

Predicting the nature of terrorist attacks

EE-558: Network Tour of Data Science

Team 34

Coen Charles-Théophile

Morel Valentin

Schumacher Cédric

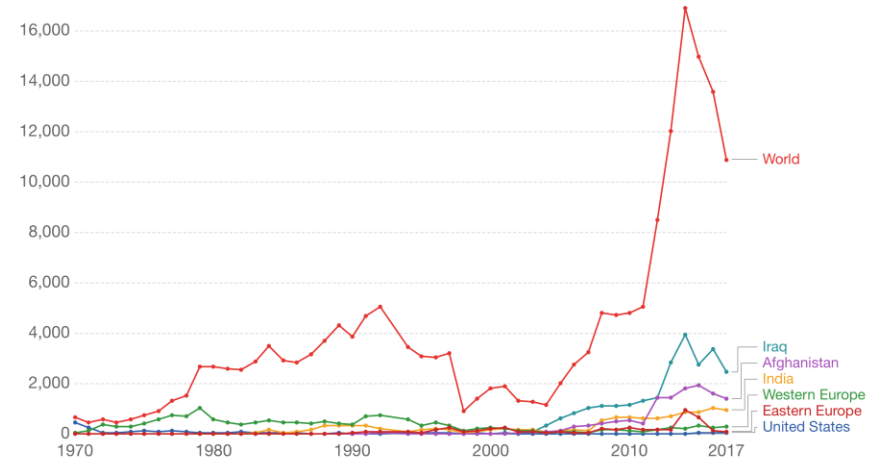
Sieber Xavier

Introduction

- Emergence of new and more radicalized terrorist group
- Necessity of analyzing, classifying and understanding these attacks
- Data science could be the answer
- Machine learning already used to detect Daesh propaganda in videos^[1]
 - Developed by Home Office, UK
 - 94% of target video detected
 - 99.99% of accuracy
 - Out of 1'000'000 videos, only 50 require additional human review

Number of terrorist incidents

The total number of terrorism-related incidents per year. The source defines a terrorist attack 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." The perpetrators of the incidents must be sub-national actors; data does not include acts of state terrorism.



Source: Global Terrorism Database (2018)

OurWorldInData.org/terrorism/ • CC BY-SA

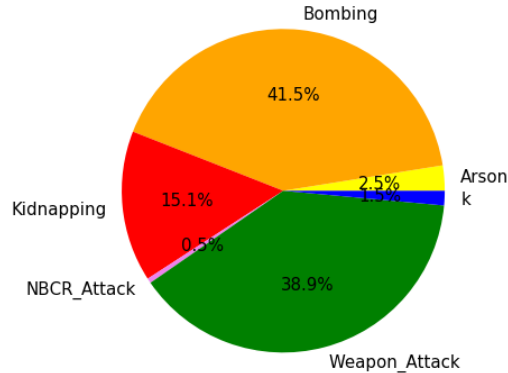
Data

- Provided by UCSC
- 1293 terrorist attacks each assigned to one of 6 labels indicating the type of the attack
 - Attack types: arson, bombing, kidnapping, NBCR attack, weapon attack or other
- Binary feature vector of dimension 106 defining the presence or absence of said feature
- Label of location created with the edge file linking collocated attacks
 - Each terrorist attack is given a specific location number
- **Focus: predicting the labels using the feature vectors**

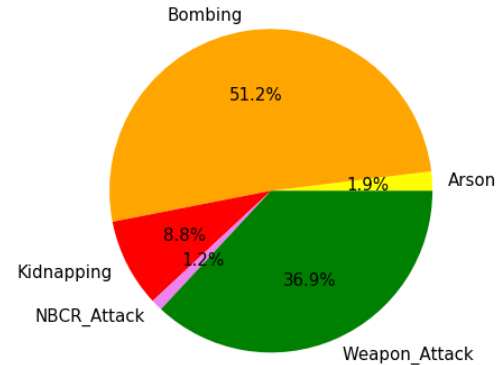
Data pre-analysis

- Pairwise correlation between the 106 features
 - Discriminate features expressing the same behaviour
 - Only 6 features show a correlation greater than 85%
- 3 attack types dominate the data with 97% of the total attacks
 - Not equally distributed

