

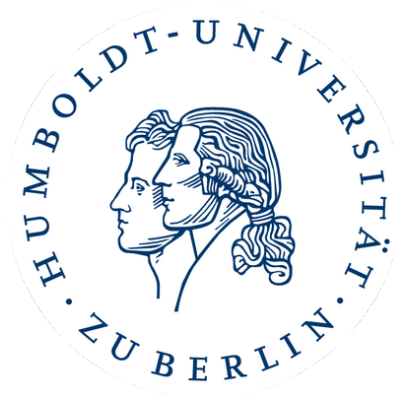
The Influence of Terrorist Hotbeds on the Global Transformation of Terrorist Tactics

Bachelor Thesis

by

Anna Franziska Bothe

(576309)



Humboldt-Universität zu Berlin

School of Business and Economics

Ladislaus von Bortkiewicz Chair of Statistics

Examiner:

Prof. Dr. Wolfgang Härdle

Prof. Dr. Weining Wang

Supervisor:

Dr. rer. nat. Sigbert Klinke

M. Litt. Jannis Jost (Kiel University (ISPK))

in partial fulfillment of the requirements

for the degree of

Bachelor of Science in Business

February 18, 2019

Abstract

The interest in studies regarding the diffusion of transnational terrorism increases constantly. Nevertheless, the existing literature is scarce. The objective of the underlying thesis is to examine the influence of hotbeds on the global transformation of tactics. The transformation and diffusion of terrorist tactics is analyzed from 1977-2017. Two cluster analysis are applied with the objective to group similar observations. Afterwards, heatmaps are generated to visualize and demonstrate the diffusion of tactics. The causality of observed patterns is discussed in detail. Moreover, definitions for the terms hotbed and tactics are given and a weighted Jaccard coefficient is introduced. Also challenges regarding the analysis of categorical data are presented and possible approaches to solve those difficulties are given.

Contents

List of Figures	iv
List of Tables	v
1 Introduction	1
2 Data	4
2.1 Global Terrorism Database	4
2.2 Data Quality	5
2.3 Manual Selection	7
2.4 Data Cleaning and Preparation	9
3 Theory	12
3.1 Terrorism	12
3.2 Hotbeds	12
3.3 Tactics	13
4 Method	14
4.1 Cluster Analysis	14
4.1.1 Proximity Measure	14
4.1.2 Hierarchical Clustering	16
4.1.3 Fuzzy C-Means	19
4.2 Contingency Tables and Heatmaps	20
5 Implementation	22
5.1 Proximity Measure	22
5.2 Cluster Analysis Difficulties	23
5.2.1 Hierarchical Clustering	23
5.2.2 Fuzzy C-Means	26
5.2.3 Other Possible Approaches	26
5.3 Hotbeds and Tactics	26

5.4 Contingency Tables and Heatmaps	28
6 Qualitative Analysis of Empirical Results	29
7 Discussion	33
8 Conclusions	35
References	37
Appendix	41

List of Figures

1	Absolute frequencies of terrorist incidents per year between 1970 and 2017	6
2	Comparison of similarity measures between the Jaccard coefficient matrix (Jac) and the weighted Jaccard coefficient matrix (Jac_w).	23
3	Dendrograms for the first period (t1) of the single and complete linkage method as well as the ward's algorithm	25
4	Armed assault with automatic or semi-automatic firearms against armed targets	29
5	Assassination with automatic or semi-automatic firearms against governmental targets	31
6	Armed assault with automatic or semi-automatic firearms against civil targets	43
7	Bombing/explosives with a grenade against an armed actor	43
8	Bombing/explosives with an unknown explosive type against an armed actor	44
9	Hostage taking (kidnapping) with an unknown gun type of civil targets	44
10	Armed assault with an unknown gun type against civil targets	45
11	Assassination with an unknown gun type against civil targets	45
12	Assassination with an unknown gun type against government targets .	46

List of Tables

1	Institutes of GTD data collection phases	5
2	Overview of inclusion criteria, (GTD Codebook, 2017).	8
3	Overview of inclusion criteria and "Doubt Terrorism Proper" variable .	8
4	Values incl. imputation of <i>vicinity</i> and <i>multiple</i> – 0 shows the absense of the variable and 1 represents the presence of the variable	10
5	The five periods that are divided by the changing collection methodolo- gies and max. time periods no longer than ten years	11
6	Overview of existing grouping algorithms with n_P , n_R and n_Q represent- ing the number of objects within the group of P, R and Q (Härdle and Simar, 2012)	17
7	Squared Euclidean distance matrix of single linkage after first merge (left), exemplary distance matrix (middle), distance matrix of complete linkage after first merge (right)	18
8	Example of variables with different factor levels coded as binary variables	22
9	The number of observations (n), with the computed elements of the symmetric distance matrix and the size of RAM that is needed for the computation for the given time periods t1-t5, specified in Table 5. . . .	24
10	Results of hotbed analysis for t1-t5	27

1 Introduction

Since 9/11, especially religiously shaped terrorism is more present to most people even though terrorism is a very old phenomenon (Richardson, 2007). Nevertheless, the total number of terrorist incidents has risen almost constantly during the last decades (Smith and Zeigler, 2017). Also, radicalization can be observed in form of an increasing lethality of attacks over time (Enders and Sandler, 2000). But how exactly do terrorism or rather terrorist tactics diffuse and change over time?

As a pioneer study Midlarsky et al. (1980) examined the spread of international terrorism within countries but also across borders, comparing Latin America and Western Europe as regions via an analysis of types of terrorism from 1968-1971 and 1973-1974. The results of the analysis show significant autocorrelation effects for bombings, kidnappings and hijackings, indicating that these tactics are contagious. A tactic is "a specific action intended to get a particular result" (Cambridge Dictionary, 2008). Midlarsky et al. (1980) also introduced the theory of hierarchies which implies that countries with a less diplomatic prominence imitate terrorist tactics of countries with a strong diplomatic presence. Even though their research proved the theory of hierarchies wrong, they demonstrated the spread of terrorist tactics.

In order to understand and interpret the results of a geographical diffusion of terrorist tactics, the social movement theory can be adapted. The social movement theory seeks to explain the reasons as well as the timing of emerging movements and its consequences on society, culture and politics (Doug et al., 1996). It is based on mobilization of resources, political opportunities and framing (Doug et al., 1996). Mobilization of resources refers to the available resources. Political opportunities arise as soon as the political environment is unstable, as for example due to corruption, violent governance or social strains, combined with a general grievance of the citizens. Framing is the justification and explanation of one's actions and is also used as a recruiting method (Beck, 2008). It is an important pillar of the social movement theory.

Beck (2008) shows that the social movement theory can be used as an approach to explain and understand the collective behaviour in established, regular political

decision-making processes, as well as in disruptive and contentious political contexts, including political radicalization like terrorism.

Hence, mobilization of resources, political opportunities and framing are also needed for terrorism to arise and proliferate. An unstable environment gives terrorist organizations the opportunity to recruit supporters and to get enough resources for terrorist acts while framing their actions with their ideologies.

According to McAdam (1983), tactics must change over time in order to hinder counter-movements from adapting and developing efficient counter-actions. As an example, from 2011-2017 the three most frequently used weapon subtypes¹ were unknown gun types, followed by unknown explosives and projectiles (rockets, mortars, RPGs, etc.). In contrast, the most used weapon types between 1980-1987 were automatic or semi-automatic rifles, handguns but also unknown explosives.² This shows that terrorist tactics, such as deciding which weapon subtypes to use, change over time.

Hardly any studies have been conducted regarding the transnational diffusion and contagion of terrorism on country level. Existing literature concerning the geography of transnational terrorism constrains their analysis on aggregated regions, such as the spread of terrorist tactics from Iraq (Stern and McBride, 2013); others focus on neighborhoods in countries with many terrorist incidents, such as the research of the transnational diffusion of terrorism on country level from 1980-1997 by Braithwaite and Li (2007).

The objective of this thesis is to examine the influence of terrorist hotbeds on the global transformation of terrorist attack tactics. In contrast to Braithwaite and Li (2007), I study transnational terrorism on a wider time frame (1977-2017) and define tactics of terrorist afflicted countries, comparing them with the global transformation of tactics. Moreover, I introduce a definition of the term *hotbed*. My hypothesis implies a change of global terrorist tactics originating from those defined hotbeds and I discuss the causality of the observed patterns presented in heatmaps.

¹ Missing values are not considered for the comparison.

² The data are drawn from the Global Terrorism Database (GTD) which is described and analyzed in detail in the up-coming chapters.

Real-life data are often categorical such as the data I am using in order to examine the hypothesis of this paper. The selected data are exclusively nominal scaled which implies unordered levels that can only be compared regarding their equality. Two approaches of cluster analysis and heatmaps as a visualization technique are applied to group different types of observations. Cluster analysis with only categorical data with different amounts of levels are also very sparse. Therefore, I introduce an adjusted similarity matrix that can be used for the analysis of categorical data with different levels.

The thesis is organized as follows. Section 2 deals with the data, the data set as well as its advantages and disadvantages are described. Also the data cleaning and preparation process is explained. Section 3 is about the general theory of terrorism and therefore gives definitions of the terms terrorism, hotbeds and tactics. The applied methods and algorithms are described in detail in section 4. Followed by the analysis of the data in section 5. In section 6 the results of the analysis are presented and evaluated. Chapter 7 includes a detailed discussion about possible bias and further research possibilities. Finally yet importantly, a conclusion regarding the stated theory and the presented results is given in section 8.

2 Data

In the following subsections, basic information are given about the development of the database, its constraints and its quality. Afterwards, the manual selection process as well as the data cleaning and preparation is described in detail.

2.1 Global Terrorism Database

The Global Terrorism Database (GTD) is the most comprehensive, unclassified open-source database of terrorist incidents. It is collected primarily via manual selection efforts by institutes and service companies in the United States of America (USA). Since 2012, also nature language processing and machine learning methods are used. It contains longitudinal data that have been collected from 1970 till the end of 2017. The data are retrieved from publicly available primary and secondary domestic as well as international data sources such as (social) media, newspaper articles, other databases and archives recording terrorist activities. Only sources are included in the database that are proven to be credible. The collection teams focus on six main categories: location, perpetrators, targets, weapons and tactics, causalities and consequences, and general information (National Consortium for the Study of Terrorism and Responses to Terrorism (START), 2018).

The responsible for the collection process have changed several times. First the Pinkerton Global Intelligence Service (PGIS) has collected the data, followed by the Center for Terrorism and Intelligence Studies (CETIS) and the Institute for the Study of Violent Groups (ISVG). Since November 2011, National Consortium for the Study of Terrorism and Responses to Terrorism (START) is collecting the data.

Not only have the responsible changed but also the data collection methodology was improved and new technologies were adapted. Data inconsistencies and possible bias resulting from the listed changes, will be described in section 2.2 and discussed in section 7.

