# CS259D: Data Mining for Cybersecurity

# Phishing
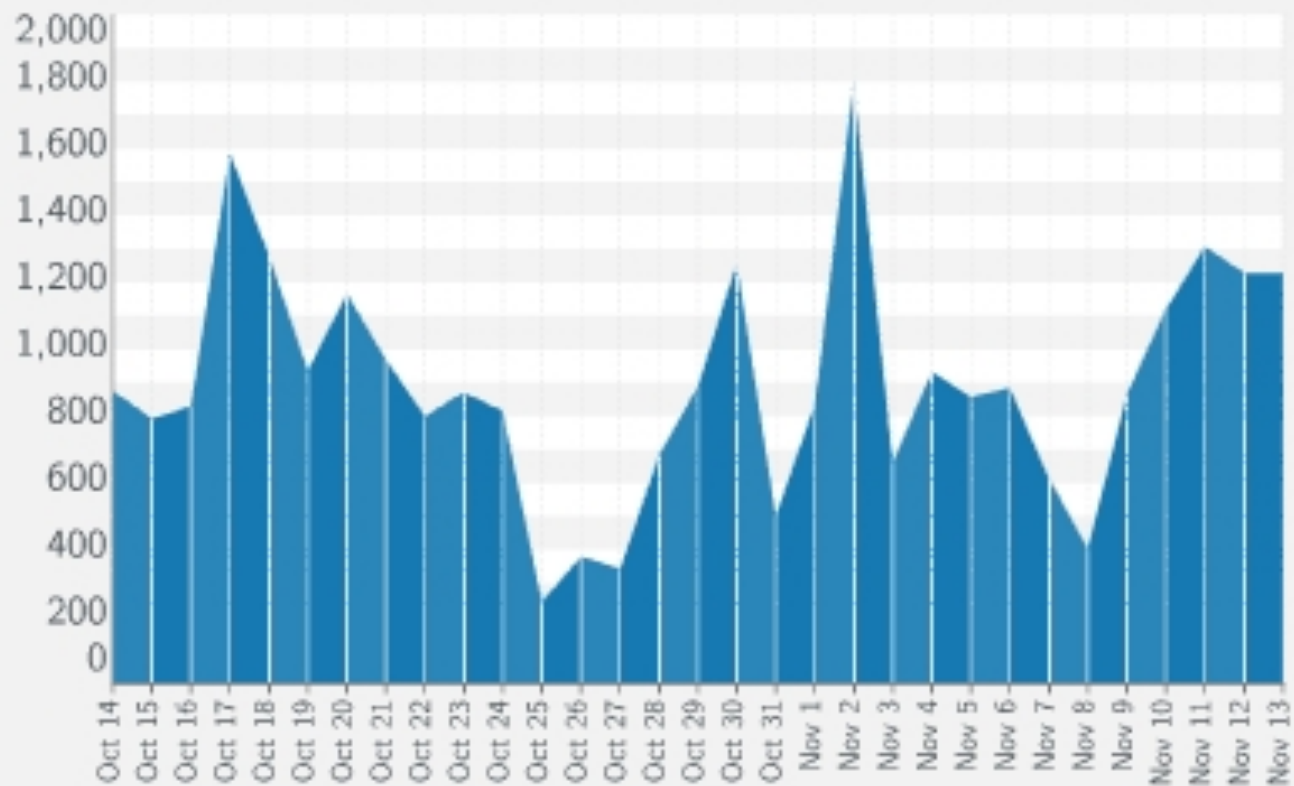
- Goal:
  - Account information
  - Logon credentials
  - Identity information
- Attack vectors:
  - Legitimate-looking emails
  - Legitimate-looking websites

# Scale of the problem



**Daily Phishes Verified**
chart created Nov 13 2014 22:50 UTC

PhishTank

# Detection

- Toolbars
  - Spoofguard
  - Netcraft
- Email filtering
  - Examples:
    - SpamAssassin
    - Spamato
  - Advantages:
    - More complete context (content, headers, etc.)
    - Completely shield user from decision-making process

