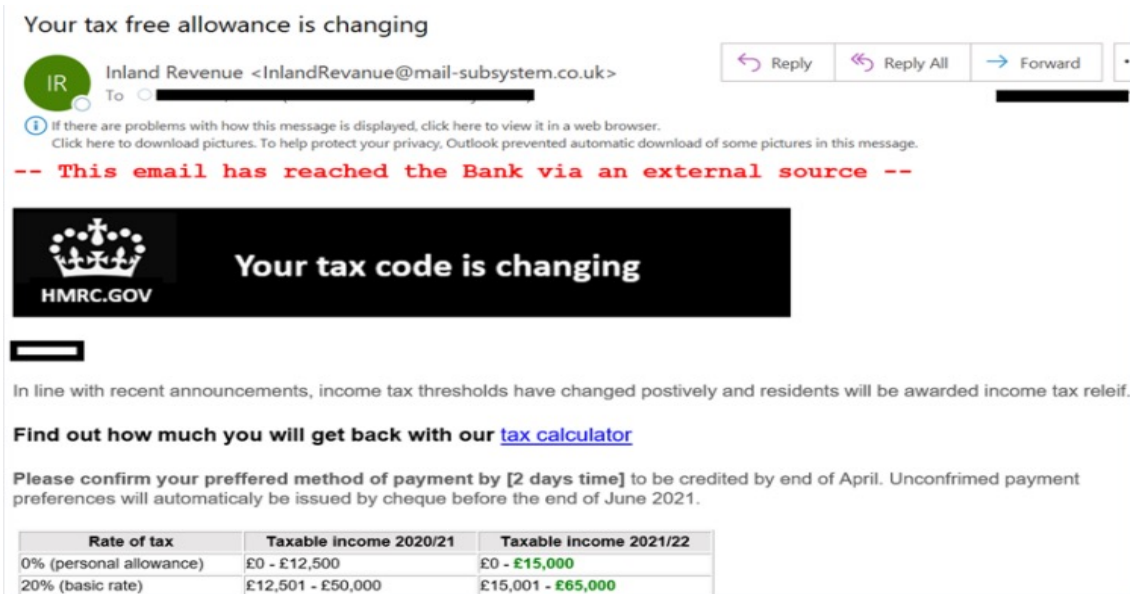


Figure C.9: Screenshot presented as fifth question at the Human-as-a-Cyber-Security-Sensor test.

“No” reasons include:

- Simple information.
Example feedback: “*Just giving details of tax.*”

“Yes” reasons include:

- Misspelling of “revenue”/Unknown sender.
Example feedback: “*email does not match what the official email would be, HMRC wouldn’t contact by email to advise of a change, spelling mistake.*”
- Content of message.
Example feedback: “*The tax office would not send such emails.*”
- Phishing email.
Example feedback: “*Its a classic case of phishing attack asking user to choose payment type.*”

Table C.5: Topic modelling results. Reasoning for text example 5

Answer	Topics obtained from LDA	Number	Frequency
No	Simple information.	1	1
Yes	Misspelling of “revenue”/Unknown sender.	19	0.13
	Content of message.	57	0.38
	Phishing email	73	0.49

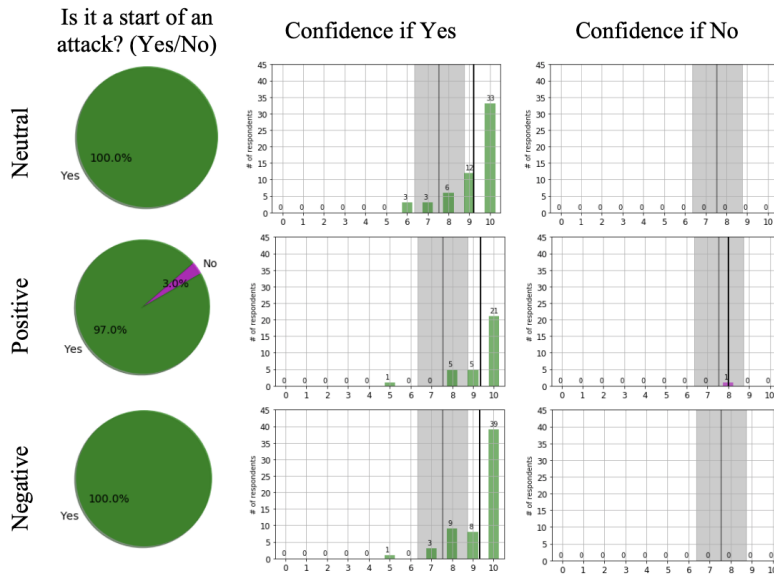


Figure C.10: Results for the fifth question of the HaaCSS test.

Exercise 6

The test example 6 is a typical example of QRishing masquerading. This case does represent the start of a cyber-attack. The screenshot presented to the respondents is observed in Figure C.11, while the results for the first two questions are presented in Figure C.12. Results for the topic modelling are presented in Table C.6.

The results show quite a bit of confusion as many employees classified this screenshot as benign. This indicates that they are not used to QR code scams and might need more training on malicious QR codes.

Does represent the start of a cyber attack? Yes

Type of cyber attack: QRishing

Table C.6: Topic modelling results. Reasoning for text example 6

Answer	Topics obtained from LDA	Number	Frequency
No	Not familiar with Instragram/Steam.	20	0.54
	Normal add/familiar with Steam ads.	5	0.14
	Nothing looks suspicious.	12	0.32
Yes	URL suspicious.	46	0.341
	Unverified account/suspicious QR code.	17	0.15
	Fake website.	50	0.44

“No” reasons include:

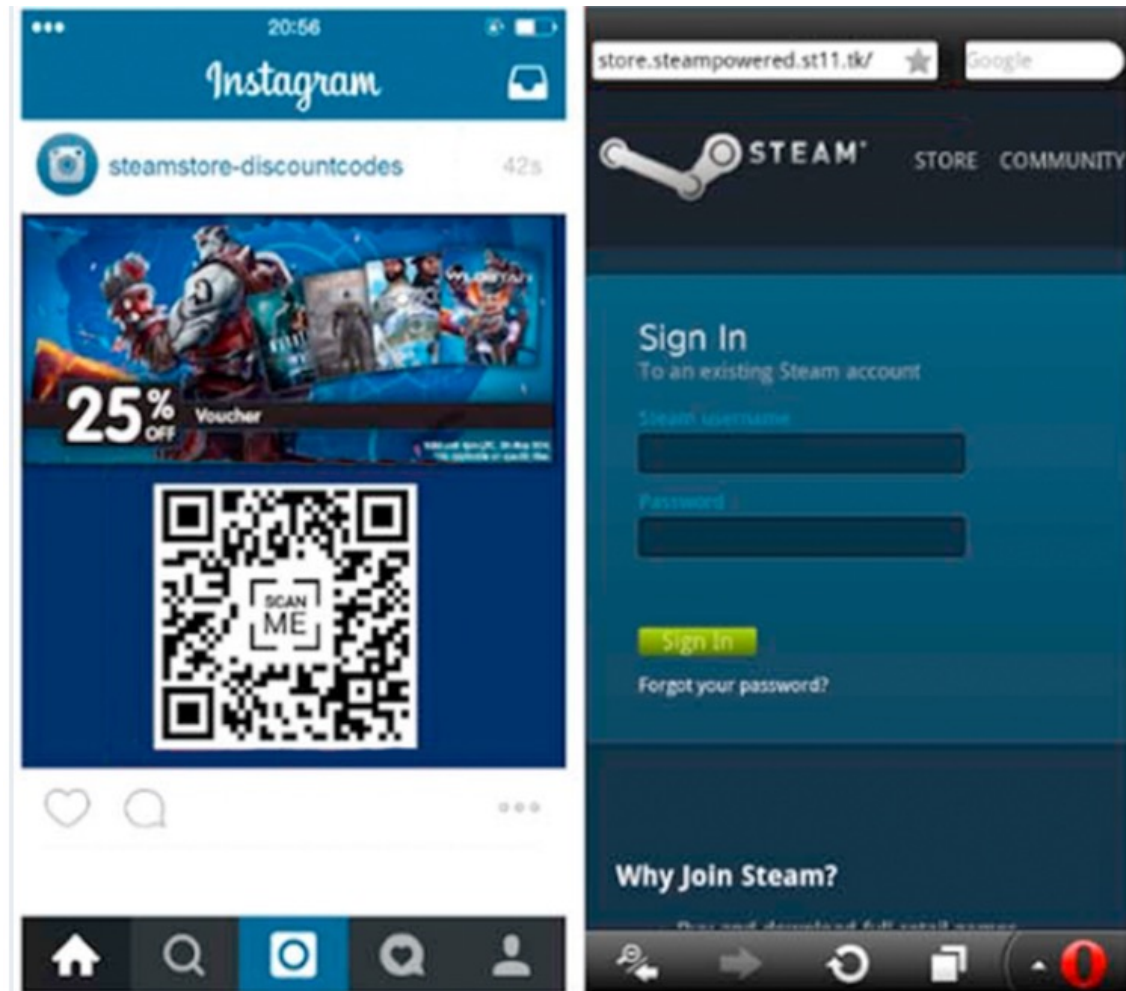


Figure C.11: Screenshot presented as sixth question at the Human-as-a-Cyber-Security-Sensor test.

