

Identifying Spam in the iOS App Store



Speaker: Rishi Chandy
Collaboration with Jay Gu

Outline

- Background
- App Store Datasets
- Observations
- Finding spam
 - supervised, unsupervised
- Future Work
- Conclusion

Background

- 2008 to now: 500k+ apps on Apple app store
- iPhone, iPad restricted to official apps only
- Popular app = millions of dollars
(or private data)



The Wrong Way: Path Uploads iOS Users' Address Books Without Permission



CHRIS VELAZCO ✓

Tuesday, February 7th, 2012

27 Comments

The Path logo, consisting of the word "Path" in a white, sans-serif font, centered on a solid red rectangular background. The letter "P" is stylized with a circular element at its top left.

The Problem: Fake Reviews

- Motivation: Fun & Profit
 - attack competitors
 - promote own apps
- Profit
 - Money
 - Private data (address book, location, ...)

Hole 9 crash every time (v3.02) ★★★★★

by Blake1775 - Version 3.0.3 - Dec 4, 2011

I want my money back! Do not buy this app.

2 out of 2 customers found this review helpful

Course updates ★★★★★

by golfer1under - Version 3.0.3 - Nov 10, 2011

Fast and accurate

2 out of 4 customers found this review helpful

Not Worth a Dime ★★★★★

by OldGuy52 - Version 3.0.3 - Jan 29, 2012

Crashes constantly, data lost, no recovery. No response from "Customer Service". I want my money back.

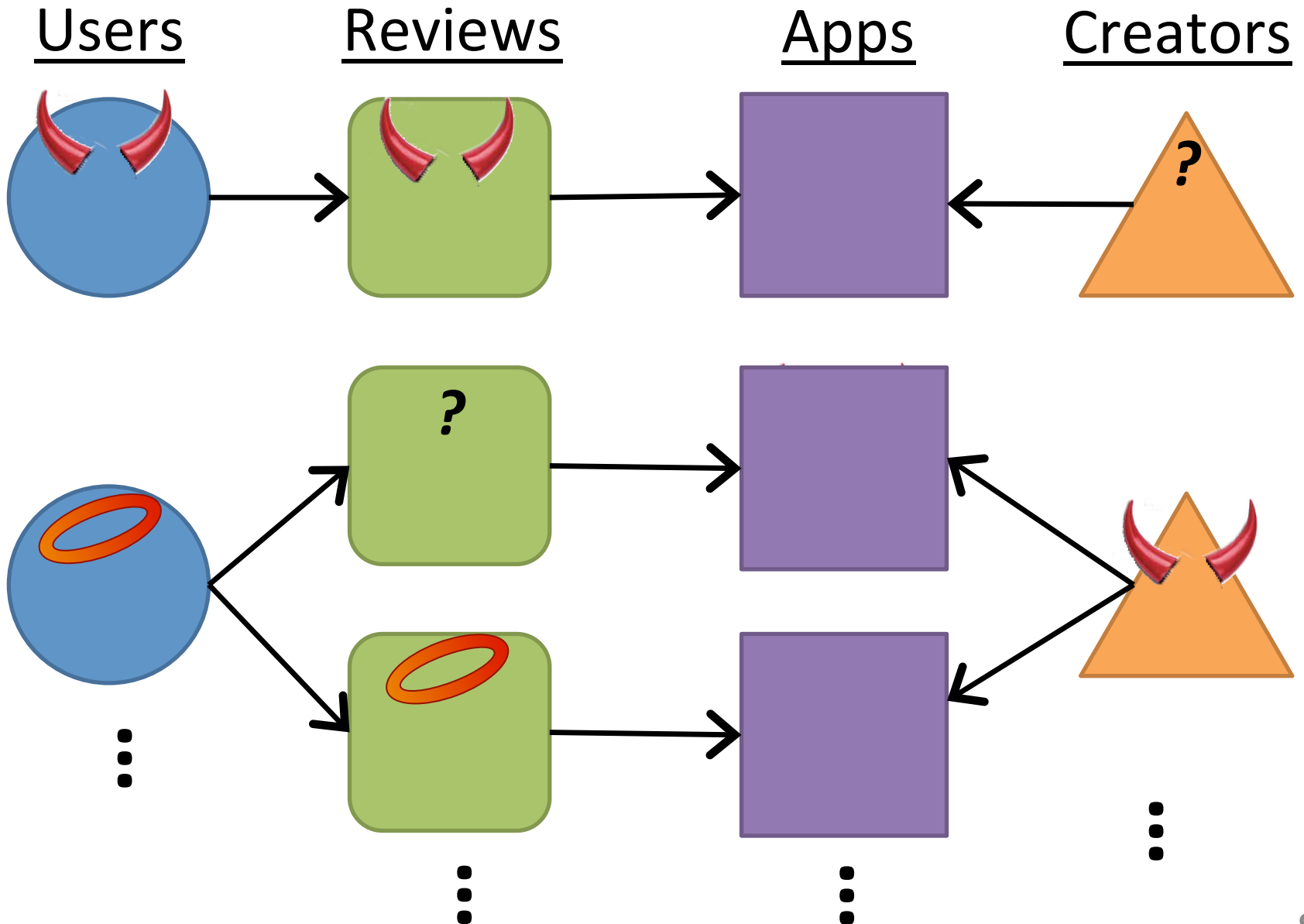
1 out of 1 customers found this review helpful

