

# PREDICTING THE SUCCESS OF CYBER-RELATED TERRORIST ATTACKS

- USING THE GLOBAL TERRORISM DATABASE PROVIDED BY START AT THE UNIVERSITY OF MARYLAND:  
[HTTP://WWW.START.UMD.EDU/GTD/](http://www.start.umd.edu/gtd/)

CITATION: NATIONAL CONSORTIUM FOR THE STUDY OF TERRORISM AND RESPONSES TO TERRORISM (START). (2018). GLOBAL TERRORISM DATABASE [DATA FILE]. RETRIEVED FROM [HTTPS://WWW.START.UMD.EDU/GTD](https://www.start.umd.edu/gtd)

# RESEARCH QUESTIONS AND OBJECTIVES

---



# BACKGROUND:

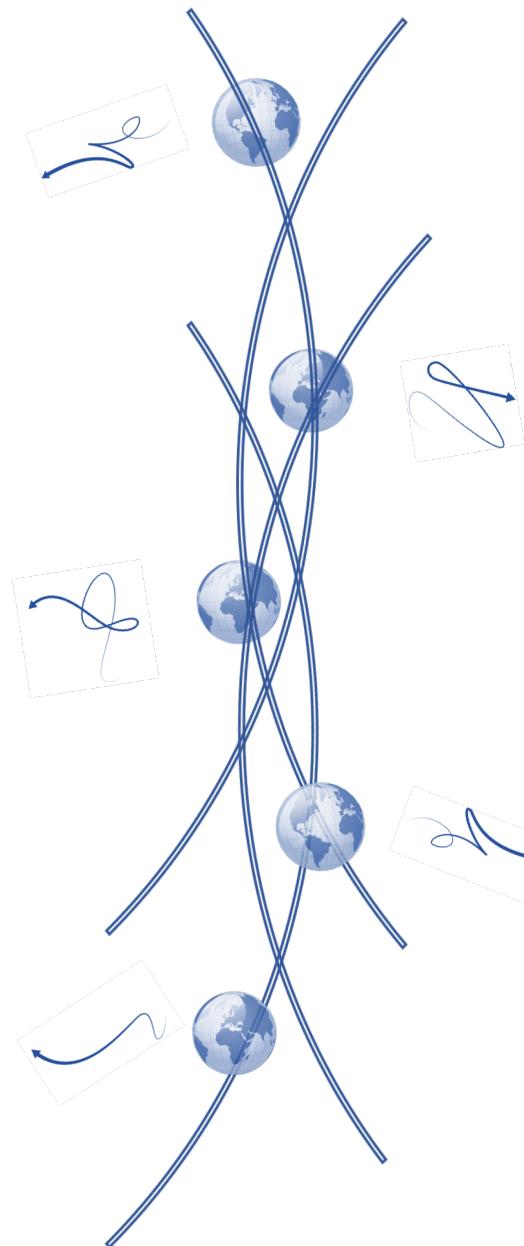
Hosted by the University of Maryland  
and funded by the US Department of  
State ([until May 2018](#)).

## What is a terrorist attack?

Page 10 of the [Codebook](#) lists that a terrorist attack must include the following characteristics: *intention, violence and non-state-related perpetrators* (i.e. those outside of the current government).

## Potential issues given the nature of this data:

Given that there are instances of terrorism which do not necessarily reach an 'enacted' state (according to this database), further researchers/people using this database (including these results) must be careful and recognize that a plethora of other instances (nearly-executed plans that went un-noticed, etc.) can and should be included for consideration when possible.



## What incidents qualify for this database?:

- "For an event to be included in the GTD, the attackers must be 'out the door,' enroute to execute the attack. Planning, reconnaissance, and acquiring supplies do not meet this threshold." p. 11 [Codebook](#).

## Regarding our output feature 'Success':

- An attack is classified as successful if noticeable consequences occur as a result of the attack.
- "The definition of a successful attack depends on the type of attack. Essentially, the key question is whether or not the attack type took place." p. 26 [Codebook](#).

# LOCATION FEATURES:



What characterizes a successful cyber attack?

## What is our output feature?

Our output feature will be the success column. In other words, this will be the information we are trying to predict. Given that we are building models to try and predict this information, we will use this column to compare our model predictions to some actual results and ascertain its accuracy before applying it to newer information (where we do not have information on the attack's success yet).

## Other features:

This data-set consists of *numerous* features. In order to stay within the time-frame of this project, we will limit our analysis to some of the more prominent and easily recognized/comprehensible features.

This selection of features could lend itself to bias and should not be considered definitive. The use of Select K-Best and PCA algorithms below will aim to reduce some of this bias.

Feature Name	Description	Object Type	Most Frequent Value or/Mean
country_txt	Country	Object	India
region_txt	Region	Object	South Asia
city	City	Object	Unknown
target_1txt	Broad target description.	Object	Private Citizens and Property
sub_targettxt	More detailed target info.	Object	Police Security Forces/Officers
corp1	Identification of targeted corporation (if applicable).	Object	'Not Applicable'
specific_target	Specific target info (if applicable).	Object	Civilians
victim_nationalitytxt	Victim nationality.	Object	India
property	Was there property damage?	Int64	-937