- Not familiar with Instagram/Steam.
Example feedback: *“Nothing forcing me to click on links, no urgency, would only be interested in this if I knew what I was and was an existing user, otherwise would have no interest.”*
- Normal add/familiar with Steam ads.
Example feedback: *“These appear to be sites for Steam Store. It seems screenshots for log in, rather than a cyber attack.”*
- Nothing looks suspicious.
Example feedback: *“Could not notice anything.”*

“Yes” reasons include:

- URL suspicious.
Example feedback: *“the .tk domain looks odd / instagram looks ok.”*

- Unverified account/suspicious QR code.
Example feedback: “Instagram - fewer controls, URL, QR code could infect phone.”
- Fake website.
Example feedback: “The website that the QR code has taken the user to ends with “.st 11.tk” which is not a legitimate suffix for Steam (this would likely be “.com”). This is a fake Steam page designed to syphon a user’s account username and password.”

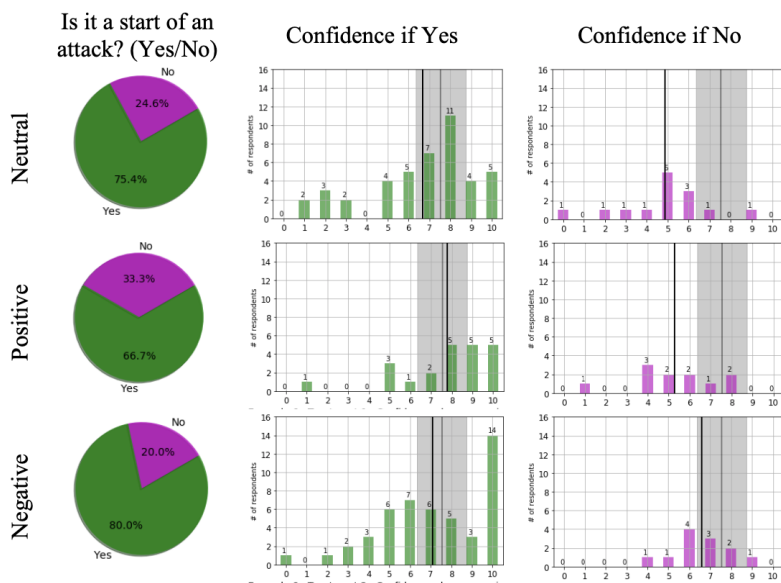


Figure C.12: Results for the sixth question of the HaaCSS test.

Appendix D

Chapter 2: Table of models

In Table D.1 we present all the different models tested. They are numbered as shown in Figure 2.3a and Figure 2.3b.

Table D.1: List of all trained models.

	Model	Loss Function	Free parameters calibrated
1	Gravity	MSE	β, γ
2	Gravity	Poisson	β, γ
3	Radiation	MSE	ρ, r
4	Radiation	Poisson	ρ, r
5	Retail	Poisson	β (travel times)
6	Retail	Poisson	β, α_1 (hospital admissions by misuse of drugs)
7	Retail	Poisson	β, α_2 (hospital admissions by poisoning of drugs)
8	Retail	Poisson	β, α_3 (police workforce)
9	Retail	Poisson	β, α_4 (knife crime events)
10	Retail	Poisson	β, α_5 (gross dispensable household income)
11	Retail	Poisson	$\beta, \alpha_1, \alpha_2$
12	Retail	Poisson	$\beta, \alpha_1, \alpha_3$
13	Retail	Poisson	$\beta, \alpha_1, \alpha_4$
14	Retail	Poisson	$\beta, \alpha_1, \alpha_5$
15	Retail	Poisson	$\beta, \alpha_2, \alpha_3$
16	Retail	Poisson	$\beta, \alpha_2, \alpha_4$
17	Retail	Poisson	$\beta, \alpha_2, \alpha_5$
18	Retail	Poisson	$\beta, \alpha_3, \alpha_4$
19	Retail	Poisson	$\beta, \alpha_3, \alpha_5$
20	Retail	Poisson	$\beta, \alpha_4, \alpha_5$
21	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_3$
22	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_4$
23	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_5$
24	Retail	Poisson	$\beta, \alpha_1, \alpha_3, \alpha_4$
25	Retail	Poisson	$\beta, \alpha_1, \alpha_3, \alpha_5$
26	Retail	Poisson	$\beta, \alpha_2, \alpha_3, \alpha_4$
27	Retail	Poisson	$\beta, \alpha_2, \alpha_3, \alpha_5$
28	Retail	Poisson	$\beta, \alpha_3, \alpha_4, \alpha_5$

29	Retail	Poisson	$\beta, \alpha_2, \alpha_4, \alpha_5$
30	Retail	Poisson	$\beta, \alpha_1, \alpha_4, \alpha_5$
31	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_3, \alpha_4$
32	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_3, \alpha_5$
33	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_4, \alpha_5$
34	Retail	Poisson	$\beta, \alpha_1, \alpha_3, \alpha_4, \alpha_5$
35	Retail	Poisson	$\beta, \alpha_2, \alpha_3, \alpha_4, \alpha_5$
36	Retail	Poisson	$\beta, \alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5$
37	Retail	MSE	β
38	Retail	MSE	β, α_1
39	Retail	MSE	β, α_2
40	Retail	MSE	β, α_3
41	Retail	MSE	β, α_4
42	Retail	MSE	β, α_5
43	Retail	MSE	$\beta, \alpha_1, \alpha_2$
44	Retail	MSE	$\beta, \alpha_1, \alpha_3$
45	Retail	MSE	$\beta, \alpha_1, \alpha_4$
46	Retail	MSE	$\beta, \alpha_1, \alpha_5$
47	Retail	MSE	$\beta, \alpha_2, \alpha_3$
48	Retail	MSE	$\beta, \alpha_2, \alpha_4$
49	Retail	MSE	$\beta, \alpha_2, \alpha_5$
50	Retail	MSE	$\beta, \alpha_3, \alpha_4$
51	Retail	MSE	$\beta, \alpha_3, \alpha_5$
52	Retail	MSE	$\beta, \alpha_4, \alpha_5$
53	Retail	MSE	$\beta, \alpha_1, \alpha_2, \alpha_3$
54	Retail	MSE	$\beta, \alpha_1, \alpha_2, \alpha_4$
55	Retail	MSE	$\beta, \alpha_1, \alpha_2, \alpha_5$
56	Retail	MSE	$\beta, \alpha_1, \alpha_3, \alpha_4$
57	Retail	MSE	$\beta, \alpha_1, \alpha_3, \alpha_5$
58	Retail	MSE	$\beta, \alpha_2, \alpha_3, \alpha_4$
59	Retail	MSE	$\beta, \alpha_2, \alpha_3, \alpha_5$
60	Retail	MSE	$\beta, \alpha_3, \alpha_4, \alpha_5$
61	Retail	MSE	$\beta, \alpha_2, \alpha_4, \alpha_5$
62	Retail	MSE	$\beta, \alpha_1, \alpha_4, \alpha_5$
63	Retail	MSE	$\beta, \alpha_1, \alpha_2, \alpha_3, \alpha_4$
64	Retail	MSE	$\beta, \alpha_1, \alpha_2, \alpha_3, \alpha_5$
65	Retail	MSE	$\beta, \alpha_1, \alpha_2, \alpha_4, \alpha_5$
66	Retail	MSE	$\beta, \alpha_1, \alpha_3, \alpha_4, \alpha_5$
67	Retail	MSE	$\beta, \alpha_2, \alpha_3, \alpha_4, \alpha_5$

Appendix E

Chapter 2: Database

The github repository github.com/LeonardoCastro/BritishDrugDynamics contains all the information presented in this section. In Section E.1 we present the different elements worth mentioning from the database. These include the different data that are included, their respective resolution, format and sources.

We also include the different territorial divisions used in Section E.1.10. These include the territorial divisions for NHS Local Authorities, Local Police Forces and English regions.

E.1 The database

All the compiled information is obtained from different British Governmental websites, being public, open and shared with the [Open Government Licence](#). Only the used travel time matrices are obtained from a third party company, which is the Maps API from Google®.

Most of the data is presented in `.csv` format and thought to be managed as *data frames* objects with packages as `pandas` for Python or `data.table` for R.

Most of the data is presented as time series. Coming from the British Government, the year steps from one data to another are not from January to December, but rather from April to March. This is because the Government takes the fiscal year as unit of time. In that sense, when a measurement reads for a “2012”, this actually means that the measurement refers to the fiscal year starting in April 1, 2011 and finishes on March 31, 2012.

In the following, we present each kind of data used in this project. Subsections of missing data are also added. This is only to acknowledge the raw data that has not been fully processed.

E.1.1 Drug related hospital admissions data

NHS hospitals present annually data about their hospital admissions. The particular set of hospital admissions related to drugs comprise three different types:

1. NHS hospital finished admissions where there was a **primary diagnosis of drug related mental health and behavioural disorders**.

2. NHS hospital finished admission episodes with a **primary or secondary diagnosis of drug related mental and behavioural disorders**.
3. NHS hospital finished admissions where a **primary diagnosis of poisoning by drugs**.

We include the time series for each kind of hospital admission from 2009 to 2019 for different geographical resolutions: England, its 9 regions, 39 police force areas and 131 counties described in Section E.1.10. We also include their respective time series for admissions by 100 thousand inhabitants. Depending of the resolution, the normalisation is done using the population of the territorial unit. That is, the measure for Northumbria in 2012 is done with the population of Northumbria in 2012.