Previous Work

- None for app store specifically
- Other contexts



App Store Review Graph



Outline

- Background
- ➔ App Store Datasets
- Observations
- Finding spam
 - supervised, unsupervised
- Future Work
- Conclusion

Datasets

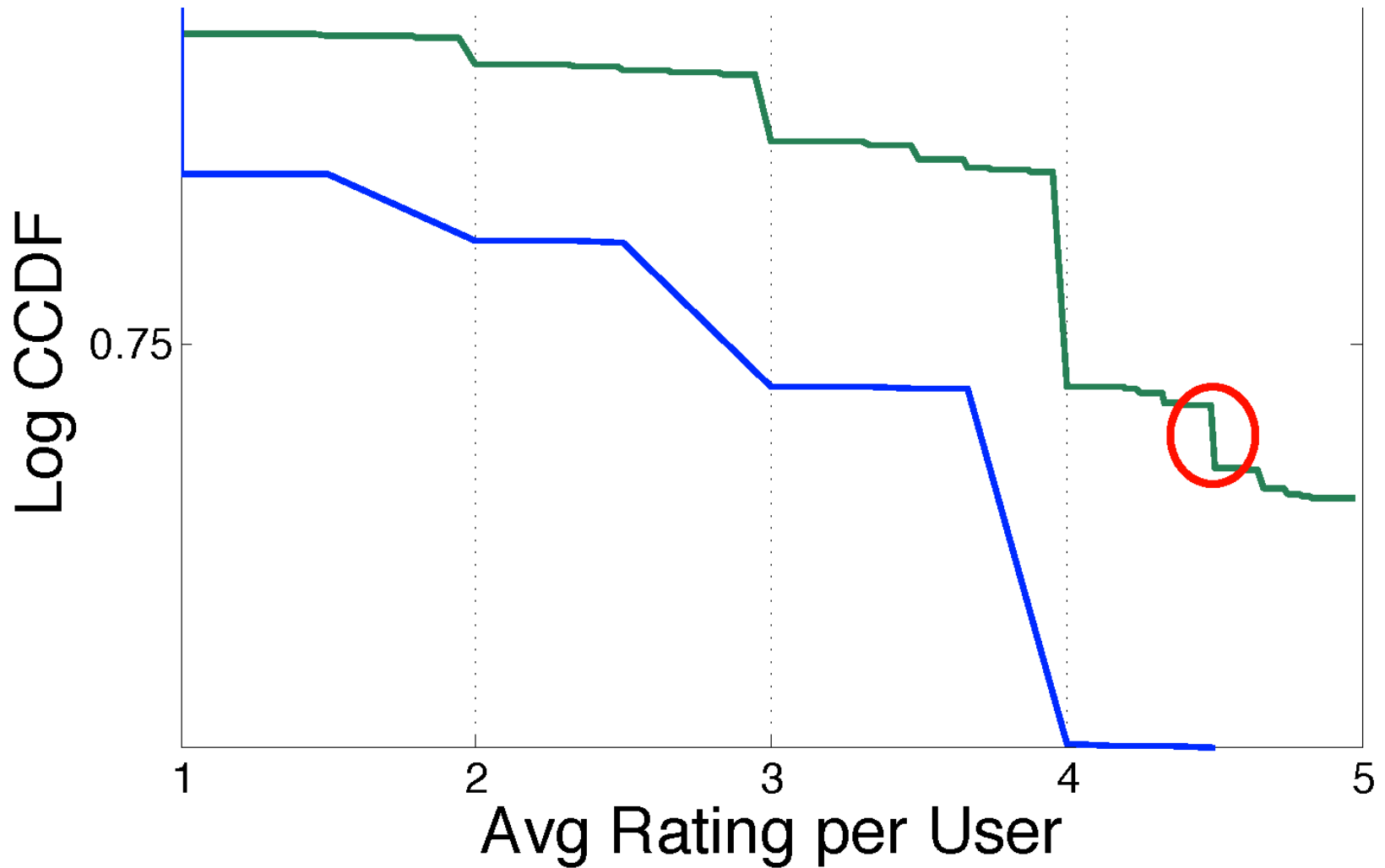
- User-Review-App-Developer relationships
- App metadata
 - name, release date, price, ...
- Reviews
 - user id, timestamp, rating, “helpfulness”

