

March 21, 2019

Last week:

- unsupervised learning
 - ↳ clustering
 - nodes into groups based on the graph
- semi-supervised learning
 - a few nodes labeled
 - try to classify the rest
- today: representation learning
 - node embeddings

Construct function f

$$f: V \rightarrow \mathbb{R}^k$$

$$k \ll n = |V|$$

low-dimensional embedding

Sparse data \Rightarrow compress
to a dense representation
that preserves structure

Why?

- convenient for "downstream"
ML tasks

Today:

02 ① Hoff et al. latent space

12 → clustering

② Henderson et. al. PolX

↳ role of nodes

16 ③ node 2 vec (Grover + Leskovec)

↳ ① + ②

Example: spectral clustering

Downstream: embedding convenient
for k-means

Latent Space Models (Hoff et al. '02)

z_1, \dots, z_n be latent positions

of n nodes in a network

$$z_i \in \mathbb{R}^k$$

Idea: nearby nodes in the latent space are likely to connect



$$P((i, j) \in E)$$

$$\frac{1}{1 + \exp(\beta \|z_i - z_j\|^2 + \alpha)}$$

Makes sense for two

social reasons

① Reciprocity

Link $i \rightarrow j$

$\Rightarrow z_i, z_j$ likely close

$\Rightarrow j \rightarrow i$ is likely

\Rightarrow reciprocation is
common

② links $i \rightarrow j, i \rightarrow k$

$\Rightarrow z_i, z_j$ likely close

z_i, z_k likely close

$\Rightarrow z_j, z_k$ likely close

$\Rightarrow j$ and k likely to link

\Rightarrow triangles / clustering

Optimize with maximum likelihood

Observe graph A

$$\max_{\{z_i\}} \Pr(A \mid z_1, \dots, z_n)$$

$$\Pr((i, j) \in E) = \frac{1}{1 + \exp(\beta \|z_i - z_j\|^2 + \alpha)}$$

$$\prod_{(i, j) \in E} \Pr((i, j) \mid z_i, z_j)$$

$$\prod_{(i, j) \in E} \Pr((i, j) \mid z_i, z_j) \cdot \prod_{(i, j) \notin E} (1 - \Pr((i, j) \mid z_i, z_j))$$

Still about clustering
(similar to spectral but
probabilistic)

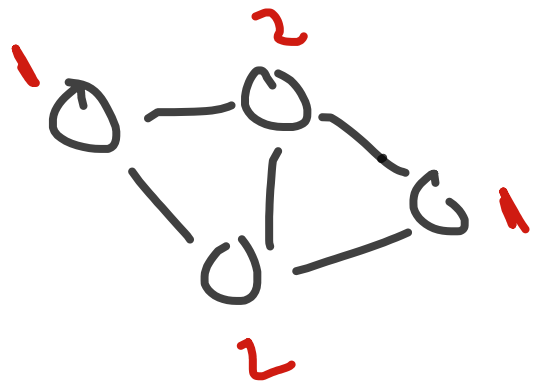
What else?

Roll X (Henderson et al. '12)

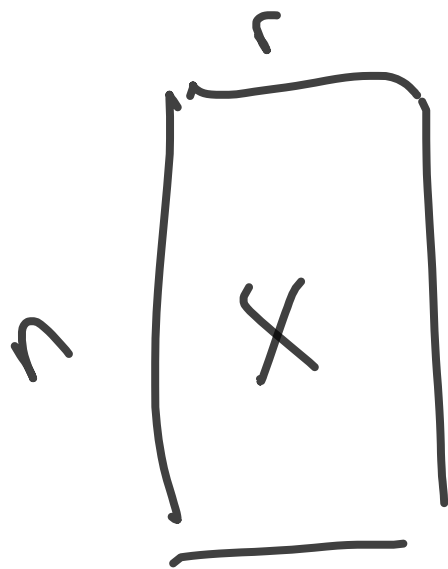
Alternative clustering approach

Collect features about nodes

- degree
- # triangles
- PageRank
- squared degree



Key idea: features of nodes



Cluster in
feature space

Dimensionality reduction of
the data matrix

- PCA
- Robust PCA
- NMF

Problem: have to manually
pick out features

Problem: clustering vs.
roles

node2vec

(Grover + Jure Leskovec 2016)

$$\max_{\{z_u\}} \sum_u \log(\Pr(N_S(u) \mid z_u))$$

sampled neighborhood

embedding

① Conditional independence

$$\Pr(N_S(u) \mid z_u) = \prod_{v \in N_S(u)} \Pr(v \mid z_u)$$

② Softmax probabilities

- $P(r | z_u)$

$$= \frac{\exp(z_v^T z_u)}{\sum_{w \neq u} \exp(z_w^T z_u)}$$

Opt. problem

$$\max_{\{z_u\}} \sum_u \left[-\log Y_u + \sum_{v \in N_s(u)} z_v^T z_u \right]$$

$$Y_u = \sum_{w \neq u} \exp(z_w^T z_u)$$

Computing Y_u can be expensive

Use "negative sampling"