At national level we include the age distribution for each kind of admission. We also include the distribution of diagnoses for admissions type 1 and 3. However, for admissions type 3 these were published only from 2013.

Admissions type 2 are not considered into the analysis and are only considered for reference. This decision is based on the lack of information from the primary and secondary diagnoses, thus being unable to detect the underlying causes of the admissions.

The main source is the [Statistics on Drug Misuse](#) published by NHS Digital annually.

E.1.2 Drug related deaths data

Statistics about drug related deaths in England and Wales are published each year. In this case, data has an inherent delay caused by the difference between the decease date and the registered date. The delay in England for 2018 had a median of 181 days, and a median delay of 172 days for 2017 according to the publishers [John, 2019]. This makes the interpretation from this statistic more difficult to handle, as some deaths registered in a particular year could have happened more than a year before.

The publishers also make the distinction between the deaths caused by poisoning of drugs and those deaths by poisoning of which were caused by misuse of drugs. This is a subtle distinction, as the drug poisoning death is defined by the WHO's [International Classification of Diseases](#). The drug misuse death is a drug poisoning death which also involves a drug abuse or dependence.

We publish the time series (2009-2018) for the different resolutions handled: England, its 9 regions, the 39 local police areas and the 131 counties described in Section E.1.10. We also include for each one of these resolutions their respective time series for deaths by 100 thousand inhabitants. Depending of the resolution, the population of the territorial unit is used. That means that for the number of deaths for each 100 thousand inhabitants in Essex for 2012, the population for Essex in 2012 is used.

In the bottom three resolutions (local authorities, police forces and regions) we only present the total number of deaths by poisoning and by misuse. However, at the national (England) level we also present the time series for underlying causes, age distribution for deaths by misuse and by poisoning, and the age distribution by drug for the total number of deaths.

Given that the data does not cover 2019 and 2020, this data was not used for the analysis of this work.

The main source of the deaths related to drugs data is the [Deaths related to drug poisoning in England and Wales: 2018 registrations](#) published annually by the Office for National Statistics.

E.1.3 Number of hospital beds data

The number of hospital beds was collected for the three different resolutions. However, we only recommend data for the regional and local police resolutions. This is due to the fact that an important number of reported hospital beds are an addition for different hospitals in different local authorities. An example of this is the Guy's and St. Thomas' Hospitals: the hospital beds are reported as an addition for both hospitals, while one is in the London Borough of Lambeth, and the other in the London Borough of Southwark. This of course is solved when counting the hospital beds for the Metropolitan Police resolution, including most of the Greater London boroughs.

This data was used to normalise the hospital admissions.

The main data is the NHS database of [hospital beds availability](#).

E.1.4 Police workforce data

The police workforce data is the only one that is published more than once a year, being published each semester. This allows to know how the workforce varies along each year. The Home Office publishes this data in different ways, considering the

number of heads working full-time and part-time jobs, and doing a conversion to the equivalent of heads working only full-time. In all years they include police officers and police staff. However, from 2012 the Home Office includes the numbers for different job titles working in the police workforces. These include community support officers, designated officers and traffic warden. In order to have a coherent database, we only consider the regular officers in the conversion to heads working only full-time.

The time series are presented for data from 2009 to 2019. They are presented for England, its 9 regions and the 39 local police forces, including the British Transport Police and the Central Service Secondments. More information about the different resolutions is found in Section E.1.10. For each resolution, workforces are given by annual mean with its correspondent standard deviation (for measurement, we take the workforce number at the beginning, middle and end of each year). We also include the same numbers for each 100k inhabitants. The normalisation is done using the respective population resolution. That means that the mean for each 100k inhabitants for Oxfordshire in 2012 is obtained using the population of the same county in 2012. In the case of the British Transport Police and the Central Service Secondments, the population of England is used.

The main source for the workforce data is the [Police workforce England and Wales statistics](#) published by the Home Office.

E.1.5 Police numbers of drug seizures data

The Home office publishes once a year the number of seizures and total quantities by drug and police force. This allows to compile a set of time series (2010-2019) for different drugs at different resolutions (England, regions and police forces). Also, at England & Wales resolution we obtain the number of seizures by weight/dosage for different drugs. The available drugs (with their dosage unit) are:

- **Class A drugs:** Cocaine (kg), Crack (kg), Ecstasy (doses), Heroin (kg), LSD (doses), Methadone (doses), Morphine (doses).
- **Class B drugs:** Herbal Cannabis (kg), Cannabis Resin (kg), Cannabis Plants (plants), Amphetamines (kg), Barbiturates (doses), Ketamine* (kg).
- **Class C drugs:** Anabolic steroids (kg), Benzodiazepines (doses), GHB (doses), Temazepam (doses).

The main source is the [Seizure of drugs in England and Wales statistics](#) published by the Home Office annually.

Given the lack of data for 2020 and the geographical resolution used, we did not use this data for the analysis of this work.

*: In the fiscal year 2014-2015, Ketamine was reclassified a Class B drug instead of a Class C drug.

E.1.6 Knife crime related data

Knife crime has been reported by the [House of Commons library](#) since 2009. The data is at the police forces resolution.

E.1.7 Disposable Income data

The Gross Disposable Household Income data is reported the [ONS](#). The data is available at county level. We aggregated the income to a police territory resolution using a weighted average using the population of each county.

E.1.8 Demographic data

As demographic data we include the time series of the population for different resolutions of England from 2009 to 2019. The different resolutions are those described in Section E.1.10, and refer to England, the 9 regions conforming England, the English Local Police Forces and the 131 Counties adopted for this project.

The main source of the demographic data is the [Estimate of the population for the UK, England, Wales, Scotland and Northern Ireland](#), from the Office for National Statistics. The Estimates are released each year.

E.1.9 Geographic data

Geographic data is analogous to demographic data. We include different `.geojson` files containing the geometries for Great Britain in different resolutions. The different resolutions are those described in Section E.1.10, and refer to England, the 9 regions conforming England, the English Local Police Forces and the 131 counties adopted for this project. We also include a fourth file for Scotland and Wales at local authority resolution.

The main source of the geographic data is the [Open Geography Portal](#) from the Office for National Statistics.

E.1.10 Territorial resolutions

In this section we present the different resolutions used along the work. We start presenting the lowest resolution, which is England and its nine regions. We then present the one used for the local police forces to then present the one used for the hospital admissions.

However, we recommend to visit the [github repository](#) as some information is replicated in a friendly way.

E.1.11 England and its 9 regions

These are the most simple and trivial resolutions. England is one of the four nations comprising the United Kingdom and shares borders with Scotland and Wales. England by itself is traditionally divided into 9 regions. These divisions by themselves are the top tier sub-national divisions, and although they do not hold governmental and administrative powers, these regions are often used for statistical and administrative means. These are: 1. East of England, 2. East Midlands, 3. London, 4. North East, 5. North West, 6. South East, 7. South West, 8. West Midlands and 9. Yorkshire and the Humber.

E.1.12 Local Police Forces

The United Kingdom has a handful of “British police forces” as the National Crime Agency, the British Transport Police or the British Borders Police are. Instead, most of the police tasks are taken by local polices acting in a limited area.

In England there are 39 local polices. Some of them act in unitary local authorities, like Lincolnshire and Northamptonshire Polices, whereas other act in metropolitan regions, like the Metropolitan Police in most of London and the Greater Manchester Police. Also, some polices act in a mixed area comprised of rural and different urban areas. Examples of these are the Thames Valley Police, acting en Oxfordshire, Reading, Milton Keynes, etc., or the Northumbria police acting in Sunderland, Newcastle upon Tyne, Northumbria, etc.

The list of the different police forces is:

1. **East of England:** Bedfordshire, Cambridgeshire, Essex, Hertfordshire, Norfolk, and Suffolk Polices.
2. **East Midlands:** Derbyshire, Leicestershire, Lincolnshire, Northamptonshire, and Nottinghamshire Polices.
3. **London:** Metropolitan Police and the City of London Police.

4. **North East:** Cleveland, Durham, and Northumbria Polices.
5. **North West:** Cheshire, Cumbria, Great Manchester, Lancashire, and Merseyside Polices.
6. **South East:** Hampshire, Kent, Surrey, Sussex, and Thames Valley Polices
7. **South West:** Avon and Somerset, Devon and Cornwall, Dorset, Gloucestershire, and Wiltshire Polices.
8. **West Midlands:** Staffordshire, Warwickshire, West Mercia, and West Midlands Polices.
9. **Yorkshire and the Humber:** Humberside, North Yorkshire, South Yorkshire and West Yorkshire Polices.

Additionally, in the [github repository](#), a list of equivalences between the police forces and merged local authorities is shown.

E.1.13 Merged local authorities

In Section E.1 we presented different statistics used throughout this work published from different Governmental offices. For most of them, mainly the ONS, the Home Office and `data.police.gov.uk`, the different resolutions used are consistent during the time interval analysed (2009-2019). However, the hospital admissions data published by NHS digital changed its lower tier territorial divisions in 2012, thus not allowing to have a coherent time series.

In order to fix this, we created our own lower tier divisions and we call these divisions as *merged local authorities*. This topology of merged local authorities is transferable to the other statistics, allowing us to have a full homogenised database.

Up until 2012, the NHS was divided in 10 different Strategic Health Areas (similar to the regions described above) and 152 Primary Care Trusts (PCTs) covering England. However, that year the British Parliament passed the Health and Social Care Act 2012, abolishing SHAs and PCTs, transferring the administrative powers to the 151 Local Authorities in England.