# Phishing classification: Features

- IP-based URLs
  - Example:
    - http://192.168.0.1/paypal.cgi?fix account
  - Compromised PCs with no DNS entries
  - Binary feature
- Age of linked-to domain names
  - Registered legitimate-sounding domain names
    - Example: playpal.com, paypal-update.com
  - Typically short life-span
    - Registered using stolen credit cards, canceled by registrar
    - Domain caught by anti-phishing monitors
    - Often lasting only ~ 48 hours
  - Obtained using a WHOIS query
  - Binary feature: Lifetime < 60 days

# Phishing classification: Features

- Non-matching URLs
  - Example: <a href="badsite.com"> paypal.com</a>
  - Binary feature: URL text different from HREF
- "Here" links to non-modal domain
  - Example: Click here to restore your account
  - Modal domain: domain most frequently linked to
  - Binary feature: link with text "link", "click", "here" that links to a domain other than modal domain

# Phishing classification: Features

- ## HTML emails
  - ◦ Binary feature: email section with MIME type text/html

- ## Number of links
  - ◦ Numeric feature: # links in HTML part(s) of email
  - ◦ Link defined by an <a> tag with href attribute
    - • Including mailto: links

# Phishing classification: Features

- Number of domains
  - Domain names for URLs starting with http/https
  - Only the main part of the domain name
    - What registrar gets paid for
      - Not necessarily same as combination of top- & 2nd-level domain
    - Example:
      - **university.edu** for www.cs.university.edu
      - **company.co.jp** for www.company.co.jp
        - Top-level: **.jp**, second-level: **.co**
  - Numeric feature: #distinct domains

# Phishing classification: Features

- Number of dots
  - Subdomains:
    http://www.my-bank.update.data.com
  - Redirection script:
    http://www.google.com/url?q=http://www.badsite.com
    - Looks to naïve user to be from google.com
    - Redirects browser to badsite.com
  - Numeric feature: Maximum number of dots in any of the links in the email

# Phishing classification: Features

- Contains javascript
  - Binary feature: string "javascript" appears in email
- Spam filter output
  - Binary feature: class assigned to email by SpamAssassin

# Evaluation

- 10-fold cross validation
- Classifier: Random forest
  - 10 decision trees
  - Each decision made on a random attribute
  - Trees pruned

# Evaluation

- SpamAssassin ham corpora
  - ~6950 non-phishing non-spam
- Publicly available phishingcorpus
  - ~ 860 phishing messages
  - Challenge with WHOIS queries
    - Only 505 domains out of 870 domains
    - Increases false negative rate

# Evaluation

| Feature | Non-Phishing Matched | Phishing Matched |
|---------|---------------------|------------------|
| Has IP link | 0.06% | 45.04% |
| Has "fresh" link | 0.98% | 12.49% |
| Has "nonmatching" URL | 0.14% | 50.64% |
| Has non-modal here link | 0.82% | 18.20% |
| Is HTML email | 5.55% | 93.47% |
| Contains JavaScript | 2.30% | 10.15% |
| SpamAssassin Output | 0.12% | 87.05% |

# Evaluation

| Feature | $\mu_{\text{phishing}}$ | $\sigma_{\text{phishing}}$ | $\mu_{\text{non-phishing}}$ | $\sigma_{\text{non-phishing}}$ |
|---|---|---|---|---|
| Number of links | 3.87 | 4.97 | 2.36 | 12.00 |
| Number of domains | 1.49 | 1.42 | 0.43 | 3.32 |
| Number of dots | 3.78 | 1.94 | 0.19 | 0.87 |

# Evaluation

| Classifier | False Positive Rate $fp$ | False Negative Rate $fn$ |
|---|---|---|
| PILFER, with S.A. feature | 0.0013 | 0.036 |
| PILFER, without S.A. feature | 0.0022 | 0.085 |
| SpamAssassin (Untrained) | 0.0014 | 0.376 |
| SpamAssassin (Trained) | 0.0012 | 0.130 |

# Review of TF-IDF

- Measure importance of word in document
- TF = frequency of word in document
- IDF = measure popularity of word in corpus
  - Log(N/#{documents having the term})
- tf-idf (t, d, D) = tf(t, d) x idf(t, D)

# Robust hyperlinks

- Lexical signatures for identifying URLs
- Signature words chosen using TF-IDF
- Experiments: 5 terms enough for unique page identification

# Observation

- Minimal changes to original page detectable via robust hyperlinks
- Phishing sites often include brand names
  - Common on brand's webpages
  - Rare on the web