The data collection process can be divided into four main phases:

	PGIS	CETIS	ISVG	START
01/1970-12/1997	x			x
01/1998-03/2008		x		x
04/2008-10/2011			x	x
11/2011-12/2017				x

Table 1: Institutes of GTD data collection phases

The data of 1993 were lost and reconstruction efforts in cooperation with the National Consortium for the Study of Terrorism and Responses to Terrorism (START) failed since many source materials were unavailable by 2008. Retrospectively, all data from January 1998 till March 2008 were collected by the Center for Terrorism and Intelligence Studies (CETIS).

The database contains 181.691 observations and 135 variables. Categories that are presented by numbers are coded as integers and text variables as factors. Variables that are ratio or interval scaled, are displayed as numeric variables. All variables regarding terrorist tactics that are used for the analysis of the underlying thesis, are nominal scaled and will be transformed to binary variables.

2.2 Data Quality

The GTD is collected by coding and collection rules defined in the codebook (GTD Codebook, 2017) which also includes all key decisions that have been made during the development of the database. The collection team aims to make the collection methodology as transparent, comprehensive and consistent as possible to maintain a high data quality.

The database is not a random sample of terrorist incidents. The collectors attempt to gather all data on terrorist attacks worldwide. However, due to technology and communication boundaries it is not always possible. The technological framework

has changed tremendously within the last decades. When the collection started in 1970, US institutes mainly collected regional data from the USA as they had limited access to international information. Their main sources were government reports, wire services and major international newspapers (National Consortium for the Study of Terrorism and Responses to Terrorism (START), 2018). According to the database, in 1970, most terrorist incidents – 468 ($\sim 72\%$) – were recorded in the United States. The second most incidents were documented in Uruguay with only 33 ($\sim 5\%$) terrorist acts. Especially within the first years of the GTD collection, the number of terrorist incidents per region does not reflect the true situation of the early 1970s. In Latin America several terrorist incidents occurred, but only about a quarter of incidents compared to the USA were recorded. Therefore, the data quality of the early 1970s, regarding generalizability, is very low. Until today, developing countries have less access than industrialized countries to the internet, being the main source of the GTD today. Hence, their incident reports are not always available. This leads to skewed information recordings in the GTD.

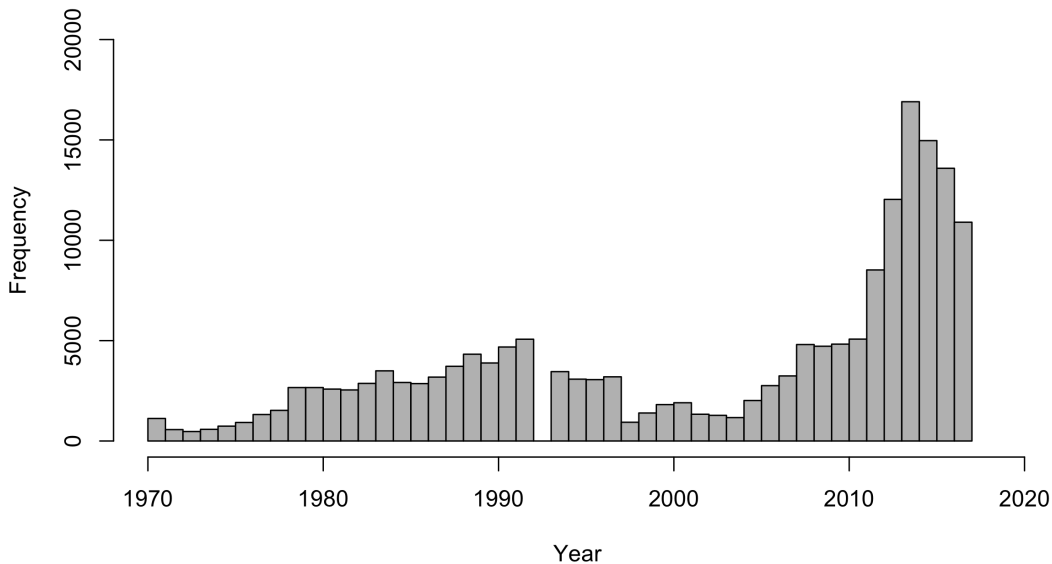


Figure 1: Absolute frequencies of terrorist incidents per year between 1970 and 2017

Additionally, the data are biased due to the different collection methods. As stated in chapter 2.1, at first the data was collected manually and later, Natural Language Processing (NLP) and machine learning algorithms were integrated. After improving the collection methodology, the total number of incidents recorded every year increased (see Figure 1).

This bias can be overcome, by subdividing the data by time periods with the same collection methodology and also, by using relative frequencies instead of absolute numbers while comparing the observations of different years or even across time periods.

Overall, the data quality is good and the data can be used by calculating relative frequencies and having the different collection methodologies in mind, when comparing absolute numbers.

2.3 Manual Selection

Firstly, the variables month, year, country and region³ are included as temporal and spatial data. Secondly, all variables that can be associated with tactics are selected. This includes the weapons (incl. weapon subtypes), the general attack type, the specific target types as well as the variables regarding the vicinity, the multiplicity of incidents within a short time period and suicide strokes. In case of variables with three different types, e.g. *attacktype1*, *attacktype2*, *attacktype3*, only the primarily used type is included in the data set to be analyzed. The second and third types are not included due to many missing values ($\emptyset \sim 97\%$), leading to a low meaningfulness.

The perpetrator information as well as the information concerning causalities and consequences are excluded from the analysis since they cannot be fully-controlled by the perpetrators. Moreover, "post-tactics" such as claiming an incident afterwards, are not analyzed since the data are not reliable because there are often competing claims and $>97\%$ missing values in the given data set.

The definition of terrorism often overlaps with other forms of violence attacks. Therefore, two filters are employed. The first filter are the inclusion criteria in order to decide whether or not an incident is added to the database. The data set includes all

³ An overview of the regional classification of the countries is given in the appendix.

terrorist events that fulfill at least two of the three inclusion criteria mentioned below. For the purpose of this analysis, all data are excluded that do not fulfill all three of the following inclusion criteria (crit1, crit2, crit3):

variable	definition
crit1	political, economic, religious, or social goal
crit2	intention to coerce, intimidate or convey to larger audience
crit3	outside international humanitarian law

Table 2: Overview of inclusion criteria, (GTD Codebook, 2017).

Secondly, a filter is included that indicates if there is a doubt regarding the terrorism property. Since the underlying thesis focuses only on terrorist incidents, all terrorist acts with a doubt if it is terrorism proper, are excluded as well. In contrast to the inclusion criteria which was implemented for the whole data set, the "Doubt Terrorism Proper" variable is only systematically available for incidents after 1997. Many of the values before 1997 are coded as unknown (-9). In the case of unknown values the decision if an observation is kept or not kept for the analysis is only made based on the inclusion criteria. Observations of contrary filter results are discarded, cp. Stern and McBride (2013) as shown in Figure 3. Only the bold marked observations are included in the analysis:

	0 (=no doubt)	1 (=doubt)	-9/NA (=unknown)
all criteria fulfilled	138.879	3.139	13.744
not all criteria fulfilled	26	25.862	41

Table 3: Overview of inclusion criteria and "Doubt Terrorism Proper" variable

Furthermore, only observations of successful terrorist attacks are included, as failed terrorist attacks are unlikely to be imitated, cp. (Braithwaite and Li, 2007).

Last but not least, all incidents that are unknown for two or all of the three following variables: attack, weapon and target type are excluded because the variables lose their meaningfulness when being "unknown". In contrast, to the variables vicinity, multiple

and suicide still contain a meaningfulness even if they are not present (0).

2.4 Data Cleaning and Preparation

Israel and Palestine are merged to *Israel/Palestine* as they are closely related in terms of terrorist incidents. Only 137 (71%) of the UN members (United Nations A/RES/67/19, 2012) accept the independence of the State of Palestine and the Oslo I Accord that was signed by Israel and Palestine in 1993 announcing Palestine as a "Palestinian interim self-government" but not as an independent Palestinian state (Shlaim, 1994). Additionally, according to the three element doctrine by Georg Jellinek, a state is defined by a state territory, state people and state power (Jellinek, 1990). In order to practically be independent, Palestine has to fulfill those criteria. Since the territories of Palestine were split in the Oslo Agreement, the Palestinian Authority just controls territory A and has only partial or no control over area B and C. Furthermore, Israel has many Israeli settlements in Palestine – 653.621 Israeli lived in the West Bank incl. Jerusalem Area J1 und J2 in 2017 (Palestinian Central Bureau of Statistics, 2018). As a consequence, Palestine does also not fulfill the criteria of its own state people. Hence, Palestine is stateless but the Palestine territories also do not belong to sovereign Israel.

Moreover, only two of five borders of Israel are internationally accepted. The Green Line, separating the West Bank from Israel, is e.g. globally not recognized (Newman, 2012).

The difficulties regarding international recognition and independence as well as the unclear borders and especially the omnipresence of the Israeli military in Palestine lead to a close relationship between the number of incidents in both countries. The merged country is named after both countries (*Israel/Palestine*) – sorted alphabetically. Other countries that are not internationally recognized as independent are not merged because they lack relevance for the later hotbed analysis.

All NAs of the weapon subtypes (*weapsubtype1.txt*) are imputed with either the value of the corresponding weapon type (*weaptype1.txt*) or with a fitting value of the weapon subtypes. As an example, no weapon subtype for the weapon type *Biological* exists, thus *Biological* is imputed as the value for the weapon subtype. If the

weapon subtype already has a category for unknown observations, the weapon subtype was coded as such. The weapon type *Explosives* already has the sub-category *Unknown Explosives Type*. The process is analogously applied to the numerical coded *weapsubtype1*.

In order to hinder the analysis of becoming too fuzzy, the target types are summed up into the following categories: civil targets, government targets, armed actors and unknown targets. The target categories are subdivided in hard and soft targets. Soft targets are unarmed and vulnerable targets (like civilians); hard targets are targets that are well protected (like governmental institutions) or that are able to fire back (like the military or other terrorist groups, thus armed actors), cp. (Berman and Laitin, 2005). All unknown target types are excluded from the soft and hard target clustering and form an own category (unknown).

The categorical variables with two levels that have missing (NA) or not available data (-9) are *vicinity* and *multiple* as shown in Table 4.

	0	1	NA/-9	imputed value
vicinity	168.932	12.724	35	0
multiple	156.658	25.032	1	0

Table 4: Values incl. imputation of *vicinity* and *multiple* – 0 shows the absense of the variable and 1 represents the presence of the variable

Since $\sim 93\%$ of the observations did not occur in the immediate vicinity of the city but rather in the city itself and $\sim 86\%$ of the incidents were not part of a multiple attack, it can be assumed that the missing observations were the same as the large majority of the observations and therefore, a 0 (non-presence) for all missing cases is imputed.

