Spotting suspicious reviews via (quasi-)clique extraction					bio
Paras Jain paras@gatech.edu	Shang-Tse Chen schen351@gatech.edu	Mozhgan A mazim003	<b>zimpourkivi</b> @cs.fiu.edu	Duen Horng Chau polo@gatech.edu	Bogdan Carbunar carbunar@cs.fiu.edu
	Overview			Results	
We construct a Rev users based on 1.6N We extract cliques a Reviewer Similarity Suspicious cliques o Large populations o • Yelp Scouts are use • Contriversial as Sco	viewer Similarity Graph o A reviews. and quasi-cliques from th A Graphs of sizes up to 11 users we of "Yelp Scouts" were disc rs paid by Yelp to review certa outs are artifical influences on	of >300K nese re found covered in venues Yelp	• Cliques up to Table 1: Co	D size 11 and quasi-cliques of (k, d)-graph 3,5 3,6 3,8 4,5 4,6 9-clique 112 152 1040 29 73 10-clique 22 25 290 3 13 11-clique 2 2 50 - 1 (k, d)-graph 6,5 6,8 7,5 7,8 8,3 9-quasiclique 144 649 94 351 42 10-quasiclique 44 315 33 134 12 11-quasiclique 7 100 4 33 - 12-quasiclique 1 20 - 4 - unts of (quasi-)cliques for various (k- larger k and smaller d values are mo	F size up to 12 were found
Motivation 16% of Yelp reviews are fraudulent (data from Yelp) Community trust and the livelihood of legitimate businesses are threatened by venue review fraud		<ul> <li>Manual examination of suspicious accounts (members of large cliques) revealed large populations of "Yelp Scouts"</li> <li>Yelp Scouts tend to associate with each other. See Figure 2 where the red nodes are closely clustered.</li> </ul>			

• Recent literature has explored finding fraud through analyzing link structures among fraudsters in eBay

## Objectives

- Generate review similarity graphs based on data from the Yelp Dataset Challenge
   Mine reviewer similarity graphs for cliques and quasi-cliques
   Chiques
   Chiques
   Chiques
   V
   V
- Study if such cliques are indeed suspicious
  - Is the (quasi-)clique structure
    - an indicator for suspicious behavior?
  - How to minimize false positives for flagged suspicious reviews?

## Reviewer Similarity Graphs

- The entire Yelp social graph is massive with 3M+ edges
  Reviewer similarity graph is created by extracting a (k,d)-graph
  - Edge represents two users



Clustering, Geo-aggregation

Manual inspection

Figure 1: Data Processing Pipeline



Figure 2: Scouts are tightly clustered and appear to associate with other Scouts. Graph of combined cliques in a weighted (6, 5)-graph. Larger nodes represent more reviews. Red nodes are Yelp Scouts while white nodes are regular users.

In the weighted (6, 5)-graph, 23% of users in a clique were Scouts
Scouts are very tightly geographically clustered. For example, Edinburg has large populations of past Yelp scouts. One reason for this might be that Yelp pays locals to seed reviews in new regions.
Reviews written by suspicious users are generally positive.



• 42,	reviewing ≥k venues within d days of each other Nodes represent Yelp users 153 venues and 1,125,458	6+ common venues within 5 days Figure 2: Review similarity graphs link users with similar reviewing patterns 8 reviews represented
	Clique Ex	traction
• The • Clie Ker	e maximum clique proble ques are extracted from t bosch Sparse real-world graphs allow time	em is NP-hard he (k,d)-graph using Bron- w clique finding in reasonable
	Quasi-clique	e Extraction
• Qu a ce • Per • Tak for	asi-cliques are sub-graph ertain threshold fect cliques are not comm ceaki Uno presented a gre quasi-cliques through gr	s with edge densities above non in real social graphs edy method for searching owing cliques

Figure 3: A geographic plot of Edinburgh with reviews (part of graph edges) from cliques in the (6, 5)-graph. The majority of reviews are located in a narrow region and are generaly positive.

## Conclusion

- Our research yielded favorable results in finding suspicious reviews.
  As a side-effect of the clique and quasi-clique extraction process, Yelp Scouts were found alongside possibly fraudulent results due to their artificial reviewing structures.
- The research has extensions to other crowd-sourced review sites, such as Amazon, Walmart.com or Tripadvisor and in fraud detection for electronic commerce platforms.
- We plan to incorporate other rich signals from the Yelp data to help with this, such as by analyzing review text, and the spatial and temporal relationships among reviewed venues.