Datasets

- Top Apps (TA), Entertainment & Lifestyle (EL)

(approx)	TA	E&L	Labeled E&L
# Apps	700	2k	114
# Reviews	6 million	37k	33k
# Users	4 million	37k	33k
# Devs.	400	2k	114

Observations



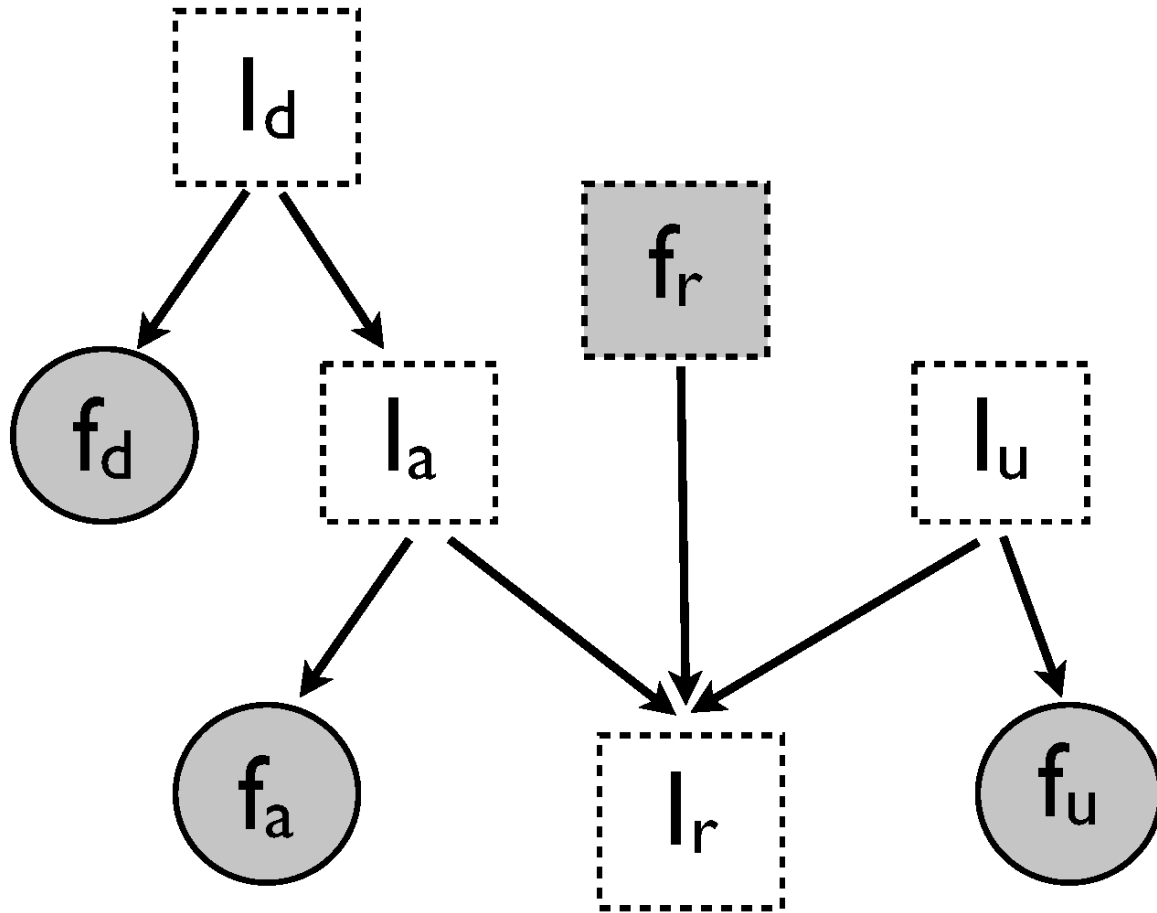
Outline

- Background
- App Store Datasets
- Observations
- ➔ Finding spam
 - supervised, unsupervised
- Future Work
- Conclusion