The data from 1970-1976 are excluded from the analysis because of the information bias described in section 2.2 as well as the few observations (<1000). The data are divided to five time periods, as shown in Table 5:

As explained in section 2.1, the year 1993 is missing completely. The loss is not

periods	time	observations
t1	01/1977-12/1987	20.722
t2	01/1988-12/1997	25.686
t3	01/1998-03/2008	14.168
t4	04/2008-10/2011	14.874
t5	11/2012-12/2017	52.114

Table 5: The five periods that are divided by the changing collection methodologies and max. time periods no longer than ten years

systematic and the data are missing completely at random (MCAR). Hence, no data were imputed for 1993.

Following the manual selection as well as the data cleaning and preparation, the data set consists of 127.564 observations and seven categorical variables regarding the tactic, two geographical variables presented by region and country as well as the year of the incident.

3 Theory

In the following chapter, the terms terrorism, hotbeds and tactic, which are relevant for the underlying thesis, are presented and defined.

3.1 Terrorism

Not every act of violence is an act of terrorism. It has to meet certain criteria in order to be categorized as a terrorist act. The term terrorism is defined by the GTD Codebook as "the threatened or actual use of illegal force and violence by a non-state actor to attain a political, economic, religious, or social goal through fear, coercion, or intimidation" (GTD Codebook, 2017). This implies that the incident has to be intentional, violent or at least a threat of violence and that the perpetrators have to be sub-national actors. The last criteria regarding the non-state actors is debatable. However, since the codebook of the GTD database as well as the previously cited papers follow the non-state actor approach, it is not be discussed in this context (see among others Beck (2008), Richardson (2007)). As the definition of terrorism often overlaps with the definition of other crimes and acts of political violence, the additional filter "doubt terrorism proper?" is added. The filter is used to exclude the following non-terrorism acts: insurgency/guerilla actions, intra/inter-group conflicts, lack of intention, state actors and other crime types (National Consortium for the Study of Terrorism and Responses to Terrorism (START), 2018).

3.2 Hotbeds

In this section, the term *hotbed* is specified by two criteria. A loose definition for a hotbed is given by Braithwaite and Li (2007): "[...] a neighborhood of countries that experiences a larger number of terrorist incidents than one would expect of an average neighborhood in the international system". This approach incorporates the foundation of both hotbed criteria. After consulting with an expert from the Institute for Security Policy Kiel University (ISPK), the amount of hotbeds per period should at least one and not more than five.

The first criterion is the amount of incidents a country experiences compared to other countries. In order to derive a threshold level for the quantity of terrorist attacks, the 95% quantile is used. Hence, the number of terrorist acts of a country has to experience must be represented in the 95% quantile in order to fulfill the first criterion of being a hotbed for the given time period.

Secondly, the escalation rate is analyzed. To achieve this, a growth rate for each country and every year is computed. This growth rate needs to be above the 75% quantile. Additionally, these growth rates should not be more than two years apart in order to represent an escalation within the country.

3.3 Tactics

As briefly described in the introduction, a tactic needs to be an intentional action that aims to achieve a particular result. This does not only include an active decision regarding the type of weapon that is used, but also the general attack type. This can be an assault, a bombing or similar, and other specifics of the incident, such as the multiplicity and vicinity of attacks. A resolution of a specific tactic does not count as a tactic yet, due to a missing action.

Most of the existing literature present single elements of incidents as tactics. This includes for example suicide attacks, bombings and hijackings, see among others Midlarsky et al. (1980), Stern and McBride (2013). Similar to Gill et al. (2013), who describes a shift of tactics, referring to the targets, the attack methods as well as the "delivery method[s]" as one tactic, I take all actively combined elements of the incident as one tactic and compare those. Therefore, the combination of the following form a tactic: attack type, weapon type and its subtype, target type and whether it was planned as a suicide stroke, as a multiple attack, in the vicinity of the city or in the city itself. When all those factors are equal in two incidents, the terrorist groups used the same tactic.

4 Method

In the following chapter different cluster analysis algorithms as well as a visualization technique are described. All of the described methods are applied on the given data set.

4.1 Cluster Analysis

Cluster analysis is an unsupervised learning method that aims to reduce the dimensions of observations by grouping observations with homogeneous properties. The algorithms minimize the intra-cluster distance and maximize the inter-cluster distance. In contrast to other classification tasks, like discriminant analysis, the groups are not known a priori. Before running the cluster algorithms two choices have to be made. A proper proximity measure and group-building algorithm have to be chosen (Härdle and Simar, 2012).

Cluster analysis is applied in many research fields such as in social science. It can be useful to analyze the identical behaviour of individuals and derive optimal treatments to target the clustered groups (Backhaus et al., 2016).

In the context of the underlying thesis at hand, the clustering approach aims to group incidents of hotbeds and non-hotbeds in order to analyze whether a trend in the spread of terrorist tactics can be observed over time. Two clustering algorithms are applied: an agglomerative hierarchical clustering and a soft clustering approach via fuzzy c-means. For both techniques, the same distance matrix was used.

4.1.1 Proximity Measure

A proximity measure computes the distance between observations in case of continuous variables and the similarity between observations in case of categorical variables of the observations (Chatfield, 2018).

To calculate a similarity matrix, pairs of observations are formed, (x_i, x_j) where $x_i^\top = (x_{i1}, \dots, x_{ip})$, $x_j^\top = (x_{j1}, \dots, x_{jp})$ and $x_{ik}, x_{jk} \in \{0, 1\}$ with p representing the number of binary variables (Hennig, 2015). All variables are compared to each other.

A high similarity of the variables results in larger values computed in the similarity matrix (Härdle and Simar, 2012). Four cases are possible:

$$x_{ik} = x_{jk} = 1,$$

$$x_{ik} = 0, x_{jk} = 1,$$

$$x_{ik} = 1, x_{jk} = 0,$$

$$x_{ik} = x_{jk} = 0.$$

Different methods for calculating a similarity matrix exist. In general, a similarity between two observations is defined as:

$$S_{ij} = \frac{a + \delta d}{a + \lambda(b + c) + \delta d} \quad (1)$$

with

$$\begin{aligned} a &= \sum_{k=1}^p I(x_{ik} = x_{jk} = 1), \\ b &= \sum_{k=1}^p I(x_{ik} = 0, x_{jk} = 1), \\ c &= \sum_{k=1}^p I(x_{ik} = 1, x_{jk} = 0), \\ d &= \sum_{k=1}^p I(x_{ik} = x_{jk} = 0). \end{aligned}$$

Similarity measures differ in their weighting of δ and λ depending on the weights given to mismatching as well as the absence or presence of the same features. In contrast, the Euclidean distance treats the expressions 0 and 1 equally (Härdle and Simar, 2012).

For the purpose of analysis, a similarity coefficient is needed which prioritizes the presence of common features and puts no weight on the absence of common features. Since the aim is to focus on tactics that were used or differently used rather than those that were not used at all.

A similarity coefficient that sets $\delta = 0$ and $\lambda = 1$, is the Jaccard similarity measure (Backhaus et al., 2016). It is computed as:

$$S_{Jac} = \frac{a}{a + b + c} \quad (2)$$

The values within the similarity matrix will always be between 0 (max. difference) and 1 (max. similarity) (Backhaus et al., 2016).

As described in Gower (1966), the similarity matrix has to be positive semi-definite and with the maximum similarity scaled as $S_{ii} = 1$ in order to correctly transform a similarity matrix into a distance matrix. This is needed for hierarchical as well as partitioning algorithms.

$$d_{ij} = \sqrt{2(1 - S_{ij})} \quad (3)$$

4.1.2 Hierarchical Clustering

Hierarchical clustering can be split into an agglomerative ("bottom-up") and a divisive ("top-down") approach. The agglomerative algorithm, which is focused, starts with as many clusters as observations and merges observations iteratively to form larger clusters until only one cluster with all observations is left. The next observations to be merged are chosen by a group-building algorithm as mentioned in section 4.1, using a linear search within the original distance matrix (Härdle and Simar, 2012).

The result of the clustering algorithm can be visualized with a dendrogram which displays the distance between clusters on the y-axis and the indices of data points on the x-axis. The dendrogram indicates how many clusters to choose. The tree should be cut when the distance between clusters is maximized. In case the presented results are ambiguous, the number of clusters to choose has to be weighted against their explanatory value (Backhaus et al., 2016).

Various grouping algorithms exist that differ mainly in their weighting factors (δ_j) which influence the way they compute the distance between an object (R) and a new cluster (P + Q) (Härdle and Simar, 2012):

$$d(R, P + Q) = \delta_1 d(R, P) + \delta_2 d(R, Q) + \delta_3 d(P, Q) + \delta_4 |d(R, P) - d(R, Q)| \quad (4)$$

with

$d(R, P)$ = distance between group R and P

$d(R, Q)$ = distance between group R and Q

$d(P, Q)$ = distance between group P and Q

Among the different grouping methods presented in Table 6, single linkage, complete linkage and Ward's distance are the most used grouping algorithms.

Grouping algorithms	δ_1	δ_2	δ_3	δ_4
Single linkage	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$
Complete linkage	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$
Average linkage (unweighted)	$\frac{1}{2}$	$\frac{1}{2}$	0	0
Average linkage (weighted)	$\frac{n_P}{n_P + n_Q}$	$\frac{n_Q}{n_P + n_Q}$	0	0
Centroid	$\frac{n_P}{n_P + n_Q}$	$\frac{n_Q}{n_P + n_Q}$	$-\frac{n_P * n_Q}{(n_P + n_Q)^2}$	0
Median	$\frac{1}{2}$	$\frac{1}{2}$	$-\frac{1}{4}$	0
Ward	$\frac{n_R + n_P}{n_R + n_P + n_Q}$	$\frac{n_R + n_Q}{n_R + n_P + n_Q}$	$-\frac{n_R}{n_R + n_P + n_Q}$	0

Table 6: Overview of existing grouping algorithms with n_P , n_R and n_Q representing the number of objects within the group of P, R and Q (Härdle and Simar, 2012)

The single linkage, also called nearest neighbour, merges the closest objects due to the distance matrix (Backhaus et al., 2016). Given the exemplary distance matrix in Table 7, B and C would be merged first and the new distance between B+C and the

other objects A and D would be calculated as shown in (5) with the weights given in Table 6. The single linkage method minimizes the distance between the merged objects and the other objects:

$$d(R, P + Q) = \min\{D(R, P), D(R, Q)\} \quad (5)$$

The resulting distance matrix is shown in Table 7. The single linkage method tends to form "chains", which can be a problematic for selecting and forming a reasonable number of clusters.