# ATTACK FEATURES:

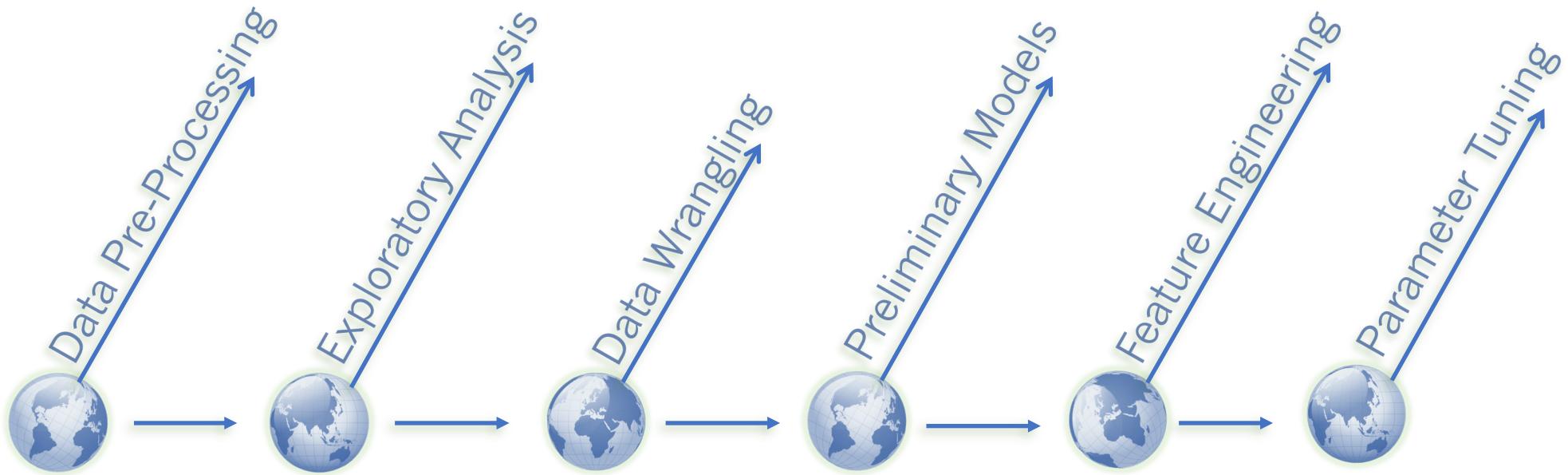


Feature Names	Description	Object Type	Most Frequent Value/Mean
crit1	“Political, Economic, Religious or Social Goal?”	Int64 / Binary	.988
crit2	“Intention to coerce, intimidate or publicize to larger audience(s)?”	Int64 / Binary	.987
crit3	“Outside international humanitarian law?”	Int64 / Binary	.948
doubt	Is there doubt regarding classifying the attack as terrorism?	Float64 / Binary	-.047
suicide	Was the attack a suicide attack?	Int64 / Binary	.040
attack_1txt	Means of attack.	Object	Bombing/Explosion
group_name	Name of perpetrating terrorist group.	Object	Unknown
group_attrib_certainty	Is group-attribution dubious?	Float64 / Binary	.190
unaffil_individ	Was perpetrator identified by name and acting as an individual?	Int64 / Binary	.005
weapon_txt	Type of weapon.	Object	Explosives
sub_weapontxt	More detailed weapon info.	Object	Unknown Gun Type
motive	Motive of attack		n/a
summary	Synopsis of attack.	Object	n/a
<b>success</b>	<b>Was the attack successful?</b>	Int64 / Binary	<b>.904</b>

Regarding our output feature ‘Success’:

- An attack is classified as successful if noticeable consequences occur as a result of the attack. (p. 26 [Codebook](#)).
- “The definition of a successful attack depends on the type of attack. Essentially, the key question is whether or not the attack type took place.” p. 26 [Codebook](#).

# WORKFLOW



# DATA PRE-PROCESSING:

## SIFTING OUT CYBER-RELATED DATA:

- Database is vast and multi-faceted.

- Method: Regex

- Issues: Regex function might not catch all instances - decreasing accuracy

## SPLITTING TRAIN/TEST DATA:

- Prevents over-fitting

- Provides realistic and reliable scoring

- Produces safer models for implementation

## REMOVING NULL VALUES:

- Not a big problem after using our Regex Filter.

- Primarily focused on replacing values with 'Unknown'

- Did not have to drop too many values completely.

## SELECTING PRIMARY FEATURES:

- Focus on understandable and actionable features

- Most translatable across potential audiences

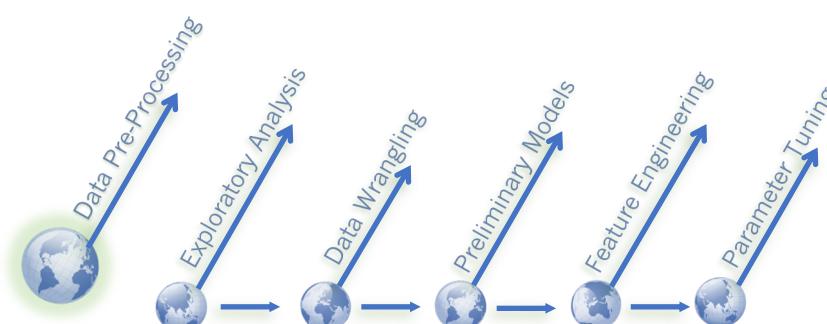
- Numerically significant (i.e. no human-assigned numeric categorizations).

## VALUE-FILTERING:

- Prepares data for later modeling

- Condenses values for feature-selection algorithm

- Keeps data comprehensible for humans while maintaining flexibility for machines/computers



# DATA PRE-PROCESSING:



## SIFTING OUT CYBER-RELATED DATA:

- Database is vast and multi-faceted (too large).
- Filter allows us to focus on cyber-instances.
- Method: Regex
- Features: *motive*, *summary* and *specific\_target*.
- Issues: Regex function might not catch all instances - decreasing accuracy.
- Issues: Could also incorporate instances that are realistically unrelated.



### Detailed Workflow:

```
1 # Regex Filter which fills na in this column with the value: 'Unknown'  
2 specific = cyber_data.specific_target.str.contains(r'(internet|cell+|radio|communic+|email+|cyber+|web|hacking+)',  
3 na = 'Unknown', flags = re.IGNORECASE)
```



```
1 # Combining the above results into a dataframe and looking at the shape:  
2  
3 cyber_data = cyber_data.loc[(cyber_data.index.isin(specific_true)) | (cyber_data.index.isin(motive_true)) |  
4 (cyber_data.index.isin(summary_true)) | (cyber_data.index.isin(specific_unknown))]  
5  
6 cyber_data.shape  
  
(13255, 24)
```



- Data Variable: 'cyber-data'
- Workflow Stage: Pre-Processing
- Sub-Stage: Filtered for Cyber-Instances
- Size: 13255 Observations x 24 Features

# DATA PRE-PROCESSING:



## SPLITTING TRAIN/TEST DATA:

- Prevents over-fitting.
- Provides realistic and reliable scoring.
- Produces safer models for implementation.
- Leakage Prevention: Split before further manipulation to maintain data-integrity.
  - i.e.: Any groupings/null-removals, etc. will be contained within each subset of data. This prevents data removal from leaving traces of its removal in another subset.