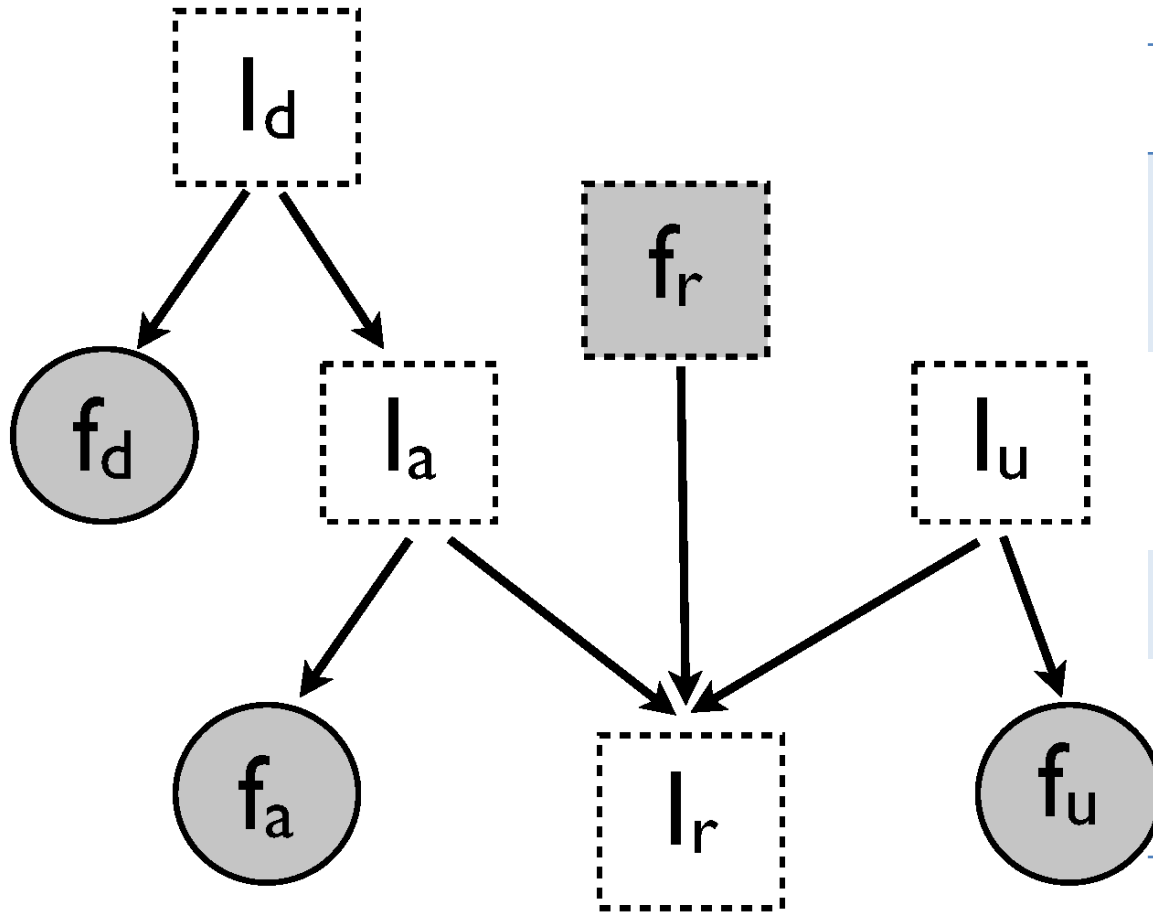
Finding Spam: Supervised

- We manually labeled 114 apps
- Decision Tree
- Simple graphical model

Simple Graphical Model



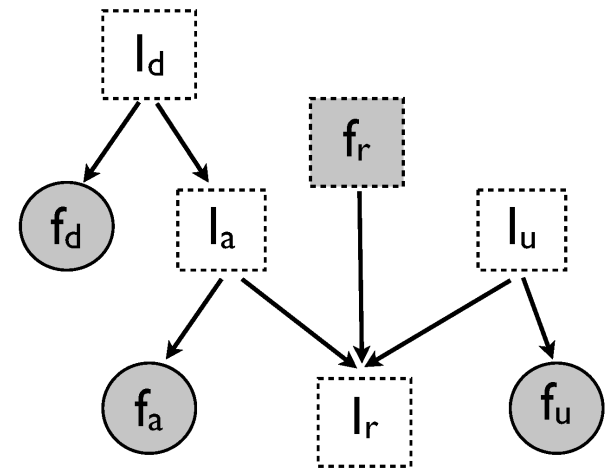
Simple Graphical Model



Node	Features