The difference between complete linkage, also known as the furthest-neighbour method, and single linkage is the calculation of the new distance between objects. Instead of minimizing the new distance matrix, the algorithm maximizes it (Härdle and Simar, 2012):

$$d(R, P + Q) = \max\{D(R, P), D(R, Q)\} \quad (6)$$

An example of the result of a merge, is shown in Table 7. The complete linkage algorithm tends to form small groups. This approach might cause problems with outliers since they are less likely to be spotted and might lead to biased results. Outliers should be eliminated before conducting the complete linkage grouping algorithm. This results in relatively homogeneous clusters which tend to be equally sized.

	A	B, C			A	B	C			A	B, C
B, C	9		⇐	B	15			⇒	B, C	15	
D	7	6		C	9	3			D	7	11
				D	7	6	11				

Table 7: Squared Euclidean distance matrix of single linkage after first merge (left), exemplary distance matrix (middle), distance matrix of complete linkage after first merge (right)

In contrast to the single and complete linkage method, the Ward's algorithm does not merge objects due to their distances, but rather due to the smallest increase of the total within-cluster variance. Furthermore, the Ward's method requires the distances between objects to be squared Euclidean distances. It can be calculated with the weighting factors in Table 6. The Ward's algorithm can be a good choice for a grouping method if outliers are absent and metrically scaled data can be guaranteed. However, it has difficulties as soon as the groups have elongated patterns or many small groups (Backhaus et al., 2016).

To sum it up, the main difference between the linkage methods and the Ward's method is that the linkage measures merge clusters based on their distance to each other and the Ward's algorithm aims to form clusters that are as homogeneous as possible. The single linkage, complete linkage and Ward's method are the most commonly used grouping algorithms (Härdle and Simar, 2012).

4.1.3 Fuzzy C-Means

Other clustering approaches are the partitioning algorithms. The initial cluster centroids are partitioned randomly and the closest observations are assigned to the cluster centroids. The number of clusters to start with has to be set in advance which is one of the main downsides of this approach as the cluster analysis is an unsupervised learning technique and therefore the number of clusters are usually not known before conducting the analysis. The starting cluster centroids are recalculated and then changed by an exchange algorithm between the groups as long as the algorithm reaches its optimum. This optimum is strongly dependent on the starting clusters. Often a local maximum is reached rather than a global maximum. Elements of clusters can be exchanged between clusters during the process, in contrast to the hierarchical algorithm where an element that forms a group with another element during the process cannot be broken up anymore (Backhaus et al., 2016).

A distinction can be made between hard and soft clustering methods. Hard clustering assigns every observation to exactly one cluster. Often problems arise because the observations in each cluster are given equal importance and the researcher can-

not tell how "typical" the observation is for the specific cluster (Keller et al., 1985). Therefore, soft clustering calculates a probability how likely every observation belongs to each cluster in the form of membership coefficients varying between 0 and 1. One soft clustering approach is called fuzzy c-means. It computes the membership degree u_{ij} of the observations x_i ($i = 1, \dots, N$) to every cluster j ($j = 1, \dots, C$) by minimizing the following function:

$$J_m = \sum_{i=1}^N \sum_{j=1}^C u_{ij}^m \|x_i - c_j\|^2 \quad (7)$$

with $1 \leq m \leq \infty$. c_j is representing the center of cluster j and $\|\cdot\|$ expresses the similarity between x_i and the cluster center c_j .

The fuzzy clustering algorithm optimizes – like the crisp c-means algorithm – the results by iterating through the clusters and updating the membership u_{ij} and the cluster centers c_j . The iteration of the fuzzy c-means algorithm stops as soon as a minimum of equation 7 is reached (Li et al., 2008) – at least a local minimum.

4.2 Contingency Tables and Heatmaps

With nominal or ordinal scaled data, many analysis methods, such as the computation of standard deviations, cannot be used for lack of meaningfulness. However, the data can be described via contingency tables and the relationship between two variables can be shown. The only requirements are two categorical variables with at least two levels each. Variable X with I categories and Y with J categories have IJ possible combinations presented as frequencies. The computed table has the size of $I \times J$ (Agresti and Kateri, 2011). A contingency table can be used in order to generate a heatmap out of the computed table.

Heatmaps are two-dimensional graphics that can be expanded on a third dimension with the help of colors which represent the values between the two variables on the axis. The name of heatmaps results from its typical color shading from either yellow to red or from blue/green to red like a temperature scale. The "warmer" the color, the higher the value. The temperature like color shading belongs to the merits of heatmaps since

one can intuitively get an impression of the underlying data. Additionally, heatmaps summarize large amounts of data quickly and trends within the data can be observed easily. A thoughtless, non-critical valuation can also lead to misinterpretations (Bojko, 2009). Possible interpretation bias will be discussed in detail in chapter 7.

5 Implementation

The following chapter is about the implementation of the theory as well as the methods. All described R libraries are listed in the R file *Pack_final.R* attached to the thesis.

5.1 Proximity Measure

In the case of the data used for the underlying thesis, only categorical data are used as described in section 2.4. Due to the binary structure, first a similarity matrix has to be calculated and afterwards transformed into a distance matrix.

Since the categorical variables are "mixed" – some with two levels and some with more than two levels – the Jaccard coefficients b and c of the factor variables with more than two levels need to be weighted. Otherwise, the variables would be double weighted. The weighted Jaccard similarity measure can be applied as follows:

$$S_{Jac_w} = \frac{a_1 + a_2}{a_1 + b_1 + c_1 + a_2 + 0.5 * b_2 + 0.5 * c_2} \quad (8)$$

with a_1, b_1, c_1 presenting variables with two levels and a_2, b_2, c_2 presenting variables with more than two levels. The variables with more than two levels need to get a weight of 0.5 in order to not falsely give double weight to the observation pairs of $x_{ik} = 0, x_{jk} = 1$ and $x_{ik} = 1, x_{jk} = 0$.

As an example: the variable suicide has two levels and the variable attack type has three levels (in the data set it has nine levels, but in terms of simplicity only three variable levels will be presented here).

	suicide	attacktype 1	attacktype 2	attacktype 3
observation 1	0	1	0	0
observation 2	1	0	0	1
observation 3	1	1	0	0

Table 8: Example of variables with different factor levels coded as binary variables

When applying the Jaccard and weighted Jaccard similarity coefficient, the following symmetric similarity matrices can be computed:

$$\begin{array}{cc} \begin{pmatrix} 1 & 0 & \frac{1}{2} \\ & 1 & \frac{1}{3} \\ & & 1 \end{pmatrix} & \begin{pmatrix} 1 & 0 & \frac{1}{2} \\ & 1 & \frac{1}{2} \\ & & 1 \end{pmatrix} \\ Jac & Jac_w \end{array}$$

Figure 2: Comparison of similarity measures between the Jaccard coefficient matrix (Jac) and the weighted Jaccard coefficient matrix (Jac_w).

As observed in the matrices, the similarities differ which is the result to a double weighting of b and c for the variables with more than two levels when computing the measures with the unweighted Jaccard coefficient. The adapted source code for the weighted Jaccard coefficient is attached to the regular R code of the thesis.

In order to transform a similarity matrix properly, the similarity matrix has to meet certain assumptions described in section 4.1.1. Figure 2 shows exemplary that the similarity matrix is positive semi-definite and that it has a maximum similarity scaled as $S_{ii} = 1$. Hence, the assumptions are met and the similarity matrix can be transformed into a distance matrix which is needed for hierarchical as well as the partitioning algorithms. The distance matrix is computed by the *dist* function.

5.2 Cluster Analysis Difficulties

The attempts of conducting a meaningful cluster analysis were not successful. The reasons are discussed in following chapters. Also, possible other approaches are provided which are not tested due to the framework of the underlying thesis.

5.2.1 Hierarchical Clustering

In order to conduct a hierarchical cluster analysis, the variables to be analyzed, have to be transformed to dummies with help of the function *createDummyFeatures* of the

package *mlr*. The produced binary data frame is used for the calculation of distance matrix which is the main problem of hierarchical cluster analysis in big data sets. Every observation has to be compared with each other (Backhaus et al. (2016)) and the computational power increases exponentially. The first time period of the analysis has about 21.000 observations. Calculating the distance matrix for this partial data set takes already more than 6 minutes to run. The distance matrix contains 214.690.281 elements. The computational performance of the used computer (MacBook Pro 2015, 2,7 GHz Intel Core i5, 8 GB RAM, macOS Mojave 10.14.2) is not sufficient enough to run the next time period which is t2 with about 26.000 observations. Also t5 with about 52.000 observations is too big to be computed. Due to the lack of time, a computer with enough RAM could not be organized.

t_x	observations (n)	matrix elements ($\frac{(n^2-n)}{2}$)	RAM in GB ($\frac{n^2*4}{1.024^3}$)
t1	20.722	214.690.281	~ 1, 60
t2	25.686	329.872.455	~ 2, 45
t3	14.168	100.359.028	~ 0, 75
t4	14.874	110.610.501	~ 0, 82
t5	52.114	1.357.908.441	~ 10, 11

Table 9: The number of observations (n), with the computed elements of the symmetric distance matrix and the size of RAM that is needed for the computation for the given time periods t1-t5, specified in Table 5.

Analyzing the results of the hierarchical cluster of t1 and the corresponding dendrograms, indicates that the cluster algorithm is not the right approach. As grouping algorithms the single and complete linkage and the ward's method have been applied, resulting in the following dendrograms:

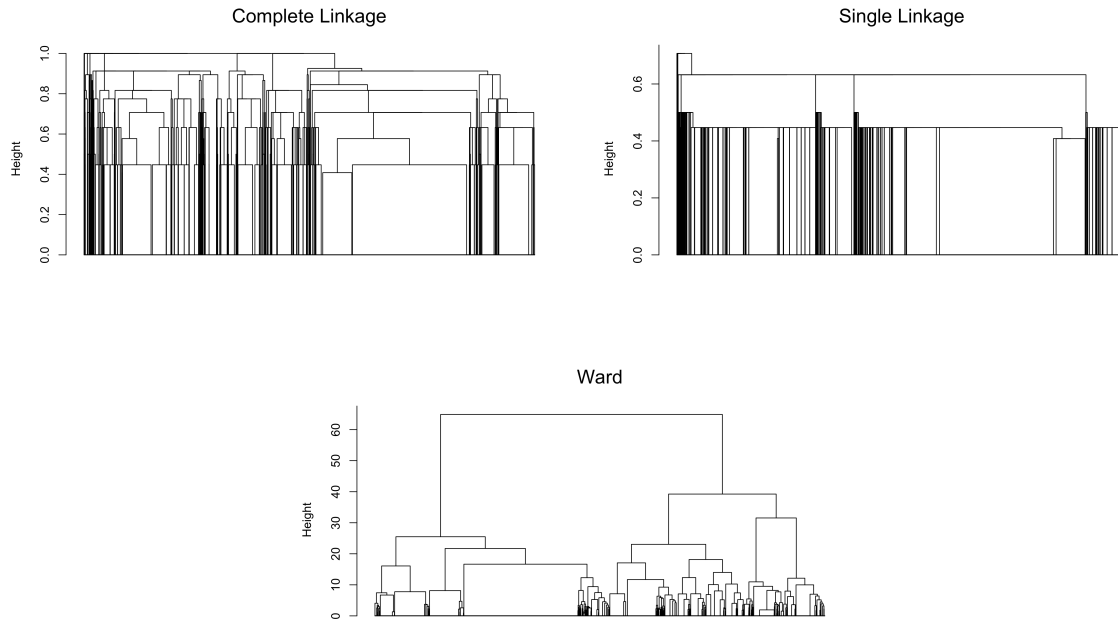


Figure 3: Dendrograms for the first period (t_1) of the single and complete linkage method as well as the ward's algorithm

The complete linkage algorithm tends to form many small groups as it can be seen in the dendrogram in Figure 3. Moreover, the grouped observations are very close to each other; a reasonable cut-off point cannot be observed. The dendrogram visualizes the results of the single linkage method. It shows a typical characteristic of this algorithm: a chain structure. An efficient cut-off can be done at least five clusters (many observations are on the left hand-side of the dendrogram but I cannot tell how many). Independently of the number of clusters, the observation almost belong always to the first cluster and the other cluster contain mostly only one to two observations- Thus, this approach does not contain much information.