# Algorithm

- Compute term TF-IDFs
- Find top 5 terms
- Submit terms as query to Google
- Check if domain is among top-N results
- Assumption: phishing pages have low pagerank

# Lowering false positives

- Include domain name in lexical signature
- Heuristic: Zero results Means Phishing

# Example

# Example

# Example

- Top terms: eBay, user, sign, help, forgot

# Other features

- Age of domain
- Known images
  - Presence of inconsistent well-known logos
  - Top-10 identified targets:  eBay, PayPal, Citibank, Bank of America, Fifth Third Bank, Barclays Bank, ANZ Bank, Chase Bank, and Wells Fargo Bank
- Suspicious URL
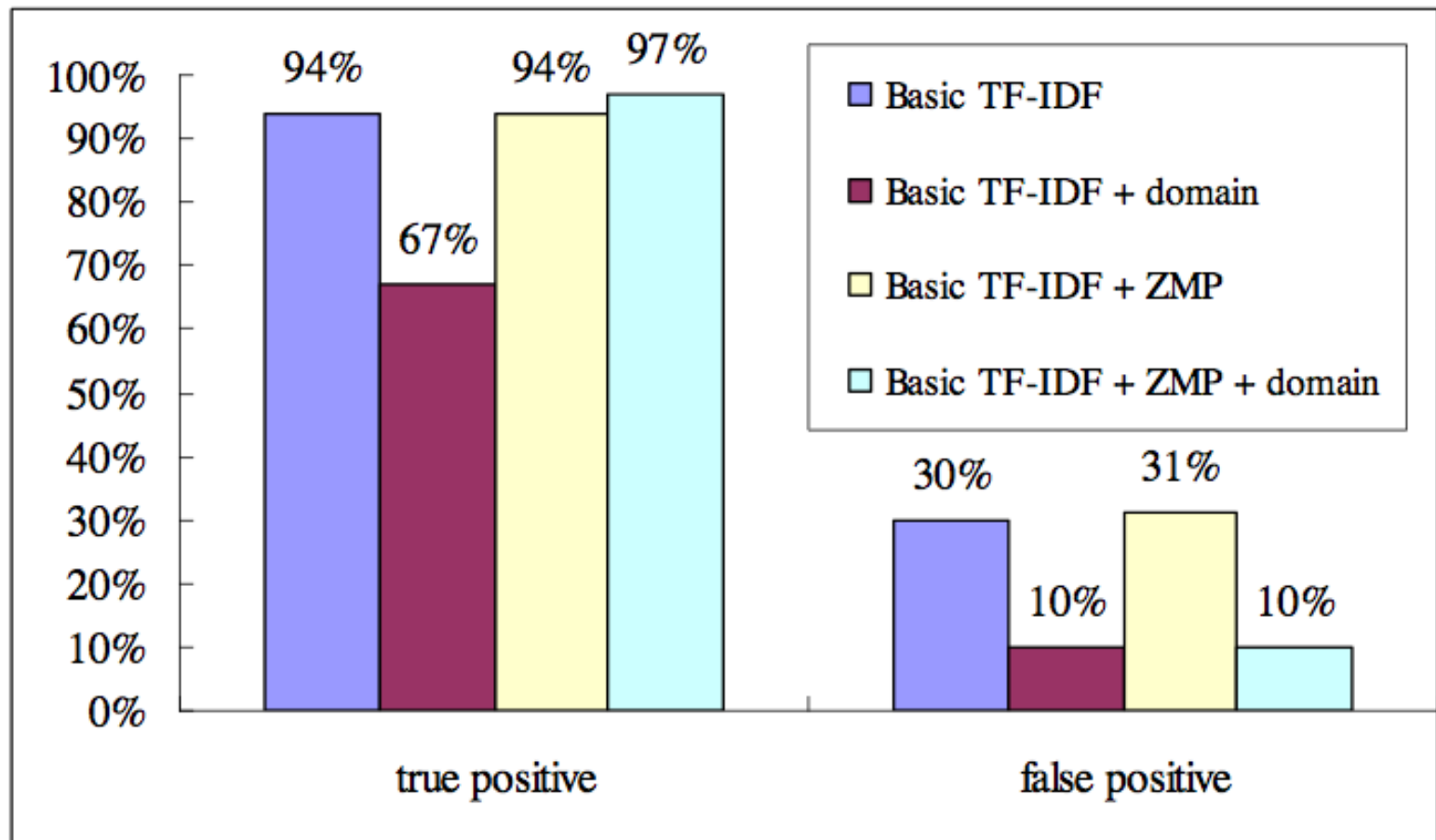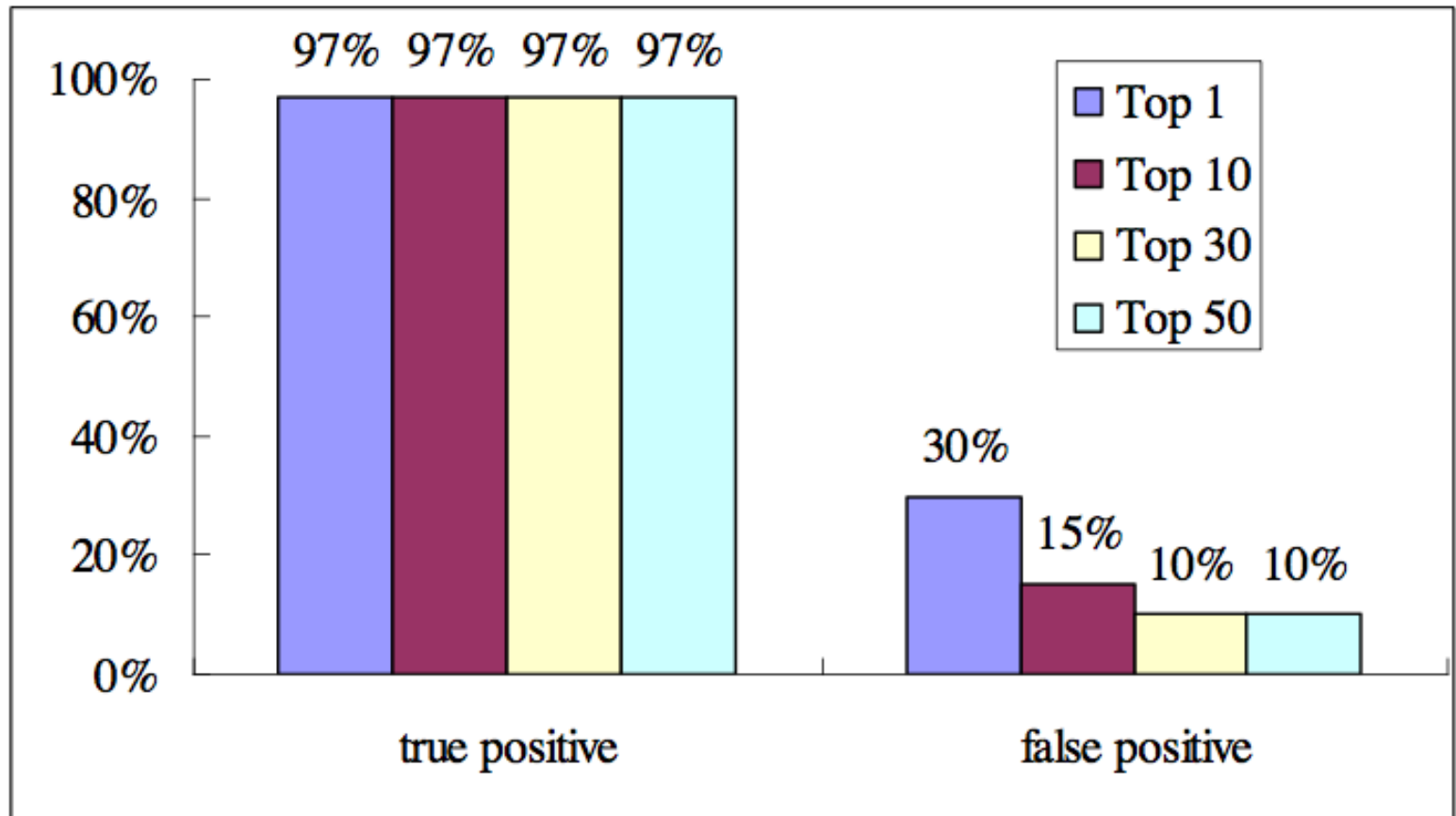  - Contains @ or – in domain name

