

Next Generation of Fraud Detection

Sihong Xie and Philip S. Yu

IEEE CIC. Oct 19, 2018. Philadelphia, PA



Frauds: *Wrongful or criminal deception intended to result in financial or personal gain*

- Review spams
- Return frauds (Amazon, Costco, other retailers)
- Search spams (click farms)
- Fake news (Facebook and Twitter)

Stories and statistics

A single couple fraudsters caused 1.2 million loss to Amazon using return fraud ¹.

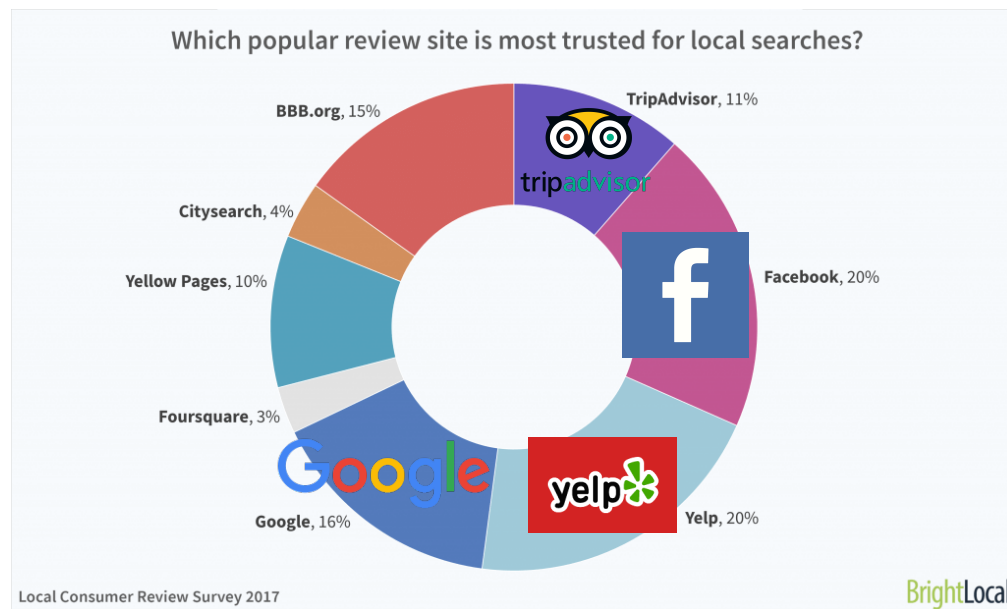
Samsung fined \$340,000 for posting fake reviews ².

1. <http://fortune.com/2018/06/05/amazon-tech-scam/>

2. <https://www.techadvisor.co.uk/feature/tech-industry/taiwans-ftc-investigating-samsung-for-defaming-htc-on-local-online-forums-3442252/>

Review frauds (spams)

Local
business
search



Review frauds:

low quality, biased,
and fake reviews
from the dishonest
brands and third-
party SEO.

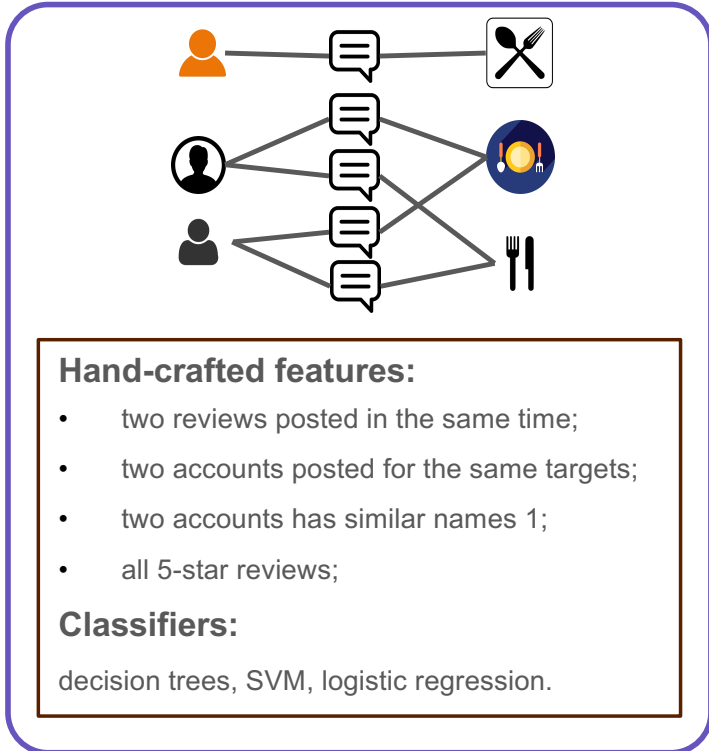
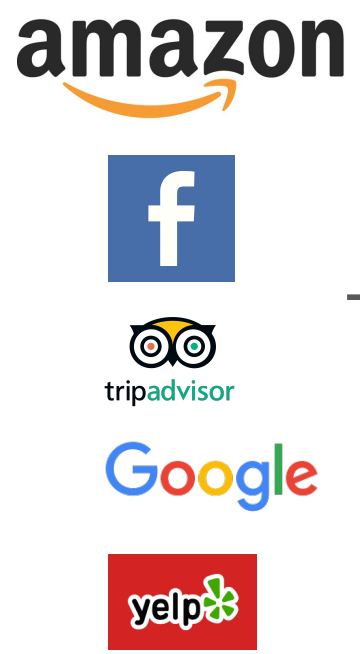
The challenges



Source: <https://www.brightlocal.com/learn/local-consumer-review-survey/>
 based on a pool of representative sample of 1,031 US-based consumers

Create a trustworthy system that spots frauds for social good.