The ward's method provides a clear cut-off visualization. The dendrogram is cut off at two clusters. Even though two types of observations seem to be very unlike, it does not contain much added value because the purpose of the analysis is to compare many different tactics and how they spread. By looking at the observations one can tell that there are more than only two types.

5.2.2 Fuzzy C-Means

For the application of the fuzzy c-means clustering algorithm, the *cluster* R package and the function *fanny* is needed. Originally, the fuzzy c-means, such as the c-means, are built for metric and not for categorical data. But since the *fanny* function also takes in a dissimilarity matrix computed by *dist*, the matrix that was calculated for the hierarchical clustering can be used. But this leads to the same big data problem as described in 5.2.1.

Before conducting the fuzzy clustering, a number of clusters have to be chosen. Since the most definite result regarding the number of clusters was shown by the dendrogram of the ward's algorithm, the analysis is conducted with two clusters. The results shows that all membership coefficients are all very close to 0.5 which implies a very fuzzy clustering. This means that the probabilities for an observation to belong to each cluster is equally distributed.

5.2.3 Other Possible Approaches

As presented in Huang and Ng (1999), a fuzzy clustering algorithm was developed that is adapted to categorical data called fuzzy c-modes. It uses the simple matching coefficient. Since a Jaccard coefficient matrix is needed in order to weight the variables properly for this analysis, this algorithm cannot be used. One could adapt the source code with the needed properties but this would overdue the frame for this analysis.

As described in 5.2.1, R has a very bad memory management because it loads everything up in memory and processes it. Thus, the data processing is limited and scaling is not efficient. A solution for this problem could be the usage of a different statistical software like e.g. Python which has full support for multithreading and does not dependent on memory (Kasson, 2018).

5.3 Hotbeds and Tactics

In order to apply the first hotbed criterion, a table with all frequencies of each country per time period (t1-t5, see Table 5) is generated and all frequencies that are equal to

0 are removed. In the second step, only countries are kept that are above the 95% quantile regarding the distribution of the frequencies.

For the second hotbed criterion, the time frames are stretched because the growth rate of the first year is always 0. In order to prevent the first year of the original time period to be 0, an extra year is added to the beginning of every time frame. Next, a table is aggregated by the year and country and their combined frequency is computed. The growth rate of the number of incidents of each country from year i to $i+1$ is calculated and added to the time frame. Afterwards, the additional year that has been added, is removed. Only the countries with a growth rate over the 75% quantile regarding the distribution of all growth rate within one period, are kept. Afterwards, all countries are removed that appear less than twice and with escalation rates further apart than 2 years.

The countries that hold the first and second criteria within one period are labeled as hotbeds in each time frame. Also, the term "hotbed" is added to the corresponding regions. The results for the hotbed analysis are as follows:

1980-1987	1988-1997	1998-2007	2008-2011	2012-2017
Chile	Algeria	Afghanistan	Afghanistan	Afghanistan
Peru	Pakistan	Pakistan	Pakistan	Iraq
El Salvador	Turkey	Iraq	India	Philippines
Guatemala		Israel/Palestine	Iraq	Somalia
		Thailand	Philippines	Nigeria

Table 10: Results of hotbed analysis for t1-t5

To combine the single incident elements to a tactic, the variables *vicinity*, *multiple*, *suicide*, *attacktype1_txt*, *weaptype1_txt*, *weapsubtype1_txt* and *targgroup* are merged to one string variable. Only tactics are kept that occur >300 times between 1977-2017. This threshold is debatable but it was chosen in order to exclude rare forms of tactics that do not occur very often and do not have a big impact. The result is a frequency table of all used tactics that are above the threshold within each period.

5.4 Contingency Tables and Heatmaps

For the preparation of the heatmaps, the time frames, presented in Table 5, are merged to t1t2, t2t3, t3t4t5, t4t5 in order to observe a wider time course with help of the heatmaps. A contingency table for each time frame is created via the function *xtabs*. The input variables are *tactics*, *iyear* and *region_txt*. Only the observations of the tactics are kept that occur at least once a year per region. This produces some white spots on the heatmap and makes it is easier to identify the actual presence of tactics. The next step is the computation of the relative frequencies of every tactics per region and per year.

For the heatmaps only the unique value of each tactic is aggregated that occurs more than 50 times in the combined time frames. A new list is created with all observations that have been executed with one of the aggregated tactics.

Finally, the heatmap is created with the help with the package *ggplot2*. The variable regarding the year is presented by the x-axis and the regions are shown on the y-axis. The colors represent the relative frequencies of each tactic per region and per year. Moreover, two histograms are generated. The upper shows the absolute number of incidents of the specific tactic during the time period and the bottom histogram shows the development of all incidents within the same time period.

6 Qualitative Analysis of Empirical Results

From all heatmaps that are generated, two are selected and the underlying plausibility of the observed pattern is discussed. A selection of other heatmaps, that shows striking patterns, is attached in the appendix. The complete selection of produced heatmaps can be found enclosed to the underlying thesis with the corresponding R Code. Before drawing a conclusion about the plausibility and causality, the following questions are discussed:

- What is the tactic?
- What is the political situation?
- What is the target and why? What is the underlying motivation?
- What weapons are used? Who supported whom?
- What are similarities to the imitating countries?

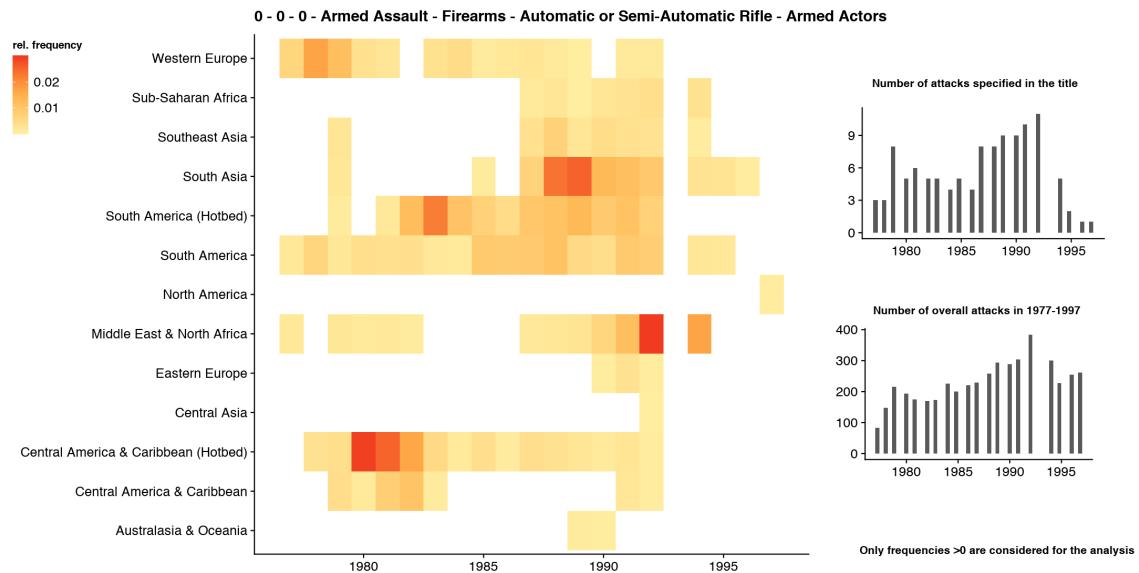


Figure 4: Armed assault with automatic or semi-automatic firearms against armed targets

The first heatmap in Figure 4 shows the relative frequencies of armed assaults with automatic or semi-automatic rifles against armed actors in 1977-1997. Armed actors can be governmental armed actors, such as the military, but also other terrorist groups. It can be observed that more than 2% of all incidents in 1980-1981 were conducted in Central American hotbeds (El Salvador and Guatemala) by using the described tactic.

In the 1970s and 1980s, the political situation in El Salvador and Guatemala was very unstable due to an underdeveloped justice system, corruptive policemen and politicians. The civil wars in both countries – Guatemala (1960-1996) and El Salvador (1980-1992) – were closely entangled with the Cold War. The largest amount of firearms has been imported by Cold War allies (Laurance and Godnick, 2001). The USA supported the government and the USSR supported the rebels (Borda, 2009). The targets are armed actors. It can be assumed that the terrorists targeted rather the military than other terrorist groups because of the contentious political situation and because only violent incidents of non-state actors are included in the GTD. Therefore, the violent behaviour of the state against terrorist groups is not included in the data.

Besides the political tension, El Salvador and Guatemala had major problems with drug trafficking. Those routes were also used for smuggling weapons to South America (Laurance and Godnick, 2001). Starting in 1982, an increase of the tactics can be observed in the hotbeds of South America (Chile and Peru) and shortly after, throughout South America. The similar political as well as the drug trafficking situation (Borda, 2009) explains the same targeting and the same attack type.

As observed in the heatmap in Figure 4, from 1987-88, the same terrorist tactic is also more frequently used in South Asia. The Soviet-Afghan War was waged by the same opponents and weapon suppliers, being the United States and the USSR, than in Central America – the Soviets supported the government and the USA supported the rebels which were also supported by foreign Muslim fighters and by Saudi Arabia and Pakistan financially. After the war was "won" by the Afghans in 1989, the different terrorist groups did not manage to establish one government and fought against each other with the weapons that were left behind by the Russians and Americans (Harvey, 2003).

In this case, at first, the armed actors are the Afghan military and later on, other terrorist groups.

From 1991-1992, an increase of the relative frequency can be observed in the Middle East and North Africa. In 1990, the Gulf War started. Iraq, who was an ally of the Soviets, invaded in Kuwait and the United States got involved (Federation of American Scientists, 2019). It can be assumed that the weapons were imported by the opponents. Automatic or semi-automatic firearms for an armed assault against armed actors were used for almost 3% of all incidents that occurred in this region in 1992.

