# Identifying Politicians in the
# PANAMA PAPERS

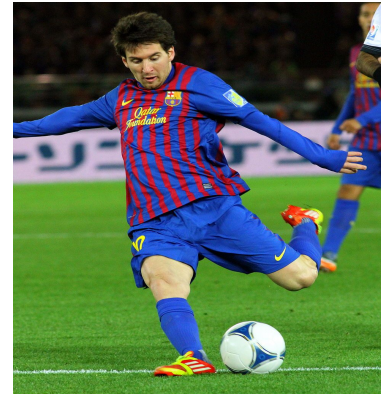Irene Beckman, Ruhi Saraf, Santhanankrishnan Ramani

# Panama Papers

Early 2016 11.5 million emails, power-of-attorney letters, and internal notes from the law firm of Mossack Fonseca were leaked to the press.
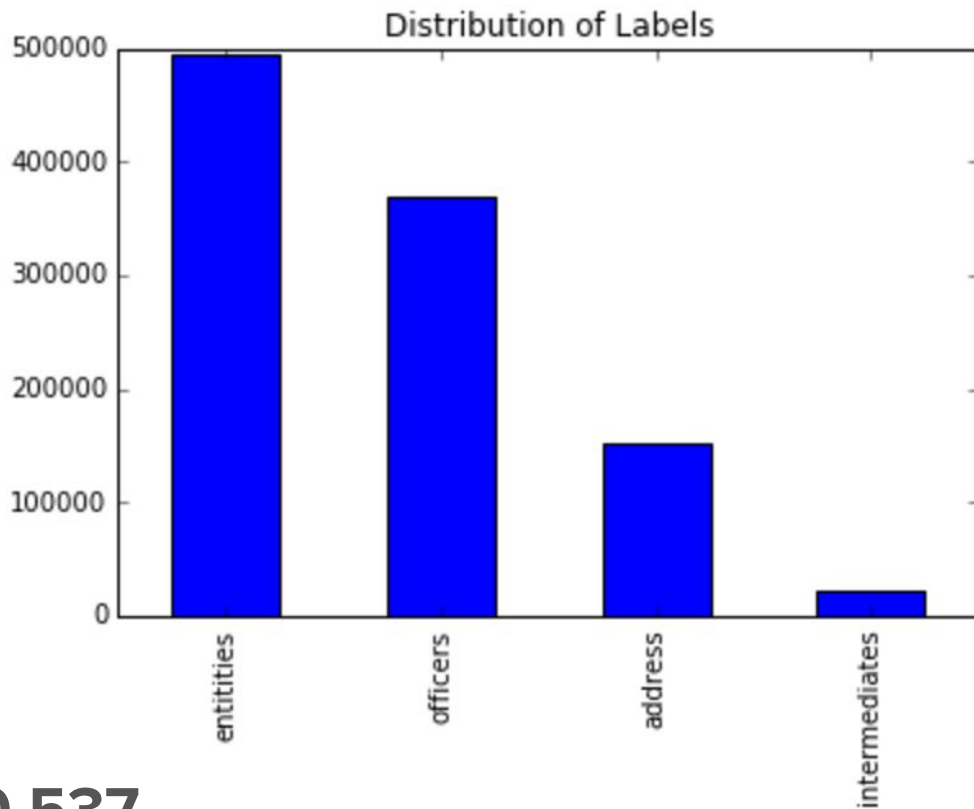




Panama is the only place in the world where one can see the sun rise on the Pacific Ocean and set on the Atlantic

