# Using Decision Trees to Identify Phishing Sites

Lherisson Medina

- Every day more and more people are susceptible to phishing attempts through email and sites they may visit.

- The affects of a successful phishing attempt are long term, and can ruin an individual's economic and social life for years.

- In May, millions of Google Gmail users were hit with a sophisticated phishing attack which took all day to identify and fix

- Estimated 85% of companies are hit with phishing attacks in recent years

- Because of its consequences, it is important to recognize when a site may be attempting to phish data

- These sites are the subject of my term project



Lherisson Medina

# The Data

- The data was downloaded form the Machine Learning Repository at UC Irvine

- Based on attributes gathered by Auckland Institute of Studies

### 9 Available Features

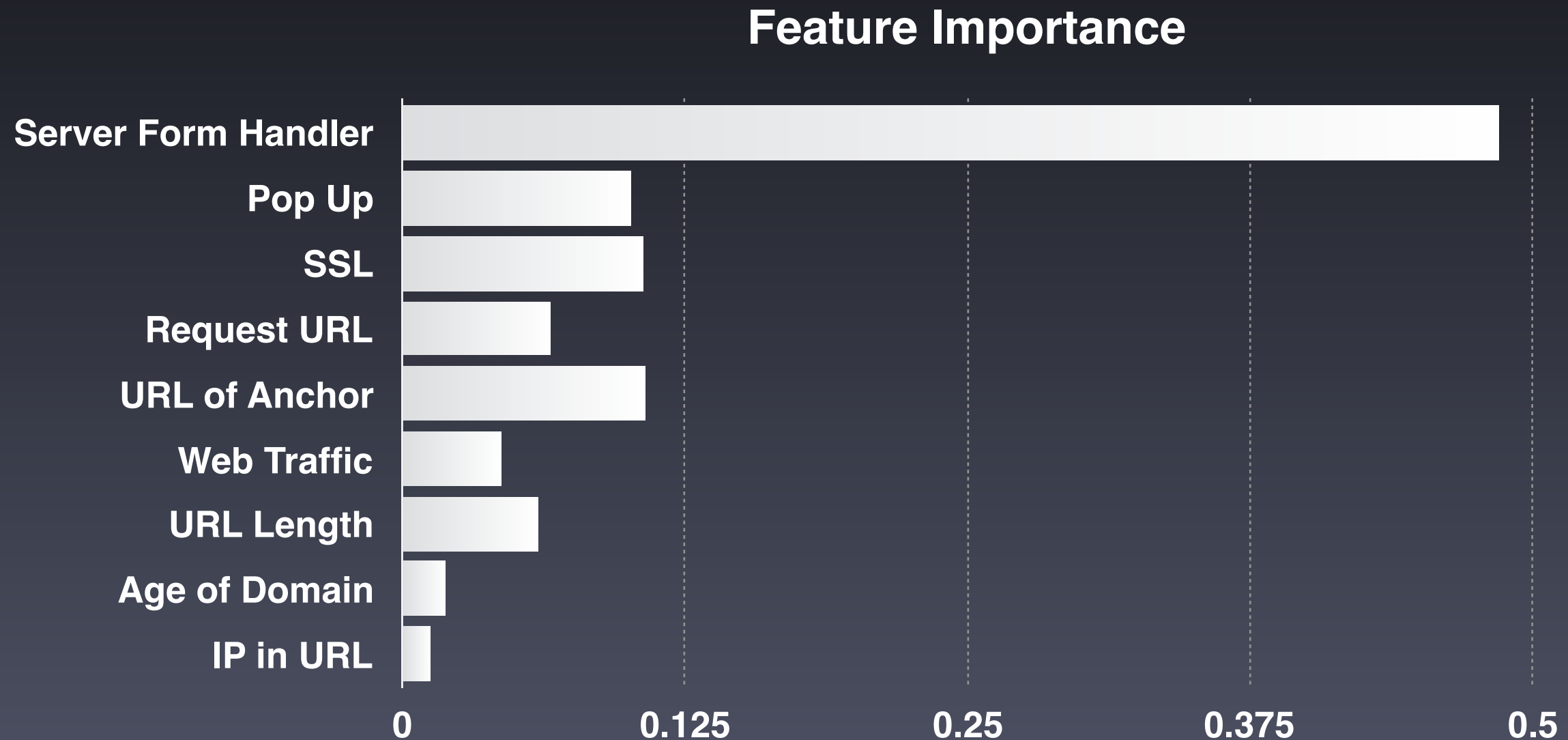| SFH | HasPopUp | SSL | Request URL | Anchor URL |
|-----|----------|-----|-------------|------------|
| Web Traffic | URL Length | Domain Age | IP Address in URL | |

- Feature values have been normalized to [-1, 0, 1] where

  -1 - Phishy
  0 - Suspicious
  1 - Legitimate

Lherisson Medina

# Approach

- I used a supervised learning approach

- Decision Trees used to create a predictive model based on the data

- Out of 1,353 Data Points, a random 25% used as test size

  - 75% used for training

# Approach (Feature Selection)

- The Server Form Handler is the most informative feature

- 'Domain age' and 'IP in URL' are the least informative

**Feature Importance**



Lherisson Medina

# Results

- With top 5 most informative features selected

  - ~87% Accuracy

- Can be improved

  - Encountered overfitting

- Regression Trees also tested



Lherisson Medina

# Further Work

- Classification Trees only one approach

- Consider comparing with other approaches

  - Neural Net, Nearest Neighbor, etc

- Optimize tree (Pruning?) to increase Accuracy and decrease overfitting

# Similar Work

- Black-Lists and White-Lists are more generally used to track Phishing and Legitimate Sites

- Norman Sadeh et al. looked at Emails and URLs to try and classify whether they are Phishing or Legitimate using Random Forest Classifier

- May 31st Google started delaying and flagging emails having predictable patterns to Phishing emails

  - No info on Algorithm

# Questions

Lherisson Medina

# Thank You