Existing efforts: reviewmeta + spofake



Outcome + Explanations

Analysis Details

FAIL Suspicious Reviewers


Take-Back Reviewers

40% Have Previously Deleted Reviews

2 of the 5 reviewers have had at least one of their past reviews for another product deleted. This is an excessively large percentage of Take-Back Reviewers which may indicate unnatural reviews.

5.0/5
From Take-Back Review

4.0/5
in Single-Day Review

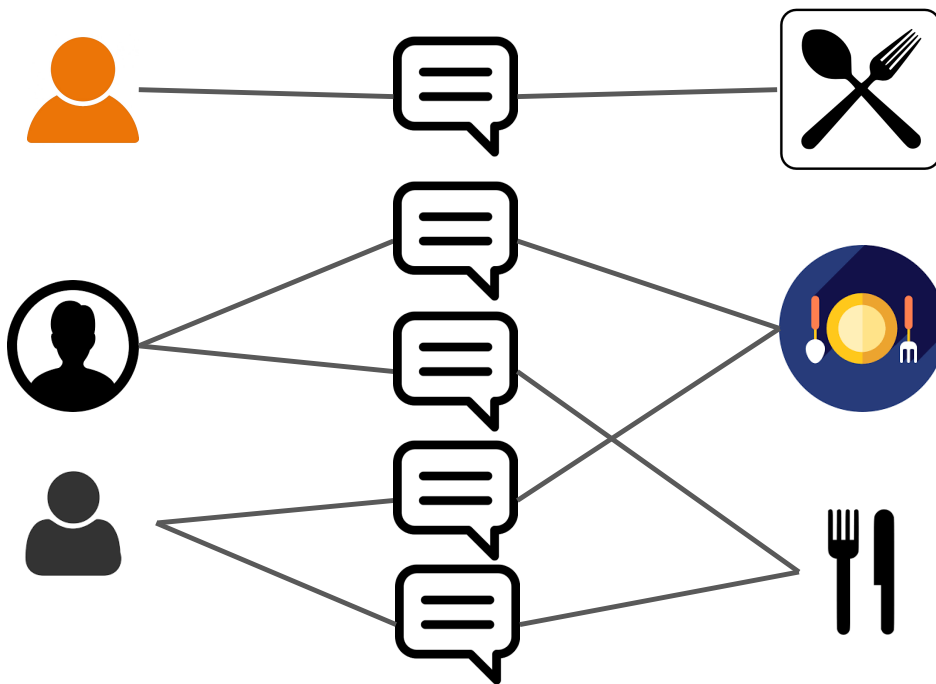


Existing efforts

Independent review fraud detectors

- <http://reviewfraud.org>
- <https://www.fakespot.com>
- <https://reviewmeta.com>

Detection pipeline

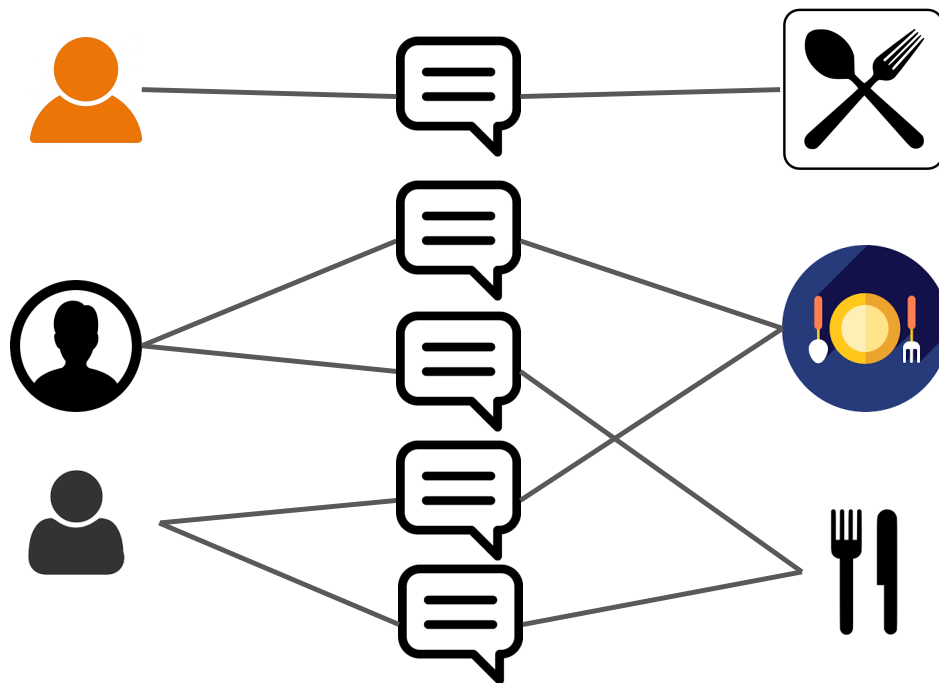


Hand-crafted features:

- two reviews posted in the same time;
- two accounts posted for the same target;
- two accounts has similar names 1;
- all 5-star reviews;
- Singleton reviews;
- near-duplicate review texts;
- near-duplicate images;
-

1. based on a true story: <http://reviewfraud.org/cloud-9-marketing-aguilar-ventures/>