The target types differ between continents. South and Central America attacked governmental armed actors. The objectives were economical (due to drug trafficking), political (due to a corruptive system) or self-defending (caused by state terrorism). In contrast, the terrorist groups in South Asia and the Middle East and North Africa followed rather religious (such as proclaiming a Jihad), closely linked to political, motivations, cp. Richardson (2007).

The second heatmap in Figure 5 is similar to the first heatmap in Figure 4.

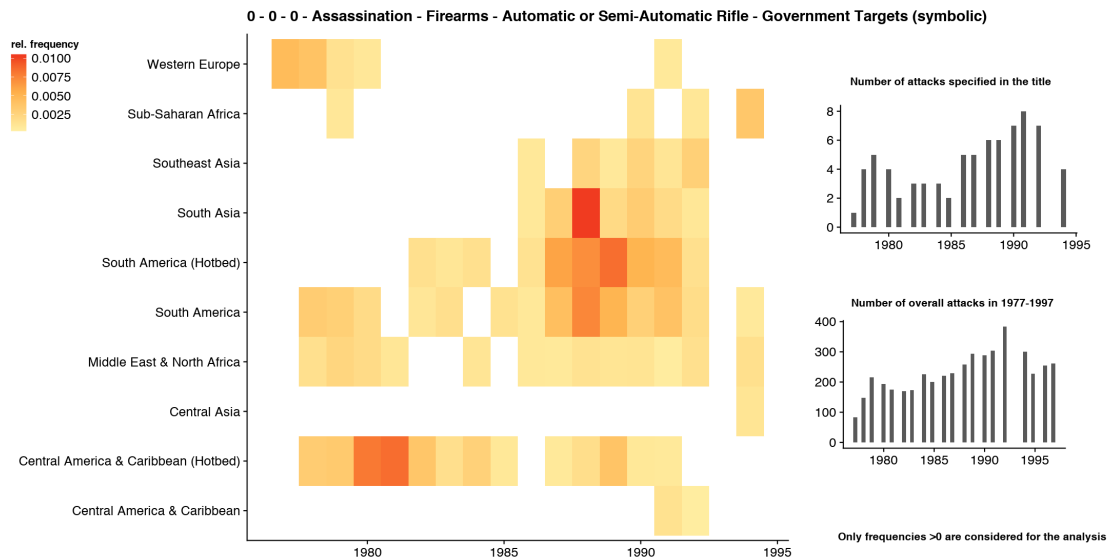


Figure 5: Assassination with automatic or semi-automatic firearms against governmental targets

The form of execution of the attack types is alike. A firearm is a controllable weapon. However, when conducting an armed assault, any death of the targeted group is a success. In contrast, an assassination has the objective to kill one or more specific individuals (GTD Codebook, 2017). The target type for this tactic are one or more individuals that are associated with the government.

In heatmap Figure 4 and Figure 5, a parallel development and "spread" of incidents, using similar tactics, can be observed, with the exception of the Middle East and North Africa. This difference can be explained by the fact that Gulf War was about resources and not against a regime (Frank, 2018). Therefore, it can be assumed that no specific governmental associates were targeted. The histograms of the heatmap in Figure 5 show that the total number of incidents increases between 1992-1993, but the total number of incidents for this specific tactic decreases between those years. At the same time, the number of incidents for Figure 4 increases approximately by the amount that Figure 5 decreases. This change could indicate a shift between the tactics but could also be random. Due to the few incidents, the statistical significance cannot be analyzed.

The causality of the spread of the assassination/armed assault against the government/armed actors with semi-automatic or automatic firearms is rather due to the American and Russian involvement and diffusion of weapons. The objective of terrorist groups in Central and South America is a different one than for the South Asia organizations. Latin America tried to get rid of the suppressing and corruptive structures and deals with major drug trafficking issues. In contrast, in South Asia different terrorist groups fight against each other with the religious fundamental of exclaiming a Jihad. Therefore, the observed spread of tactics is not an imitation of the hotbeds but rather influenced by the weapon suppliers.

7 Discussion

The discussion is split into three parts: data - theory - method. In the first part, the database and its assemblage are examined, secondly my decisions regarding the analysis are debated and lastly the visualization and its potential misleading structure is discussed.

One major factor that influences the validity as well as the reliability of the GTD is the collection process. Besides different collectors, the data have been collected in real-time but also retrospectively. Also, some data have been lost complete, like the data from 1993. Additionally, different methods have been used. The codebook itself states: "Users should note that differences in levels of attacks and casualties before and after 1997, 2008, and 2012 may be at least partially explained by differences in data collection; researchers should adjust for these differences when modeling the data." (GTD Codebook, 2017). As a countermeasure I introduced the different time frames and relative frequencies, but the analysis is still influenced by inconsistencies in the collection methodology within the database.

Other possible bias arise from the sourcing. Sources are only added if the source is found to be credible. International media sources from developing countries or countries with politically repressive systems might not be credible in this matter and therefore terror attacks might not be listed. The codebook states: "Note that particular scarcity of high- quality sources in certain geographic areas results in conservative documentation of attacks in those areas in the GTD." (GTD Codebook, 2017).

As observed and described in 6, it sometimes is difficult to separate terrorist incidents from non-terrorist incidents. According to the database codebook all of the non-terrorist incidents are excluded but the incidents in Latin America were often rather acts of guerrillas than terrorists. Moreover, state terrorism is not included in the GTD. But state terrorism can explain the behaviour or attacks of violent actors, it can also be argued that a violent act against a suppressing, illegitimate regime is not an act of terrorism but rather self-defense. Looking back in history, terrorist were often retrospectively declared to be freedom fighters (Richardson, 2007). Therefore,

the GTD misses important information such as a variable about the political situation indicating e.g. state terrorism.

In the data preparation phase but also during the analysis, I had to decide on several thresholds. The chosen thresholds are based on interviews with an expert but also on common sense. Thus, further research could be to vary the thresholds and examine the impact on the results.

Regarding the heatmaps, the downside of this visualization technique is that it is easily misleading. The heatmap is heavily influenced by outliers. For example if one relative frequency is really high compared to all others, this observation will be presented in red, all the other observations in contrast will be presented by a light yellow. Even if they also varying, it will hardly be observed. This leads to a major loss of information. As a countermeasure, I print the histograms of the absolute frequency of the tactics within the time frame next to the heatmap. Nevertheless, the results of the heatmap has to be evaluated carefully.

8 Conclusions

A partitioning as well as a hierarchical clustering algorithm has been implemented with the objective to group similar observations. The clustering approaches were not successful which is mainly due to the large data set combined with the insufficient memory management of R and the scarce amount of research, which has been done for cluster analysis on categorical data. As one approach to contribute something to the categorical clustering research, a weighted Jaccard coefficient has been introduced. As an attribution to the terrorism literature, the terms tactic and hotbed have been defined. The results have been visualized with heatmaps.

Even though the generated heatmaps have to be considered with caution, many indicate a relationship between hotbeds and the global diffusion of terrorist tactics, originating from hotbeds. The observed geographical spread of terrorism can be explained with the social movement theory. All named countries had contentious politics in common and are resourcing weapons from their allies. The suppressing governments and the strong influence by foreign powers, such as Russia and the United States, could have been used for framing the actions of the terrorist groups.

According to Midlarsky et al. (1980), tactics might be contagious not only because they are easy to conduct but also because they attract simultaneously a lot of attention. Even international attention can be achieved when taking for example hostages from countries abroad or causing deaths of foreigners due to bombings. After analyzing the data, I assume that tactics spread because similar resources are available due to similar circumstances, such as having the same ally.

The presented results only show attacks with expensive, professionally produced weapons. I believe that the causality of the spread is not originating from the hotbeds but rather from other underlying factors, such as foreign powers financing wars by delivering weapons to political contentious countries. The analysis of Figure 4 and 5 in chapter 6 supports this theory.

To put it in a nutshell, the causality of the influence of terrorist hotbeds on the global transformation of terrorist tactics could not be proven.

As an further research approach, a chi-squared test could be conducted in order to examine if the variables are independent from each other. My hypothesis is that certain weapon types are dependent on the attack type. As an example, an assassination is more likely conducted with firearms than with a bomb because the target is a specific individual and it is more likely to be successful with a weapon that can fully be controlled. All produced heatmaps could be analyzed in detail and compared. The tactics could also be aggregated due to the results of the chi-squared test. Furthermore, the parameters, such as the thresholds, the time frames, the variables forming a tactic or the geographical classification, could be varied.

References

- AGRESTI, A. AND M. KATERI (2011): *Categorical data analysis*, Springer.
- BACKHAUS, K., B. ERICHSON, W. PLINKE, AND R. WEIBER (2016): *Multivariate Analysemethoden*, Springer.
- BECK, C. J. (2008): “The contribution of social movement theory to understanding terrorism,” *Sociology Compass*, 2, 1565–1581.
- BERMAN, E. AND D. LAITIN (2005): “Hard targets: Theory and evidence on suicide attacks,” Tech. rep., National Bureau of Economic Research.
- BOJKO, A. A. (2009): “Informative or misleading? Heatmaps deconstructed,” 30–39.
- BORDA, S. P. (2009): “The internationalization of domestic conflicts: a comparative study of Colombia, El Salvador and Guatemala.” .
- BRAITHWAITE, A. AND Q. LI (2007): “Transnational Terrorism Hot Spots: Identification and Impact Evaluation,” *Conflict Management and Peace Science*, 24, 281–296.
- CAMBRIDGE DICTIONARY (2008): “Cambridge advanced learners dictionary,” *PONS-Wörterbücher, Klett Ernst Verlag GmbH*.
- CHATFIELD, C. (2018): *Introduction to multivariate analysis*, Routledge.
- DOUG, M., J. D. MCCARTHY, AND M. N. ZALD (1996): “Introduction: Opportunities, mobilizing structures, and framing processes Toward a synthetic, comparative perspective on social movements,” *Comparative perspectives on social movements: Political opportunities, mobilizing structures, and cultural framings*, ed. D. McAdam, J. McCarthy, and M. Zald, 1–20.
- ENDERS, W. AND T. SANDLER (2000): “Is transnational terrorism becoming more threatening? A time-series investigation,” *Journal of Conflict Resolution*, 44, 307–332.