Distribution of types of attack for locations with one attack



Distribution of types of attack for locations with more than one attack



Collocation prediction

- Is there any links between the features and the location of the attacks
 - Could underline the modus operandi of organizations
- Clustering on the location of attack
 - Only location with more than one attack
- Results: globally poor thus location may not be related to the 106 features

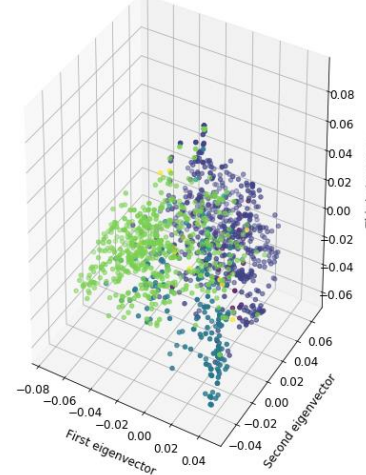
Techniques	AMI score
Birch	0.32
HDBscan	0.0
Kmeans	0.31
Spectral clustering	0.25

Prediction on the nature of the attack

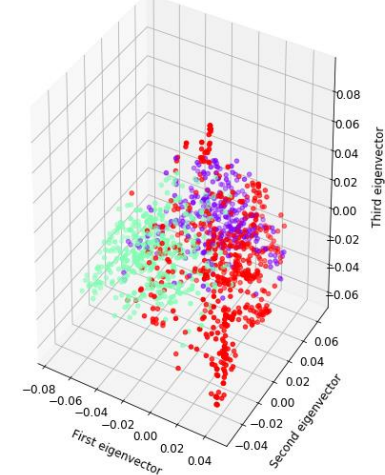
- Another important information:
 - The attack type
- Is there any links between the type of the attack and the 106 features
- Clustering on the nature of attack
 - 4 different type of clustering
 - 2 different measure for feature distance
 - 3 or 6 clusters
- Results: globally poor, the data needs more complex analysis

Cluster number	Euclidean		Cosine	
	3	6	3	6
Birch	0.24	0.26	0.24	0.26
HDBScan	0.14	0.13	0.14	0.13
KMeans	0.24	0.23	0.23	0.23
Spectral Clustering	0.27	0.21	0.25	0.16

Laplacian eigenmap with ground truth.



Laplacian eigenmap with cluster assignment.



Hybrid approach to clustering

- Feature vectors approach failed
- Projection into a lower dimensional space

$$\tilde{y} = yP \quad P \in \mathbb{R}^{106 \times 6}$$

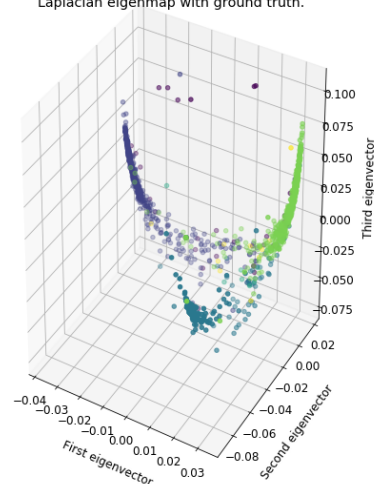
- Find the optimal projection

$$\arg \min_P \|yP - y_t\|_2^2$$

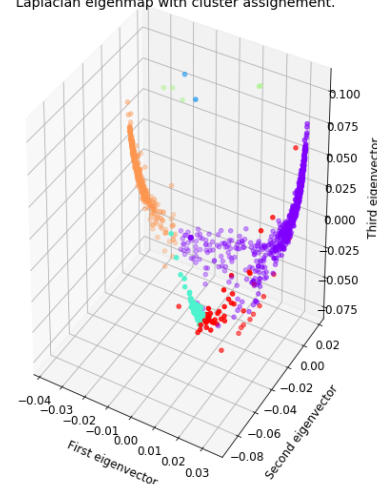
- Spectral clustering was performed on the results

	Euclidean		Cosine	
Cluster number	3	6	3	6
Birch	0	0	0	0
HDBScan	0.20	0.34	0.20	0.34
KMeans	0.53	0.43	0.53	0.43
Spectral Clustering	0.53	0.46	0.55	0.58

Laplacian eigenmap with ground truth.