Our homogenisation process involves the detection of the local authorities comprising each PCT, and the detection of PCTs comprising each local authorities. Once done that, the largest number of merged local authorities comprising the 152 PCTs are chosen. The result is a list of 131 merged local authorities. The equivalence between these, the pre-2012 PCTs and the current local authorities can be found in the [github repository](#).

Appendix F

Chapter 3

F.1 Formation of sectors for both datasets

In this section we present how the collected Forbes Fortune databases were divided so our analysis could be made for different company sectors. First, we present the divisions for the US sample (Table F.1), to then present the divisions for the Global sample (Table F.2).

Table F.1: Industry sectors for US-based companies used with respect to the industry labels at the Forbes Fortune database.

Industry		
Sector	Industry label	Examples of companies in industry label
Finance (N=39)	Banking, Financial Services (N=32)	J.P. Morgan, Wells Fargo, MetLife
	Insurance (N=7)	State Farm Insurance, Allstate, New York Life Insurance
Food and Wholesale (N=16)	Food (N=10)	PepsiCo, Coca-Cola, Starbucks
	Tobacco (N=1)	Phillip Morris International
	Wholesale (N=5)	Sysco, US Foods Holding, United Natural Foods
Health Care and Pharmaceutical (N=10) (N=26)	Health Care	CVS Health, United Health Group, McKesson
	Medical Equipment (N=5)	Abbot Laboratories, Thermo Fischer Scientific, Danaher
	Pharmaceutical (N=11)	AmerisourceBergen, Walgreens Boots Alliance, Johnson & Johnson
Energy (N=23)	Energy (N=8)	World Fuel Services, Exelon, Duke Energy
	Oil and gas (N=8)	Exxon Mobil, Chevron, Marathon Petroleum
	Petroleum (N=6)	Energy Transfer, ConocoPhillips, Baker Hughes

F.1. FORMATION OF SECTORS FOR BOTH DATASETS

Table F.1: Industry sectors for US-based companies used with respect to the industry labels at the Forbes Fortune database.

Industry		
Sector	Industry label	Examples of companies in industry label
	Metals (N=1)	Nucor
Chemicals (N=12)	Chemicals (N=7)	Dow, 3M, DuPont
	Consumer Goods (N=5)	Procter & Gamble, Kimberley-Clark, Colgate-Palmolive
Heavy Industry (N=49)	Automotive (N=6)	Ford Motor, General Motors, Tesla
	Semiconductors (N=7)	Intel, Jabil, Qualcomm
	Heavy Equipment (N=3)	Caterpillar, Deere, Paccar
	Conglomerate (N=7)	Berkshire Hathaway, General Electric, Honeywell International
	Advanced Tech (N=1)	Lockheed Martin
	Aerospace (N=3)	Raytheon Technologies, Boeing, General Dynamics
	Defence (N=3)	Northrop Grumman, Raytheon, Howmet Aerospace
	Security (N=1)	BlackRock
	Computing and Conglomerate (N=18)	Amazon, Apple, Alphabet
	Services (N=13)	Advertising and Marketing (N=2)
Social Media (N=1)		Facebook
Media and Entertainment (N=2)		Walt Disney, Netflix
Telecommunications (N=5)		AT&T, Verizon Communications, Comcast
Networking (N=1)		Cisco Systems

F.1. FORMATION OF SECTORS FOR BOTH DATASETS

Table F.1: Industry sectors for US-based companies used with respect to the industry labels at the Forbes Fortune database.

Industry		
Sector	Industry label	Examples of companies in industry label
	Outsourcing (N=2)	ManpowerGroup, Aramark
Retail (N=23)	Retail (N=17) Apparel (N=3) Automotive Services (N=3)	Walmart, Costco Wholesale, Kroger Nike, Gap, Ross Stores Penske Automotive Group, Autonation, Car Max
Miscellaneous (N=16)	Real Estate (N=1) Construction (N=5) Hospitality (N=1) Logistics (N=2) Railroads (N=1) Air Transport (N=4) Courier (N=2)	CBRE Group Lennar, AECOM, Fluor Marriott International XPO Logistics, Waste Management Union Pacific Delta Airlines, American Airlines Group, United Airlines Holdings FedEx

Table F.2: Industry sectors for Global companies used with respect to the industry labels at the Forbes Global database.

Industry		
Sector	Industry label	Examples of company in industry label
Finance (N=24)	Banking, Financial Services (N=22) Trading (N=2)	Axa (France), Bank of China (China), Allianz (Germany) Itochu (Japan), China Minmetals (China)
Health Care and Pharmaceutical (N=14)	Health Care (N=7) Pharmaceutical (N=6)	CVS Health (USA), United Health Group (USA), McKesson (USA) Johnson & Johnson (USA), Roche Group (Switzerland), Bayer (Germany)

F.1. FORMATION OF SECTORS FOR BOTH DATASETS

Table F.2: Industry sectors for Global companies used with respect to the industry labels at the Forbes Global database.

Industry		
Sector	Industry label	Examples of company in industry label
	Food (N=1)	Nestlé (Switzerland)
Energy (N=17)	Oil and Gas (N=14)	Gazprom (Russia), Total (France), BP (UK)
	Commodity (N=2)	Glencore (Switzerland), Trafigura Group (Switzerland)
	Petroleum Refining (N=1)	Sinopec Group (China)
	Automotive (N=13)	Automotive (N=13)
Computing (N=14)	Computing (N=1)	Apple (USA)
	Computing and Conglomerate (N=4)	Amazon (USA), Microsoft (USA), Alphabet (USA)
	Electronics (N=3)	Samsung Electronics (South Korea), Huawei Investment and Holding (China)
	Telecommunications (N=6)	AT&T (USA), Nippon Telegraph and Telephone (Japan), Deutsche Telekom (Germany)
Conglomerate (N=15)	Conglomerate (N=9)	Berkshire Hathaway (USA), Japan Post Holdings (Japan), Siemens (Germany)
	Courier (N=1)	China Post Group (China)
	Retail (N=5)	Walmart (USA), Kroger (USA), Carrefour (France)
Construction (N=12)	Construction (N=5)	China State Construction Engineering (China), China Railway Construction (China), Vinci (UK)
	Electric Utility (N=5)	State Grid (China), Enel (Italy), Electricité de France (France)
	Holding Company (N=2)	Exor Group (Netherlands), China Railway Engineering Group (China)

F.2 Topic Modelling

F.2.1 Perplexity plots to find optimal number of topics

Topic Modelling using Latent Dirichlet Allocation (LDA) is a generalised method in literature to extract information about the different subjects discussed in a given collection of texts. The aim of the method is to obtain K number of topics, each being formed by a set of words extracted from the vocabulary used in the documents. LDA works in a reverse way, assuming the existence of K topics which generate the collection of texts. The algorithm then finds the best available set of words for each of the topics. However, finding the optimal number of topics K is not a trivial task, as a large number can come to redundant topics, while a small number might omit important information in the collection of documents.

A way to infer the optimal number of topics is by computing the perplexity of the topics over the collection. Perplexity is a measure of how well a probabilistic distribution (the topics in this case) predicts a sample (the collection of texts). If D is the collection of documents and p_k is the probability distribution defined by the set of k topics, then the perplexity is defined as

$$\log_2 PP_k = \sum_{d \in D} p_k(d) \log_2 p_k(d). \quad (\text{F.1})$$

Results for both samples (US and Global) are found in Figures F.1a and F.1b. We ran 20 simulations for each number of topics $k=1, \dots, 15$. In that way we can obtain a measurement with its standard deviation, to know which number k is the optimal for each sample.

F.2.2 Insight about resulting topics – Most important words in each of them

In this section we present the resulting topics for each of the two databases. The number of topics is chosen with respect to the perplexity analysis done in Section F.2.1.

F.2. TOPIC MODELLING

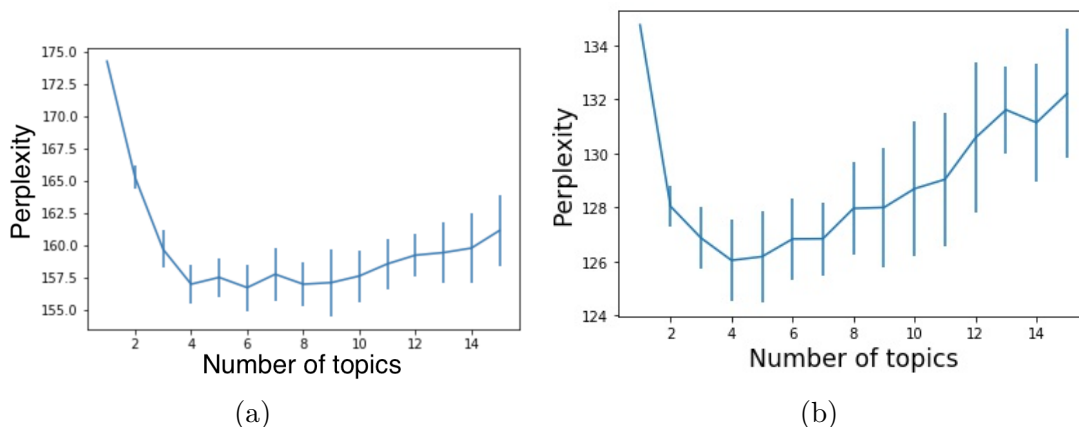


Figure F.1: Perplexity plots for (a) US-based companies and for (b) Global companies.

