

domain knowledge

gaurav and suriyan

Why predict the category of content hosted by a domain?

Why predict the category of content hosted by a domain?

- malware and phishing

Why predict the category of content hosted by a domain?

- malware and phishing
- adult content

Why predict the category of content hosted by a domain?

- malware and phishing
- adult content
- personalization

Why predict the category of content hosted by a domain?

- malware and phishing
- adult content
- personalization
  - you can have your *cookie* and eat it too!

But if we were to learn from domain names, how would we?

But if we were to learn from domain names, how would we?

- manual labeling—PhishTank, DM0Z, etc.



But if we were to learn from domain names, how would we?

- manual labeling—PhishTank, DMOZ, etc.
- meta tags—self reports

But if we were to learn from domain names, how would we?

- manual labeling—PhishTank, DMOZ, etc.
- meta tags—self reports
- metadata: how the pages are laid out, how many images, etc.

But if we were to learn from domain names, how would we?

- manual labeling—PhishTank, DMOZ, etc.
- meta tags—self reports
- metadata: how the pages are laid out, how many images, etc.
- content understanding: text, images, etc.

But if we were to learn from domain names, how would we?

- manual labeling—PhishTank, DM0Z, etc.
- meta tags—self reports
- metadata: how the pages are laid out, how many images, etc.
- content understanding: text, images, etc.
- signal in the url

But if we were to learn from domain names, how would we?

- manual labeling—PhishTank, DM0Z, etc.
- meta tags—self reports
- metadata: how the pages are laid out, how many images, etc.
- content understanding: text, images, etc.
- signal in the url

## Considerations

## Considerations

- computational burden, speed

## Considerations

- computational burden, speed  
    > 2B domains
- privacy



## Considerations

- computational burden, speed  
    > 2B domains
- privacy
- accuracy

Intuition

Intuition

- Patterns in names.

## Intuition

- Patterns in names.

words like *xxx*, *porn*, *adult*, *sex* etc. common in domains carrying pornographic content

## Intuition

- Patterns in names.

words like *xxx*, *porn*, *adult*, *sex* etc. common in domains carrying pornographic content

- > 200 PhishTank URLs have the word *paypal* in them.

## Intuition

- Patterns in names.

words like *xxx*, *porn*, *adult*, *sex* etc. common in domains carrying pornographic content

- > 200 PhishTank URLs have the word *paypal* in them.

compare with Alexa 1M to build classifiers that can tell those apart

## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier

## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.



## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.
  - Use a large corpus

## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.
  - Use a large corpus
  - Learn the context well

## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.
  - Use a large corpus
  - Learn the context well
  - Preserve tens of hundreds of vectors and pass them to a model

## How to Classify Text

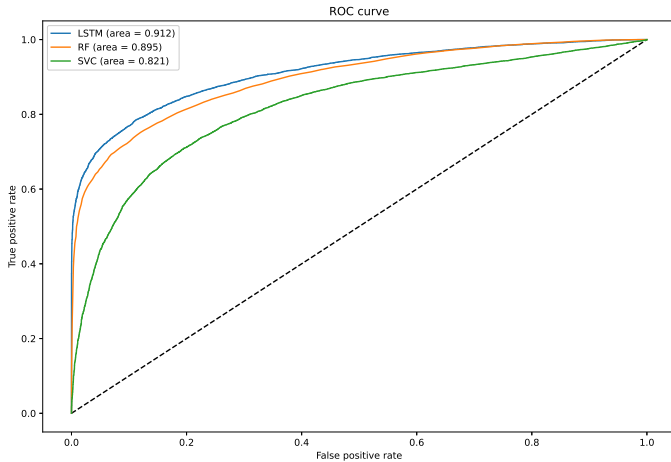
- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.
  - Use a large corpus
  - Learn the context well
  - Preserve tens of hundreds of vectors and pass them to a model
- In Our Case:

## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.
  - Use a large corpus
  - Learn the context well
  - Preserve tens of hundreds of vectors and pass them to a model
- In Our Case:
  - Embeddings of common bi-chars

## How to Classify Text

- Text  $\rightsquigarrow$  Embeddings  $\rightsquigarrow$  Classifier
  - Embeddings leverage the adage:  
You are the company you keep.
  - Use a large corpus
  - Learn the context well
  - Preserve tens of hundreds of vectors and pass them to a model
- In Our Case:
  - Embeddings of common bi-chars
  - LSTM



PhishTank

Application



## Application

- Are disadvantaged people most at risk online?