f_u	avg rating, num reviews
f_a	avg rating, num reviews
f_r	$l(\text{score})$
f_d	avg rating, num apps

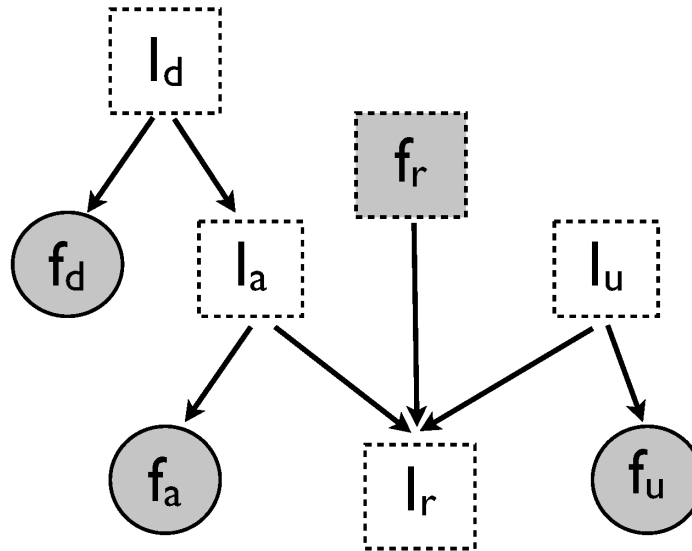
Finding Spam: Supervised

- We manually labeled 114 apps
- Decision Tree
 - 41% error
- Simple graphical model
 - 27% error

