



# Detecting Fraud with Data Science

...

David Hogue, Jeff Grana, Neil Bezdek, Ryan Busby

galvanize

# Overview

Objective: Examine Eventbrite data for fraudulent events

Dataset: 43 features + 1 fraud label

Workflow: Feature engineering - 37 features

Model building - Random Forest, Gradient Boost, ADABOOST

Profit Curve Analysis

Deliverables: Dashboard to predict fraud from live stream

# Examine fraud cases to guide feature engineering

**"BEST EVENT EVER --  
BELIEVE ME"**

Events with high percentage of  
capitalized text more likely to be  
fraud

**\$\$\$**

Fraudulent events more likely to  
promote high-end "luxury"  
products



# Feature Engineering

## Dropped:

'object\_id'  
'sale\_duration'  
'user\_created'  
'user\_type'  
'ticket\_type'  
'venue\_address'  
'venue\_latitude'  
'venue\_longitude'  
'venue\_country'  
'venue\_name'  
'venue\_state'  
'payout\_type'  
'email\_domain'  
'description'  
'approx\_payout\_date'  
'event\_start'  
'event\_end'  
'event\_published'  
'event\_created'  
'currency'  
'country'  
'delivery\_method'  
'acct\_type'  
'name'

## Kept:

'body\_length'  
'gts'  
'has\_logo'  
'num\_order'  
'show\_map'  
'channels'  
'has\_analytics'  
'num\_payouts'  
'user\_age'  
'name\_length'  
'sale\_duration\_2'  
'fb\_published'

## Edited:

'has\_header'  
'org\_desc'  
'listed'  
'org\_facebook'  
'org\_twitter'  
'payee\_name'  
'org\_name'  
'event\_published'  
'blank\_venue\_name'  
'time\_to\_pay'  
'time\_to\_pay2'  
'planning\_time'  
'has\_pub\_date'

## Engineered:

'avg\_ticket\_price'  
'max\_ticket\_price'  
'ticket\_tiers'  
'total\_available'  
'pct\_caps'  
'previous\_payouts'  
'blank\_venue\_address'  
'same\_countries'  
'org\_facebook\_nan'  
'nan\_venue'  
'org\_twitter\_nan'  
'duration'

# Feature Importance

(u'previous\_payouts', 0.33717854008808879),  
(u'time\_to\_pay2', 0.11508351431934936),  
(u'sale\_duration2', 0.086194935447564056),  
(u'gts', 0.067157648378076376),  
(u'user\_age', 0.057314953453766054),  
(u'num\_order', 0.053524360465484186),  
(u'pct\_caps', 0.029795932631783963),  
(u'num\_payouts', 0.027944121385088033),  
(u'body\_length', 0.025926269426975432),  
(u'name\_length', 0.022668285447869457),  
(u'avg\_ticket\_price', 0.02188740770758547),  
(u'max\_ticket\_price', 0.021178139089864909),  
(u'org\_facebook', 0.01733205830640748),  
(u'org\_twitter', 0.012126223484773999),  
(u'duration', 0.011166827233625016),  
(u'time\_to\_pay', 0.0109335820590898),  
(u'planning\_time', 0.010801764808993627)

Engineered: df['previous\_payouts'].apply(len)

Engineered: df['approx\_payout\_date'] - df['event\_start']

Kept

Kept

Kept

Kept

Engineered: percent capitals from df['name']

Kept

Kept

Kept

Engineered: Weighted Average (# tickets \* cost @ each tier)

Engineered: Max price of all ticket types

Edited: df['org\_facebook'].fillna(0.)