### Detailed Workflow:

```
3 # Defining our Input and Output data:  
4 # (Cleaning afterwards to prevent leakage)  
5  
6 X = cyber_data.drop(['event_id', 'success', 'summary', 'motive'], axis = 1)  
7 Y = pd.DataFrame(cyber_data['success'])  
8  
9 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = .25)
```



```
1 print(X_train.shape)  
2 print(X_test.shape)  
  
(9941, 20)  
(3314, 20)
```



- Data Variable: 'X\_train/X\_test/Y\_train/Y\_test'
- Workflow Stage: Pre-Processing
- Sub-Stage: Splitting training/test sets.
- Train-Size: 9941/3314 Observations x 20 Features

# DATA PRE-PROCESSING:

## REPLACING NULL VALUES:

- Not a big problem after using our Regex Filter.
- Primarily focused on replacing values in text columns with 'Unknown':
  - This maintains the rest of the row's information in the dataset without providing misleading information (such as an n/a or '0' value) for this specific column
- Only numeric column was *group\_attrib\_certainty*.
  - This was a very low number. Dropping these values before modeling will still maintain data integrity.



### Detailed Workflow:

Select comprehensible text features:



	0		0
city	26	city	6
sub_targettxt	585	sub_targettxt	202
corp1	1671	corp1	578
specific_target	478	specific_target	158
victim_nationalitytxt	103	victim_nationalitytxt	24
group_attrib_crtainty	40	group_attrib_crtainty	9
sub_weapontxt	1051	sub_weapontxt	344



Fill object columns with 'Unknown.'



Drop numeric nulls just before modeling.



```
1 print(cyber_train_X.shape)
2 print(cyber_test_X.shape)
```

(9941, 19)  
(3314, 19)



- Data Variable: 'cyber\_train/test\_X, cyber\_train/test\_Y'
- Workflow Stage: Pre-Processing
- Sub-Stage: Replacing nulls and selecting meaningful features.
- Size: 9941/3314 Observations x 19 Features

## SELECTING PRIMARY FEATURES:

- Focus on understandable and actionable features.
- Most translatable across potential audiences.
- Numerically significant (i.e. no human-assigned numeric categorizations).

# DATA PRE-PROCESSING:



## VALUE-FILTERING:

- Prepares data for later modeling.
- Condenses values for feature-selection algorithm.
- Keeps data comprehensible for humans while maintaining flexibility for machines/computers.
- Method: Pandas .replace()
- Features: city, group\_name, specific\_target, sub\_target

### Issues:

- Not safe from human bias.
- Highly dependent on careful research.

```
1 # Here is the function we will use to take our list of values and replace it with specific group/category names:  
2 # We will do this for each of the test and training sets:  
3  
4 def magic_value_replacer(df, column, variable, string):  
5  
6     df[column] = df[column].replace(variable, string)  
7
```



### sub\_targettxt:

```
32 internet_comm_information = ['Newspaper Journalist/Staff/Facility', 'Radio Journalist/Staff/Facility',  
33 'Television Journalist/Staff/Facility', 'Other (including online news agencies)',  
34 'Radio', 'Internet Infrastructure', 'Television', 'Electricity',  
35 'Telephone/Telegraph']
```



### specific\_target:

```
1 comm_related = ['Cell tower', 'Cell phone tower', 'Cell Tower', 'Cell Phone Tower', 'Cell Phone Shop',  
2 'Telecommunication Tower', 'Telecommunications Tower', 'Radio Stations', 'Radio Station',  
3 'Radio station', 'radio station', 'Radio station antenna', 'A mobile phone tower',  
4 'A mobile tower was targeted in the attack.', 'A Globe Telecom cell site', 'Internet Cafe',  
5 'Telecommunications office', 'Telecommunication Institute', 'Communications Tower',  
6 'Telecommunications Mast', 'An internet cafe']  
7 polling_areas = ['Polling Station', 'Polling Center', 'Polling Stations', 'Polling Booth']
```



### group\_name:

```
1 palestinian_separatists = ['Hamas (Islamic Resistance Movement)', 'Palestinian Islamic Jihad (PIJ)',  
2 'Popular Front for the Liberation of Palestine (PFLP)', 'Popular Resistance Committees',  
3 'Al-Fatah']
```



### city:

```
1 sunni_cities = ['Mosul', 'Kirkuk', 'Sanandaj', 'Ramadi', 'Trabzone', 'Diarbekir',  
2 'Damascus', 'Gwadar', 'zahedan', 'Kandahar', 'Khiva', 'Fallujah',  
3 'Dakhla', 'Tajura', 'Sabrata', 'Azizia', 'Kasabat', 'Misrata', 'Tripoli',  
4 'Takrit', 'Tikrit']
```