Table F.3: Topics and 10 most important keywords for the US-based companies database.

	Topic	10 most important keywords
1	<i>Help and support to families, and donations to healthcare organizations</i>	"fund", "foundation", "donate", "relief", "organizations", "food", "clients", "workers", "small", "families"
2	<i>Office protocols for workers</i>	"store", "associate", "wear", "mask", "order", "return", "office", "team members", "protocols", "hours"
3	<i>Office protocols for clients and suppliers</i>	"travel", "mask", "manufacture", "site", "suppliers", "flight", "protocols", "colleagues", "best", "equipment"
4	<i>Enabling Work and Services from Home Protocols</i>	"patients", "virtual", "organizations", "tool", "remote", "test", "network", "healthcare", "technology", "connect"
5	<i>Financial statement for investors and markets</i>	"statements", "market", "net", "result", "forward looking", "cash", "billion", "year", "cost", "income"
6	<i>Economic waivers and insurance coverage</i>	"policy", "insurance", "state", "coverage", "visit", "question", "payment", "account", "leave", "benefit"

Table F.4: Topics and 10 most important keywords for the Global companies database.

	Topic	10 most important keywords
1	<i>Ensure production chain and equipment at work</i>	"store", "network", "return", "order", "associate", "members", "equipment", "address", "remote", "production"
2	<i>cost reduction, focus on income and production</i>	"income", "net", "cost", "reduce", "end", "fee", "enterprises", "production", "share", "current"
3	<i>Financial statement for investors and markets</i>	"loan", "bank", "market", "finance", "clients", "bond", "mortgage", "economic", "capital", "coverage"
4	<i>Work from Home policy</i>	"employee", "suspend", "japan", "items", "hours", "stakeholders", "production", "vehicle", "july", "members"

F.2.3 Frequency tables

Given the size of the tables, these are found at the end of the document (Table F.9 and F.10).

The sum of the proportions for each of companies is close to 1, but not exactly equal to 1. This is simply because we are omitting those topics whose distribution in a company's document is smaller than 0.001.

F.3 LIWC

F.3.1 Dictionary of words for each dimension

(LIWC forbids to publish all the dictionary, but allows to publish a small set)

DO NOT PUBLISH

The LIWC framework is optimised to search for different psychological dimensions with respect to the different dictionaries. Over this work we are interested only in a set of available dimensions, these being: "Positive emotions", "negative emotions", "Social", "Cognitive processes", "Work", "Drives", "Home", "Past-oriented",

“Present-oriented”, “Future-Oriented”. In Table F.5 we present 20 words present in the dictionaries for each of the dimensions.

****NOTE TO AUTHORS/REVIEWERS:** a larger sample is attached as a separate pdf document (**NOT TO PUBLISH**) ******

Table F.5: Sample of words from the LIWC dictionaries for each of the dimensions used.

Dimension	Sample of dictionary
Positive emotions	Aggregable, beautiful, care, determined, excellent, fearless, giving, happy, interest, joy, kind, loyal, magnific, nicely, optimism, peaceful, respect, safe, thankfully, useful, valued, warmest.
Negative emotions	Adversity, bitterness, complain, danger, empty, fear, gloomy, hopeless, ignore, jealous, kill, lonely, missing, nervous, overwhelming, pain, rejection, sadness, tension, unfair, victim, weak, yelling
Social	Adult, baby, citizen, dance, e-mail, family, gather, help, interrupt, kid, listening, maternity, name, owner, partner, question, relationship, speak, team, visits, whoever, yourselves
Cognitive processes	Admitting, believing, complicated, deciding, exception, finding, generating, hoping, imagination, justify, knowledge, launching, maybe, normally, opposite, perspective, question, realizing, solving, thought, unacceptable, virtually, wanting.
Work	Agent, broker, collaborator, delegate, employer, finance, government, hardworking, investment, legal, manager, negotiation, officer, project, qualification, regulation, scholar, taxes, unemployed, worker.
Drives	Achievable, better, capable, destruction, excitement, family, great, highest, inferior, joining, kids, love, management, nomination, overcame, partner, quit, respect, shared, together, unacceptable, victim, wager.
Home	Address, bed, couch, domestic, family, garden, home, kitchen, lease, mortgage, neighbor, oven, pet, renovation, studio, tenant, window.

Table F.5: Sample of words from the LIWC dictionaries for each of the dimensions used.

Dimension	Sample of dictionary
Past-oriented	Accepted, bought, carried, denied, ended, formerly, guessed, hoped, included, joined, left, mastered, overcame, provided, questioned, remembering, supposed, taught, undid, viewed, wanted.
Present-oriented	Add, begins, commit, determines, enters, forbids, happens, include, join, know, leave, manage, nowadays, organizes, passes, runs, searches, tweets, walks.
Future-oriented	Anticipate, coming, eventually, foresee, going, henceforth, imminent, looming, might, onward, plan, someday, tomorrow, upcoming, wish.

F.3.2 Tables about averages of each dimension in both datasets

The LIWC framework allows to compute how much a collection of texts includes a given dimension. In Tables F.6, F.7 and F.8 we present the proportions for each dimension for both databases.

Table F.6: Proportion of dimensions in collection of texts for the US-based companies database per industrial sector.

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home
Base (all US-based industries)	3.81	0.8	10.66	8.28	14.05	1.36	7.9	1.19	8.77	0.62
Finance	3.45	0.84	9.48	8.64	13.19	1.37	7.78	1.23	8.86	0.73
Food & Wholesale	3.31	0.86	9.89	7.79	13.25	1.18	7.51	1.11	8.63	0.59
Health Care & Pharma	3.74	0.93	9.57	8.48	13.22	1.24	7.39	1.31	9.02	0.47
Energy	2.8	0.94	7.68	9.25	12.59	1.62	6.8	1.27	9.18	0.53
Chemicals	3.88	0.97	9.54	9.03	13.95	1.54	7.29	1	8.8	0.68
Industry	3.49	0.75	10.7	8.31	13.71	1.3	7.68	1.3	8.63	0.64

F.3. LIWC

Table F.6: Proportion of dimensions in collection of texts for the US-based companies database per industrial sector.

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home
Services	4.18	0.68	9.72	7.78	12.78	1.29	6.68	1.07	8.16	0.54
Retail	3.67	0.86	11.7	8.67	14.37	1.26	8.57	1.16	8.87	0.76
Miscellaneous	3.79	0.65	11.26	9.32	12.96	1.65	8.36	1.4	8.49	0.43

Table F.7: Proportion of dimensions in collection of texts for the global companies database per industrial sector.

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home
Base (all US-based industries)	3.26	0.8	9.03	7.18	12.71	1.77	6.65	1.1	8.56	0.58
Automotive	3.08	0.67	9.05	7.33	13.36	2	6.71	1.63	8.41	0.63
Finance	3.46	0.98	7.63	6.32	12.15	1.7	5.59	1	10.03	0.8
Energy	3.34	0.81	8.5	6.36	13	1.88	5.98	0.89	8.74	0.46
Computing & Electronics	3.24	0.84	10.99	8.49	13.07	1.25	8.08	1.09	8.41	0.93
Construction	2.99	1	8.88	7.45	13.08	2.95	6.21	0.88	7.97	0.46
Conglomerate	2.32	0.69	7.64	8.27	10.54	1.93	6.93	1.07	7.11	0.41
Health Care & Pharma	3.9	1.01	8.05	7.35	12.57	1.43	6.67	1.13	8.89	0.52

Table F.8: Proportion of dimensions in collection of texts for the global companies database per region.

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home
Base (all global industries)	3.26	0.8	9.03	7.18	12.71	1.77	6.65	1.1	8.56	0.58
China	2.9	1.1	5.05	6.07	11.05	2.1	4.09	0.82	9.44	0.44
USA	3.59	0.87	9.99	7.74	13.3	1.45	7.35	1.14	8.92	0.82

F.4. RESULTS

Table F.8: Proportion of dimensions in collection of texts for the global companies database per region.

	Positive emotions	Negative emotions	Social	Cognitive processes	Drives	Past-oriented	Present-oriented	Future-oriented	Work	Home
Western Europe	3.2	0.84	9.1	7.4	12.52	1.69	7.31	0.95	8.63	0.47
Asia	2.89	0.64	8.71	7.24	11.8	2.22	6.32	1.49	7.8	0.8

F.4 Results

Statistically significant correlations ($p < 0.05$) between the 2 steps shown in Figure 3.2 of the main manuscripts are shown here as a Sankey diagram. To know the values of the correlations, please contact the authors.