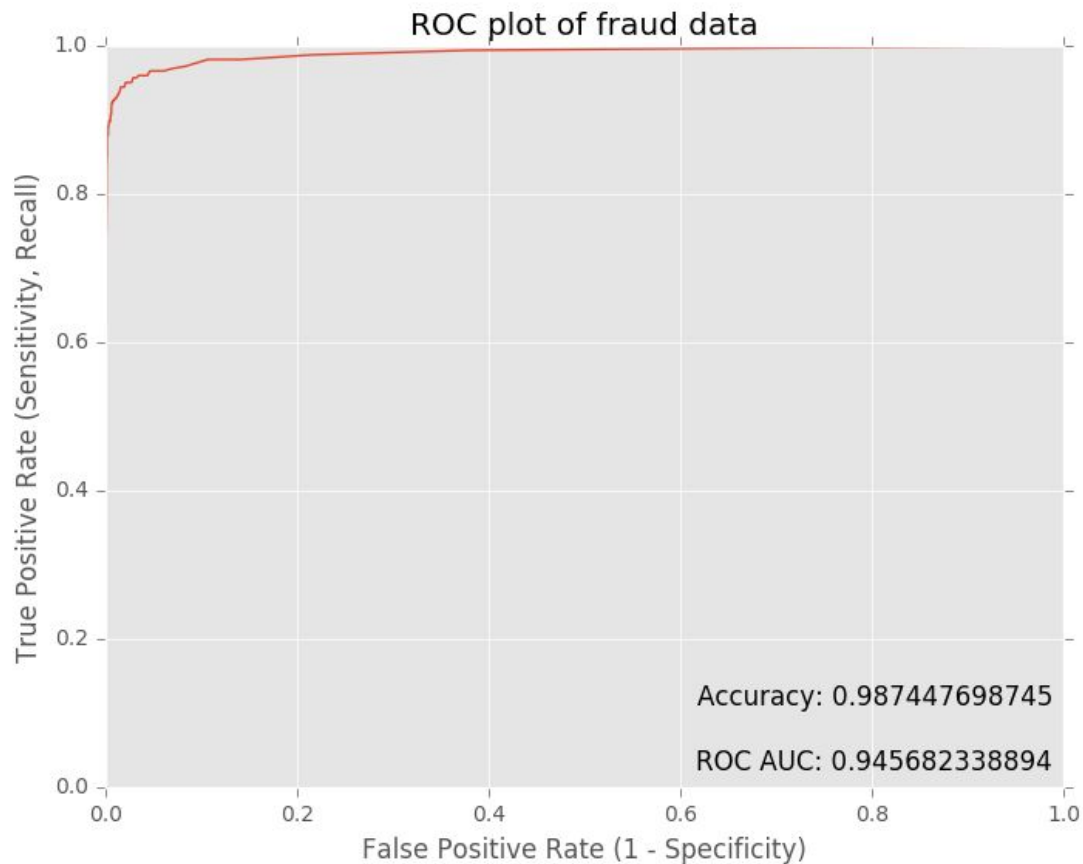
Edited: df['org\_twitter'].fillna(0.)

Engineered: df['event\_end'] - df['event\_start']

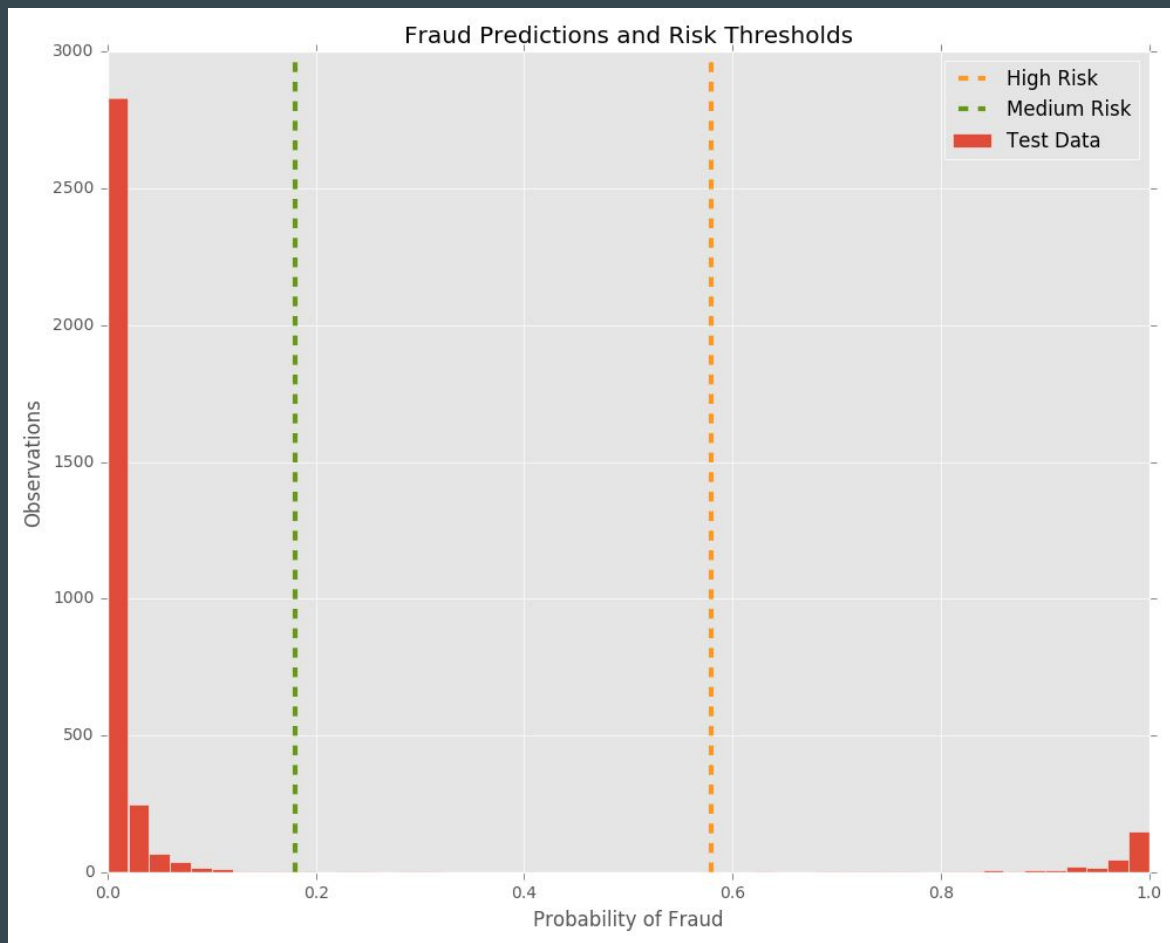
Engineered: df['approx\_payout\_date'] - df['event\_start']