# EXPLORATORY ANALYSIS:

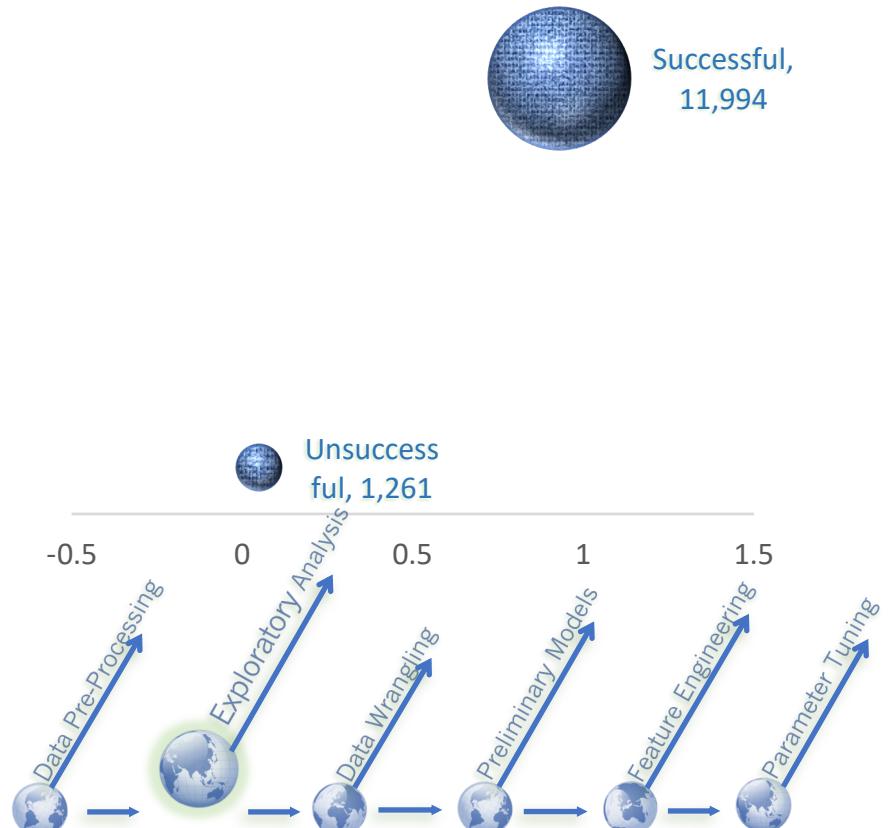


## CLASS IMBALANCE:

Can lead to a warped perspective:

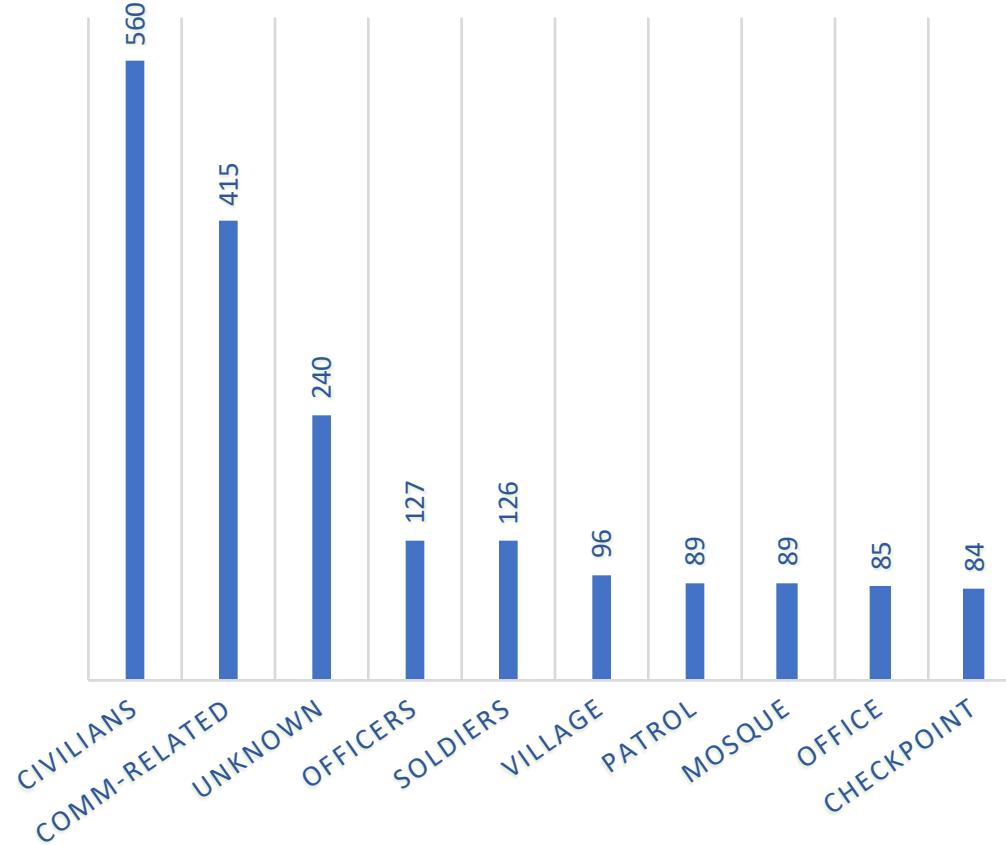
- Lots of information about successful attacks.
  - Not enough information on unsuccessful attacks.
- 
- Keeping an eye on our false positive predictions will be necessary. Most of our predictions will be positive. Limiting our incorrectly negative predictions will tell us how strong our model is.
  - Using a balanced accuracy score will help our evaluations take class imbalance into account.

