# An Information Security Approach to Feature Engineering

RELIAQUEST

CAMLIS

# Brian P. Murphy
## ReliaQuest Chief Architect

RELIAQUEST

- Driving technical vision on the GreyMatter platform. Graduate of the University of Limerick in Ireland.

- Worked at Splunk, Elasticsearch, Loggly.

- ReliaQuest partners with Splunk and other Fortune 1000 customers in managing the Security Model across their entire organization.

# Feature Engineering

# Feature Engineering
## Definition

- Feature engineering is the process of using domain knowledge of the data to create features that make machine learning algorithms work.

- "Coming up with features is difficult, time-consuming, requires expert knowledge. 'Applied machine learning' is basically feature engineering."
  — Andrew Ng, Machine Learning and AI via Brain simulations

# Feature Engineering
## Outcome Desired

- Create encoders that can represent individual elements in a set of related values while also maintaining their relationships.

- URIs / Domains / Log messages
  - Classify and detect outlier and/or malicious activity

- Geo locations
  - Typically similar attacks will come from the same geo regions, but not necessarily the same countries.

RELIAQUEST

# Feature Engineering
## Machine Learning

- Most algorithms require columnar numerical values

- Models suffer from GIGO (Garbage In/Garbage Out)

- Memory constraints normally impossible to hold the entire data set in RAM

# Feature Engineering
## Common Options

- Feature hashing

- One hot encoding

- Ordinal/Label encoding

# Feature Engineering
## Issues

- Feature hashing
  - Locality of information is lost

- One hot encoding
  - Can result in feature explosion
  - Dealing with new values not seen in training set

- Ordinal/Label encoding
  - Locality is not always obvious
  - Dealing with new values is hard also

# Feature Engineering
## Proposed solution

- Min Hash Shingle

- Min Hash ngram

- Geo Hash

# Feature Engineering

## ngram/shingling

- Ngrams are contiguous sequences of characters from events.
  - Camlis -> 3gram -> ["cam","aml","mli","lis"]


- Shingles are to tokens what ngrams are to characters.
  - Camlis is really awesome -> 3shingle -> ["camlisisreally", "isreallyawesome"]

# Feature Engineering
## Min Hashing

- Jaccard Similarity

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$

  - Very accurate.
  - But need to compare every element of a set with every other element of the set.

- Min Hash
  - Approximates Jaccard by creating 'k' hash functions and hashing each ngram or shingle, then finding the min value for each hash.
  - Events that share min hashes are similar.

Vassilvitskii, Sergey - COMS 6998-12

# Feature Engineering
## Min Hash Benefits

- Encodes all values into a known number of columns.

- New values are handled and if they are similar to fitted values they will be similarly encoded.

- Linear time event by event processing.

# Feature Engineering
## Geo Hash Approach

- Partitions the globe into a hierarchical NxN (32x32) grid, using a z-curve.

# Feature Engineering

## Geo Hash Approach

- Calculate a bivariate (3D) normal distribution with the location at the peak of this distribution.

- Based on distance from the grid peak assign values across the grid.

- Results in NxN columns representing a decaying weight from the detected action.

# Feature Engineering
## Geo Hash Benefits

- Allows the encoding of all possible geolocations into known number of columns.

- New locations are allowed for in the encoding.

- Maintains relationships between geolocations.

- Lookups and values can be precomputed to accelerate encoding.

# Feature Engineering
## Demo

# New Search

Save As ▾     Close

    index=_internal uri=* | minhashngram field=uri  | head 10000 | apply urikmeans3 | search cluster=4          Last 24 hours ▾     🔍

✓ 133 events (10/25/19 2:00:00.000 PM to 10/26/19 2:04:54.000 PM)     No Event Sampling ▾          Job ▾  ⏸ ⏹ ↗ 🔖 ⬇     ☰ Verbose Mode ▾

**Events (133)**     Patterns     Statistics     Visualization

Format Timeline ▾          — Zoom Out          + Zoom to Selection          ✕ Deselect          1 hour per column

List ▾     ✎ Format     20 Per Page ▾          ‹ Prev   **1**   2   3   4   5   6   7   Next ›

| ↓ | Time | Event |
|---|------|-------|
| < Hide Fields     ☰ All Fields | > | 10/26/19<br>2:03:58.023 PM | 127.0.0.1 - admin [26/Oct/2019:14:03:58.023 -0400] "GET /servicesNS/nobody/Splunk_ML_Toolkit/properties/mlspl/defau<br>ltToutput_mode=json&count=-1 HTTP/1.0" 200 3545 - - - 1ms |

SELECTED FIELDS
# bytes 5
# host 1
# source 2
# sourcetype 2

bytes = **3545**     host = **ubuntu**     source = /home/gwelladh/FEDemo/splunk/var/log/splunk/splunkd_access.log
sourcetype = splunkd_access

> 10/26/19     127.0.0.1 - admin [26/Oct/2019:14:03:56.808 -0400] "GET /servicesNS/nobody/Splunk_ML_Toolkit/properties/mlspl/defau
2:03:56.808 PM     ltToutput_mode=json&count=-1 HTTP/1.0" 200 3545 - - - 1ms

INTERESTING FIELDS
# clientip 1
# cluster 1

bytes = **3545**     host = **ubuntu**     source = /home/gwelladh/FEDemo/splunk/var/log/splunk/splunkd_access.log
sourcetype = splunkd_access