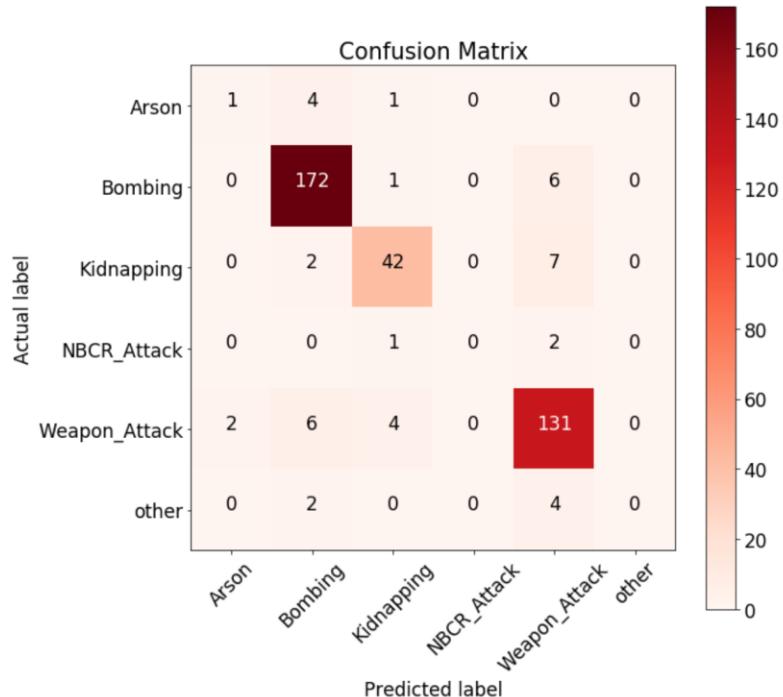
Laplacian eigenmap with cluster assignment.



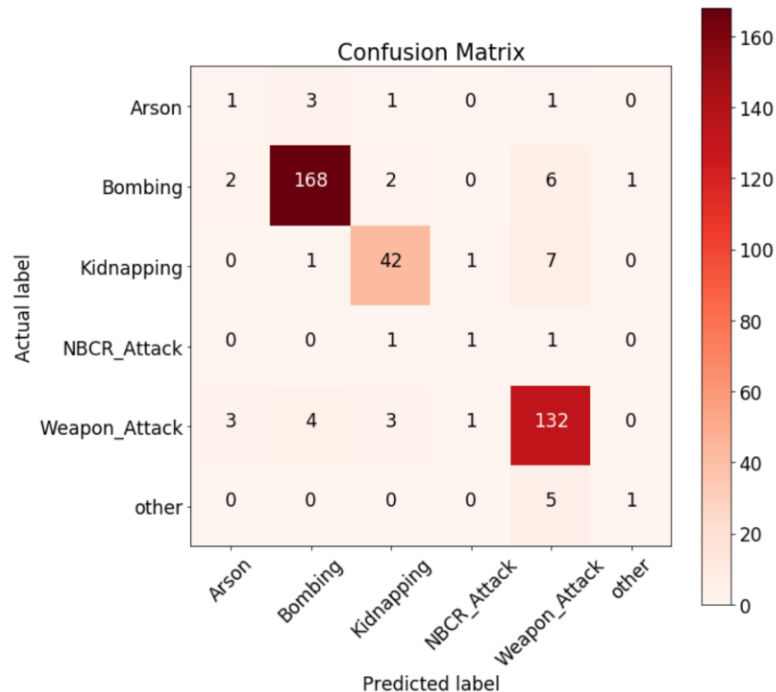
Classification (1)

- Support vector machine classifier and Gradient boosting tree classifier
- 70% of the 1293 attacks to train / 30% to test
- Gridsearch with 5-fold crossvalidation on hyperparameters
- Both yield good results: Attack nature is predictable
- Arson and “other” attack type are more difficult to classify
 - Possible reason: only few data

Classification (2)