NUMBER OF SUCCESSFUL V.  
UNSUCCESSFUL ATTACKS

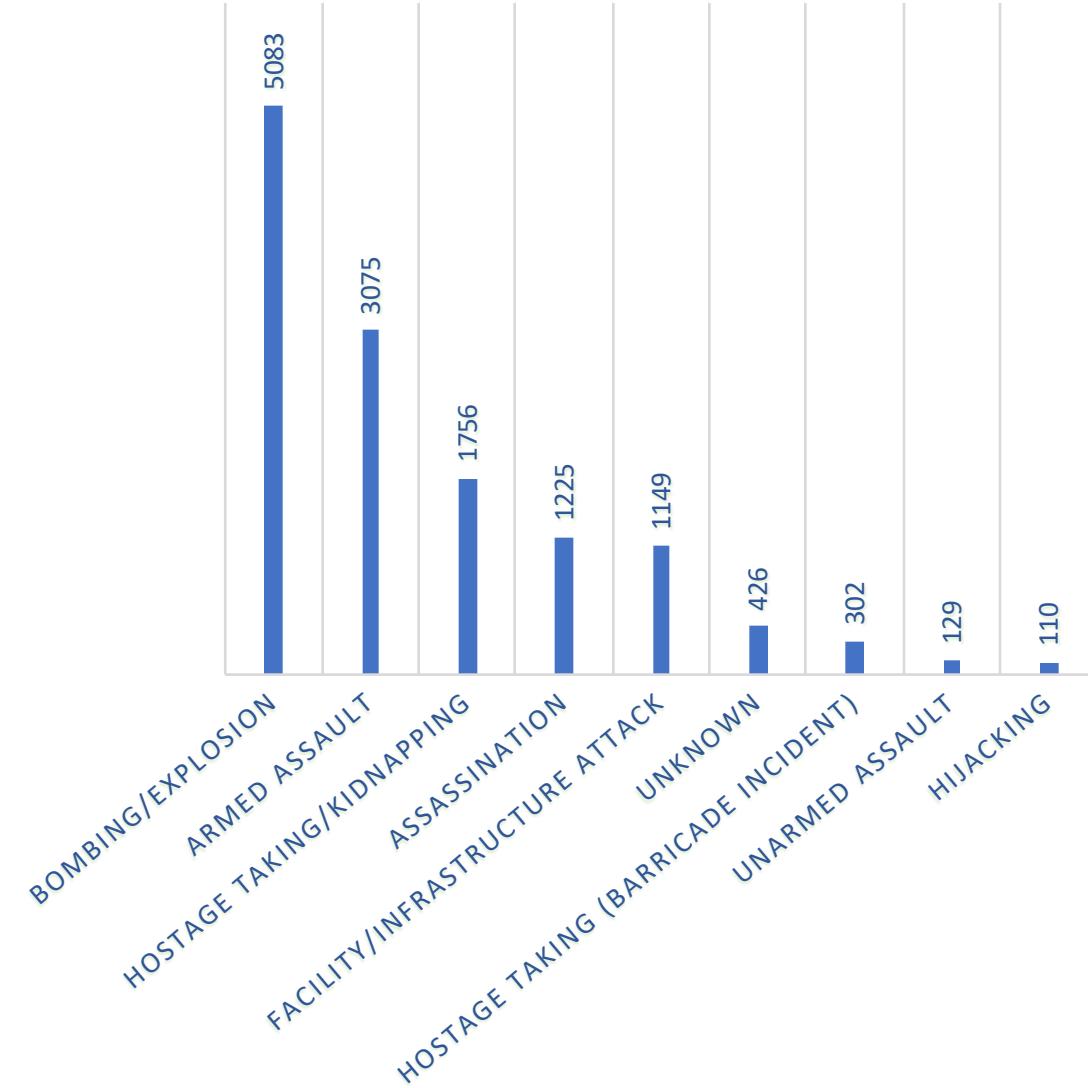


# FEATURE DISTRIBUTIONS:

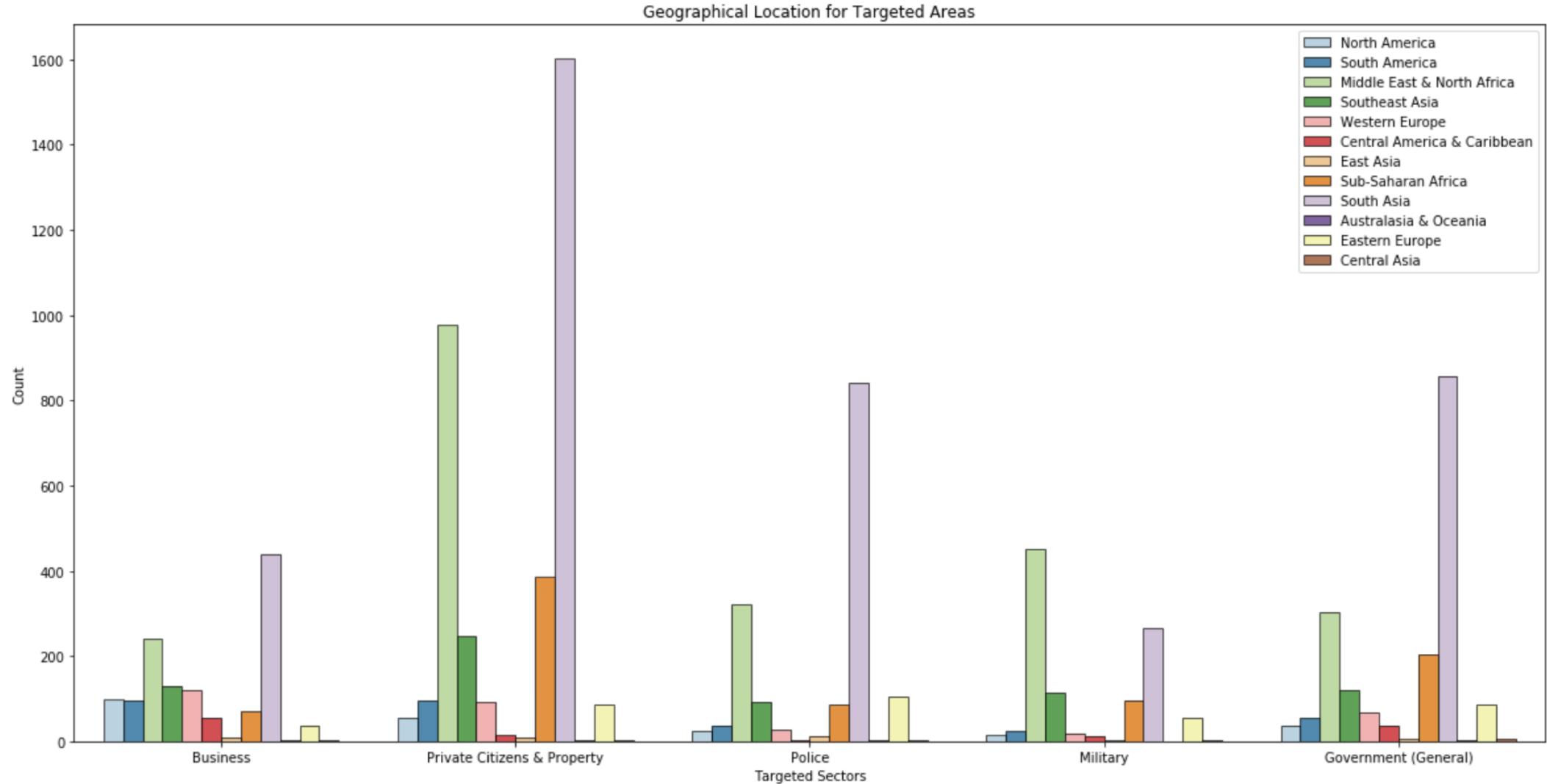
## TOP 10 SPECIFIC TARGETS



## TYPE OF ATTACKS

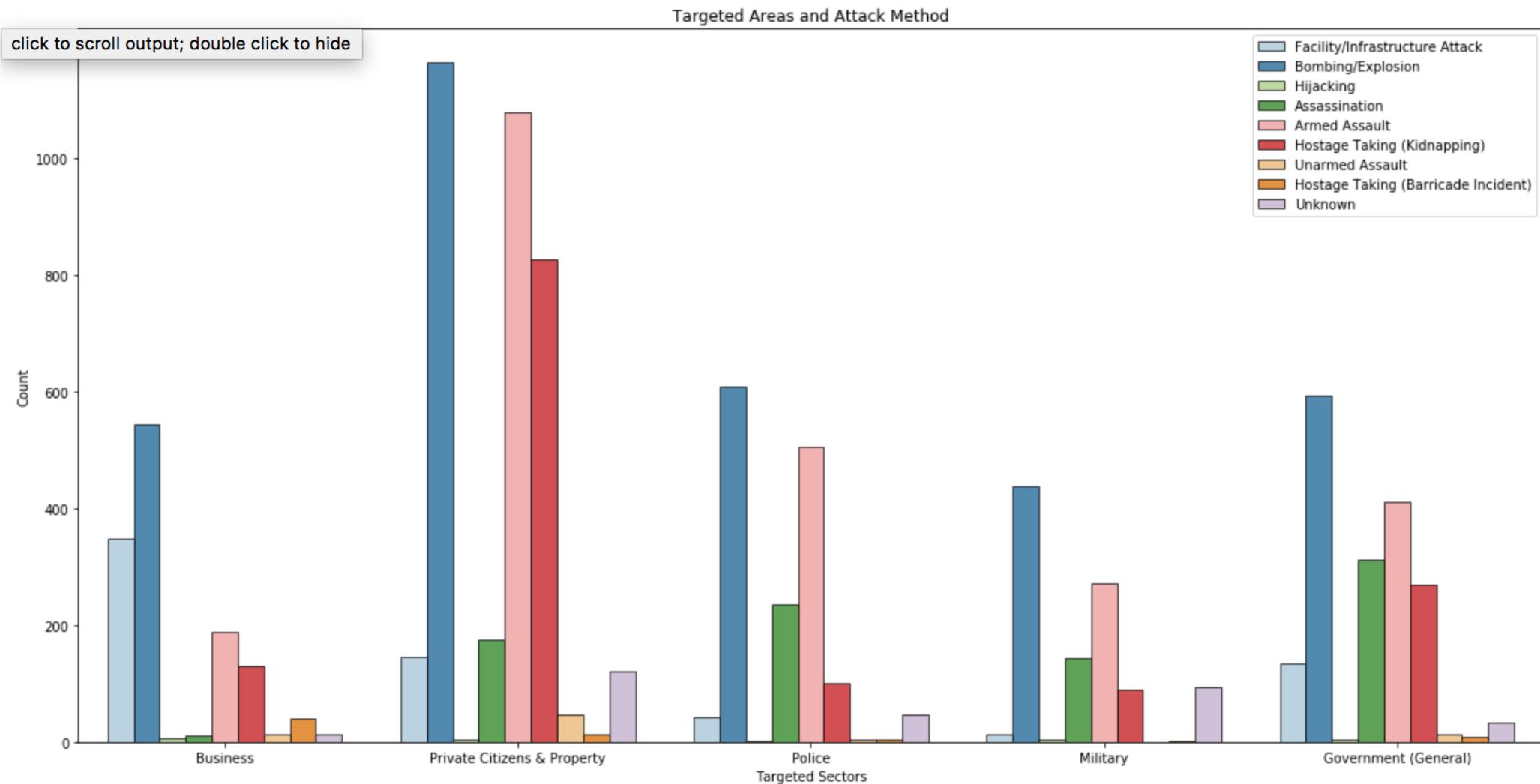


# PRELIMINARY GRAPHS:



This graph displays the distribution of 'sectors' which were targeted and how these numbers are distributed throughout the main regions of the world.

# PRELIMINARY GRAPHS:



This graph displays the distribution of 'sectors' which were targeted and how frequently these attacks were conducted with different methods.

# DATA WRANGLING / FEATURE SELECTION

## SELECT K-BEST:

- This algorithm uses statistical measures to highlight features which have the most promising correlation to our success columns.
