

10/30/2020

① Missing Node Labels.

② Missing Edges.

Why do we have missing data?

- In Medicine, diff<sup>↑</sup> collection standards.
- If "Optional" field on a survey...  
<sup>data</sup> (some nodes  
opt to not  
disclose label).
- Classified Data - unreleased (privacy issues)
- Sampling → got network, but not all the labels.
  - Partial visibility
- Data Cleaning, Wrangling Issues.

Heuristic to fill in the gaps.



M.L. models that we call  
"semi supervised"

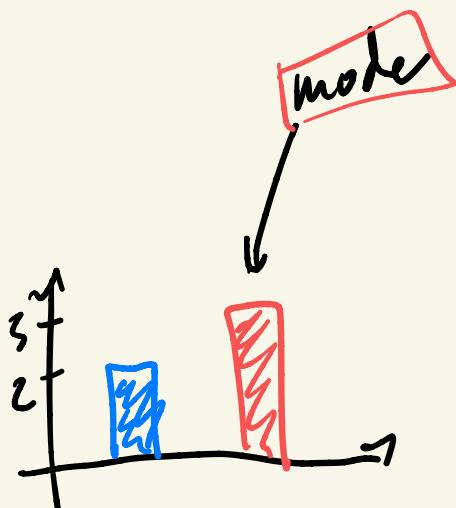
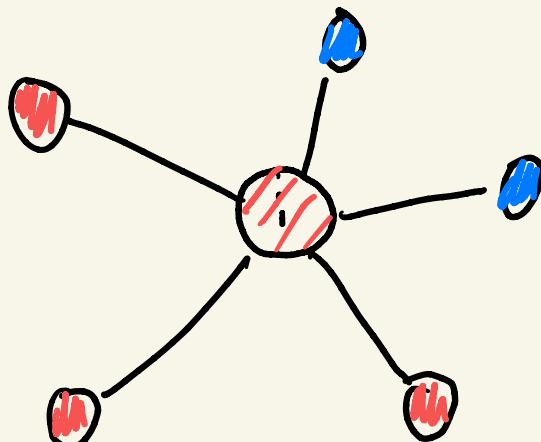
Goal: impute / generalize  
known labels → unkno labels.

# "Guilt By Association"

GBA:

- For a node  $i$  with no label, "guess" that its missing label is the mode of its neighbors' (known) labels.
- Break ties U.A.R.

Local Neighborhood.



- When is this likely to work well?  
[Network Gerrymandering??]
- Could be tough if there are too many unlabeled nodes.
- Assumes Homophily.

Test Bed for How this heuristic performs:

- ① take a network w/ full labels.
- ② chuck out U.A.R. fraction  $1-f$  of labels ( $f$  remains)
- ③ Apply heuristic
- ④ Evaluate my prediction
- ⑤ Repeat, redrawing the fraction.

## Monophily:

Tendency to link with others of only one type.

Altenberger  
+  
Ugander.

## Link Prediction

### Test Bed:

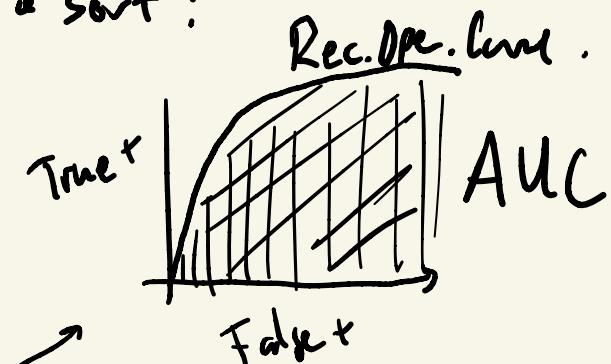
- ① Reveal only  $f$  of links
- ② Apply alg'm or heuristic:

↓  
 Ranked list, starting from most likely non-edge that we think should be an edge, and decreasing.

- ③ Evaluate by moving down the list... correct vs incorrect as a function of how far down list...

magic (3 examples,  
next slide)

- score every  $(i, j)$  in network where  $A_{ij} = 0$  (non-edge)
- sort !



## 3 methods for Link Prediction.

score(i, j)

① Chung Lu / C.M.

Degree - Product.

$$\text{score}(i, j) = k_i \cdot k_j$$

③

Shortest Path.

$$\text{score}(i, j) = \frac{1}{\sigma(i, j)}$$

length of shortest path  
between i and j.

② Normalized Common Nbrs.

Let  $N_i$  be neighbors of  $i$ :

$$\text{score}(i, j) = \frac{|N_i \cap N_j|}{|N_i \cup N_j|}$$

↑ cap                      ↑ cup

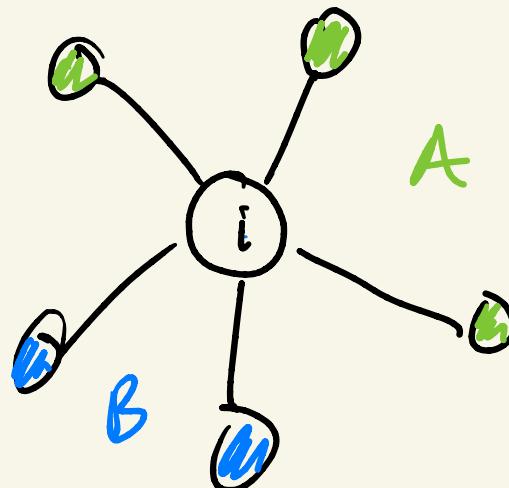
Jaccard Coeff. → Botany. Triadic Closure

- E.R. Network.  $G(n,p)$

- 2 labels.

- Fraction class I =  $\beta$

- f fraction labels



$\Pr(i \text{ gets label } A)$

$$= \Pr(k_A > k_B) + \frac{1}{2} \Pr(k_A = k_B) = \underbrace{\quad}_{\text{Sum}} \quad \underbrace{\quad}_{? \geq \beta}$$

Degree  $k_i = \text{Bin}(n, p)$

$$k_A = \text{Bin}(k_i, \beta)$$

$$k_B = \text{Bin}(k_i, 1 - \beta)$$

$$\Pr(k_A > k_B) = \Pr(k_A - k_B > 0) \leftarrow \text{sym.}$$

$$\Pr(k_A - k_B > 0) = \sum_{+} \Pr(k_A - t > 0) \Pr(k_B = +)$$

$$\sum_{+} \Pr(k_A - k_B > 0 | k_B = +) \Pr(k_B = +)$$

$$\Pr(k_A = k_B) = \sum_s [?] ?$$