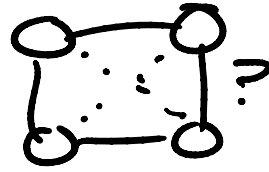
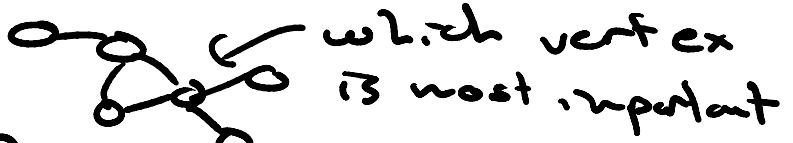


Graph mining:

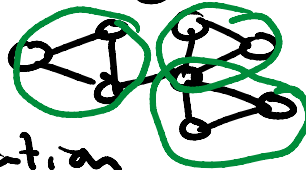
1. Link prediction



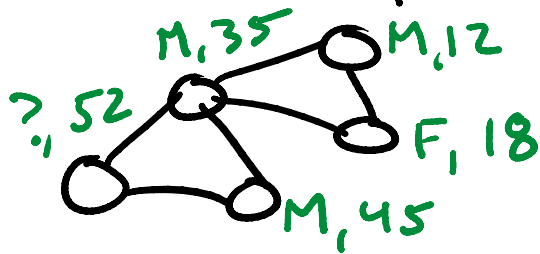
2. Centrality



3. clustering / CO



4. vertex classification and label prediction



Vertex labeling problem

Given graph $G = (V, E, w, Y)$

V = vertex set

E = edge set

w = edge weights

Y = vertex labels

Example: social network

V = people

E = relationships

w = strength of relationship

w = strength of relationship
 Y = demographic information

The problem:

Given $G = (V, E, w, Y_e)$

Predict Y_u

where Y_e = labeled data

Y_u = unlabeled data

Our approach: iterative classification

Features

Age of v

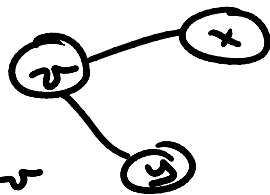
Gender of v

Politics of v

Avg. age of $N(v)$

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

Proportions L/R/C of $N(v)$



features of v

- metadata of v

- metadata of $N(v)$

Basic classification problem:

- we have pieces of data
- we construct features from it
- train a classifier on known classes
- use to predict unknown classes

- use to predict unknown classes

Our iterative classification algorithm:

Construct Φ_e, Φ_u from $G = (V, E, w, Y_e)$
 ← features for labeled/unlabeled

train f from Φ_e (classifier)
 $(\min_{v_e \in V} \|f(\Phi_e) - Y_e\|)$

min error output from classifier relative to ground truth

for some # iterations

predict $Y_u = f(\Phi_u)$

update Φ_u ← updating features

return Y_u

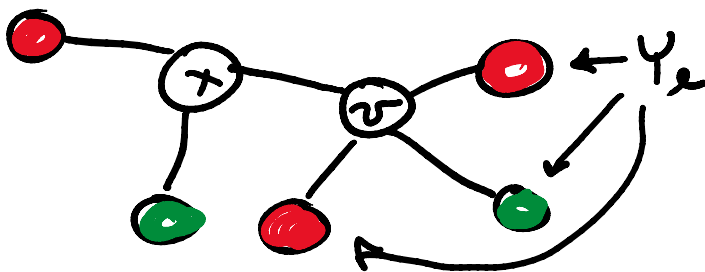
Classifiers

→ recall label propagation

Y_e = ground truth for $v_e \in V$

ϕ = labels of neighbors

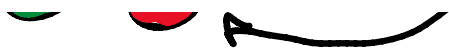
f = return label w/ max count
 (ties broken randomly)



ϕ_v = counts for each label in $N(v)$

$\phi_v = \{ \bullet : 2, \bullet : 1, \dots \}$

$f(\phi_v) = \bullet$



$$f(\phi_v) = \bullet$$

Naive Bayes Classifier

\bar{X} = features

\bar{X}_v = features for vertex v

$$\bar{X}_v = (x_{v_1}, x_{v_2}, x_{v_3}, \dots, x_{v_n})$$

To classify some v as class C_i

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

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

Two things first:

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

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

for more than two events

$$P(A_1 \cap A_2 \dots A_n) = P(A_1 | A_2 \dots A_n) P(A_2 \dots A_n)$$

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

~~$P(\bar{X}_v)$~~

Note: this is constant ~~$P(x_{v_1})$~~

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

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

chain rule again

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

apply 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

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

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

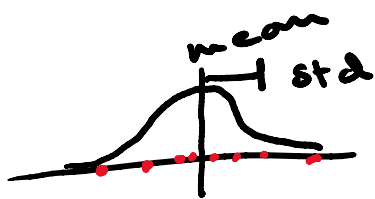
$$Y_u(\bar{x}_v) = \max_{i \in C} P(C_i) \prod_{j=1}^k P(x_{v_j} | C_i)$$

$P(C_i)$ = ratio of C_i labels

$P(C_i) = \text{ratio of } C_i \text{ labels}$

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

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



we need to assume a
distribution for numeric
feature values

we can assume a normal
distribution associated with
each feature

what we're calculating for $P(x_{v_j} | C_i)$
is the probability that feature
value x_{v_j} is sampled from the
distribution (mean, variance) calculated
from known labeled x_{v_j} features

Random walks

Idea: the probability of some
 $v \in V_u$ assuming label C_i is the

--
 $v \in V_u$ assuming label C_i is the probability of a random walk from v ends on some $u \in V_l$ with label C_i

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

P = transition prob. matrix

D = diagonal degree matrix

W = weighted adjacency matrix

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

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

Note: PP gives us a matrix defining 2-hop distribution

Note x2: we want to stop our walk if we land on a labeled vert.

Let's consider such a modification to P

$$P_i = \overset{\leftarrow \text{identity}}{e_i} \text{ if } i \in V_l$$

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

We can order vertices from labeled

We can order vertices from labeled first then unlabeled second

we stop at l ← don't go from $l \rightarrow u$

$$P = \begin{pmatrix} P_{ll} & P_{lu} \\ P_{ul} & P_{uu} \end{pmatrix} = \begin{pmatrix} I & 0 \\ P_{ul} & P_{uu} \end{pmatrix}$$

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

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

going to zero as we don't stop at and visit u

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

Y_l = probability distribution over some class label

vertex in labeled set

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

← i th column

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 prediction

→ so to assign same label to v_i

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

Iterative approach to calculate this

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

via our original P matrix

Is this equivalent?

$$Y_u^{t+1} = \sum_{i=1}^{t+1} P_{uu}^{i-1} P_{ue} Y_e + P_{uu}^t Y_u$$

as $t \rightarrow \infty$

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

Yes