- These were the top columns selected from our dataset:

	0
0	country_txt_India
1	country_txt_Nepal
2	country_txt_Uruguay
3	country_txt_Yemen
4	attack_1txt_Armed Assault
5	attack_1txt_Assassination
6	attack_1txt_Hostage Taking (Kidnapping)
7	target_1txt_Government (General)
8	target_1txt_Private Citizens & Property
9	target_1txt_Unknown
10	sub_targettxt_Unknown



### Detailed Workflow:

Ran Select-KBest and PCA on the whole dataset



Ran both algorithms on subsets of the data which were isolated according to the following features:

*group\_name,  
city,  
specific\_target*

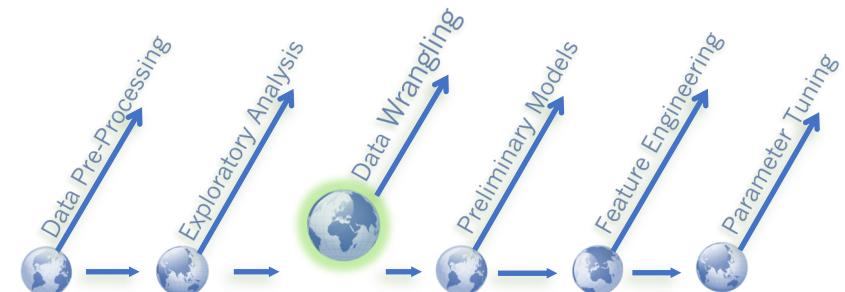


Combined results from the four steps above into two specific data-frames to test with our models:

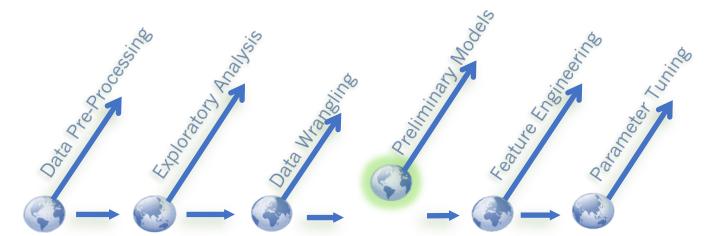
- KBest\_train/test\_X/Y
- PCA\_train/test\_X/Y

## PRINCIPLE COMPONENT ANALYSIS:

- Takes a look at where all of the data-point lies.
  - Places vectors through areas with more variance.
  - Outputs these vectors as features to describe the data's variance.
- Difficult for humans to understand what these selected features/results mean (unlike Select K-Best).
- Can return the amount of 'explained variance,' which can be helpful.



# PRELIMINARY MODELS:



Here, we can see the count of the errors in each trial model prediction as well as its Balanced-Accuracy.

## Why Balanced-Accuracy?

Since we have an output feature ('Success') which is highly imbalanced, we will need an accuracy score that takes class-imbalance into account (i.e. SKLearn's Balanced Accuracy metric).

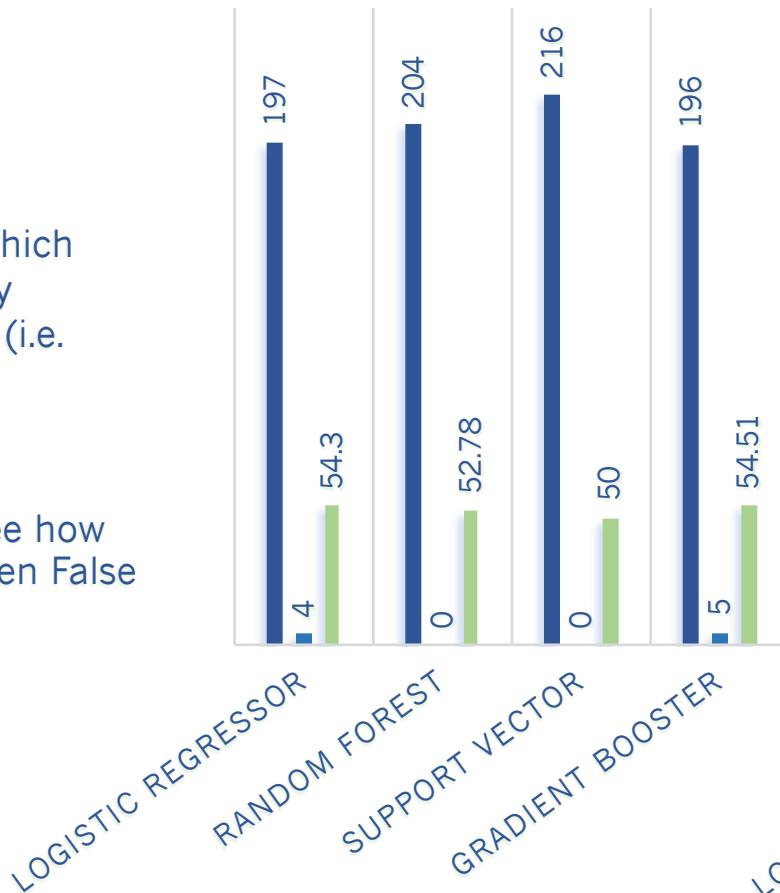
This metric gives us an additional matrix to see how our two types of errors balance out (i.e. between False Positives and False Negatives).

*Winner: Random Forest with K-Best.*

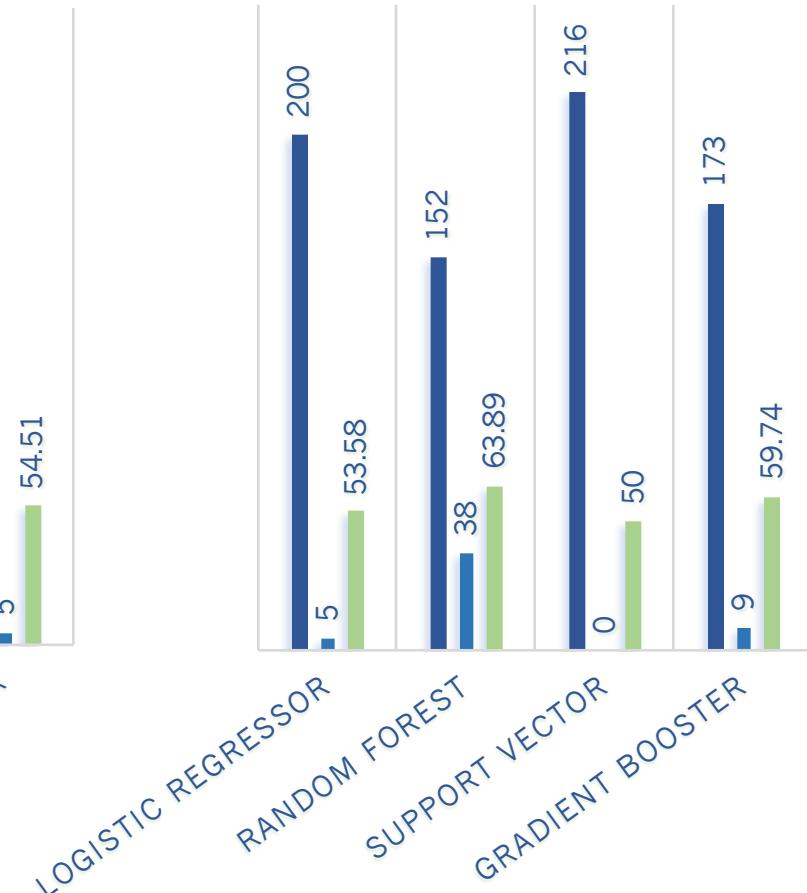
Starting with the least False Negatives is ideal for saving lives.

Relatively low False Positives will help increase our Balanced Accuracy.