## Application

- Are disadvantaged people most at risk online?

## Application

- Are disadvantaged people most at risk online?

## Data

## Application

- Are disadvantaged people most at risk online?

## Data

- comScore panel

## Application

- Are disadvantaged people most at risk online?

## Data

- comScore panel
- domain level data for a machine in a household + household attributes

# Concerns

## Concerns

- Domain  $\neq$  URL

## Concerns

- Domain  $\neq$  URL
- Net error is a function of FP - FN.



## Concerns

- Domain  $\neq$  URL
- Net error is a function of FP - FN.
- The larger the number of domains, the larger the error in their score

## Concerns

- Domain  $\neq$  URL
- Net error is a function of FP - FN.
- The larger the number of domains, the larger the error in their score
- Visits are right-skewed. Errors in the right tail are very costly.

## Concerns

- Domain  $\neq$  URL
- Net error is a function of FP - FN.
- The larger the number of domains, the larger the error in their score
- Visits are right-skewed. Errors in the right tail are very costly.

## Solutions

## Concerns

- Domain  $\neq$  URL
- Net error is a function of  $FP - FN$ .
- The larger the number of domains, the larger the error in their score
- Visits are right-skewed. Errors in the right tail are very costly.

## Solutions

- Calibration  $\rightsquigarrow FP = FN$  locally

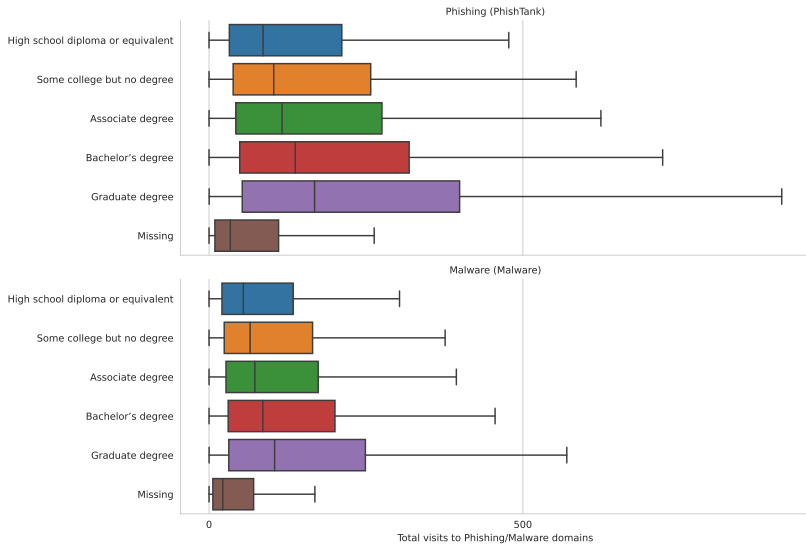
## Concerns

- Domain  $\neq$  URL
- Net error is a function of FP - FN.
- The larger the number of domains, the larger the error in their score
- Visits are right-skewed. Errors in the right tail are very costly.

## Solutions

- Calibration  $\rightsquigarrow FP = FN$  locally
- Using curated lists (manual coding) for popular domains

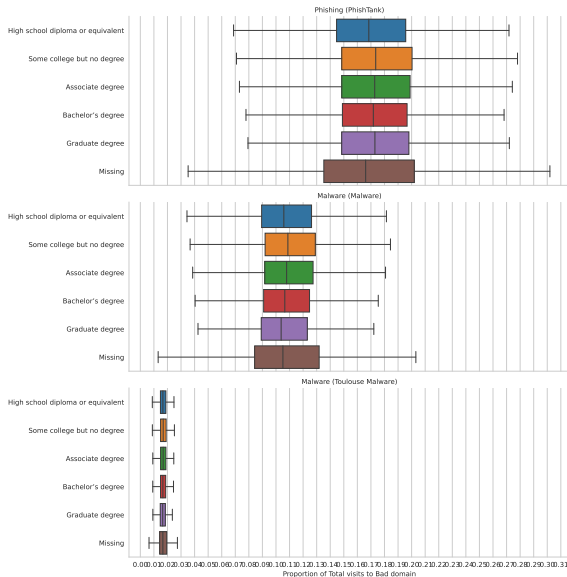
The better educated visit phishing/malware domains most often



Visits to phishing/malware domains

But the better educated don't choose worse





Proportion of visits to phishing/malware domains