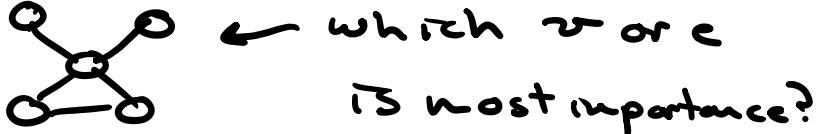


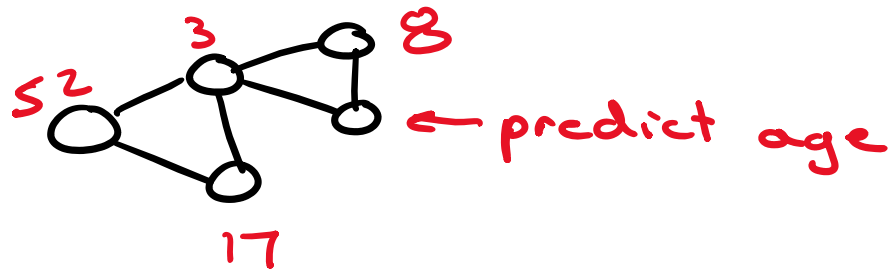
Graph Mining classic topics

1. Link prediction 

2. Centrality 

3. Clustering/CO 

4. Vertex (edge) classification



Vertex Labeling Problem

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

V = vertex set

E = edge set

W = weights (edges)

Y = vertex labels

\mathcal{Y} = vertex labels

Ex: social network

U = people

E = friendships / relationships

W = strength of relationships

\mathcal{Y} = demographic information

The Problem:

Given $G = \{U, E, W, \mathcal{Y}_e\}$

↑
existing labels

Predict \mathcal{Y}_u ← unlabeled data

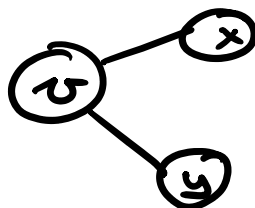
Approach: iterative classification

↳ we update our predictions using
our current predictions

Features/labels (\mathcal{Y}_e)

Age of v

Gender



Gender
Politics
etc.



Average age $N(w)$

Proportion of M/F/NB/etc. of $N(w)$

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

Basic classification problem:

- We have Y_e , pieces of data
- We construct features
- Train classifier using Y_e /features
- Predict Y_u , unlabeled data

Formalize:

Construct Φ_e, Φ_u from \mathcal{G}
(features)

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

For same # iterations:

Predict $Y_u = f(\Phi_u)$

update Φ_u (or Φ_e)

(potentially retrain f)

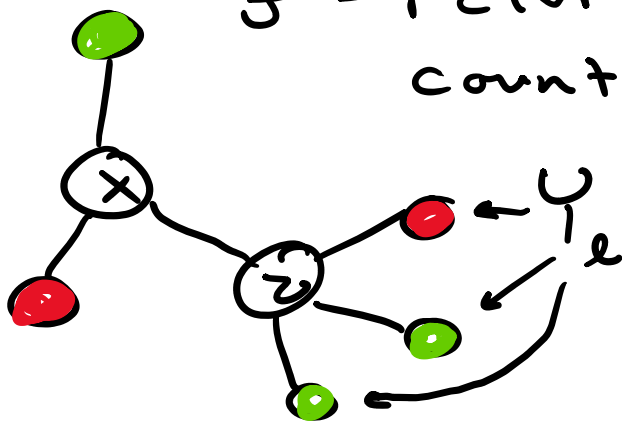
Return final Y_u

Classifier = Label propagation

$Y_e =$ ground truth labels

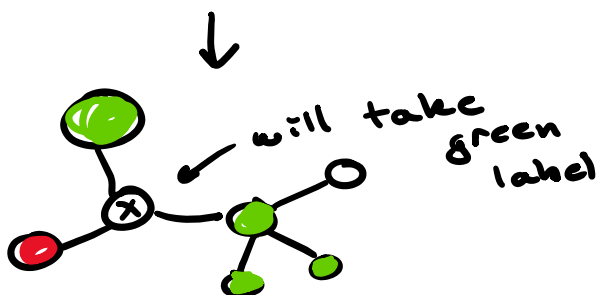
$\Phi_u =$ labels in $N(u)$

$f =$ return label with max count in $N(u) \mid \Phi$



$$\Phi_z = \{ \bullet = 2, \bullet = 1 \}$$

$$f(\Phi_z) = \{ \bullet \}$$



(same as for z)

$$\Phi_x = \{ \bullet = 2, \bullet = 1 \}$$

$$f(\Phi_x) = \{ \bullet \}$$

Naive Bayes Classifier

\bar{x} = features

\bar{x}_v = features for v

$\bar{x}_v = \{x_{v_1}, x_{v_2}, \dots, x_{v_k}\}$
(probabilities)

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

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

for 3+ "events"

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

Classification:

trivially to calc.
given labels

what we use to
learn

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

(Bayes)

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

↳ constant

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

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

chain rule again

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

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



recurse chain rule

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

Naive Bayes' assumption

↳ all x_j are independent

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

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

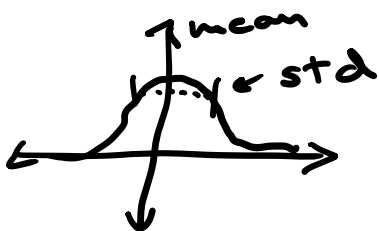
$$Y_u(\bar{x}_u) = \max_{i \in C} P(C_i) \prod_{j=1}^k P(x_{v_j} | C_i)$$

$P(C_i)$ = ratio of class labels

$$P(C_i) = \frac{|C_i|}{|V_e|} \leftarrow \# \text{verts in class } i$$

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

discrete = easy
 continuous = need to assume a given distribution



$P(x_{v_j} | C_i)$ = the probability that the feature value x_{v_j} is sampled from some distribution (mean, var.) calculated from known

calculated from known
labeled x_{uj} features
 $u \in U_2$

Random Walks

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

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

P = transitional prob. matrix

D = diagonal degree matrix (sum of weights)

W = weighted adj matrix

P_{ij} = prob. of a random walk from $i \rightarrow j$

$P = D^{-1}W \rightarrow$ a single step

$\lim_{t \rightarrow \infty} P^t =$ steady state solution

However: we only care about stopping on labeled vertices

We can modify our P

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

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

We label vertices from V_e first and then V_u second

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

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 at a V_u vertex

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

$$r = \left((I - P_{uu})^{-1} P_{ue} \bar{0} \right)$$

$Y_e =$ prob. distribution over same class label

Vertex in labeled set

$$v_i \in V_e, Y_{e_i} = [0 \dots 1 \dots 0]$$

← i^{th} column

unlabeled

$$v_j \in V_u, Y_{u_j} = [0.1 \dots 0.2 \dots 0.05 \dots 0.0001]$$

(arbitrary example)

Prediction matrix

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

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

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

Iteratively (to calculate)

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

CODE MOOE

Discogs → musical artists (album database)

we can construct bipartite graphs

artist → genre (what we'll classify
as singular)

artist → style (multiple)

↳ used to construct graph

we'll construct an overlapping
community graph

artists → overlapping communities
based on style

⇒ Predict genre