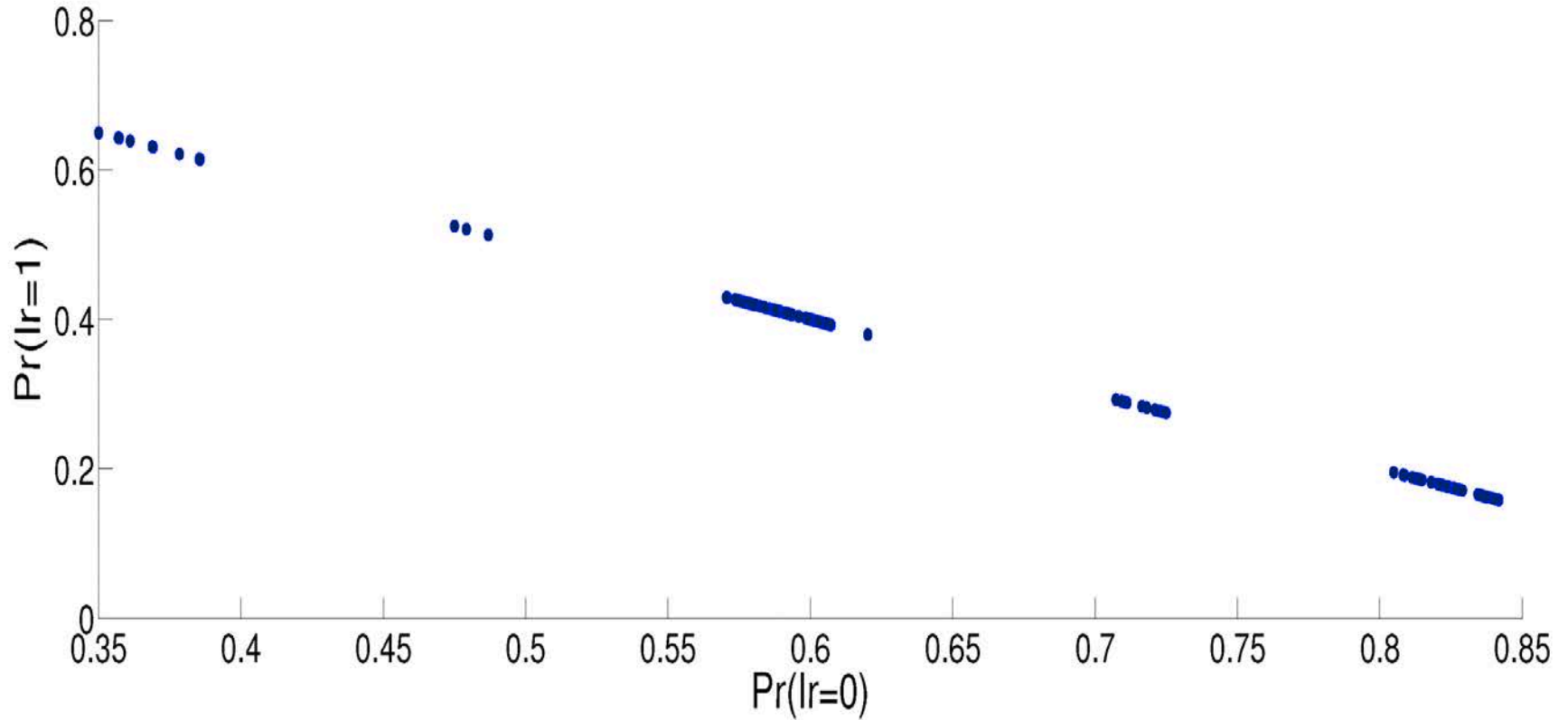


Finding Spam: Unsupervised