- The leak shed light on the financial dealing and details of:
  - World political leaders
  - Fraudsters
  - Drug traffickers
  - Billionaires
  - Celebrities
  - Sports stars and more
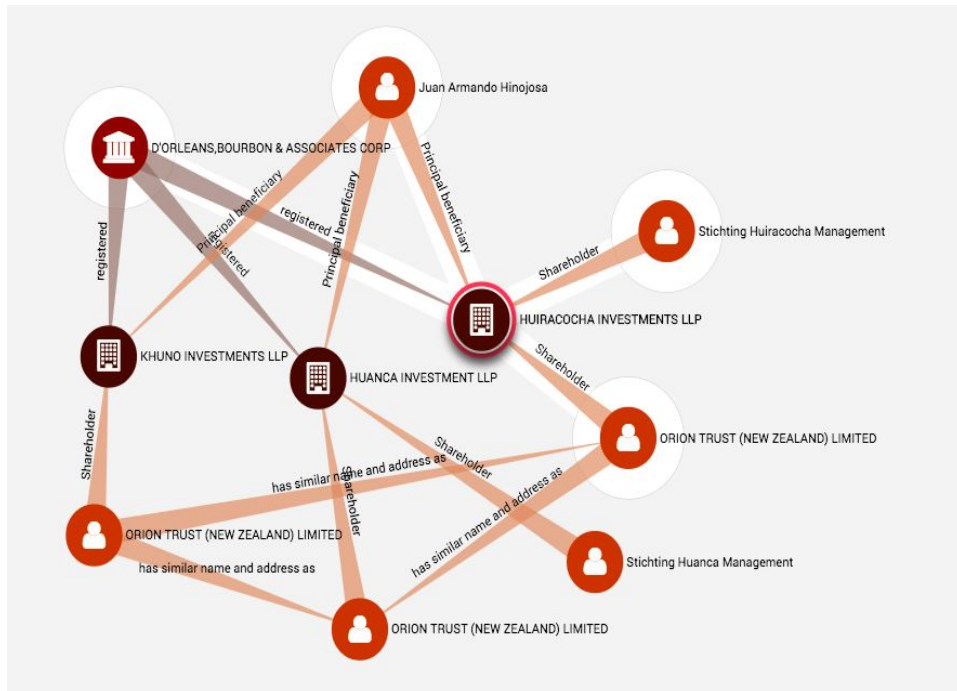- In some cases, brought about new political scandals

# The Dataset:  Nodes

❖ Entities
❖ Officers
  ➢ People: ≅200,000
  ➢ Politicians:  102
❖ Intermediaries
❖ Addresses



Distribution of Labels

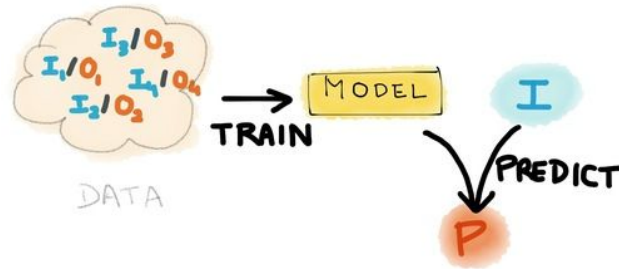Total number of nodes: **1,040,537**

# The Dataset: Edges



- Any affiliation between the nodes
  - beneficiary
  - shareholder
  - registered
  - similar name
  - director

Total number of edges: **4,429,023**

# Hypothesis

*Can we identify people officer nodes with political affiliation by modeling this as a binary classification problem?*

# Data Parsing and Filtering

❖ *Problem*:  Large data with distinct nodes types made it hard to pinpoint a subset within officer types

➢ *Solution*: Reducing the data set to only include people from officers list

❖ *Problem*: No label for politicians and affiliates

➢ *Solution*: manually labeled all 102 known politicians

❖ *Problem*: Duplicate names with a variety of spelling for same node

➢ *Solution*:  Due to the unknown extent of this issue, we could not fix it globally and so left the data as it came

# Feature Engineering

- ❖ Vertex Level Measures
- ❖ Ego Networks
- ❖ Use of directed and undirected representation of graphs
- ❖ Motifs

# Vertex Level Measures



Page Rank

Mean Page Rank ( 1 = Politician, 0 = Non-Politician)

In Degree

Mean In-Degree ( 1 = Politician, 0 = Non-Politician)

Out Degree

Mean Out-Degree ( 1 = Politician, 0 = Non-Politician)

# Vertex Level Measures

# Ego Network Measures

# Ego Network Measures



Rich Club Coefficient

Mean Rich Coefficient ( 1 = Politician, 0 = Non-Politician)