Detection pipeline

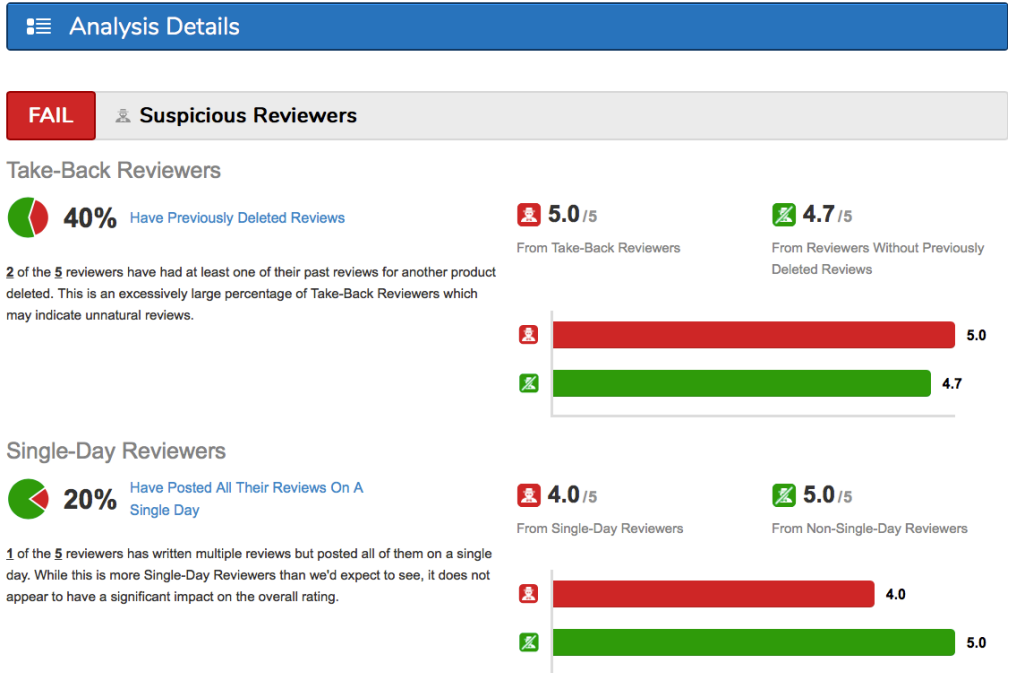


- **Supervised:**
 - decision trees,
 - SVM
 - logistic regression.
- **Unsupervised:**
 - feature histogram
 - graph pattern
 - burst detection
- **rules:**

Detection pipeline

Explain the working and outcomes

- End-users deserve to know the fact;
- To grow trustworthiness among users;
- Developers need to debug the models.



Challenges

1. **Accuracy vs. Explainability.**
2. Reactive Detection vs. Active Fraudsters.
3. Explainability vs. Security.

Click to add header

Review data

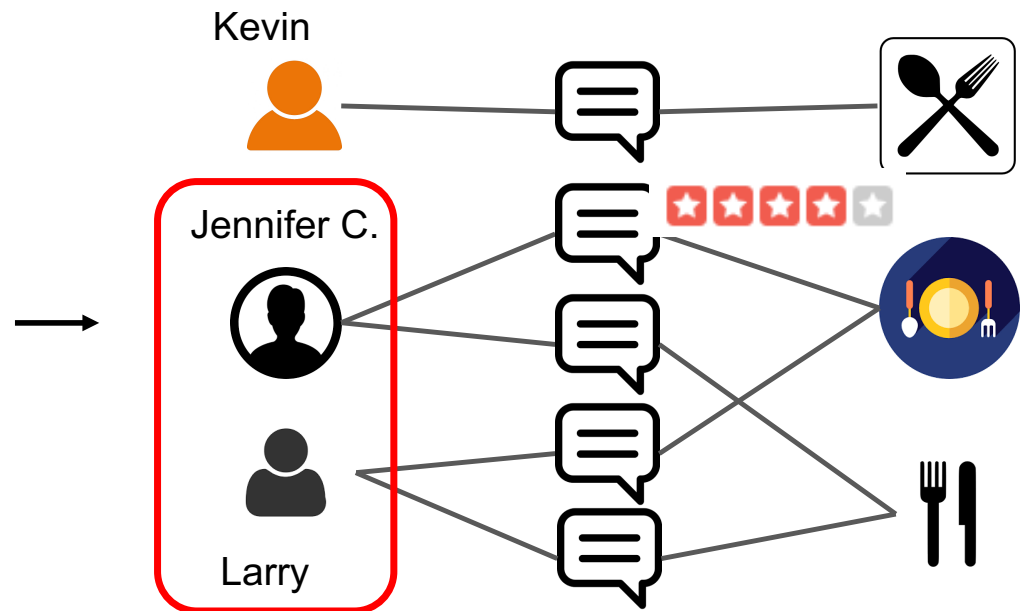


Pasta Prego

\$\$ · Italian
1502 Main St
Napa, CA 94559

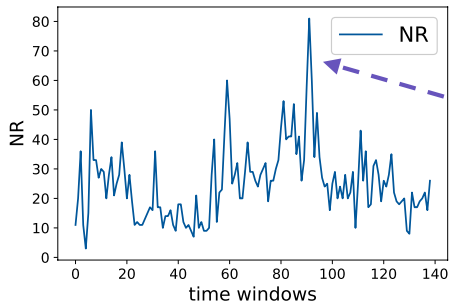
★★★★☆ 10/14/2018

This was a cozy and friendly pasta place in Napa. I loved the penne pasta which came with smoked chicken, mozzarella, basil and tomato sauce. Was pretty good although a bit salty. Everything blended in well together.



Spam detection

burst of number of reviews?