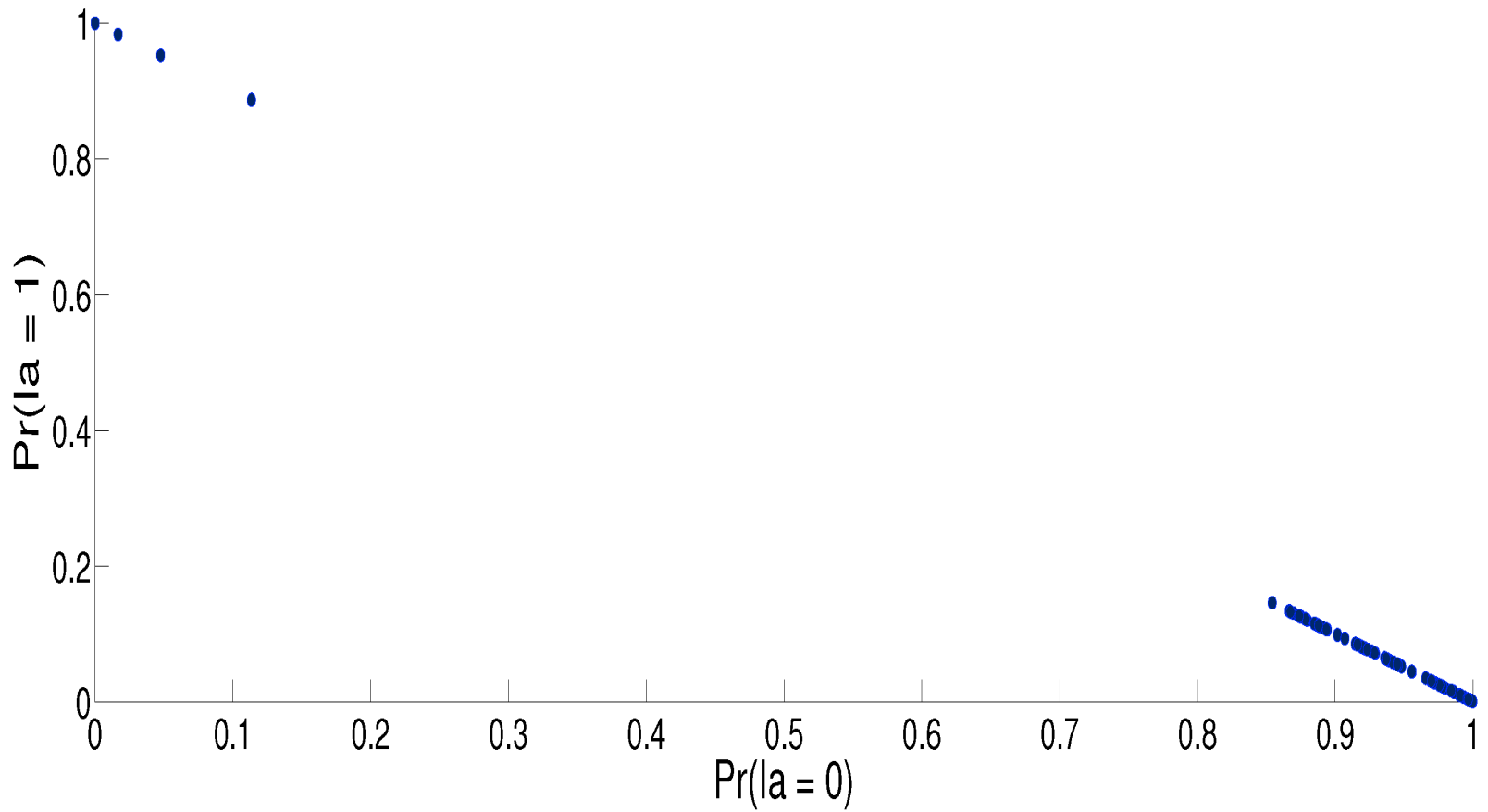
- Unlabeled!
- Cluster using our simple model



