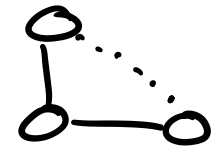
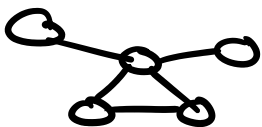
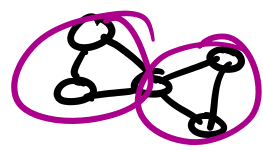
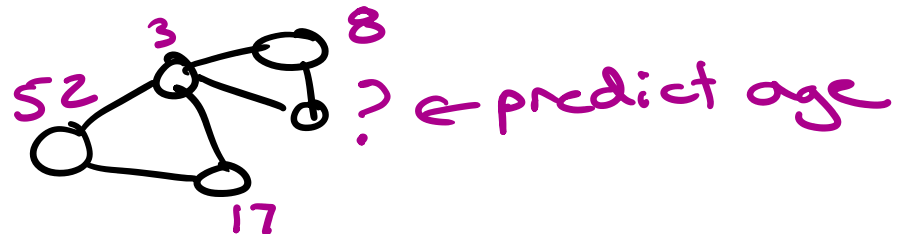


Graph mining classic topics:

1. Link prediction 
2. Centrality  ← which vertex is important
3. clustering / CO 
4. Vertex classification  ? ← predict age

Vertex Labeling problem

Given graph $G = \{V, E, W, Y\}$

V = vertex set

E = edge set

W = edge weights

Y = vertex labels

Ex: social networks

V = people

$V = \text{people}$

$E = \text{friendships}$

$W = \text{strength of friendships}$

$Y = \text{demographic information}$

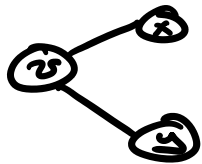
The Problem:

Given $G = \{V, E, W, Y_e\}$

Predict $Y_u \leftarrow \text{unlabeled data}$
 \uparrow existing labels

Approach: iterative classification

Features



Age of v

Gender of v

Politics of v

Avg. age of $N(v)$

Proportions of M/F/NB/etc. of $N(v)$

Generally: metadata of $v, N(v)$

Basic classification problem:

1. ... Y ... at ...

Classification problem

We have Y_e , pieces of data

We can construct features

Train classifier using Y_e , features

Predict Y_u , unlabeled data

Our general:

Construct Φ_e, Φ_u from G
(features)

Train f on Φ_e ($\min_{\phi_e \in V} \|f(\phi_e) - Y_e\|$)
(classifier)

For some # iterations:

Predict $Y_u = f(\Phi_u)$

update Φ_u

Return final Y_u

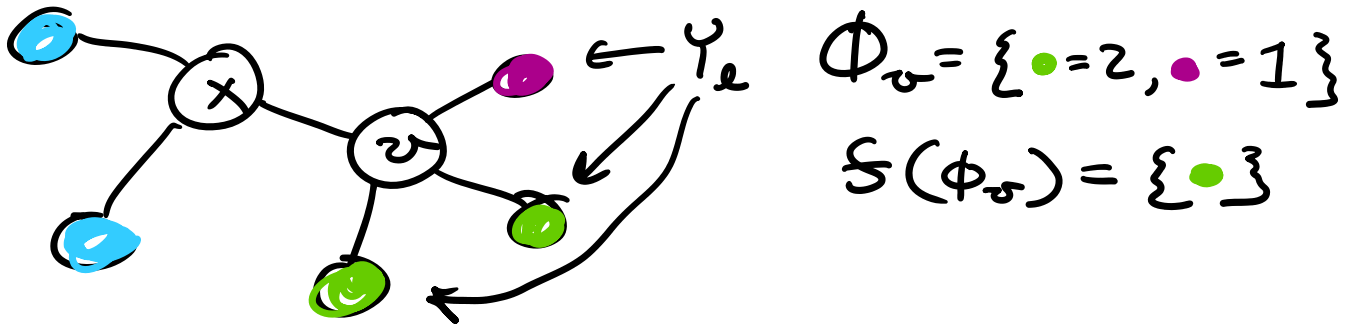
Classifier = Label Propagation

Y_e = ground truth labels

Φ_v = labels in $N(v)$

$\Phi_v = \text{labels in } N(v)$

$S = \text{return label with max count in } N(v) | \Phi$



Naive Bayes Classifier

$\bar{X} = \text{features}$

$\bar{X}_v = \text{features for vertex } v$

$\bar{X}_v = \{x_{v_1}, x_{v_2}, \dots, x_{v_k}\}$

To classify some v as class C_i

$\max_{i \in C} P(C_i | \bar{X}_v)$

↑ highest probability over all C_i given features \bar{X}_v of v

Two things first:

Bayes' Theorem $\rightarrow P(A|B) = \frac{P(A)P(B|A)}{P(B)}$

$$\text{Bayes' Theorem} \rightarrow P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$

$$\text{Chain rule} \rightarrow P(A \cap B) = P(A) P(B|A)$$

for 3+ events

$$P(A_1 \cap A_2 \cap \dots \cap A_k) = P(A_1 | A_2 \cap \dots \cap A_k) \\ * P(A_2 \cap \dots \cap A_k)$$

$$P(C_i | \bar{x}_v) = \frac{P(C_i) P(\bar{x}_v | C_i)}{P(\bar{x}_v)}$$

Note: constant

$$P(C_i) P(\bar{x}_v | C_i) = P(C_i \cap \bar{x}_v)$$

$$P(C_i \cap \bar{x}_v) = P(C_i \cap x_{v_1} \cap x_{v_2} \dots x_{v_k})$$

$$\hookrightarrow P(x_{v_1} \cap x_{v_2} \cap \dots \cap x_{v_k} \cap C_i)$$

chain rule again

$$= P(x_{v_1} | x_{v_2} \dots x_{v_k} \cap C_i) P(x_{v_2} \dots x_{v_k} \cap C_i)$$

←
Chain rule recursively

$$= P(x_{v_1} | x_{v_2} \dots C_i) P(x_{v_2} | x_{v_3} \dots C_i) \dots P(C_i)$$

Naive Bayes assumption:

→ all x_j are independent

$$\rightarrow P(x_{v_1} | C_i) P(x_{v_2} | C_i) \dots P(C_i)$$

$$= P(C_i) \prod_{j=1}^k P(x_{v_j} | C_i)$$

$$Y_u(\bar{x}_v) = \max_{i \in C} P(C_i) \prod_{j=1}^k P(x_{v_j} | C_i)$$

$P(C_i)$ = ratio of C_i labels

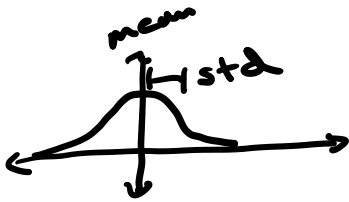
$$P(C_i) = \frac{|C_i|}{|V_e|}$$

$P(x_{v_j} | C_i)$ = over all C_i labels,
how often feature
 x_{v_j} shows up

→ if discrete → easy
if continuous

if continuous

we need to assume a given distribution of feature values



We can assume e.g., a normal distribution for a given feature

$P(x_{uj} | C_i)$ = the probability that the feature value x_{uj} is sampled from some distribution (mean, variance) calculated from known labeled x_{uj} features $u \in U_i$

Random Walks

Idea: the probability of some $v \in U_1$ assuming some label C_i is the probability of a random walk from v ends on some $u \in U_2$ with label C_i

$$P = D^{-1}W$$

P = transition prob. matrix

D = diagonal degree matrix

W = weighted adjacency matrix

$P_{ij} \in P \rightarrow$ prob. of random walk from $i \rightarrow j$

$P = D^{-1}W$ only gives one step of a random walk

$\lim_{t \rightarrow \infty} P^t =$ gives us steady-state walk probability distribution

However: what if we want to stop our walk if we land on a labeled vertex?

We can modify our P

$$P_i = e_i \quad \leftarrow \text{identity} \quad \text{if } i \in V_L$$

$$P_j = (D^{-1}W)_j \quad \text{if } j \in V_U$$

We can order vertices from V_L

first and then V_U second

we stop here forever

we don't do this

$\left(\begin{array}{c} \text{ } \\ \text{ } \end{array} \right)$

$$P = \begin{pmatrix} P_{ee} & P_{eu} \\ P_{ue} & P_{uu} \end{pmatrix} = \begin{pmatrix} I & 0 \\ P_{ue} & P_{uu} \end{pmatrix}$$

we stay here forever (pointing to P_{ee}) *we don't do this* (pointing to P_{eu})

we still $\lim_{t \rightarrow \infty} P^t$

$$P^\infty = \begin{pmatrix} I & 0 \\ (I - P_{uu})^{-1} P_{ue} & P_{uu}^\infty \end{pmatrix}$$

↑ goes to zero as we don't stop a v_u

$$P^\infty = \begin{pmatrix} I & 0 \\ (I - P_{uu})^{-1} P_{ue} & 0 \end{pmatrix}$$

Y_e = probability distribution over same class label

vertex in labeled set

$$v_i \rightarrow Y_{e_i} = [0 \dots 1 \dots 0]$$

i-th column (pointing to the 1)

unlabeled

$$v_j \rightarrow Y_{u_j} = [0.1 \dots 0.2 \dots 0.05]$$

Prediction matrix

$$Y_u = (I - P_{uu})^{-1} P_{ue} Y_e$$

we want to find a max. prob.
for our prediction

$$\operatorname{argmax}_{C_i \in C} Y_{u_i}$$

Can also do this iteratively

$$Y_u^{t+1} = P_{ue} Y_e + P_{uu} Y_u^t$$