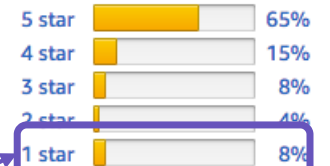


similar texts and images?
suspicious linguistic patterns?

Customer reviews

★★★★☆ 2,341

4.2 out of 5 stars



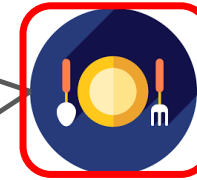
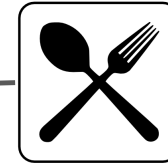
[See all 2,341 customer reviews](#)

extreme rating?

Jennifer a spammer?



Larry a spammer?



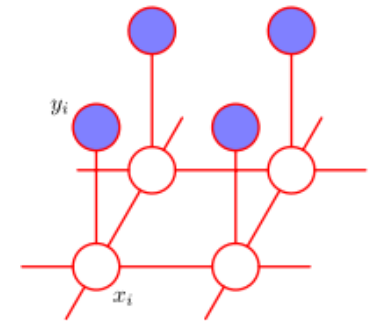
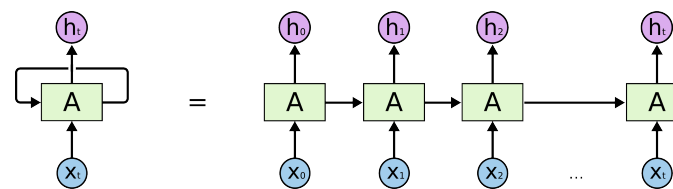
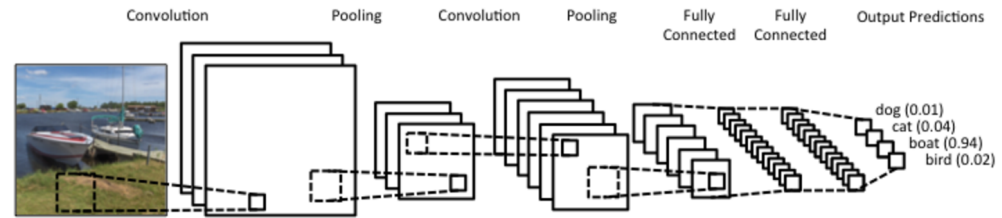
similar connectivities?

Committed spams?

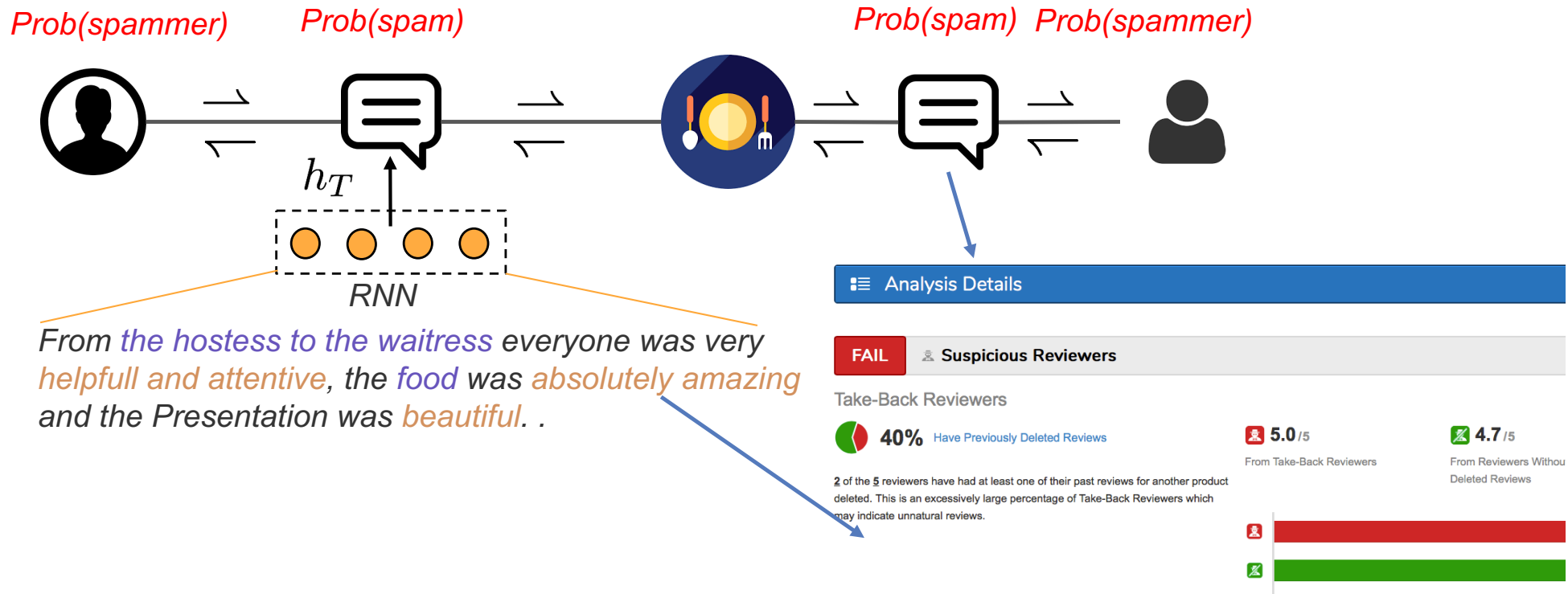
Advanced models are desired

Features that matters:

- Text and image similarity;
- Time series patterns;
- Graph connection patterns;