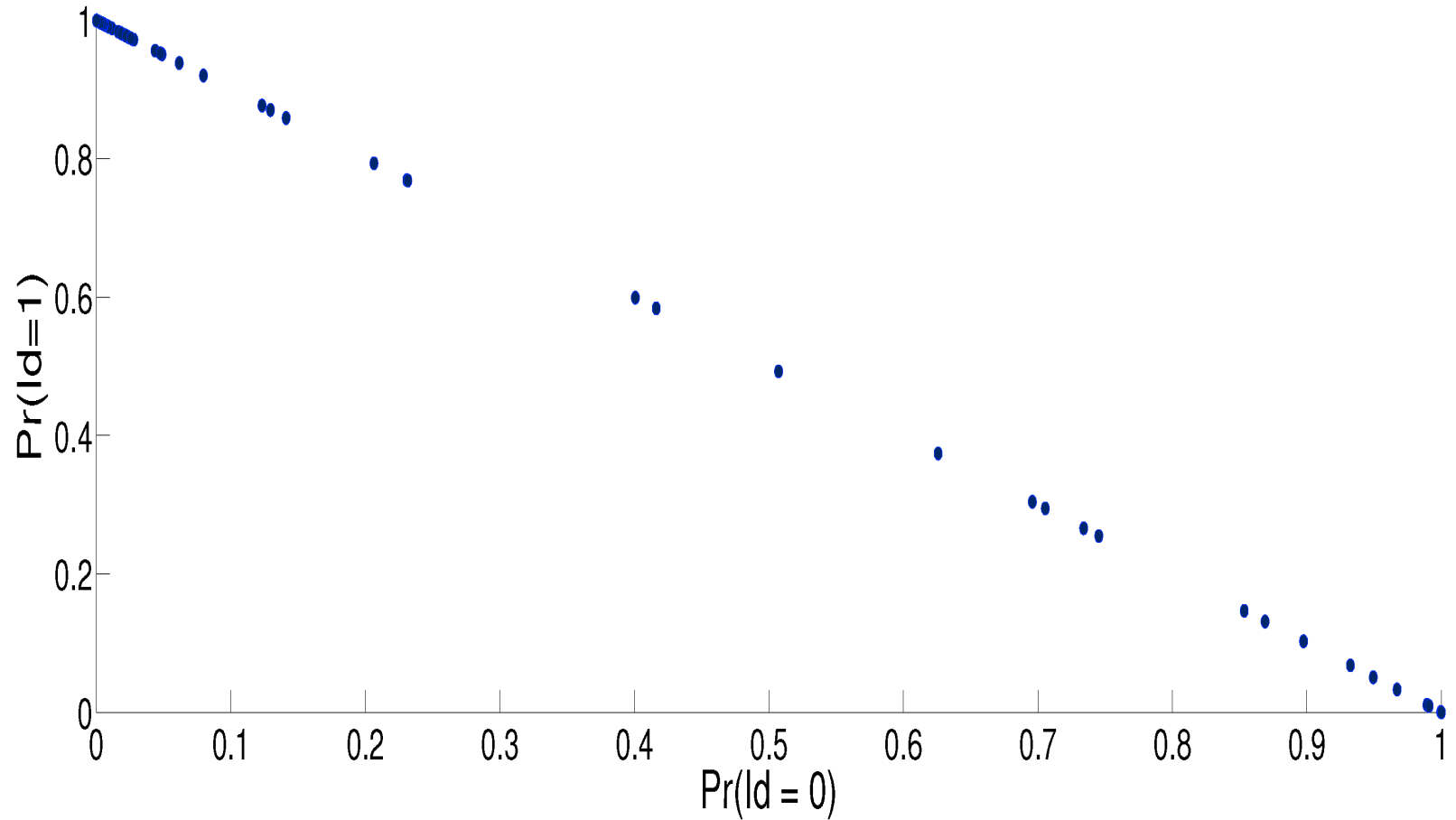
Clustering Reviews



Clustering Apps



Clustering Developers



Outline

- Background
- App Store Datasets
- Observations
- Finding spam
 - supervised, unsupervised

 Future Work

- Conclusion

Future Work

- Larger dataset
- LDA on review text
- Temporal analysis

Conclusion

- App store spam: important problem
- Baseline methods have decent performance
- Huge opportunities for future work

