

Prediction

Analysis of complex interconnected data







Common prediction tasks

- Link Prediction
- Node Classification
- Graph Classification

What is unsupervised node classification?

Examples:

https://paperswithcode.com/task/link-prediction

https://paperswithcode.com/task/node-classification

https://paperswithcode.com/task/graph-classification



° (* 1

Link Prediction

Given a graph G(V, E) predict the missing (future) links between nodes

- Modelling of the network evolution
- Predict likely interactions, not explicitly observed (e.g. terrorist network monitoring)
- Link recommendation: "friend" suggestion in social networks



° (200

Predicting Missing Links

Only a subset of edges are observed, and graph is sparse \Rightarrow searching for O(n) needles in a $\Theta(n^2)$ haystack and a random baseline would have 1/n accuracy

Define a **Score(i , j)** over unconnected pairs

local topological predictors

- Number of common neighbours
- Number of shortest paths
- Product of degree
- Same cluster
- ...

Which one corresponds to the common neighbours?



People in the Higher Education industry you may know

Slides based on lecture notes from Caluset

° (200

See all

Local Neighbourhood Scores

• Common Neighbours

$$n_{ij} = \sum_{k} A_{ik} A_{kj} = |\mathcal{N}(i) \cap \mathcal{N}(j)|$$
Neighbours of node j

 $n_{ij} = 3$ $J_{ij} = 3/6$

Example:

• Jaccard similarity

$$J_{ij} = \frac{\sum_{k} A_{ik} A_{kj}}{d_i + d_j - \sum_{k} A_{ik} A_{kj}} = \frac{n_{ij}}{d_i + d_j - n_{ij}} = \frac{|\mathcal{N}(i) \cap \mathcal{N}(j)|}{|\mathcal{N}(i) \cup \mathcal{N}(j)|}$$



Link Prediction with Jaccard Similarity

$$J_{ij} = \frac{\sum_{k} A_{ik} A_{kj}}{d_i + d_j - \sum_{k} A_{ik} A_{kj}} = \frac{n_{ij}}{d_i + d_j - n_{ij}} = \frac{|\mathcal{N}(i) \cap \mathcal{N}(j)|}{|\mathcal{N}(i) \cup \mathcal{N}(j)|}$$

We can add a small noise to break the ties randomly: $score(i, j) = J_{ij} + \epsilon$

What happens to the network if we add edges based on this? Closes triangles

What is J_{ii} for this example ?





0.091 (1/11)





| | | Jaccard |
|----------|---|-----------------------------|
| i | j | $\operatorname{score}(i,j)$ |
| 4 | 5 | 1/2 + r |
| 2 | 3 | 1/2 + r |
| 3 | 6 | 1/3 + r |
| 1 | 4 | 1/3 + r |
| 1 | 5 | r |
| 1 | 6 | r |
| 2 | 6 | r |
| 2 | 4 | r |
| 2 | 5 | r |
| 3 | 5 | r |

° (_____



Link prediction: comparing different scores

Consider the example below and two predictors, Jaccard and Degree product



Based on Preferential attachment

4 + r

4 + r

4 + r

2 + r

2+r

2 + r

2 + r

2 + r

2 + r

1 + r

Reminder.

Measuring performance for binary classification

Defined based on four values: TP (positive in both results and ground-truth), FP (positive in results but not the ground truth), FN (negative in the results but positive in the ground-truth), and TN (negative in results and ground truth)

```
Accuracy = (TP+TN) / (P+N)
Precision = TP / RP
Recall = TP / P
                                     {also called sensitivity}
F1 score = 2. Precision x Recall / (Precision + Recall)
                                                     {Harmonic mean}
Miss rate = FN / P
Fallout = FP / N
                                {also called false positive rate}
False discovery rate = FP / RP
Selectivity = TN / N
                                      {also called specificity}
False omission rate = FN / RN
Negative predictive value = TN / RN
```

 $\begin{tabular}{|c|c|c|c|} \hline & $Truth$ & Σ \\ \hline TP & FP & RP \\ \hline $Results$ & FN & TN & RN \\ \hline Σ & P & N \\ \hline \end{tabular}$

These measures depend on the threshold The tradeoff between precision and recall

> Max threshold: everything negative Perfect Fallout & Precision

Min threshold: everything positive Perfect Recall

Threshold Invariant Metrics: ROC & AUC

Receiver Operating Characteristic (ROC) as a function of prediction threshold

TPR(t) = TP(t)/P (recall, sensitivity at t) FPR(t) = FP(t)/N (fallout, false alarm at t)

Area Under the Curve (AUC) of ROC

gives the probability of ranking a random positive edge higher than a random negative edge

For **unbalanced** data use the Precision Recall Curve instead (when the number of negatives and positives differ significantly)





9

° (_____

Reminder.

