

Fraud Detection by Dense Subgraph Detection

Tong Zhao

Dense Subgraph Detection

- Given a graph $G = (V, E)$ with vertices V and edges $E \subseteq V \times V$.
- Find a subgraph S such that $d(S)$ is maximized.
- Edge density (average degree): $d(S) = \frac{|E(S)|}{|S|}$

Charikar's greedy algorithm (2000) [1]

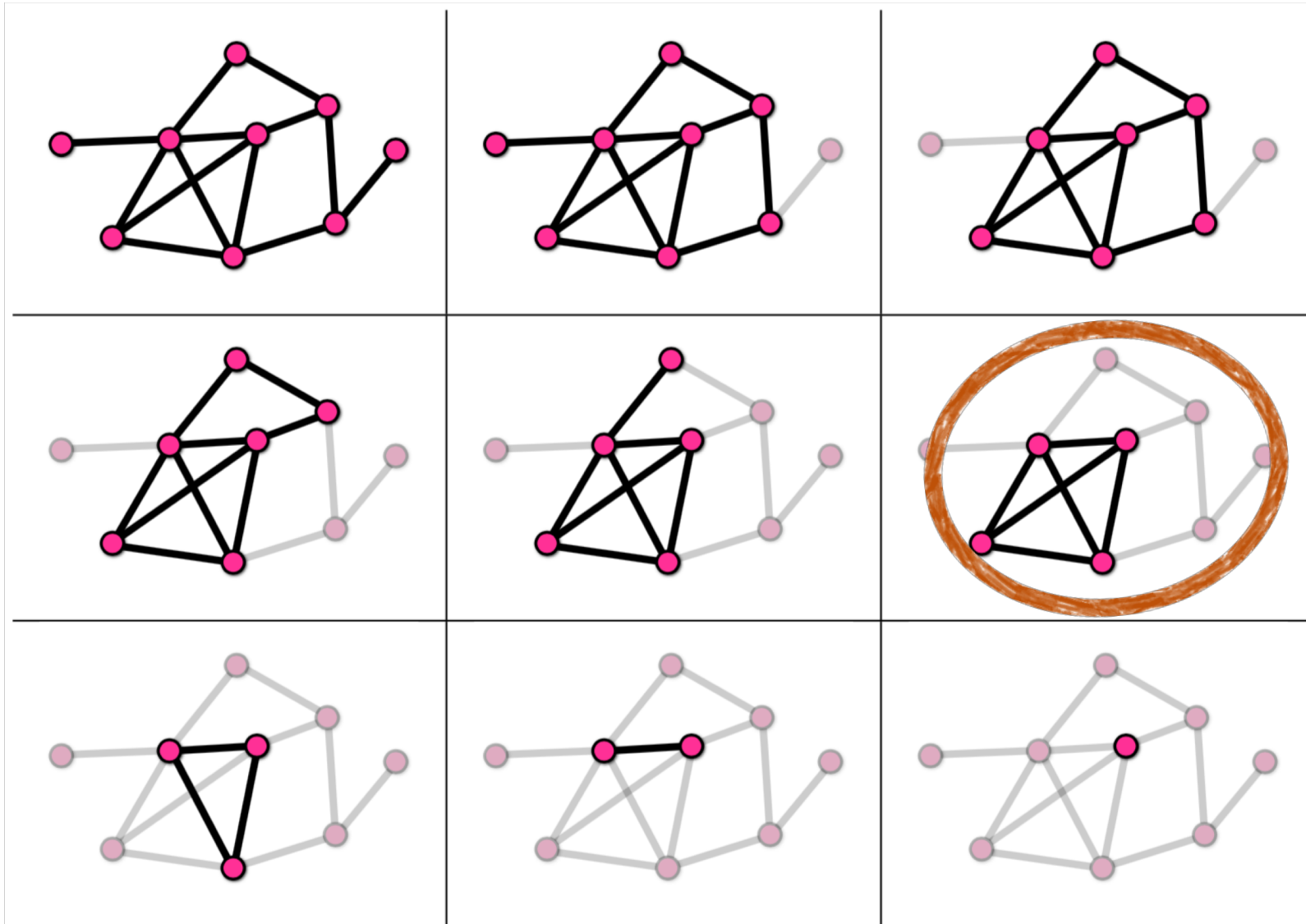


Figure from [3].

2-approximation Guarantee

- Density of the result is theoretically guaranteed.
 - Charikar's algorithm is a provable 2-approximation algorithm.

$$d(S') \geq \frac{1}{2} d(S_{opt})$$

- S' denotes the result subgraph by Charikar's algorithm.
 - S_{opt} denotes the optimal solution.
- Usually gives close-to-optimal result in real life graphs.