- FEDERATION OF AMERICAN SCIENTISTS (2019): “Patterns of Global Terrorism: 1990 - Middle East Overview,” https://fas.org/irp/threat/terror_90/mideast.html, accessed: 2019-02-17.
- FRANK, A. G. (2018): “A third-world war: A political economy of the Persian Gulf War and the new world order,” in *Triumph Of The Image*, Routledge, 3–21.
- GILL, P., J. HORGAN, S. T. HUNTER, AND L. D. CUSHENBERY (2013): “Malevolent creativity in terrorist organizations,” *The Journal of Creative Behavior*, 47, 125–151.
- GOWER, J. C. (1966): “Some distance properties of latent root and vector methods used in multivariate analysis,” *Biometrika*, 53, 325–338.
- GTD CODEBOOK (2017): “Global Terrorism Database. Codebook: Inclusion Criteria and Variables,” <https://www.start.umd.edu/gtd/downloads/Codebook.pdf>, accessed: 2019-02-14.
- HÄRDLE, W. AND L. SIMAR (2012): *Applied multivariate statistical analysis*.
- HARVEY, K. (2003): “Afghanistan, The United States, and the Legacy of Afghanistans Civil War,” .
- HENNIG, C. (2015): “Clustering strategy and method selection,” *arXiv preprint arXiv:1503.02059*.
- HUANG, J. Z. AND M. K. NG (1999): “A fuzzy k-modes algorithm for clustering categorical data,” *IEEE Trans. Fuzzy Systems*, 7, 446–452.
- JELLINEK, G. (1990): “Allgemeine Staatslehre (= Recht des modernen Staates, Bd. 1),” *Berlin (2. Auflage 1905 (Digitalisat)*, 3.
- KASSON, E. G. (2018): “Is Python Edging Out R in the Data Science Space?” *retrieved from <https://insights.dice.com/2018/02/16/python-edging-r-data-science-space/>*.
- KELLER, J. M., M. R. GRAY, AND J. A. GIVENS (1985): “A fuzzy k-nearest neighbor algorithm,” *IEEE transactions on systems, man, and cybernetics*, 580–585.

- LAURANCE, E. AND W. GODNICK (2001): “Weapons Collection in Central America: El Salvador and Guatemala,” *Managing the Remnants of War: Weapons Collection and Disposal as an Element of Peace-Building*.
- LI, M. J., M. K. NG, Y.-M. CHEUNG, AND J. Z. HUANG (2008): “Agglomerative fuzzy k-means clustering algorithm with selection of number of clusters,” *IEEE transactions on knowledge and data engineering*, 20, 1519–1534.
- MCADAM, D. (1983): “Tactical innovation and the pace of insurgency,” *American Sociological Review*, 735–754.
- MIDLARSKY, M. I., M. CRENSHAW, AND F. YOSHIDA (1980): “Why violence spreads: The contagion of international terrorism,” *International Studies Quarterly*, 24, 262–298.
- NATIONAL CONSORTIUM FOR THE STUDY OF TERRORISM AND RESPONSES TO TERRORISM (START) (2018): “Global Terrorism Database,” *retrieved from* <https://www.start.umd.edu/gtd>.
- NEWMAN, D. (2012): “Borders and conflict resolution,” *A companion to border studies*, 249–265.
- PALESTINIAN CENTRAL BUREAU OF STATISTICS (2018): “Number of Settlers in the Settlements in the West Bank by Governorate and Type of Settlement, 2017,” *retrieved from* [http : //www.pcbs.gov.ps/Portals/Rainbow/Documents/SETT8E – 2017.html](http://www.pcbs.gov.ps/Portals/Rainbow/Documents/SETT8E-2017.html).
- RICHARDSON, L. (2007): *Was Terroristen wollen: Die Ursachen der Gewalt und wie wir sie bekämpfen können*, Campus Verlag.
- SHLAIM, A. (1994): “The Oslo Accord,” *Journal of Palestine Studies*, 23, 24–40.
- SMITH, M. AND S. M. ZEIGLER (2017): “Terrorism before and after 9/11—a more dangerous world?” *Research & Politics*, 4, 2053168017739757.

STERN, J. AND M. MCBRIDE (2013): “Terrorism after the 2003 invasion of Iraq,”
Group Brown University, Eisenhower Study Google Scholar.

UNITED NATIONS A/RES/67/19 (2012): “General Assembly resolution 67/97. Status of Palestine in the United Nations,” *retrieved from*
<http://undocs.org/A/RES/67/19>.

Appendix

Division of Regions (GTD Codebook, 2017)

North America

Canada, Mexico, United States

Central America & Caribbean

Antigua and Barbuda, Bahamas, Barbados, Belize, Cayman Islands, Costa Rica, Cuba, Dominica, Dominican Republic, El Salvador, Grenada, Guadeloupe, Guatemala, Haiti, Honduras, Jamaica, Martinique, Nicaragua, Panama, St. Kitts and Nevis, St. Lucia, Trinidad and Tobago

South America

Argentina, Bolivia, Brazil, Chile, Colombia, Ecuador, Falkland Islands, French Guiana, Guyana, Paraguay, Peru, Suriname, Uruguay, Venezuela

East Asia

China, Hong Kong, Japan, Macau, North Korea, South Korea, Taiwan

Southeast Asia

Brunei, Cambodia, East Timor, Indonesia, Laos, Malaysia, Myanmar, Philippines, Singapore, South Vietnam, Thailand, Vietnam

South Asia

Afghanistan, Bangladesh, Bhutan, India, Maldives, Mauritius, Nepal, Pakistan, Sri Lanka

Central Asia

Armenia, Azerbaijan, Georgia, Kazakhstan, Kyrgyzstan, Tajikistan, Turkmenistan,

Uzbekistan

Western Europe

Andorra, Austria, Belgium, Cyprus, Denmark, Finland, France, Germany, Gibraltar, Greece, Iceland, Ireland, Italy, Luxembourg, Malta, Netherlands, Norway, Portugal, Spain, Sweden, Switzerland, United Kingdom, Vatican City, West Germany (FRG)

Eastern Europe

Albania, Belarus, Bosnia-Herzegovina, Bulgaria, Croatia, Czech Republic, Czechoslovakia, East Germany (GDR), Estonia, Hungary, Kosovo, Latvia, Lithuania, Macedonia, Moldova, Montenegro, Poland, Romania, Russia, Serbia, Serbia-Montenegro, Slovak Republic, Slovenia, Soviet Union, Ukraine, Yugoslavia

Middle East & North Africa

Algeria, Bahrain, Egypt, Iran, Iraq, Israel, Jordan, Kuwait, Lebanon, Libya, Morocco, North Yemen, Qatar, Saudi Arabia, South Yemen, Syria, Tunisia, Turkey, United Arab Emirates, West Bank and Gaza Strip, Western Sahara, Yemen

Sub-Saharan Africa

Angola, Benin, Botswana, Burkina Faso, Burundi, Cameroon, Central African Republic, Chad, Comoros, Democratic Republic of the Congo, Djibouti, Equatorial Guinea, Eritrea, Ethiopia, Gabon, Gambia, Ghana, Guinea, Guinea-Bissau, Ivory Coast, Kenya, Lesotho, Liberia, Madagascar, Malawi, Mali, Mauritania, Mozambique, Namibia, Niger, Nigeria, People's Republic of the Congo, Republic of the Congo, Rhodesia, Rwanda, Senegal, Seychelles, Sierra Leone, Somalia, South Africa, South Sudan, Sudan, Swaziland, Tanzania, Togo, Uganda, Zaire, Zambia, Zimbabwe

Australasia & Oceania Australia

Fiji, French Polynesia, New Caledonia, New Hebrides, New Zealand, Papua New Guinea, Solomon Islands, Vanuatu, Wallis and Futuna

Exemplary Heatmaps

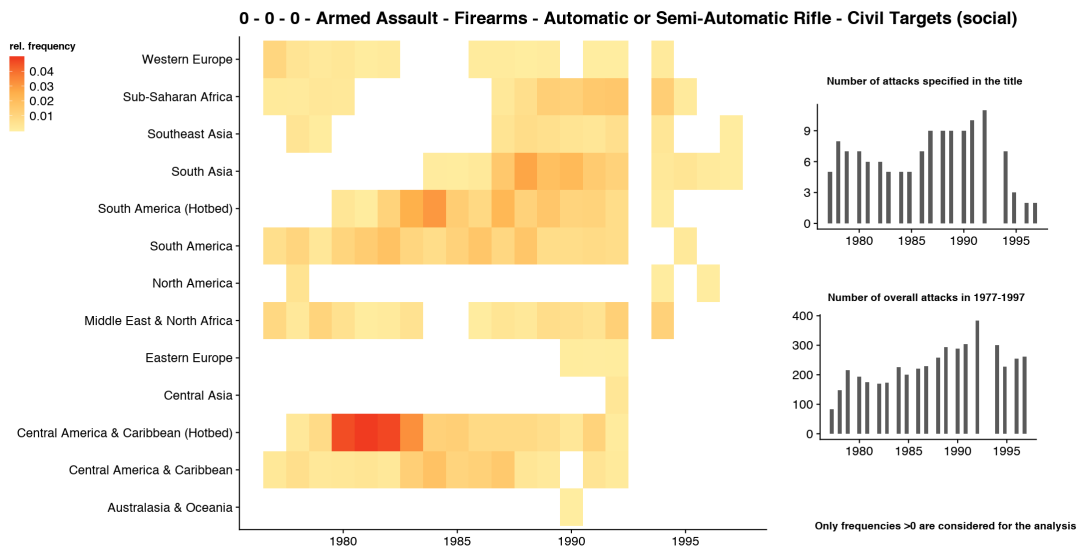


Figure 6: Armed assault with automatic or semi-automatic firearms against civil targets

