

Lecture 5-a

Feb 13, 2024

Empirical Data $G^* = (V^*, E^*)$ is often incomplete.

~~Let x_i be the attributes of node i .~~

We have four types of predictions.

- *₁ missing node attributes
 - missing link attributes
 - *₂ missing links
 - missing nodes
- } interpolation (imputation)
- } extrapolation
- ↓ increasing difficulty.

*₁ Lecture 5 (associative mining)
 *₂ Lecture 6 (next week)

Start from ① & ②

Then go here!

①

Vershynin (2018)

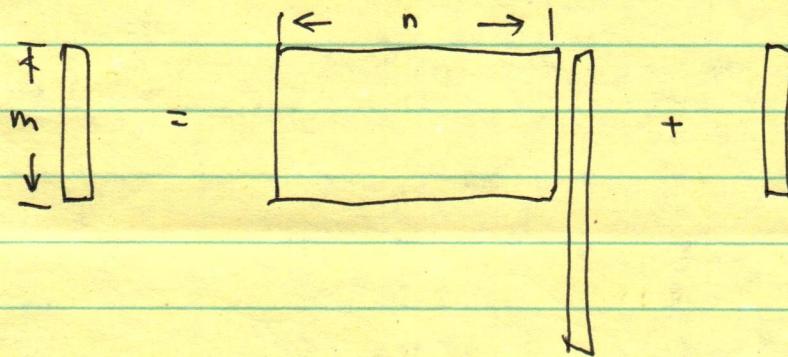
General intuition from "high-dimensional signal recovery problems".

unknown signal $x \in \mathbb{R}^n$

\equiv random, linear, noisy measurements of x , $y \in \mathbb{R}^m$

$$y = Ax + w$$

$\left\{ \begin{array}{l} A: \text{unknown measurement matrix} \\ w: \text{unknown noise vector.} \end{array} \right.$



$y \quad A \quad x \quad w$

→ Goal: Recover x from A and y as accurately as possible.

Example:

- Audio sampling

- Linear regression

freq.

~~Y~~ \rightarrow time

$$Y = X\beta + w$$

$X \in \mathbb{R}^{m \times n}$ sample of predictor

X_{ij} = expression of gene j in patient i

Variables.

Y_i : whether (or to what extent) patient has disease.

$m = |Y| \sim 100$ # patients while $n = |\beta| \sim 10000$ # genes.

(2)

Problem :

① If $m \ll n$, even if $w=0$, problem is ill-posed!

② ~~the~~ soln lives in a space of $\dim \geq n-m$

Soln (common assumption)

$$x \in T \subset \mathbb{R}^n$$



but T is informed by some prior information about the signal x .

Something ^{we} know, believe, or enforce.

/ want to

Benefit :

small $T \Rightarrow$ signal recovery is possible when $m \ll n$.

(3)

1. Predict missing node attributes

$$\left\{ \begin{array}{l} G^* = (V^*, E^*) \\ X = X^o \cup X^* \end{array} \right. \quad \begin{array}{l} \text{actual data (ground truth)} \\ \text{node metadata} \end{array}$$

$$X^o = \{\phi\} \quad k_i \in X = \left\{ \begin{array}{ll} \text{categorical} \\ \text{scalar} \end{array} \right. \quad \mathbb{R}, \mathcal{Z}$$

Goal: Recover X^o from X and G^* as accurately as possible.
(Predict)

2. Predict missing link attributes

$$\left\{ \begin{array}{l} G^* = (V^*, E^*) \\ w = w^o \cup w^*, \quad w^o = \{\phi\} \end{array} \right.$$

Goal: Recover w^o from w^* and G^* as accurately as possible.
(Predict)

3. Predict missing links

$$\left\{ \begin{array}{l} G^* = (V^*, E^*) \\ E = E^o \cup E^*, \quad E^* = V^* \times V^* - E^* \\ \text{so } |E^o| \text{ is big!} \end{array} \right.$$

Goal: Predict E^o from G^* as accurately as possible.
(i.e., V^* and E^*)

(4)

$\left\{ \begin{array}{l} \text{Predict missing nodes (extrapolation)} \end{array} \right.$

$$\left\{ \begin{array}{l} G = (V^*, E^*) \end{array} \right.$$

$$\left\{ \begin{array}{l} V = V^\circ \cup V^*, E = E^\circ \cup E^*, \text{ where } E^\circ \in V \times V - E^* \end{array} \right.$$

new nodes!

Goal: Predict V° (and E°) from V^* and E^* as accurately as possible.

- less explored!

- may resort to growth models. (week 13)

network

network sampling / growth.