Possible enhancements

- Dense subgraph detection for larger graphs.
- Dense subgraph detection for dynamic graphs.

Scalability

- Observation:
- In social media graphs: $|E| \gg |V|$
- Can we only store the vertices?

Charikar's greedy algorithm (2000) [1]

```
1: procedure DENSEST-SUBGRAPH( $G$ )
2:   Input: Undirected graph  $G = (V, E)$ .
3:   Output: Dense subgraph  $S$  of  $G$ .
4:    $n \leftarrow |V|$ 
5:    $G_n \leftarrow G$ 
6:   for  $k \leftarrow n$  down to 1 do
7:      $v \leftarrow$  the vertex with smallest degree in  $G_k$ 
8:     Delete all edges incident on  $v$ .
9:     Delete all vertices with 0 degree.
10:     $G_{k-1} \leftarrow$  the remaining of graph  $G_k$ 
11:   return The subgraph with maximum density among  $G_1, G_2, \dots, G_n$ .
```

Scalability

- Most intuitive approach:
- Store only the vertices with their degrees in RAM.
- $O(|V|)$ passes.

Enhanced Algorithm [2]

- Remove a set of vertices each time.

Algorithm 1 Densest subgraph for undirected graphs.

Require: $G = (V, E)$ and $\epsilon > 0$

```
1:  $\tilde{S}, S \leftarrow V$ 
2: while  $S \neq \emptyset$  do
3:    $A(S) \leftarrow \{i \in S \mid \deg_S(i) \leq 2(1 + \epsilon)\rho(S)\}$ 
4:    $S \leftarrow S \setminus A(S)$ 
5:   if  $\rho(S) > \rho(\tilde{S})$  then
6:      $\tilde{S} \leftarrow S$ 
7:   end if
8: end while
9: return  $\tilde{S}$ 
```

Enhanced Algorithm

- For any $\epsilon > 0$, this algorithm has
 - $O(\log_{1+\epsilon} |V|)$ passes.
 - $(2 + 2\epsilon)$ -approximation guarantee.
- This algorithm can be parallelized or distributed.
 - Originally implemented in MapReduce.

Implementation

- Written in Python 3.
- ~100 lines.
- No paradigm was used.
- multiprocessing for parallelization.

Implementation

```
62     def converge(self):
63         best_subgraph = []
64         best_density = 0.0
65         while len(self.subgraph) > 0:
66             self.tmp_counter = list(self.subgraph.items())
67             q = Queue()
68             ranges = self.getST(len(self.tmp_counter))
69             processes = []
70             for i in range(10):
71                 p = Process(target=self.getBadVertices, args=(q, ranges[i]))
72                 p.start()
73                 processes.append(p)
74             for p in processes:
75                 p.join()
76             while not q.empty():
77                 del self.subgraph[q.get()]
78             self.updateDegrees()
79             current_density = self.getDensity(self.subgraph)
80             if current_density > best_density:
81                 best_density = current_density
82                 best_subgraph = list(self.subgraph.elements())
83             print('The final result subgraph contains {} vertices with density of {}.'. \
84                   format(len(best_subgraph), best_density))
85             return best_subgraph, best_density
```

Dataset

- Twitter dataset.
 - 41.7 million users (vertices).
 - 1.47 billion follows (edges).
 - 25Gb.
- Very slow due to the I/O part.
 - $O(\log_{1+\epsilon} |V|)$ passes. (43 with $\epsilon = 0.5$)

Dataset

- Twitter dataset.
 - 41.7 million users (vertices).
 - 1.47 billion follows (edges).
 - 25Gb.
- Very slow due to the I/O part.
 - $O(\log_{1+\epsilon} |V|)$ passes. (43 with $\epsilon = 0.5$)
- It works.
 - Large improvement from MemeroyError.

References

- [1] Charikar, Moses. "Greedy approximation algorithms for finding dense components in a graph." *International Workshop on Approximation Algorithms for Combinatorial Optimization*. Springer, Berlin, Heidelberg, 2000.
- [2] Bahmani, Bahman, Ravi Kumar, and Sergei Vassilvitskii. "Densest subgraph in streaming and mapreduce." *Proceedings of the VLDB Endowment* 5.5 (2012): 454-465.
- [3] Gionis, Aristides, and Charalampos E. Tsourakakis. "Dense subgraph discovery: Kdd 2015 tutorial." *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2015.