Example

every candidate (i, j) is either a missing link or not

TPR = TP/P (**recall**, sensitivity) Links above threshold / Total Positives

FPR = FP/N (**fallout**, false alarm) Non-links above threshold / Total Negatives





actual G = (V, E)



AUC = \sum_{t} TPR(t) (FPR (t)- FPR (t-1))

box-rule approximation

| | | | | Jaccard | | | | degree | $\operatorname{product}$ |
|-------------------------|--|---------------|----------|---------|--------|--|----------|--------------------|--------------------------|
| | | r i | $j \mid$ | score | e(i,j) | i | j | score | $\mathrm{e}(i,j)$ |
| | 돈 | 4 | 5 | 1/2 | +r | 1 | 4 | 4 - | +r |
| | - Ei | 1 2 | 3 | 1/2 | + r | 1 | 6 | 4 - | +r |
| three | hold | 3 | 6 | 1/3 | +r | 3 | _6_ | 4 - | +r |
| unco | | ۲ 1 | 4 | 1/3 | + r | 1 | 5 | 2 - | +r |
| | | 1 | 5 | 1 | r | 2 | 3 | 2 - | +r |
| | <u>-</u> | 1 | 6 | r | | 2 | 6 | 2 - | +r |
| | - II - | 2 | 6 | 1 | r | 2 | 4 | 2 - | +r |
| | loh | 2 | 4 | 1 | r | 3 | 5 | 2 - | +r |
| | F | 2 | 5 | 1 | r | 4 | 5 | 2 - | +r |
| | | 3 | $5 \mid$ | 1 | r | 2 | 5 | 1 . | +r |
| | | | | | | | | | |
| | | Truth | | | Σ | Tru | th | | Σ |
| | results | TP = 2 | FI | P = 2 | RP = 4 | TP = | 0 | FP = 4 | RP = 4 |
| | resurts | FN = 0 | T | N = 6 | RN = 6 | FN = | 2 | TN= 4 | RN = 6 |
| | Σ | P = 2 | N | I = 8 | | P = 2 | 2 | N = 8 | |
| | TPR = $2/2 = 1$, FPR = $2/8 = 0.25$ TPR = $0/2 = 0$, FPR = $4/8 = 0.5$ | | | | | | | | |
| | | | | | | | | | |
| rue positive rate (TPR) | 10 0 0 .8 .6 .4 .2 | • • • • • • • | | | | Tue positive rate 700 - 200 - | | | |
| 1- | 0 0.2 0.4 0.6 0.8 1 0 0.2 0.4 0.6 0.8 False positive rate (FPR) False positive rate (FPR) | | | | | | | 6 0.8 1 e (FPR) | |

Jaccard coefficient, AUC = 1.00



degree product, AUC = 0.31

- 6

Topological Link Predictors

| ′D) |
|-----|
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |

There are many **alternatives** to give a score to a given pair of nodes (*i*,*j*)

Many also used as similarity measures since homophily is a strong force in link creation

See A Survey of Link Prediction in Complex Networks, 2015 & Link prediction in complex networks: A survey, 2011

Which one to use?

no one predictor is best, or worst, across all realistic inputs

See: Ghasemian A, Hosseinmardi H, Galstyan A, Airoldi EM, Clauset A. Stacking Models for Nearly Optimal Link Prediction in Complex Networks. PNAS 2020.,

Common Neighbour is a Strong Predictor



Adamic/Adar: a notable measure

Sums over size inverse of log of degree of common neighbours

$$AA(i,j) = \sum_{k \in \mathcal{N}(i) \cap \mathcal{N}(j)} \frac{1}{|\mathcal{N}(k)|} = \sum_{k} \frac{A_{ik}A_{kj}}{d_k}$$

Liben-Nowell D, Kleinberg J. <u>The link-prediction problem for social networks</u>. Journal of the American society for information science and technology. 2007

මා

° (200

Common prediction tasks

Link Prediction

Classic example: $z(i, j) = |\mathcal{N}(i) \cap \mathcal{N}(j)|$



e.g. how likely it is for them to become friend?

• Node Classification

What can be a simple predictor?

$$z(i) = \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} z(j)$$



z(i)

e.g. what is the age of a user based on his friends ages?

Label Propagation Algorithm

proposed originally for semi-supervised classification of iid data by defining a fully connected distance graph

Transition matrix
$$M = AD^{-1}, M_{ij} = \frac{A_{ij}}{\sum_k A_{kj}}$$

row normalize
$$\hat{M}_{ij} = \frac{M_{ij}}{\sum_k M_{ik}}$$

Repeat: Set the labeled data in Y and propagate labels with $Y = \hat{M}Y$

What is the main assumption here?