Rich Coefficient

isPolit



Radiality Centrality  (w/ ego)

Mean Radiality Centrality ( 1 = Politician, 0 = Non-Politician)

Radiality Centrality

isPolit

# Structure of Network

# Motifs

❖ Node/Edge Motifs

❖ Ego Network Motifs

❖ Neighbors Edge Motifs

# Sampling and Class Imbalance
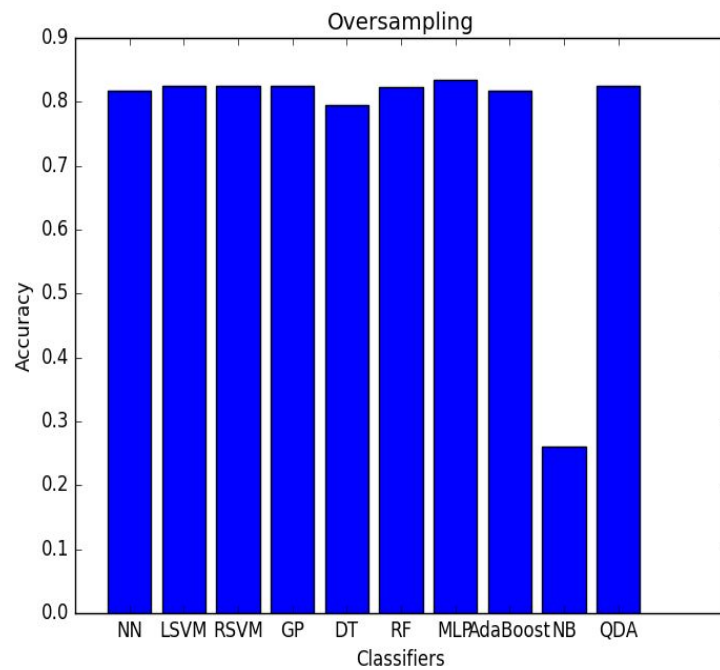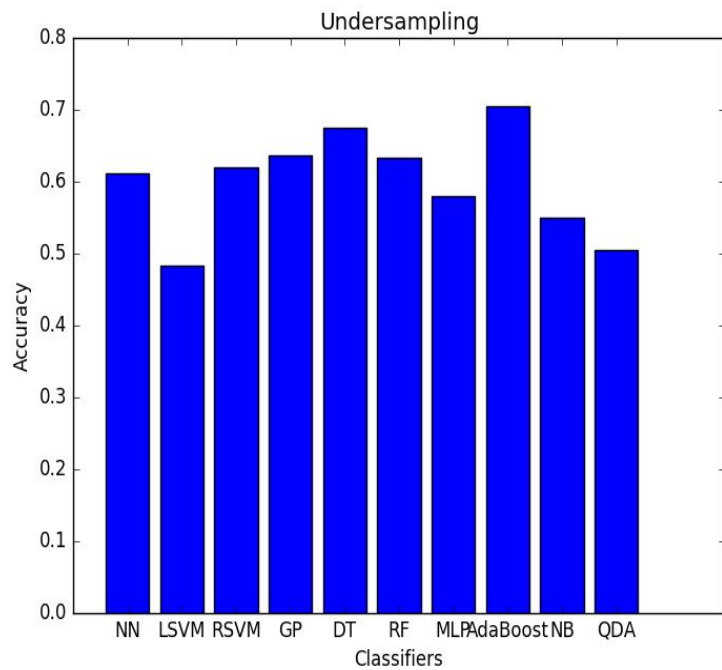
Issue: Combating Imbalanced Classes

What we tried:

1. Resampling: Undersampling, Oversampling, Unequal Ratios
2. Changed Performance Metric

# Method

❖ Various Classifiers ( Random Forest, SVC, Neural Net, AdaBoost )

❖ Parameter Optimization using Grid Search

❖ 10-fold testing to compute average accuracy

# Results

# Future Work:

❖ Feature Selection

❖ Try penalized models

❖ Ensemble Classification using Multiple subsets

❖ Balanced Cascade

❖ Incorporate country/region data manually
   (meta data associated with the nodes)

# Questions?