USA sample by sector

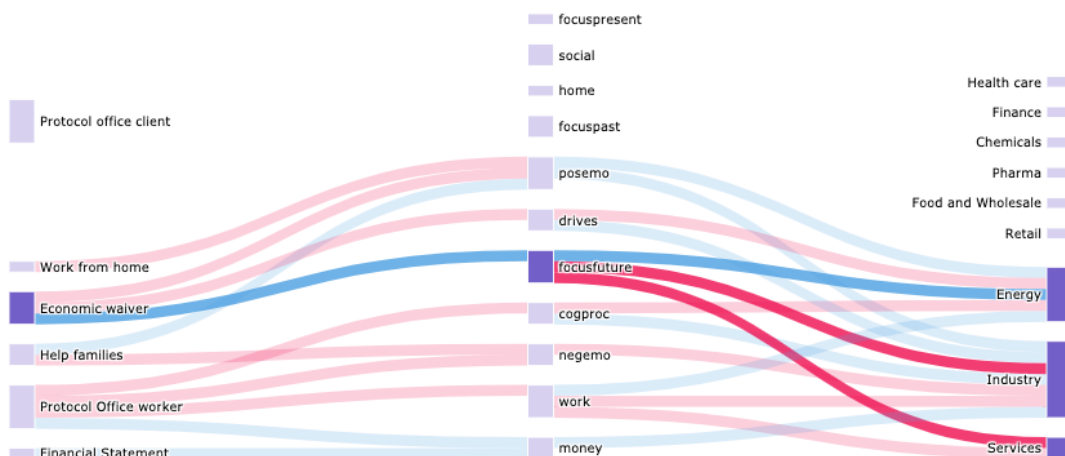


Figure F.2: Sankey diagram representing statistically significant correlations between internal factors (topics – on the left), the attitudes, values and beliefs (psychological dimensions – centre) by industrial sectors for the US sample. Blue links represent positive correlations, while red links represent negative correlations. We highlight those links presented in the main manuscript.

F.4. RESULTS

Global sample by sector

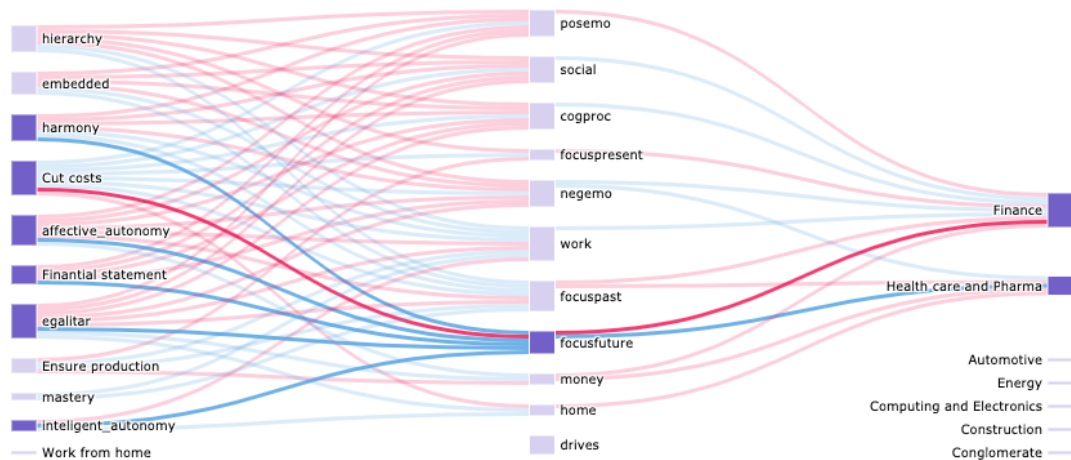


Figure F.3: Sankey diagram representing statistically significant correlations between internal and external factors (topics, cultural value orientations, GCI – on the left), the attitudes, values and beliefs (psychological dimensions – centre) by industrial sectors for the Global sample. Blue links represent positive correlations, while red links represent negative correlations. We highlight those links presented in the main manuscript.

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
1 Walmart	0.801	0	0.14	0	0.055	0
2 Amazon.com	0.516	0	0.302	0.179	0	0
3 Exxon Mobil	0.69	0	0	0	0	0.304
4 Apple	0.356	0	0.156	0.33	0	0.158
5 CVS Health	0.606	0.234	0	0.125	0.034	0
6 Berkshire Hathaway	0	0.831	0	0	0.029	0.138
7 UnitedHealth Group	0.349	0.21	0.104	0.28	0.055	0
8 McKesson	0.619	0	0.1	0.275	0	0
9 AT&T	0	0.429	0	0.453	0.117	0

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
10 AmerisourceBergen	0.411	0	0	0.585	0	0
11 Alphabet	0	0	0.443	0.456	0	0.098
12 Ford Motor	0.093	0	0	0.027	0	0.88
13 Cigna	0.203	0.376	0	0.199	0.22	0
14 Costco Wholesale	0.998	0	0	0	0	0
15 Chevron	0	0	0.251	0.342	0	0.404
16 Cardinal Health	0.096	0	0	0.547	0	0.355
17 JPMorgan Chase	0.038	0.16	0.611	0.119	0.072	0
18 General Motors	0.11	0.473	0.297	0	0	0.118
19 Walgreens Boots Alliance	0.182	0.202	0	0.369	0	0.244
20 Verizon Communications	0.715	0.15	0	0.134	0	0
21 Microsoft	0	0	0.128	0.871	0	0
22 Marathon Petroleum	0.254	0.104	0.201	0.07	0	0.37
23 Kroger	0.858	0	0.136	0	0	0
24 Fannie Mae	0	0.998	0	0	0	0
25 Bank of America	0	0.751	0.247	0	0	0
26 Home Depot	0.743	0.157	0.098	0	0	0
27 Phillips 66	0	0.039	0.096	0.248	0	0.615
28 Comcast	0	0.346	0.246	0.405	0	0
29 Anthem	0	0.05	0.034	0	0.916	0
30 Wells Fargo	0	0.82	0.177	0	0	0
31 Citigroup	0	0.146	0.286	0	0.567	0
32 Valero Energy	0	0.213	0.269	0	0	0.513
33 General Electric	0	0	0	0.632	0	0.362

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company		Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
		<i>Office protocols for clients & suppliers</i>	<i>Office protocols for workers</i>	<i>Economic waivers & insurance coverage</i>	<i>Help & support to families. Donations to healthcare organizations</i>	<i>Enabling Work & Services from Home Protocols</i>	<i>Financial statement for investors & markets</i>
34	Dell Technologies	0	0	0	0.993	0	0
35	Johnson & Johnson	0	0.093	0.162	0.577	0	0.167
36	State Farm Insurance	0	0.422	0.355	0	0.22	0
37	Target	0.808	0	0.191	0	0	0
38	International Business Machines	0	0	0	1	0	0
39	Raytheon Technologies	0.103	0	0.477	0.173	0	0.245
40	Boeing	0	0.222	0	0	0	0.777
41	Freddie Mac	0	0.996	0	0	0	0
42	Centene	0	0	0	0.128	0.871	0
43	United Parcel Service	0.488	0	0	0	0	0.504
44	Lowe's	0.76	0	0.238	0	0	0
45	Intel	0	0	0	0.996	0	0
46	Facebook	0.068	0.123	0.231	0.526	0	0.05
47	FedEx	0.123	0.084	0	0	0.105	0.687
48	MetLife	0.317	0.263	0	0	0	0.418
49	Walt Disney	0.662	0.141	0	0	0	0.196
50	Procter & Gamble	0	0.064	0.399	0.235	0	0.302
51	PepsiCo	0.096	0	0.849	0	0	0.054
52	Humana	0.535	0.296	0	0.057	0.11	0
53	Prudential Financial	0	0.51	0.313	0.176	0	0

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
54 Archer Daniels Midland	0.051	0	0	0	0	0.946
55 Albertsons	0.995	0	0	0	0	0
56 Sysco	0.753	0	0	0	0	0.245
57 Lockheed Martin	0	0.067	0.527	0.253	0	0.151
58 HP	0	0.21	0	0.654	0.134	0
59 Energy Transfer	0.089	0.379	0	0.396	0	0.133
60 Goldman Sachs Group	0	0	0.996	0	0	0
61 Morgan Stanley	0	0	0.672	0	0.327	0
62 Caterpillar	0.026	0	0.526	0.259	0	0.186
63 Cisco Systems	0	0	0	0.998	0	0
64 Pfizer	0	0	0.074	0.777	0.032	0.116
65 HCA Healthcare	0.096	0	0.187	0.309	0.374	0.033
66 American International Group	0	0.999	0	0	0	0
67 American Express	0	0.998	0	0	0	0
68 Delta Air Lines	0	0	0	0	0	0.998
69 Merck	0.036	0.158	0	0.735	0	0.07
70 American Airlines Group	0	0	0	0	0	0.998
71 Charter Communications	0	0.406	0.189	0.38	0	0.024
72 Allstate	0	0	0.993	0	0	0
73 New York Life Insurance	0	0.888	0.108	0	0	0
74 Nationwide	0	0.057	0.114	0.203	0.536	0.089

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
75 Best Buy	0.083	0.044	0.053	0	0.656	0.165
76 United Airlines Holdings	0	0	0	0	0.868	0.131
77 Liberty Mutual Insurance Group	0	0.387	0.406	0.202	0	0
78 Dow	0	0	0	0.377	0	0.613
79 Tyson Foods	0.374	0	0.62	0	0	0
80 TJX	0.754	0	0	0	0	0.241
81 TIAA	0	0.431	0.142	0	0.425	0
82 Oracle	0	0.09	0	0.902	0	0
83 General Dynamics	0.797	0	0	0	0	0.196
84 Deere	0.25	0.104	0	0	0	0.643
85 Nike	0.049	0.07	0.525	0.198	0	0.158
86 Progressive	0	0.989	0	0	0	0
87 Publix Super Markets	0.72	0	0	0	0	0.268
88 Coca-Cola	0.167	0	0.307	0	0	0.524
89 Massachusetts Mutual Life Insurance	0	0.35	0.647	0	0	0
90 Tech Data	0	0	0	0	0.226	0.773
91 World Fuel Services	0	0	0.241	0	0	0.756
92 Honeywell International	0	0	0.339	0	0	0.654
93 ConocoPhillips	0	0	0.589	0.113	0	0.295
94 USAA	0.118	0.232	0.645	0	0	0
95 Exelon	0.148	0.358	0	0.249	0	0.244