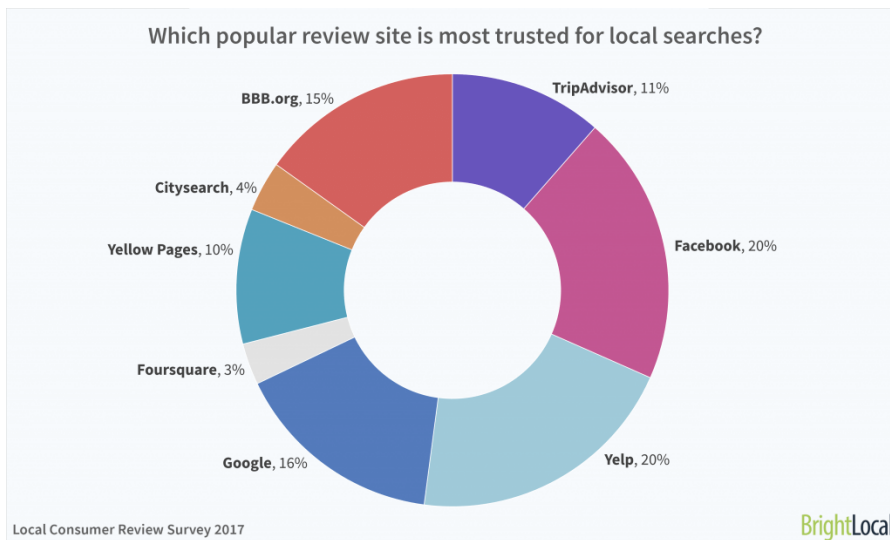


Deep structured prediction



Click to add header

Explaining complex detectors



Multiple sources of supervision

<http://reviewfraud.org>

<https://www.fakespot.com>

<https://reviewmeta.com>

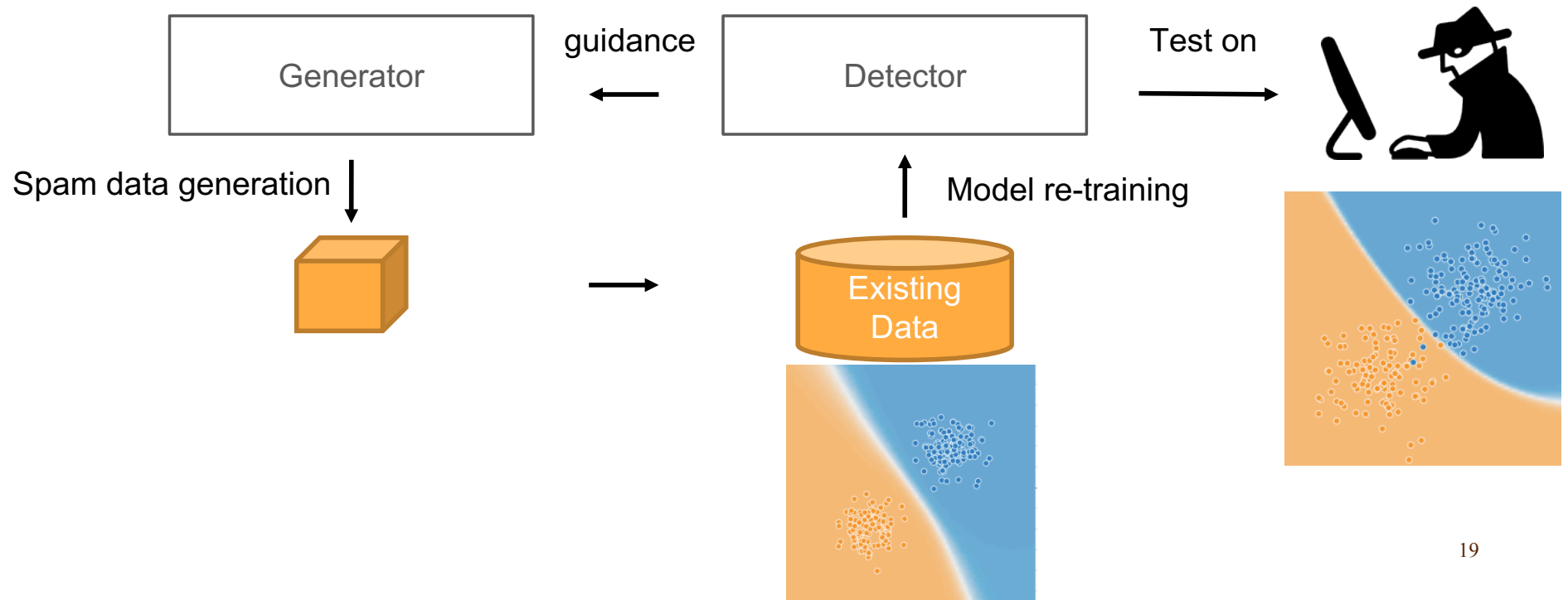
Challenge 2

Dealing with active fraudsters – it is too late when it happens.

Proactive detection is widely deployed in computer softwares and networks, auction networks.

Much more difficult in review fraud detection systems.