# Other features

- Suspicious links
  - Same as suspicious URLs
- IP address as domain
- Dots in URL
  - Binary: #Dots > 5
- Forms
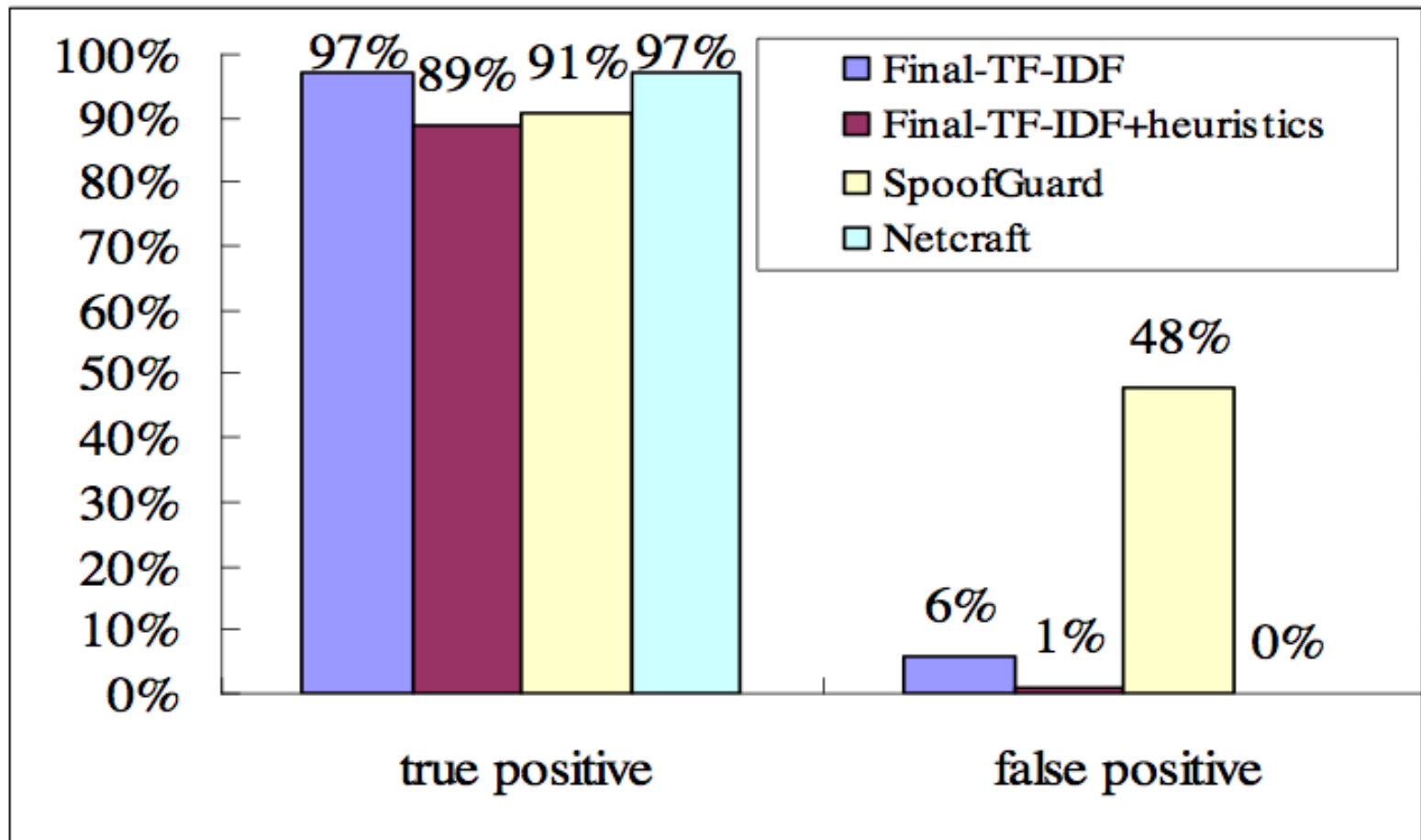  - HTML <input> tag, with text such as "credit card", "password"

# Evaluation

# Evaluation

# Evaluation

# References

- "Learning to Detect Phishing Emails", Fette et al, 2007

- "Cantina: A content-based approach to detecting phishing websites", Zhang et al, 2007