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
96 Northrop Grumman	0.042	0.033	0.241	0.68	0	0
97 Capital One Financial	0.114	0.877	0	0	0	0
98 Plains GP Holdings	0	0	0	0	0.998	0
99 AbbVie	0	0.045	0.333	0.561	0	0.059
100 StoneX	0	0	0.248	0.467	0.283	0
101 Enterprise Products Partners	0	0.041	0.15	0.275	0.434	0.099
102 Northwestern Mutual	0	0.511	0.091	0.384	0	0
103 3M	0	0	0.248	0.145	0	0.604
104 Abbott Laboratories	0.024	0.063	0.147	0.678	0	0.088
105 CHS	0	0.136	0	0.621	0	0.24
106 Travelers	0.26	0.375	0.363	0	0	0
107 Philip Morris International	0.134	0	0.259	0.604	0	0
108 Raytheon	0	0	0.321	0.307	0.369	0
109 Hewlett Packard Enterprise	0	0	0.24	0.759	0	0
110 Arrow Electronics	0	0	0	0.411	0	0.584
111 Dollar General	0.778	0.079	0.142	0	0	0
112 Starbucks	0.61	0.08	0.248	0	0	0.061
113 Bristol-Myers Squibb	0	0	0	0.994	0	0

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
114 US Foods Holding	0.811	0	0	0	0	0.187
115 Mondelez International	0	0	0.386	0.121	0.026	0.466
116 Paccar	0	0	0	0	0.766	0.229
117 Thermo Fisher Scientific	0.57	0	0	0	0	0.429
118 Macy's	0.634	0	0	0	0	0.357
119 Jabil	0.045	0	0	0.514	0	0.44
120 Kraft Heinz	0	0	0	0	0.999	0
121 Duke Energy	0	0.991	0	0	0	0
122 Tesla	0.24	0.371	0	0	0	0.387
123 Qualcomm	0	0	0.552	0.443	0	0
124 CBRE Group	0	0	0	0.85	0	0.123
125 Baker Hughes	0	0	0.121	0.262	0	0.613
126 Synnex	0.229	0.132	0.072	0.241	0	0.326
127 Dollar Tree	0.999	0	0	0	0	0
128 Cummins	0.124	0.163	0.083	0.626	0	0
129 United Natural Foods	0.765	0	0.174	0	0	0.06
130 Micron Technology	0.289	0	0	0	0	0.71
131 Amgen	0	0.246	0	0.729	0	0.011
132 Penske Automotive Group	0.038	0	0.132	0	0.726	0.103
133 Visa	0.033	0.481	0.131	0.324	0.03	0
134 Broadcom	0	0.546	0	0.443	0	0
135 Nucor	0.999	0	0	0	0	0
136 Gilead Sciences	0	0.123	0	0.871	0	0

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company		Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
		<i>Office protocols for clients & suppliers</i>	<i>Office protocols for workers</i>	<i>Economic waivers & insurance coverage</i>	<i>Help & support to families. Donations to healthcare organizations</i>	<i>Enabling Work & Services from Home Protocols</i>	<i>Financial statement for investors & markets</i>
137	Southwest Airlines	0	0	0	0	0	0.998
138	Halliburton	0	0	0	0	0	0.995
139	CenturyLink	0.159	0.43	0	0.272	0	0.138
140	International Paper	0.098	0	0.383	0	0	0.513
141	Eli Lilly	0	0.12	0.095	0.78	0	0
142	Aflac	0.097	0.901	0	0	0	0
143	Lennar	0.382	0.61	0	0	0	0
144	Occidental Petroleum	0.089	0	0.163	0.375	0	0.371
145	Union Pacific	0.822	0	0	0	0	0.177
146	Rite Aid	0.915	0	0.08	0	0	0
147	PNC Financial Services Group	0.708	0.139	0	0	0	0.151
148	DuPont	0	0	0	0.295	0	0.699
149	Southern	0	0.323	0.32	0.354	0	0
150	AutoNation	0.996	0	0	0	0	0
151	DXC Technology	0	0	0	0.981	0	0
152	McDonald's	0.581	0	0	0	0	0.417
153	Marriott International	0.996	0	0	0	0	0
154	ManpowerGroup	0	0	0	0.91	0	0.088
155	Bank of New York Mellon	0	0	0.143	0.616	0	0.238
156	Hartford Financial Services Group	0	0.225	0.533	0	0	0.24

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
157 Danaher	0.411	0	0	0.386	0	0.197
158 Whirlpool	0.283	0	0	0	0	0.713
159 AECOM	0.015	0.012	0	0.935	0	0.022
160 Netflix	0	0	0.823	0	0	0.171
161 Kohl's	0.998	0	0	0	0	0
162 Lear	0.328	0.209	0	0	0	0.46
163 Performance Food Group	0.444	0.029	0	0.157	0	0.369
164 Avnet	0	0	0	0.178	0	0.813
165 Synchrony Financial	0.02	0.525	0.069	0.385	0	0
166 Genuine Parts	0	0	0	0	0.998	0
167 NextEra Energy	0	0.356	0.512	0	0	0.128
168 CarMax	0.919	0	0.078	0	0	0
169 Tenet Health-care	0.452	0	0	0.286	0	0.258
170 Kimberly-Clark	0.228	0	0.557	0	0	0.211
171 Emerson Electric	0	0	0.208	0.593	0	0.197
172 WestRock	0	0.149	0	0.211	0	0.636
173 CDW	0.304	0	0	0	0.273	0.417
174 Sherwin-Williams	0.217	0.082	0.104	0.137	0	0.46
175 Fluor	0	0	0.159	0.322	0	0.515
176 PayPal Holdings	0	0	0.713	0.265	0	0
177 D.R. Horton	0	0.437	0	0	0	0.559
178 HollyFrontier	0	0	0	0	0.998	0
179 Tenneco	0	0	0	0	0.686	0.313

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company		Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
		<i>Office protocols for clients & suppliers</i>	<i>Office protocols for workers</i>	<i>Economic waivers & insurance coverage</i>	<i>Help & support to families. Donations to healthcare organizations</i>	<i>Enabling Work & Services from Home Protocols</i>	<i>Financial statement for investors & markets</i>
180	Becton Dickinson	0	0	0.126	0.34	0.092	0.44
181	Lincoln National	0	0.999	0	0	0	0
182	PG&E	0	0.634	0	0	0	0.362
183	salesforce.com	0.031	0	0.296	0.51	0	0.161
184	Mastercard	0	0.088	0.22	0.575	0	0.116
185	General Mills	0	0	0.994	0	0	0
186	Cognizant Technology Solutions	0	0	0.135	0	0.864	0
187	Marsh & McLennan	0.057	0.25	0.142	0.29	0.121	0.14
188	XPO Logistics	0.091	0.234	0.074	0.09	0	0.51
189	Dominion Energy	0	0.32	0.122	0.323	0	0.235
190	Western Digital	0	0	0.227	0.51	0	0.261
191	Gap	0.627	0	0.071	0.147	0	0.153
192	Aramark	0.102	0	0.138	0.619	0	0.14
193	Principal Financial	0	0.529	0.117	0.088	0.263	0
194	Ross Stores	0.157	0	0	0	0.827	0.015
195	Colgate-Palmolive	0.073	0	0.734	0.19	0	0
196	American Electric Power	0.176	0.449	0.094	0	0	0.28
197	Nordstrom	0.213	0.405	0.199	0	0.173	0
198	Jacobs Engineering Group	0	0	0.328	0.496	0	0.174
199	Waste Management	0.161	0.657	0	0	0	0.179

F.4. RESULTS

Table F.9: Distribution of topics found in the companies of the US database.

Company	Topic 1 <i>Office protocols for clients & suppliers</i>	Topic 2 <i>Office protocols for workers</i>	Topic 3 <i>Economic waivers & insurance coverage</i>	Topic 4 <i>Help & support to families. Donations to healthcare organizations</i>	Topic 5 <i>Enabling Work & Services from Home Protocols</i>	Topic 6 <i>Financial statement for investors & markets</i>
200 C.H. Robinson Worldwide	0	0	0	0.49	0.139	0.364
201 PPG Industries	0	0	0.996	0	0	0
202 Omnicom Group	0	0	0.147	0	0.685	0.164
203 Loews	0.094	0.081	0.457	0.252	0	0.115
204 Ecolab	0.211	0	0	0.378	0	0.409
205 Stryker	0	0.095	0.415	0.485	0	0
206 Estee Lauder	0.041	0	0.321	0	0.355	0.282
207 Goodyear Tire & Rubber	0.798	0	0	0	0	0.197
208 Truist Financial	0	0	0.583	0.41	0	0
209 Applied Materials	0	0	0	0.997	0	0
210 BlackRock	0	0	0.997	0	0	0
211 Stanley Black & Decker	0	0	0.144	0.522	0.056	0.277
212 Freeport-McMoRan	0	0	0	0.339	0.105	0.551
213 Texas Instruments	0	0.143	0	0.852	0	0
214 Biogen	0.115	0	0	0.792	0.085	0
215 Parker-Hannifin	0	0	0.211	0.615	0	0.169
216 Reinsurance Group of America	0	0.076	0.26	0.517	0	0.143
217 Howmet Aerospace	0	0	0	0	0.697	0.298