Homophily, similar nodes are connected



Zhu X, Ghahramani Z. Learning from Labeled and Unlabeled Data with Label Propagation.



Feedback Loop of Relations & Characteristics

There is an interplay between characteristics of nodes and their relations, a positive feedback loop derived by two social theories:

- Social Selection
 - Similarity of individuals' characteristics motivates them to form relations
 Similarity of node's attributes is a link predictors in addition to structure proximity
- Social Influence
 - Characteristics of individuals may be affected by the characteristics of their relations
 - \Rightarrow Your neighbours' attributes can reveal yours

Assortativity & Homophily

Similar nodes tend to link to each other

How to measure age homophily in a given friendship graph when you know age of every node? Similar to degree assortativity, measure the correlation of age across all edges

How to measure occupation homophily?

birds of the same feather flock together





Assortativity & Homophily: Mixing Matrix

Mixing matrix measures the **ratio** of edges connecting each pair of attribute values What indicates homophily in this matrix? How does homophily look like? **dominant diagonal**

Assortativity index

$$r = \frac{\sum_{i} e_{ii} - \sum_{i} e_{i.}e_{.i}}{1 - \sum_{i} e_{i.}e_{.i}} = \frac{Tr[E] - ||E^{2}||}{1 - ||E^{2}||}$$

This is a normalized Q-modularity assuming the characteristic partitions the graph

$$Q = \sum_{i} e_{ii} - \sum_{i} e_{i.}e_{.i} = Tr[E] - ||E^{2}||$$
$$E_{kl} = \sum_{ij} \frac{A_{ij}}{2m} \delta(C(i), k)\delta(C(j), l) = \sum_{i \in c_{k}, j \in c_{l}} \frac{A_{ij}}{2m}$$

Is there other mixing patterns?





Mixing Patterns & Structural Correlation

Mixing matrix shows the number of edges connecting each pair of attribute values

Other mixing patterns, such as **heterophily** can also be reflected in the mixing matrix

How does heterophily look like in the mixing matrix? How does correlation look like in the mixing matrix?

Low dispersion, a **dominant cell** in each row and column

How to quantify the overall correlation Does it resemble anything?

confusion matrix for pairwise cluster overlaps but with edges between pair of values

measure the total dispersion similar to clustering agreement indexes





° (200

"inclination or predisposition toward a particular thing" Proclivity

Can also be used to measure cross-correlation: correlation between your income and occupations of your friends

Example:

Facebook friendship network of a US college Color: distinct attribute values Placement : connections

| | major | gender | year | status | dorm | highschool | minor |
|------------|-------|--------|------|--------|------|------------|-------|
| major | 0.01 | 0.00 | 0.01 | 0.00 | 0.01 | 0.05 | 0.01 |
| gender | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| year | 0.01 | 0.00 | 0.25 | 0.07 | 0.07 | 0.07 | 0.01 |
| status | 0.00 | 0.00 | 0.07 | 0.09 | 0.02 | 0.02 | 0.00 |
| dorm | 0.01 | 0.00 | 0.07 | 0.02 | 0.16 | 0.10 | 0.02 |
| highschool | 0.05 | 0.00 | 0.07 | 0.02 | 0.10 | 0.31 | 0.07 |
| minor | 0.01 | 0.00 | 0.01 | 0.00 | 0.02 | 0.07 | 0.01 |

Comp 599: Network Science

correlation between the entrance year of a student and the high school his friends are from, the dormitory they are in and most strongly their entrance year



dorm major 62(76) values 23(25) values 9.94% missing 48.2% missing



gender 2(2) values 5.87% missing

status 5(6) values 0.03% missing

9(20) values 12% missing

highschool 198(2881) values 13.7% missing

Even if you don't put your information online, that information can be inferred/predicted based on what your friends reveal about themselves

| Tov | | (i) self-proclivity | (ii) cross-proclivity | | | |
|---------------------------|-------------|---------------------------------------|--|-------------|----------------|--|
| Toy Examplasi | Given: node | e's color; Guess: neigh | Given: node's color; Guess: neighbor's shape | | | |
| Examples: | i.1 | i.2 | i.3 | ii.1 | ii.2 | |
| | | | | | | |
| | homophily | heterophily | random | correlation | no correlation | |
| Q [16] | 0.75 | -0.25 | -0.25 | X | X | |
| r [17] | 1.0 | -0.33 | -0.33 | × | X | |
| PRONE _l | 1.0 | 1.0 | 0.21 | 0.67 | 0.0 | |
| PRONE ₂ | 1.0 | 1.0 | 0.11 | 0.5 | 0.0 | |
| PRONE ₃ | 1.0 | 1.0 | 0.05 | 0.33 | 0.0 | |





year



