

Intro

Class imbalance is prevalent across all problems.

→ How to do anomaly detection.

- Medical
- Net. Security
-

CARLETON DATA SCIENCE

Focus on data science

↳ for Masters Dgr.

2 New Programs of D.S

→ w- PHD program

Data Day on Mar 31st!

BENJAMIN FUNG PRESENTATION (on Authorship & Malware Analysis).

Data-mine & ML for

funded by CFI.

Today focus on National Defense Sector

Authorship Analysis

feature-engineering w-embedding

→ can be used for different problems
- deep learning not so reusable.

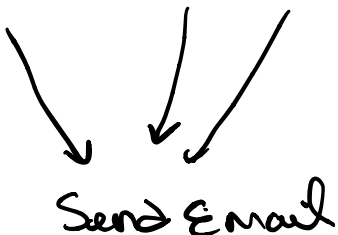
Cheque Scam

- ① Send cheque to someone
- ② Ask for return on some money
- ③ Bounce the cheque

How to prove the scammer
is the author of the
email?

Authorship Identification

Candidate X C₁ C₂



There is networking proof
but no link between
keyboards & typer

Handcrafted Stylometric Features

- ▣ lexical (caps/no caps)
- ▣ syntactic (spaces, commas)
- ▣ structural
- ▣ content-specific
- ▣ idiosyncratic features

why Models don't work:

Naive Bayes

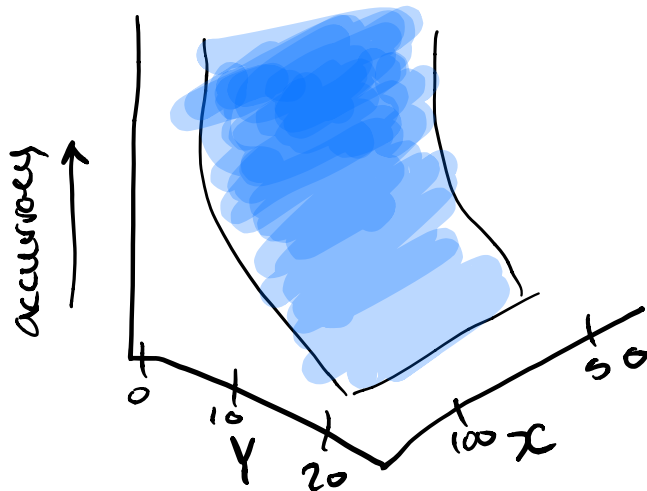
- low classification accuracy.

Decision Trees

- accuracy \approx SVM

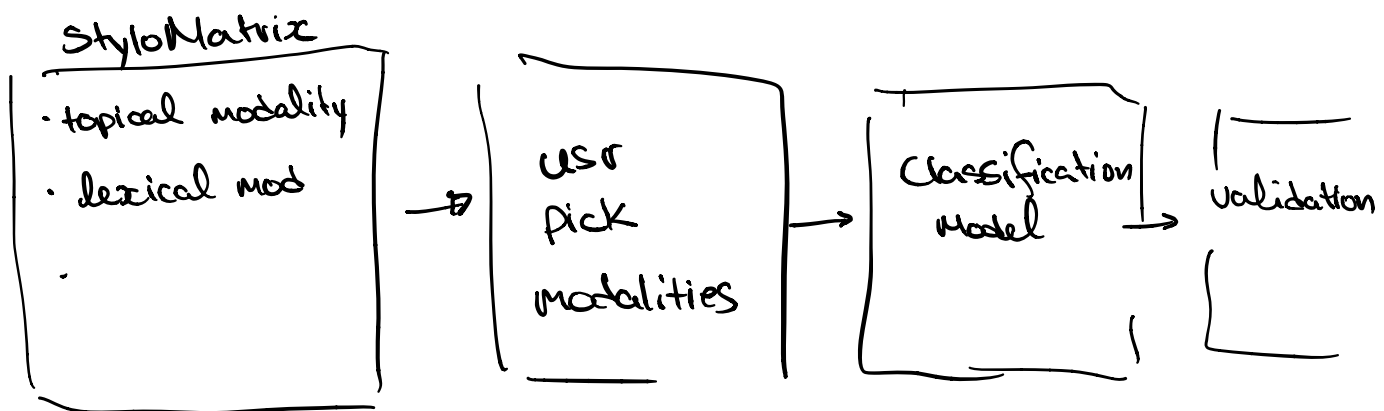
Support Vector Machine

- small data



X: # writing samples

Y: number of candidates



topic-lexic-2-vec

1. find topic
 2. find context
 3. find bias
- } find in emails.