F.4. RESULTS

Global database by region

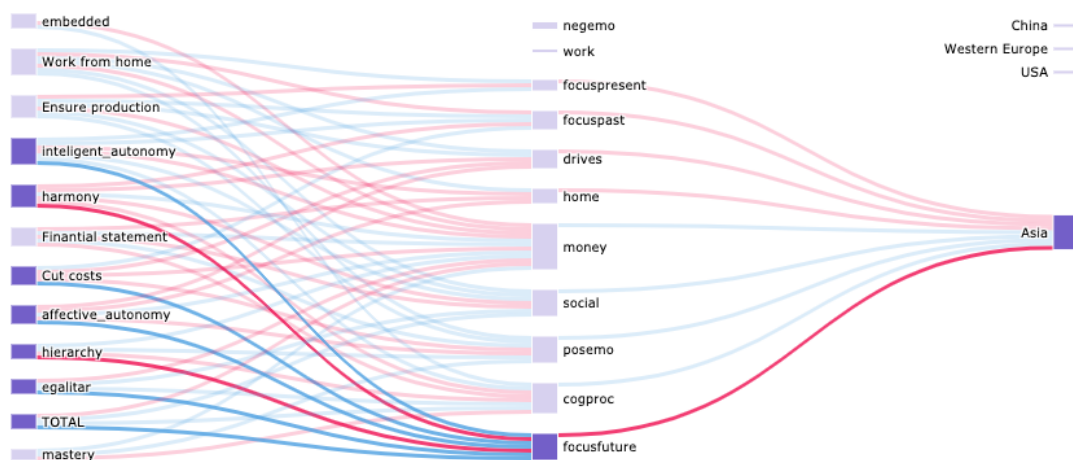


Figure F.4: Sankey diagram representing statistically significant correlations between internal and external factors (topics, cultural value orientations, GCI – on the left), the attitudes, values and beliefs (psychological dimensions – centre) by region for the Global sample. Blue links represent positive correlations, while red links represent negative correlations. We highlight those links presented in the main manuscript.

Table F.10: Distribution of topics found in the companies of the Global database.

Company	Topic 1 <i>Ensure production chain & equipment at work</i>	Topic 2 <i>Cost reduction, focus on income & production</i>	Topic 3 <i>Financial statement for investors & markets</i>	Topic 4 <i>Work from Home policy</i>
1 Walmart	0.08	0	0	0.917
2 Sinopec Group	0	0	0	0.99
3 State Grid	0.99	0	0	0
4 China National Petroleum	0	0.053	0	0.938
5 Royal Dutch Shell	0	0	0.193	0.804
6 Saudi Aramco	0	0	0	0.994
7 Volkswagen	0	0	0	0.997
8 BP	0	0	0.026	0.972
9 Amazon.com	0.052	0.113	0	0.834
10 Toyota Motor	0	0	0.245	0.746
11 Exxon Mobil	0	0	0	0.993
12 Apple	0	0	0	0.998
13 CVS Health	0.012	0.017	0	0.971

F.4. RESULTS

Table F.10: Distribution of topics found in the companies of the Global database.

Company			Topic 1	Topic 2	Topic 3	Topic 4
			<i>Ensure production chain & equipment at work</i>	<i>Cost reduction, focus on income & production</i>	<i>Financial statement for investors & markets</i>	<i>Work from Home policy</i>
14	Berkshire Hathaway		0.124	0.871	0	0
15	UnitedHealth Group		0.117	0.042	0	0.84
16	McKesson		0	0.036	0	0.959
17	Glencore		0	0.085	0	0.91
18	China State Construction Engineering		0	0	0	0.984
19	Samsung Electronics		0	0	0	0.996
20	Daimler		0.127	0	0.231	0.638
21	Ping An Insurance		0.412	0.544	0	0.043
22	AT&T		0.011	0	0	0.988
23	AmerisourceBergen		0	0.072	0	0.925
24	Industrial & Commercial Bank of China		0.584	0.163	0.104	0.149
25	Total		0	0	0	0.992
26	Hon Hai Precision Industry		0.01	0	0	0.977
27	Trafigura Group		0	0	0	0.988
28	EXOR Group		0.253	0.09	0	0.657
29	Alphabet		0	0.31	0	0.687
30	China Construction Bank		0	0.999	0	0
31	Ford Motor		0.018	0	0	0.981
32	Cigna		0.255	0	0	0.743
33	Costco Wholesale		0	0	0.997	0
34	AXA		0	0.145	0	0.853
35	Agricultural Bank of China		0.283	0.716	0	0

F.4. RESULTS

Table F.10: Distribution of topics found in the companies of the Global database.

Company	Topic 1 <i>Ensure production chain & equipment at work</i>	Topic 2 <i>Cost reduction, focus on income & production</i>	Topic 3 <i>Financial statement for investors & markets</i>	Topic 4 <i>Work from Home policy</i>
36 Chevron	0	0	0.095	0.902
37 Cardinal Health	0	0	0	0.996
38 JPMorgan Chase	0.188	0.363	0	0.449
39 Honda Motor	0	0.023	0.142	0.832
40 General Motors	0	0.245	0.312	0.441
41 Walgreens Boots Alliance	0	0	0	0.994
42 Mitsubishi	0	0	0.99	0
43 Bank of China	0	0	0	0.994
44 Verizon Communications	0	0	0.037	0.962
45 China Life Insurance	0	0.989	0	0
46 Allianz	0.827	0	0.169	0
47 Microsoft	0	0	0	0.999
48 Marathon Petroleum	0	0.082	0	0.916
49 Huawei Investment & Holding	0	0	0	0.998
50 China Railway Engineering Group	0	0	0.126	0.862
51 Kroger	0	0	0	0.994
52 SAIC Motor	0	0.044	0	0.954
53 Fannie Mae	0	0.998	0	0
54 China Railway Construction	0	0	0	0.992
55 Gazprom	0	0	0	0.994
56 BMW Group	0	0	0.239	0.755
57 Lukoil	0	0	0	0.998
58 Bank of America	0	0.593	0	0.405
59 Home Depot	0	0	0.072	0.925

F.4. RESULTS

Table F.10: Distribution of topics found in the companies of the Global database.

	Company	Topic 1 <i>Ensure production chain & equipment at work</i>	Topic 2 <i>Cost reduction, focus on income & production</i>	Topic 3 <i>Financial statement for investors & markets</i>	Topic 4 <i>Work from Home policy</i>
60	Japan Post Holdings	0	0	0.994	0
61	Phillips 66	0.1	0	0	0.898
62	Nippon Telegraph and Telephone	0	0	0	0.998
63	Comcast	0.231	0	0	0.767
64	China National Offshore Oil	0.999	0	0	0
65	China Mobile Communications	0	0	0	0.989
66	Assicurazioni Generali	0.713	0	0	0.285
67	Credit Agricole	0	0.997	0	0
68	Anthem	0.999	0	0	0
69	Wells Fargo	0	0.442	0	0.557
70	Citigroup	0.475	0.519	0	0
71	Valero Energy	0	0	0	0.993
72	Itochu	0.201	0.204	0.023	0.572
73	HSBC Holdings	0.034	0.731	0	0.235
74	Siemens	0	0	0	0.998
75	Pacific Construction Group	0.089	0.56	0	0.35
76	Rosneft Oil	0	0	0	0.998
77	General Electric	0	0	0	0.994
78	China Communications Construction	0.338	0.498	0	0.162
79	China Resources	0	0	0	0.998
80	Prudential	0.997	0	0	0
81	Dell Technologies	0	0	0	0.995
82	Nestle	0	0.076	0	0.921
83	Nissan Motor	0	0	0.998	0
84	Hyundai Motor	0	0	0.996	0

F.4. RESULTS

Table F.10: Distribution of topics found in the companies of the Global database.

Company	Topic 1 <i>Ensure production chain & equipment at work</i>	Topic 2 <i>Cost reduction, focus on income & production</i>	Topic 3 <i>Financial statement for investors & markets</i>	Topic 4 <i>Work from Home policy</i>
85 Legal & General Group	0	0.132	0	0.866
86 Deutsche Telekom	0	0	0	0.998
87 Enel	0.148	0	0	0.85
88 Aviva	0	0.277	0	0.72
89 China FAW Group	0	0	0	0.994
90 China Post Group	0	0	0	0.998
91 Amer International Group	0	0.991	0	0
92 China Minmetals	0.98	0.019	0	0
93 Banco Santander	0.055	0.604	0.169	0.173
94 SoftBank Group	0.211	0	0	0.785
95 Bosch Group	0	0	0.994	0
96 Reliance Industries	0	0.193	0	0.805
97 SK Holdings	0	0	0.051	0.947
98 Carrefour	0.025	0.018	0.014	0.943
99 BNP Paribas	0.047	0.949	0	0
100 Dongfeng Motor	0	0	0	0.995
101 Johnson & Johnson	0	0	0	0.991
102 China Southern Power Grid	0.141	0	0	0.856
103 Electricité de France	0	0	0	0.994
104 Centene	0.248	0	0	0.748
105 Humana	0	0	0	0.994
106 Roche Group	0	0	0	0.996
107 Tokyo Electric Power	0.46	0	0	0.526
108 Vinci	0	0	0	0.985
109 Bayer	0	0	0	0.995