Proactive detection via retraining

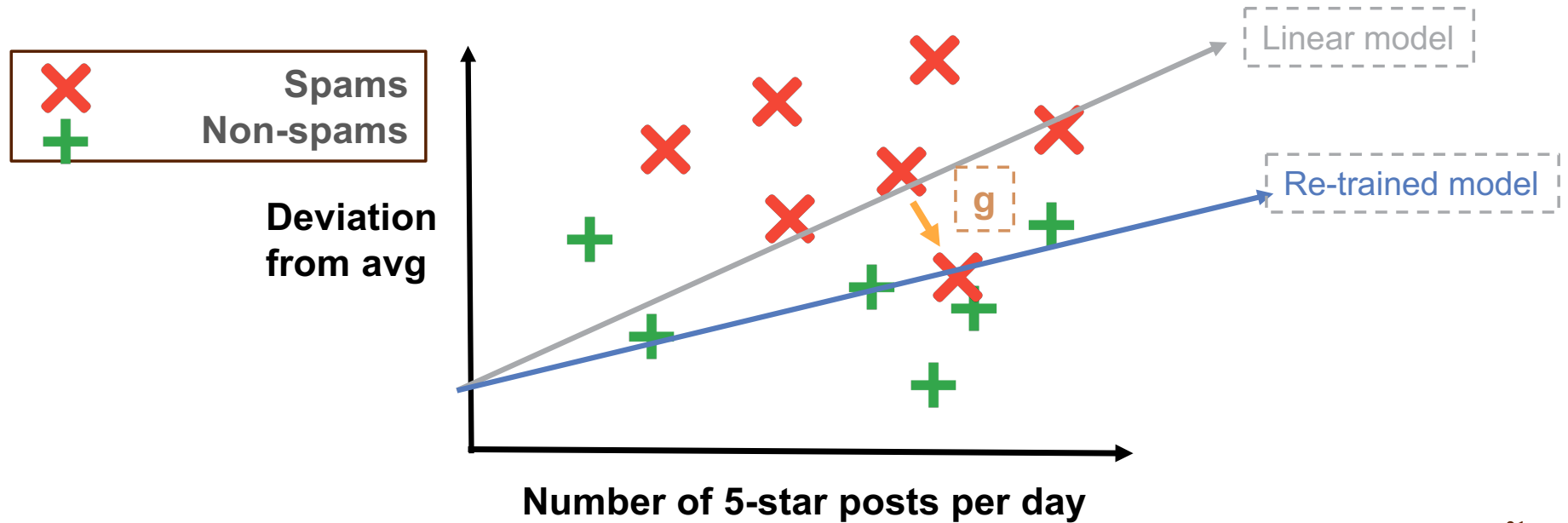


Challenges

1. Accuracy vs. Explainability.
2. **Reactive Detection vs. Active Fraudsters.**
3. Explainability vs. Security.

Challenge 2

Proactive detection via gradient attack.



Generate spams in the input space

Proactive detection via attack simulation.

- When to post a spam?
- Ratings of spams?
- Which account to post a spam?
- What contents to put down in a spam?

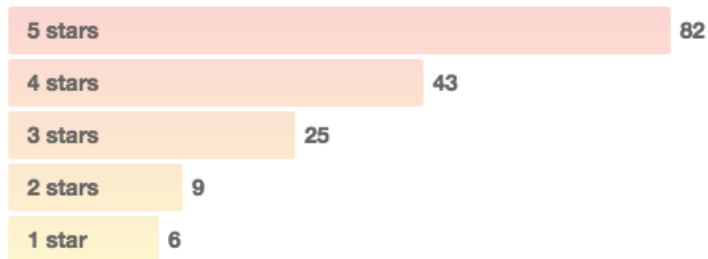
Partial solutions

How to generate spam data?

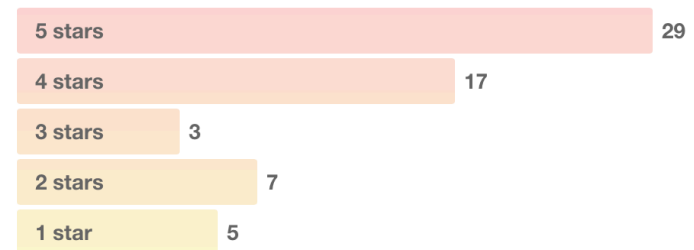
1. Maximum entropy to find the attack rating distribution.
2. Burst-avoiding techniques for attack timing.
3. Graph-based attack.
4. Review text generation.

Find evasive rating distribution

P: spammer target distribution



Q: a normal rating distribution

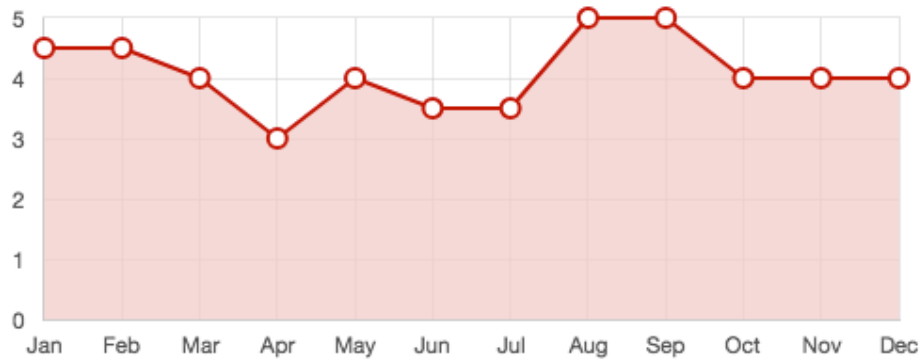


\max_P **similarity (P, Q)**
subject to **some constraints**

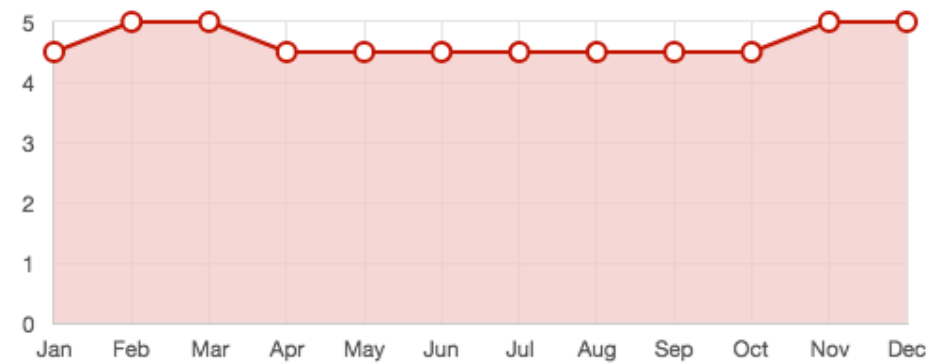
Find evasive posting frequency

Burst-avoiding techniques for attack timing.