K-BEST SELECTED FEATURES



PCA SELECTED FEATURES

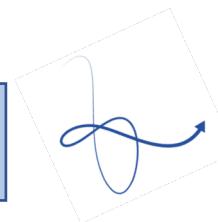


■ FALSE POSITIVES ■ FALSE NEGATIVES ■ BALANCED ACCURACY

# FEATURE ENGINEERING:



attack\_1txt: Assassination



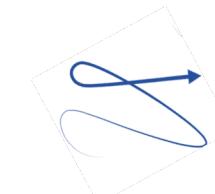
Entire Data-set Data Grouped by Assassination

	success
1	11994
0	1261

	success
1	606
0	310

Our previous class imbalance decreases drastically.

success: 0



	success
0	27
1	16

Success rate of assassination attempts within ‘Sunni\_Cities.’

Our imbalance is now inverted.

city: top 5:



city	
Sunni_Cities	27
Unknown	16
Mogadishu	14
Split_Cities	11
Benghazi	6

Any instance in our dataset meeting these criteria (i.e. any failed assassination attempt in any of these top-5 cities), will receive a 1 in the *Barometer* column.

All else will receive a 0.

Barometer: 0/1 (Binary)

NB: Inserting other location related features into the filtering code for this engineered feature can also be done – depending on which features appear in our KBest feature set whenever we run our models.

# PARAMETER TUNING:

*n\_estimators* (range: 40-280):  
The number of 'trees' helping us find correct classifications.

- Smaller parameters:
- Decreased accuracy.
  - Decreased runtime

- Larger parameters:
- Increased accuracy.
  - Increased runtime
  - Potential over-fitting.

NB: By looking at some more of the location features included our KBest-feature set and adding them to our filtering function for the engineered feature (see previous slide), we are able to increase our accuracy further:

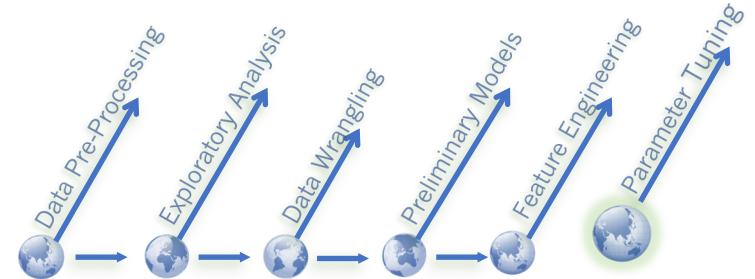
- Random Forest-Original Features (+ params): 0.5837
- Random Forest-Original Features (- params): 0.5807
- Random Forest-Engineered Features (+params): 0.5933
- Random Forest-Engineered Features (- params): 0.5914

## Parameters tuned:

*max\_depth* (range: 3-15):  
The 'size' of the trees in our forest.

- Smaller parameters:
- Insufficient accuracy.
  - Decreased runtime

- Larger parameters:
- Considerable accuracy.
  - Increased runtime
  - Considerable over-fitting



Tuned Model Results	Max Depth	Balanced Accuracy	# of Estimators
Random Forest-Original Features (+ params)	15	56.369	80
Random Forest-Original Features (- params)	15	56.348	40
Random Forest-Engineered Features (+ params)	15	56.3	120
Random Forest-Engineered Features (- params)	15	56.46	80

# FINAL SUMMARY:

---



## Potential Weaknesses:

- More Research: The above filters are merely a start and represent a base amount of research. It could certainly serve as a starting point for someone with expertise in the field.
- Bias: the filter itself will not be able to catch all instances of cyber-related terrorism conclusively; more likely than not it will miss one certain aspect or another. Increasing the number of eyes on the model and data will help with this.

## Explanatory v. Predictive Power:

- Our goal : *Prediction*.
  - We want to take *new/incoming data* and gain a prediction of whether it will be successful or not before we know the results.
  - As little variation as possible in this prediction and its accuracy score will be necessary for reliable results.
- Explanatory models: more interest to scholars or academic institutions
  - Reveals behaviors and patterns of observations that have already occurred
  - No bearing on future observations.
  - No suggestions on what these patterns might say regarding incoming data (and in this case, future terror attacks).
- Caveat: making sure that our model adapts well to new input (or test) data.
  - Allows for the variation between training and test results to be as minimal as possible (i.e. so the difference between the model's prediction and reality is minimal).
  - Predictive analysis, especially in a terror-related context, involves considerably more risk than an explanatory model, and should be handled with an intense attention to detail and accuracy.

## FURTHER RESEARCH:

---

Given the richness of the variables within this dataset, there are plenty of avenues for further research:

- One could use Regex to search the summary and motive columns for further detail and insight into the nature of these cyber-related tasks. The filter above is relatively broad and potentially encapsulates instances that many might not consider related to 'cyber' events.
- It would be interesting to create a more intricate regex filter which could give us a more detailed understanding of the 'cyber' aspect of these events:
  - How specifically are they related and in what manners?
  - What geographical locations tend to be hotbeds for such activity.
  - What targets are specified and why?
  - The 'motive' feature, in particular, could have extensive benefits with prediction, depending on the vocabulary used by those conducting the study and entering the data.

# FINAL SUMMARY:

---

## General References:

Referenced the following sites for honing my knowledge of the models, python, etc:

<https://www.analyticsvidhya.com/blog/2017/06/a-comprehensive-guide-for-linear-ridge-and-lasso-regression/>

<https://machinelearningmastery.com/feature-selection-machine-learning-python/>

<https://medium.com/@pushkarmandot/what-is-the-significance-of-c-value-in-support-vector-machine-28224e852c5a>

A big shout out to Mike Swirsky for pointing out this link to me (submitted by the user 'Abdou' on Stack Overflow: <https://stackoverflow.com/questions/40993626/list-memory-usage-in-ipython-and-jupyter>)

For a look at the multiple resources used for filtering our data-set in the beginning stages of the work, please see the original [Jupyter Notebook](#) here.

We also used Scikit-Learn for our algorithms and modeling:

[Scikit-learn: Machine Learning in Python](#), Pedregosa et al., JMLR 12, pp. 2825-2830, 2011.