Engineered: f['approx\_payout\_date'] - df['event\_start']

# ROC Curve



# Model Probabilities



# Profit Curve

TP 450	FP -50
-----------	-----------

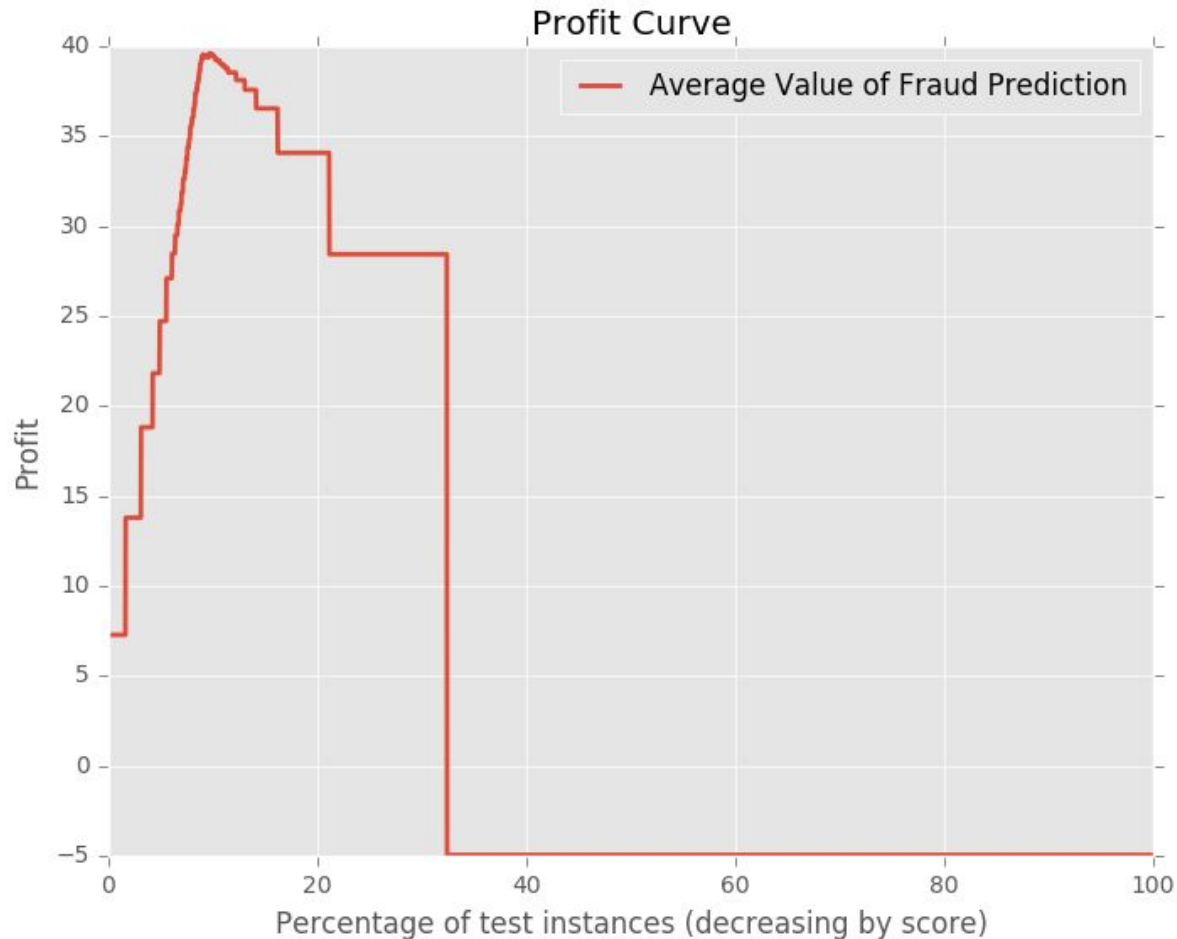
0 FN	0 TN
---------	---------

Threshold = 0.18  
Profit = 39.59

TP 450	FP -1000
-----------	-------------

0 FN	0 TN
---------	---------

Threshold = 0.58  
Profit = 38.05





# Conclusions

EDA and feature engineering contributed to model accuracy. Use common sense before feeding numbers to your model!

Random Forest model was most accurate model in this case.

Application of model was informed by analysis of cost function. In this case, classification changed from .5 to .1 as a result.

We are better at data science than web development or network engineering.