Also is **Authorship Verification**, Auth. Characterization.

↓ age ↓ political orientation → sex.

Assembly - Reverse Engineering

Given binary code - how to understand its behaviour.

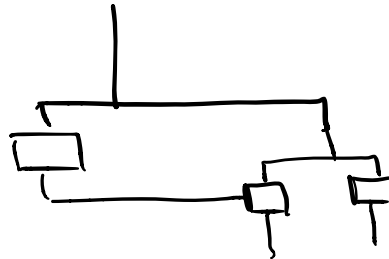
push eax

test ...

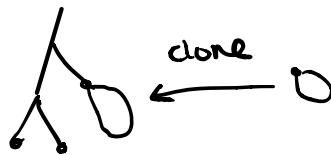
jz ----



Control Flow Graph



PROBLEM - how to find subgraph clone



Approach

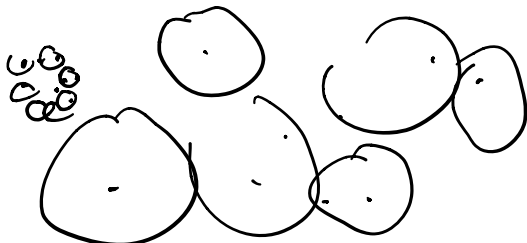
① when same
Hash code to compare

② diff operators
normalize & hash

③ almost same idiosyncratic

can't hash → do subgraph clone w- search engine.

Adaptive Locality Sensitive Hashing

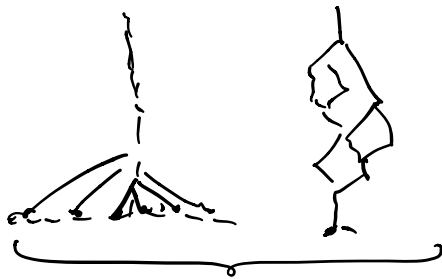


shrink scope for dense distrib.

locality sensing
hashing

- find circular reference
- expand circ. ref as much possible.

Obfuscation & Optimization

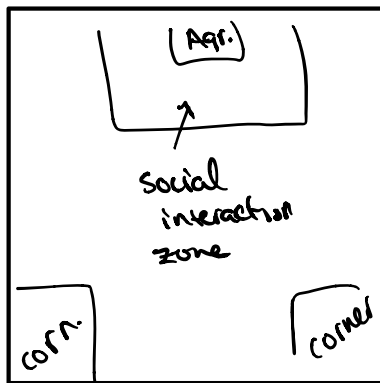


Same Code - how detect?

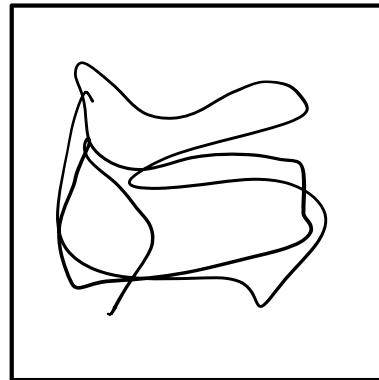
→ use ML/AI (Kamryn@).

Fighting Rats

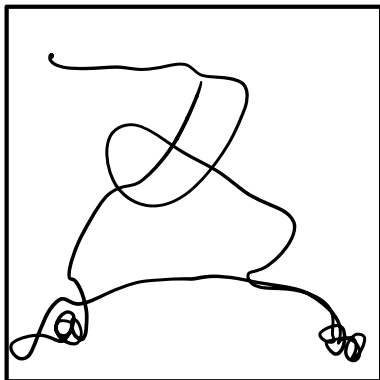
Origin



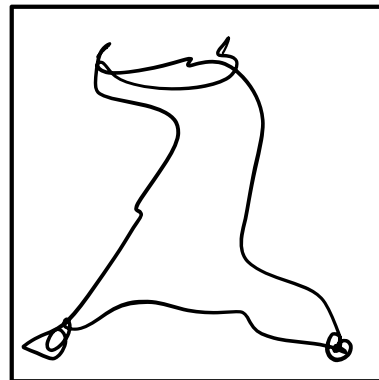
Session #1



Session #2



Session #n



(*) Could classify political sentiment using 2013 twitter data.

→ StyloMatriz. App on GitHub