Easy to catch



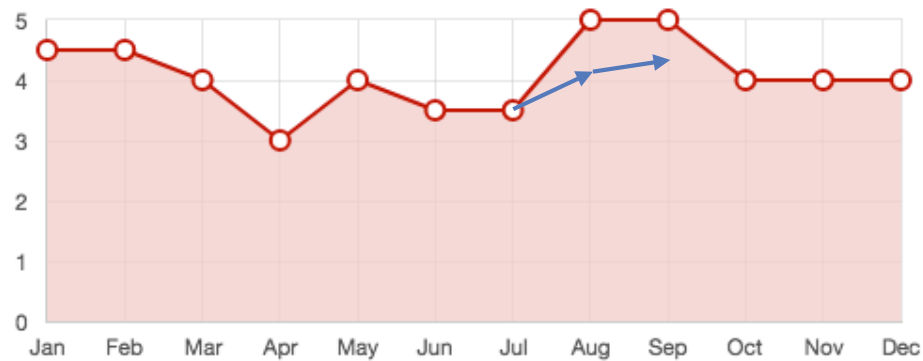
Looks more like a normal one



Find evasive rating distribution

Burst-avoiding techniques for attack timing.

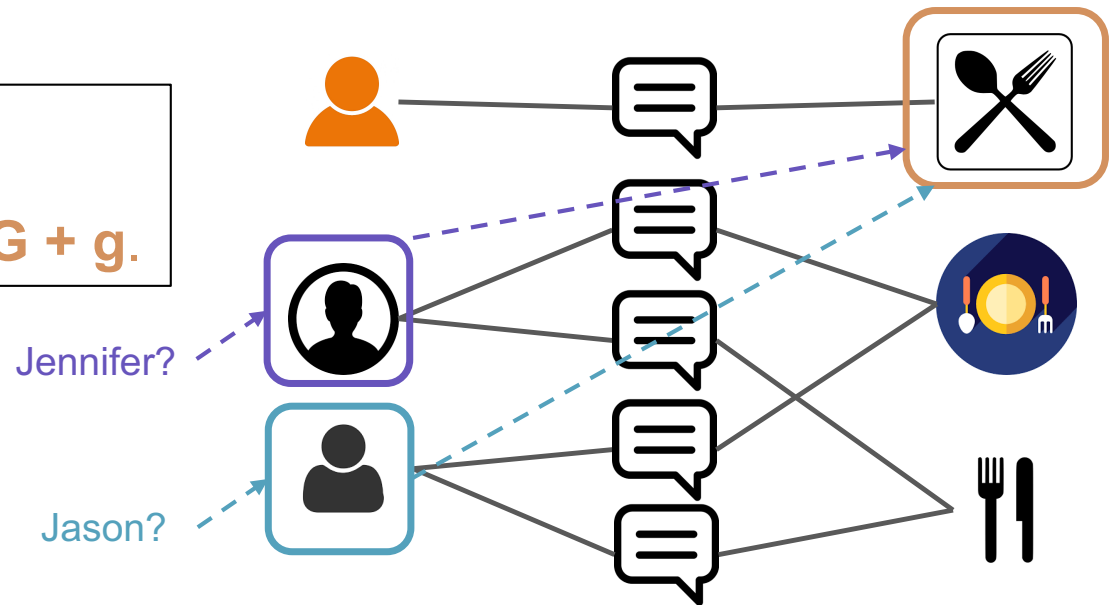
Abnormal rating dynamics



max current + future promotion
subject to smoothness constraint

Find attacking accounts

$\max_g f(G + g)$
subject to constraints over $G + g$.



Generating fake review texts

- **Crowdturfing: fraudsters are evolving to adopt more natural sounding templates and writing** ¹