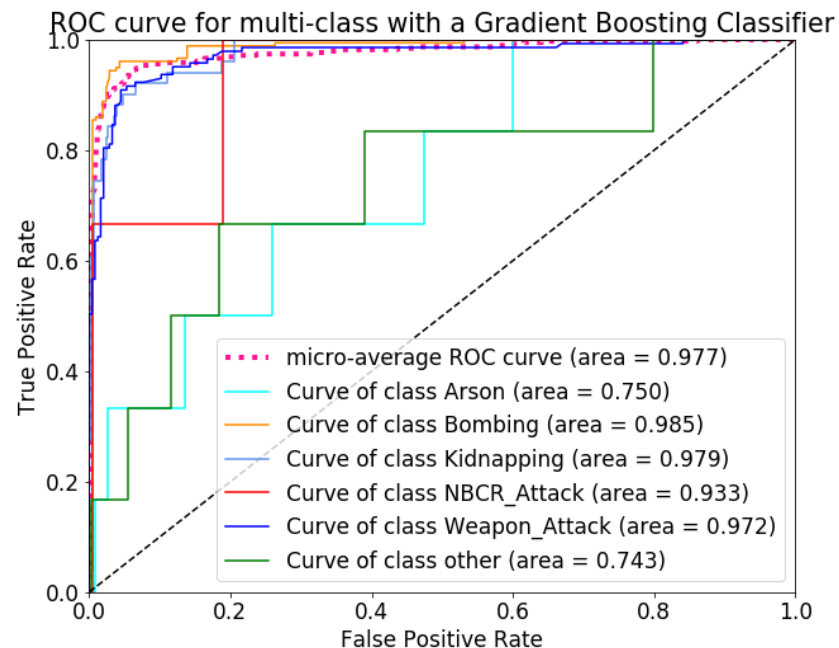
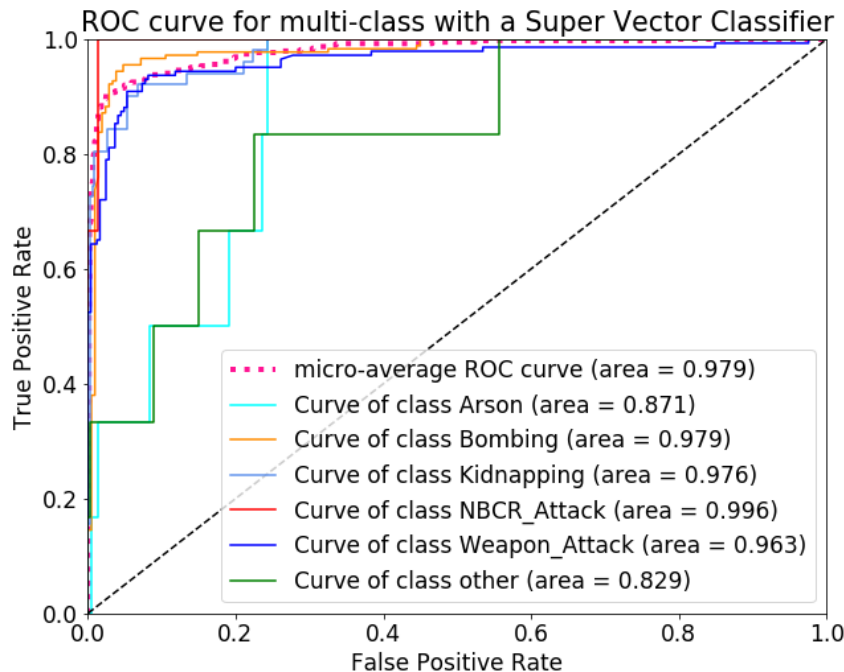


SVM



Gradient boosting classifier

Classification (3)



Classifier	F1-score	weighted AUC
Super Vector Machine Classification	0.875	0.969
Gradient Boosting Tree Classification	0.872	0.973

Conclusion

- Limited performance of pure network approach as feature similarities do not reflect label similarities directly
 - Indicates complex interaction between features and labels
 - Way around can be found by optimized projection
- Classical machine learning approach yields best performance
- Lack of information on features makes interpretation difficult
- Our solutions allows future attacks to be classified by their type using the features without other prior knowledge

Overview

- Enhancing the database
 - Known features
 - Exact location

=> Could yield better tangible implications
- Other labels could be chosen for clustering
 - Terrorist groups
 - Date of the attack
- Understanding the machine learning algorithm could be used to understand the underlying relations between terrorist attacks
 - Understanding the choice of support vectors in SVM could highlight the difference between specific attack types