# cluster_distance 5
# count 1
# date_hour 5

> 10/26/19     127.0.0.1 - admin [26/Oct/2019:14:03:55.910 -0400] "GET /servicesNS/nobody/Splunk_ML_Toolkit/properties/mlspl/defau
2:03:55.910 PM     ltToutput_mode=json&count=-1 HTTP/1.0" 200 3545 - - - 1ms

# date_mday 2

bytes = **3545**     host = **ubuntu**     source = /home/gwelladh/FEDemo/splunk/var/log/splunk/splunkd_access.log

# New Search

Save As ▾     Close

```
index=_internal uri==
| minhashngram field=uri
| head 10000
| apply urikmeans3
| search cluster=4
| stats first(uri), count by cluster_distance
```

Last 24 hours ▾     🔍

✓ 139 events (10/25/19 2:00:00.000 PM to 10/26/19 2:06:41.000 PM)     No Event Sampling ▾                                      Job ▾   ‖  ▣  ↗  🖨  ⬇    📖 Verbose Mode ▾

Events (139)      Patterns      **Statistics (5)**      Visualization

50 Per Page ▾      ✎ Format      Preview ▾

| cluster_distance ⇅  ✎ | first(uri) ⇅                                                                                                    ✎ | count ⬆ ✎ |
|---|---|---|
| 62268186098824.42 | /servicesNS/nobody/Splunk_ML_Toolkit/properties/mlspl/default?output_mode=json&count=-1 | 128 |
| 1.1785543172030556e+16 | /en-US/splunkd/__raw/servicesNS/admin/user-prefs/data/user-prefs/general | 5 |
| 3.3741051419322736e+16 | /servicesNS/admin/Splunk_ML_Toolkit/data/lookup-table-files/__mlspl_example_hard_drives_StandardScaler_1.csv?output_mode=json | 3 |
| 3.3824262151073404e+16 | /servicesNS/admin/Splunk_ML_Toolkit/data/lookup-table-files/__mlspl_example_hard_drives_StandardScaler_1.csv | 2 |
| 1.0854925036232244e+16 | /en-US/splunkd/__raw/services/search/timeparser?output_mode=json&time=-24h&_=1572092886531 | 1 |

splunk>enterprise    App: Splunk Machine Learning To... ▾    ⚠ Administrator ▾    🔵 Messages ▾    Settings ▾    Activity ▾    Help ▾    Find   🔍

Showcase    Experiments    **Search**    Models    Classic ▾    Settings    Docs 🗗    Video Tutorials 🗗    🎧 Splunk Machine Learning Toolkit

# New Search

Save As ▾    Close

```
source="port_443_external.json.out host="ubuntu" index="dotconf" dst_geo.latitude=* dst_geo.longitude=*,
| euclid latfield=dst_geo.latitude lonfield=dst_geo.longitude
| head 10000
| apply groeuclidkmeans
| stats count, values(dst_geo.country_name) as countries, values(dst_geo.city_name) as cities by cluster
```

All time ▾   🔍

✓ 10,000 events (before 10/26/19 8:30:10.000 AM)    No Event Sampling ▾        Job ▾   ⏸ ⏹ ↗ 🖨 ⬇    ⬛ Verbose Mode ▾

Events (10,000)    Patterns    **Statistics (20)**    Visualization

50 Per Page ▾    ✎ Format    Preview ▾

| cluster ⇅ ✎ | count ⇅ ✎ | countries ⇅ ✎ | cities ⇅ ✎ |
|---|---|---|---|
| 0 | 964 | Canada<br>Japan<br>United States | Boydton<br>Lansing<br>Port Coquitlam<br>Tokyo |
| 1 | 2003 | United States | Ashburn<br>Chantilly<br>Macon<br>Reston |
| 10 | 123 | Hong Kong<br>United States | Ann Arbor<br>Central<br>Clifton<br>Parsippany |
| 11 | 252 | France | Austin |

```
| fields + dst_geo.city_name, dst_geo.country_name, dst_geo.latitude, dst_geo.longitude
| head 10000
| geohash latfield=dst_geo.latitude lonfield=dst_geo.longitude
| apply geokmeans2
| stats count, values(dst_geo.country_name) as countries, values(dst_geo.city_name) as cities by cluster
```

✓ 10,000 events (before 10/25/19 5:58:50.000 PM)    No Event Sampling ▾                                    Job ▾   ⏸  ⏹  ↗  🔥  ⬇    ⊞ Verbose Mode ▾

Events (10,000)    Patterns    **Statistics (19)**    Visualization

50 Per Page ▾    ✎ Format    Preview ▾

| cluster ⇅ ✎ | count ⇅ ✎ | countries ⇅ ✎ | cities ⇅ ✎ |
|---|---|---|---|
| 0 | 3301 | United States | |
| 1 | 4125 | United States | Ashburn<br>Beltsville<br>Boydton<br>Chantilly<br>Clifton<br>New York<br>North Bergen<br>Parsippany<br>Reston<br>Secaucus<br>Staten Island<br>Washington |
| 10 | 277 | United States | Austin<br>Dallas<br>Irving<br>San Antonio |

# Feature Engineering
## Demo

# Thank You

https://github.com/GaelTadh/rq_feature_engineering
(coming soon)

pmurphy@reliaquest.com