Linguistics-based detectors: < 70%



- **Cheap automatic text generators can fool linguistic-based detectors** ²

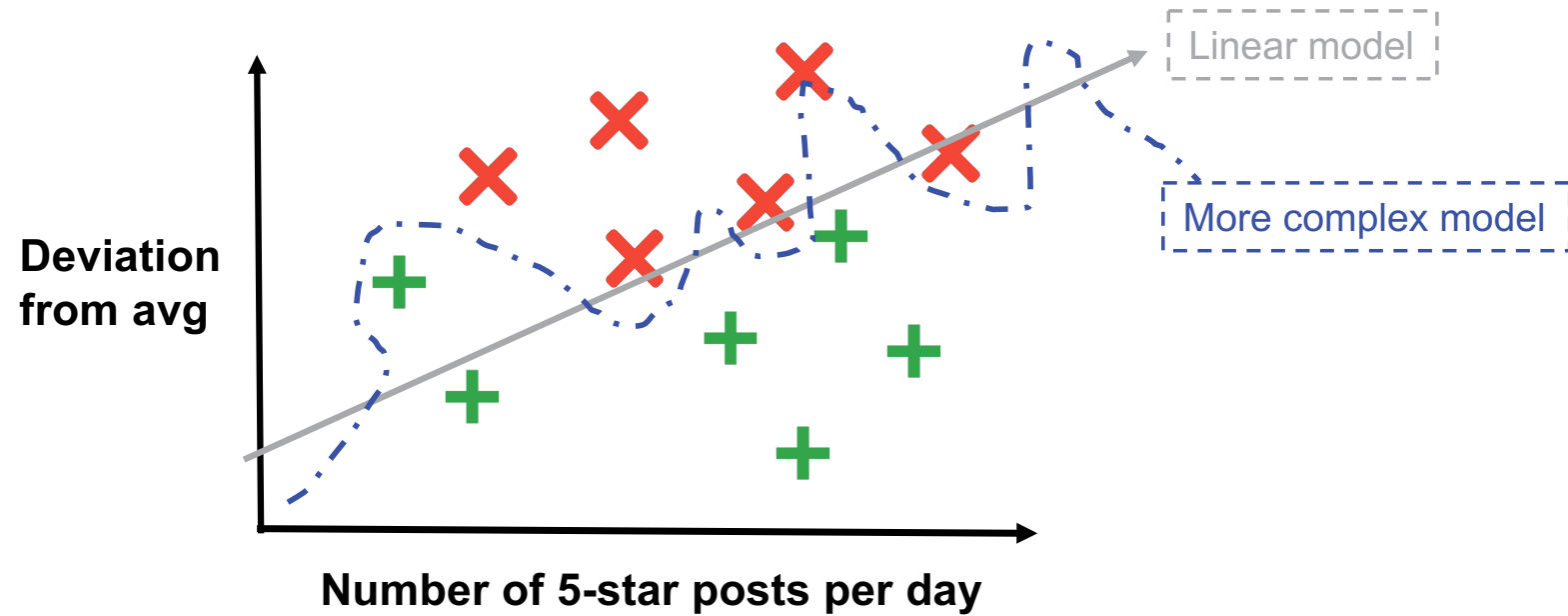
RNN is practical for short texts: 30% human detectors, 40% machine detectors F1-score

1. What Yelp Fake Review Filter Might Be Doing? ICWSM, 2013
2. Automated Crowdturfing Attacks and Defenses in Online Review Systems, CCS, 2017
3. Maximum-Likelihood Augmented Discrete Generative Adversarial Networks, ICML, 2017

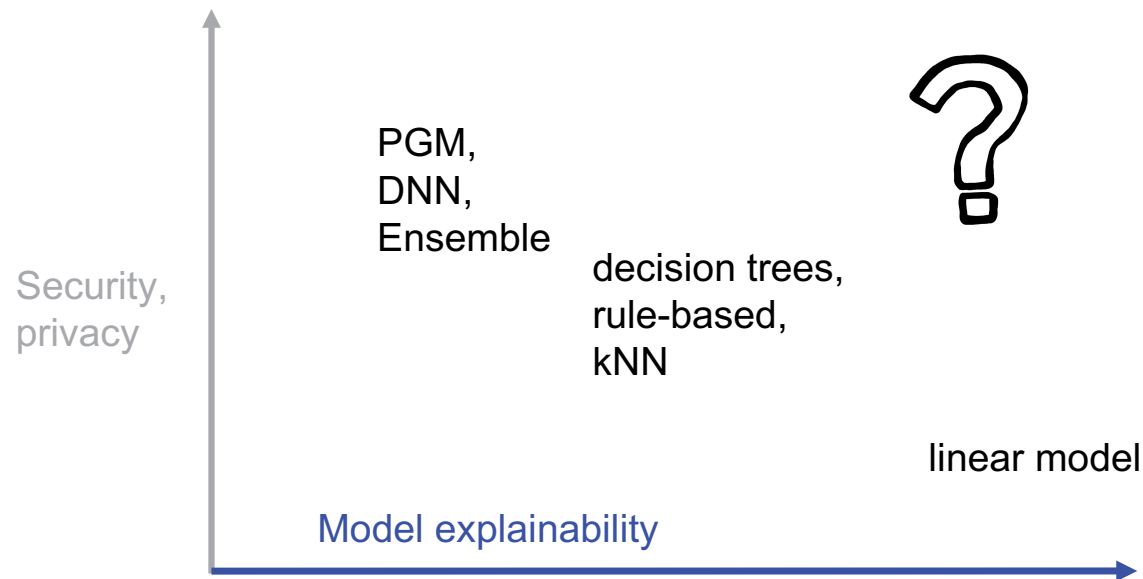
Challenges

1. Accuracy vs. Explainability.
2. Reactive Detection vs. Active Fraudsters.
3. **Explainability vs. Security.**

Explainability vs. security



Explainability vs. security




1. Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy. Explaining and harnessing adversarial examples. ICLR. 2015.
2. Florian Tramèr, et al. Stealing machine learning models via prediction apis. USENIX. 2016.

Click to add header

The reality

Spotfake.com

ReviewMeta.com
FAKESPOT REVIEW AN



[More details](#)

Least Trusted Reviews

S 5/5 **Awesome unit. Small footprint with huge quality audio**

B Awesome unit. Small footprint with huge quality audio. Now want ... [\[Go to full review\]](#)

S Jan 30, 2018

C

lr

0% TRUST

- Unverified Purchaser
- Created on a high volume day

Reviewer: David R Fleig

- Never-Verified Reviewer
- Single-Day Reviewer (posted all reviews on Jan 30, 2018)
- Easy Grader (avg. rating: 5.0)

Music

1

Thank you