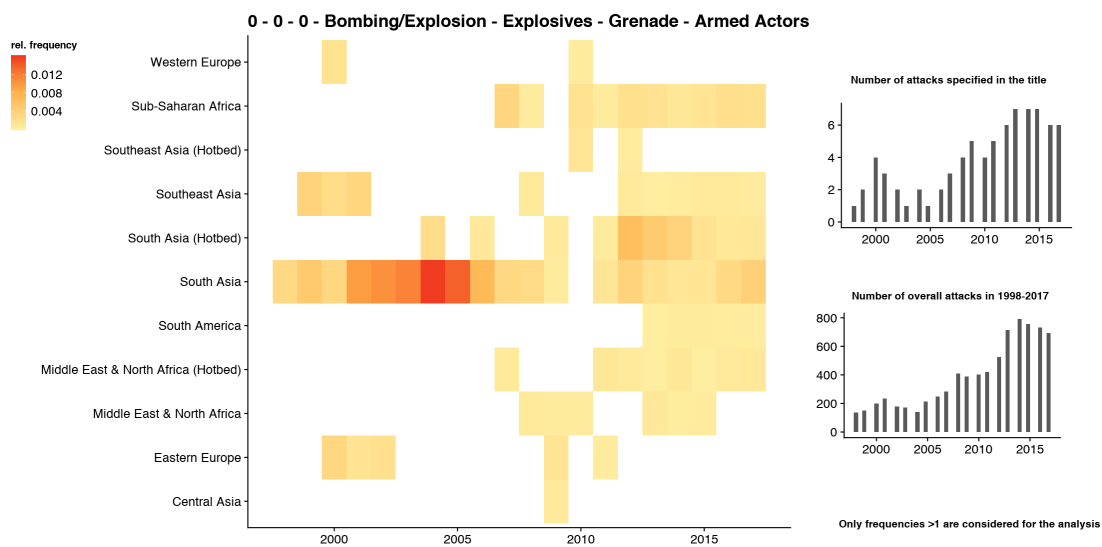


Figure 7: Bombing/explosives with a grenade against an armed actor

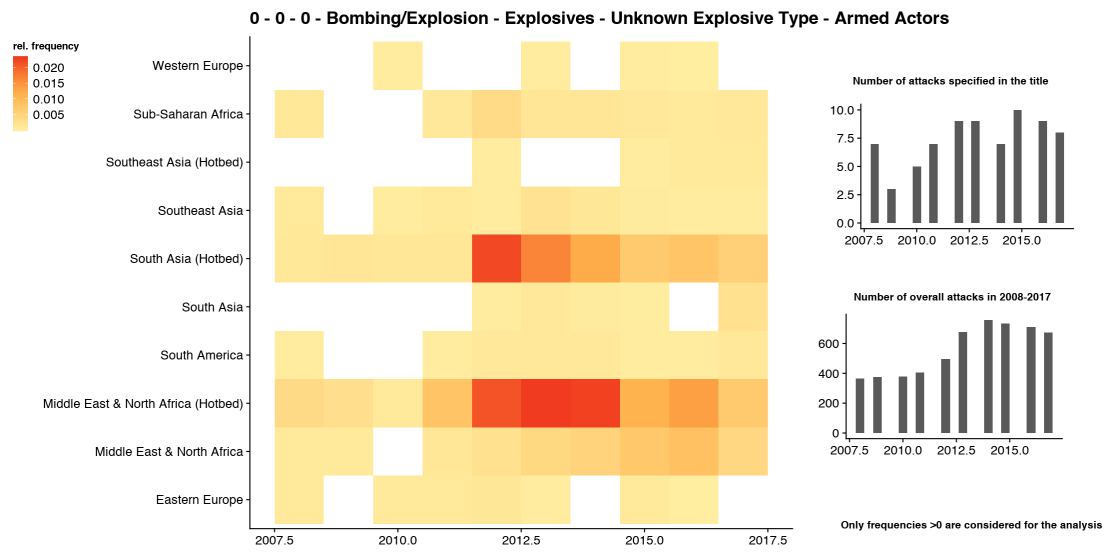


Figure 8: Bombing/explosives with an unknown explosive type against an armed actor

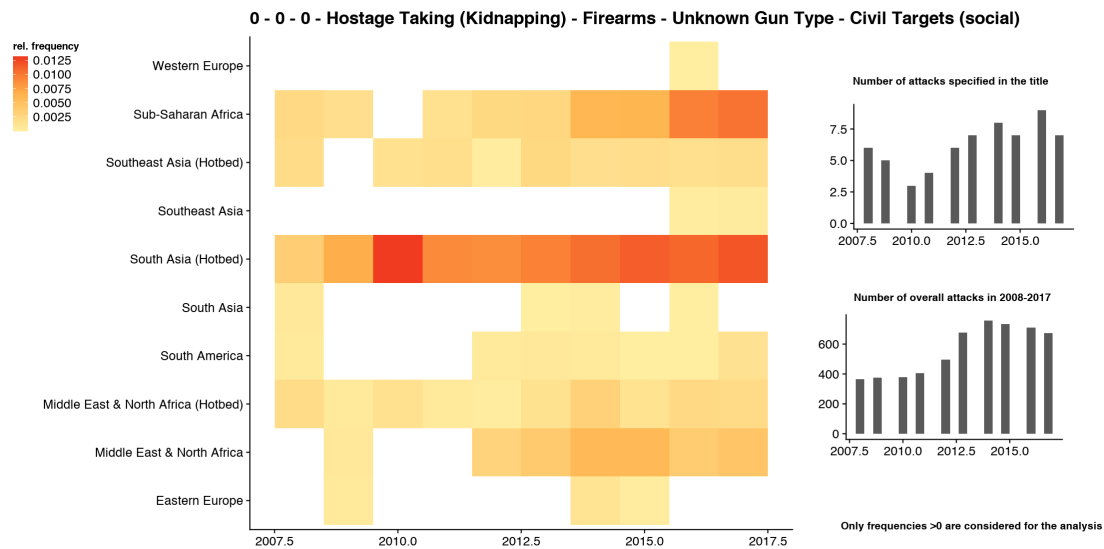


Figure 9: Hostage taking (kidnapping) with an unknown gun type of civil targets

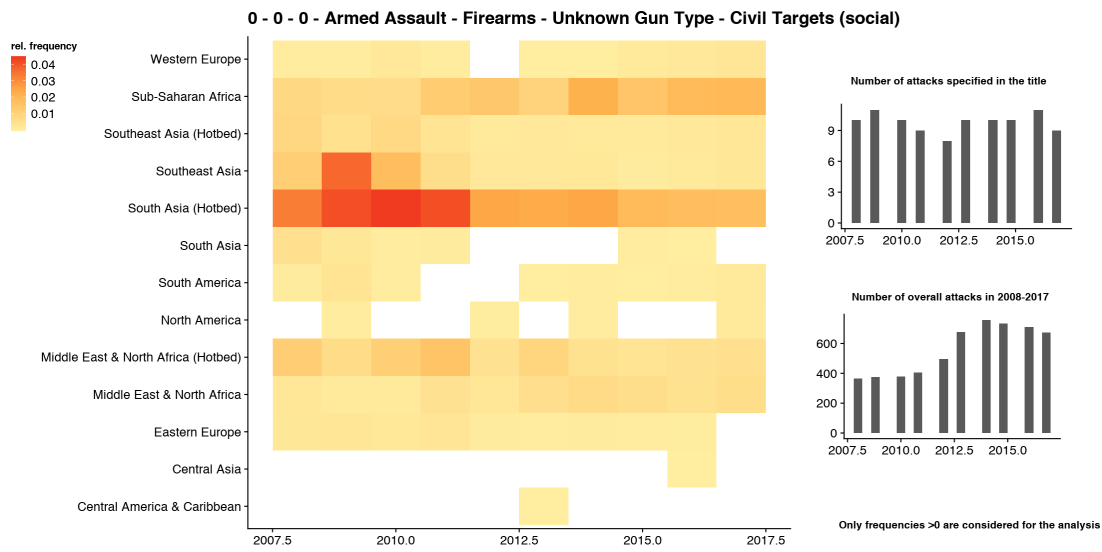


Figure 10: Armed assault with an unknown gun type against civil targets

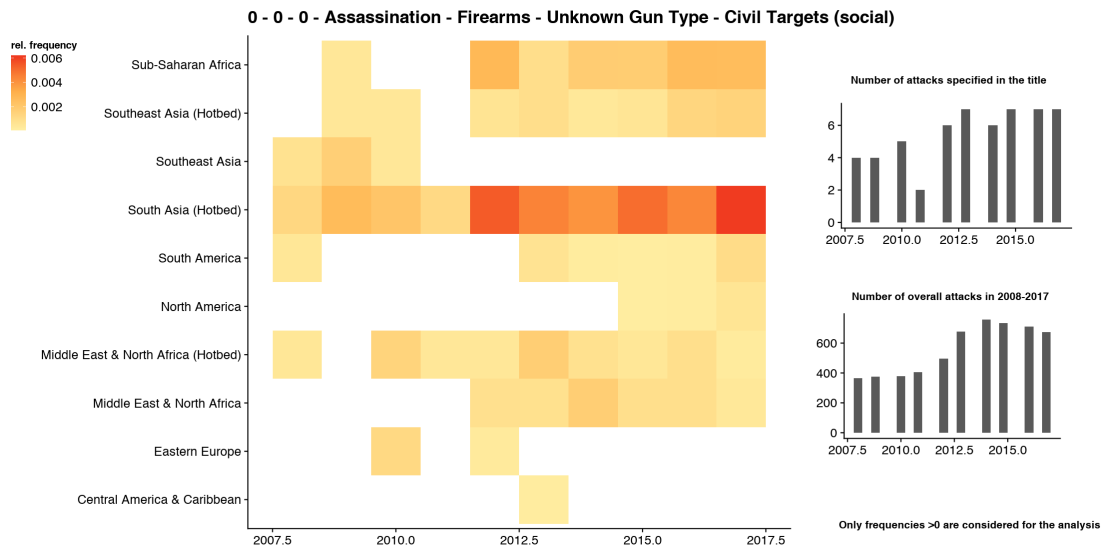


Figure 11: Assassination with an unknown gun type against civil targets

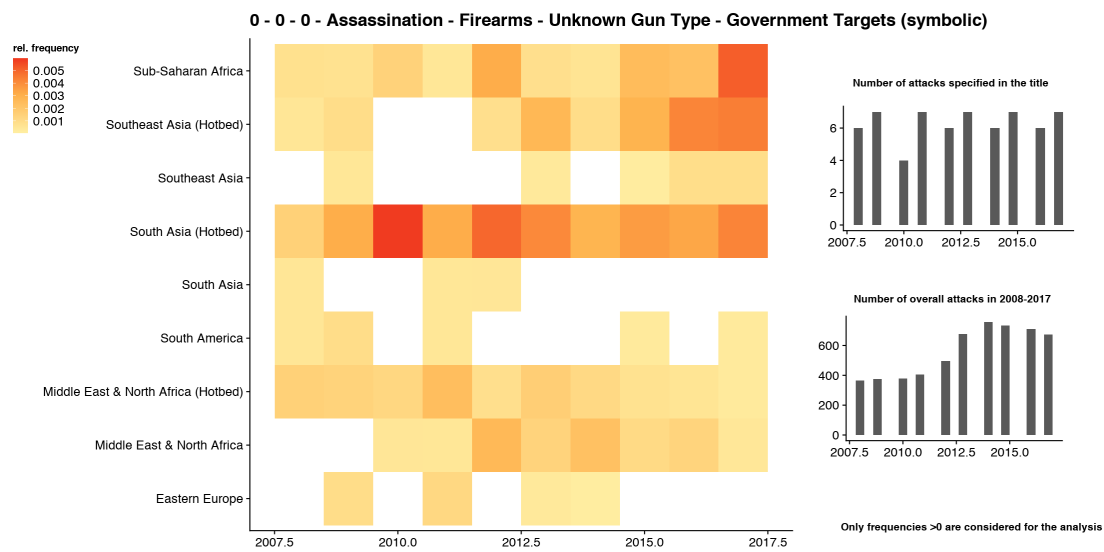


Figure 12: Assassination with an unknown gun type against government targets

Declaration of Authorship

I hereby confirm that I have authored this Bachelor's thesis independently and without use of others than the indicated sources. All passages which are literally or in general matter taken out of publications or other sources are marked as such.

Berlin, February 18, 2019

Anna Franziska Bothe