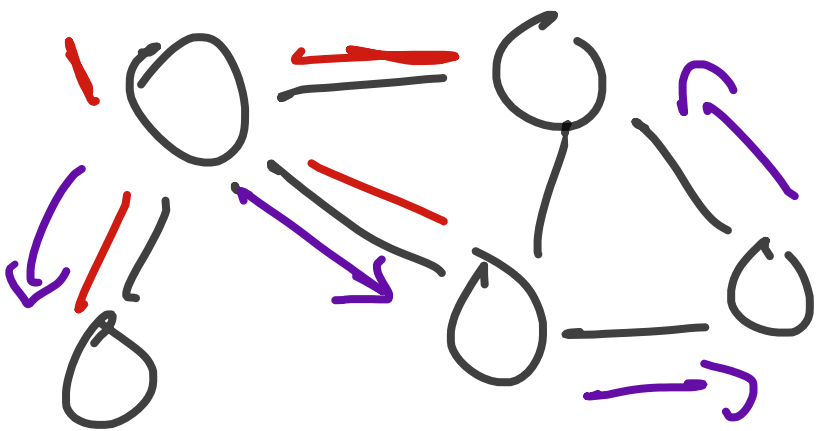
Approximating γ_u with a small random sample $W \subseteq V$

Remaining part is how we actually get $N_S(u)$?

Two sampling strategies

① BFS

② DFS



Basic idea: run BFS or DFS for a little while (truncate)

also take multiple samples

BFS - nearby nodes first

captures homophily (similar nodes are friends), similar to clustering

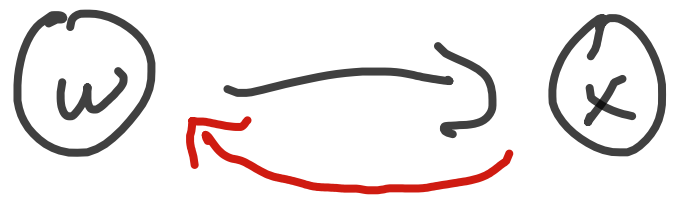
DFS - captures more macro-scale network properties

closer to role discovery?

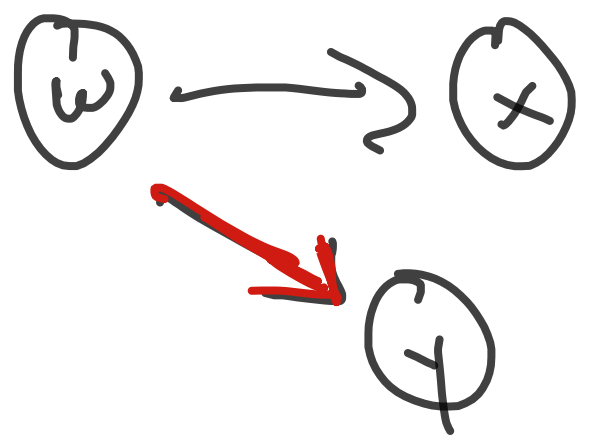
Actual strategy

Start: 

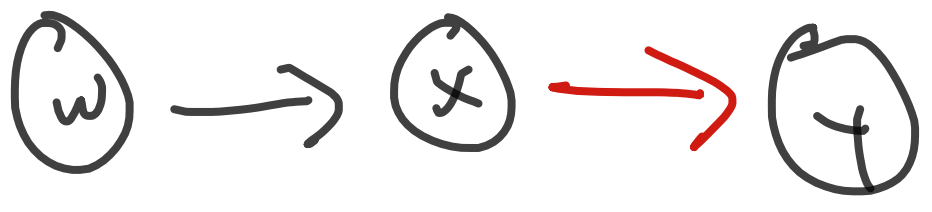
with prob. p ; return to w



with prob. q , BFS



with prob. $1-p-q$, DFS



A lot of hyperparameters

p, q control BFS / DFS

also how long "walks" or
samples should be

also how many samples
per node?

Last week / this week:

ML for graphs

Rest of networks:

- ① Small subgraph patterns
- ② Centrality / ranking
graphs &
hypergraphs