In ensuing two weeks we focus on $\left\{ \begin{array}{l} 1 \\ 3 \end{array} \right.$.

(*) $\left(\begin{array}{l} \rightarrow \text{"baseline algorithm"} \\ \rightarrow \text{refined ones that either use local / global information.} \end{array} \right)$

(Δ) $\left(\begin{array}{l} \rightarrow \text{measuring performance. (confusion matrix c, AUC)} \end{array} \right)$

area under

ROC curve.

b) W5-a (*) : up to assortative mixing / local smoothing
 - b (Δ) + advanced methods (semi-supervised ML)

W6-a (*) on "link prediction"

- b Depth dive — Ghasemian et al. (PNAS, 2018)

and Peixoto & Kirkley (PRE, 2023)

(5)

{ 1 missing node attributes.

(global)

- baseline algorithm

~~take~~ x_i^* from Uniform $(\bar{X} - \phi)$

take
(impute)

side: why is this "baseline"?

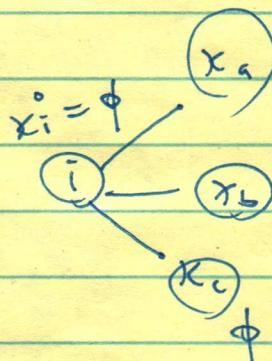
ans: ignores G.

→ note: this works whether X categorical or in \mathbb{Z} or \mathbb{R} .

- local smoothing algorithm

take x_i^* from $\begin{cases} \text{mean } (\bar{X}_{nb(i)} - \phi) & \text{if } x \in \mathbb{R} \\ \text{mode } (\bar{X}_{nb(i)} - \phi) & \text{if } x \in \{1, \dots, c\} \end{cases}$

The predictor.



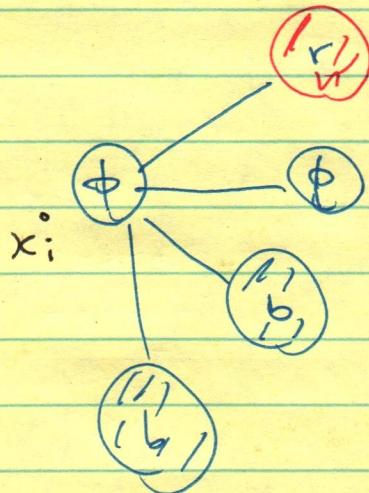
$\bar{X}_{nb(i)}$ set of ~~is~~ neighbor attributes.
(attrib. of neighbors of i)
mode data

(used G's info)

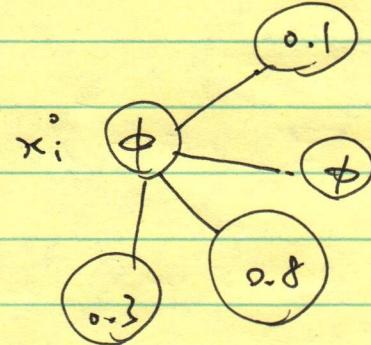
- synchronous updating (don't eat your predictors)
- what if $\{\bar{X}_{nb}\} = \{\phi\}$
 - make $nb(i)$ bigger $\rightarrow nb$'s of nb's.
 - use baseline predictor.

6

ex:



compute x_i with $\underline{b} =$



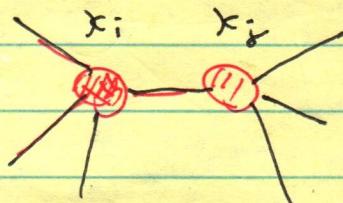
compute x_i with $\frac{0.1 + 0.8 + 0.3}{3} = 0.6$

Version of "label propagation"
(whole class of algorithms !)

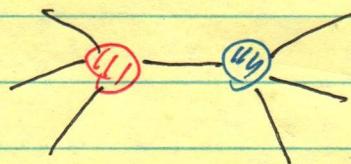
⑦

Mixing patterns

~~Phenomenon~~ : Assortative mixing (homophily)
Disassortative -



A-



Disa-

x_i : belief, age, social status.

economic/ecological #'

producer/consumers

Given $(i, j) \in E$, how similar are x_i, x_j ?

null models!

EXAMPLES:

- Ad health # : assortative mixing by age and race,
Newman & Clauset (Nature but not by gender!
Communication 2016)

- Weddell Sea food web : disassortative mixing by
same paper as above.
by feed type ; and
by metabolic category

- Peel et al. (PNAS 2018)
mixing pattern can be locally heterogeneous.

Caution
for future classes.

: Even after showing the slides, we still ran 10 minutes early!

①

Wk - 5b

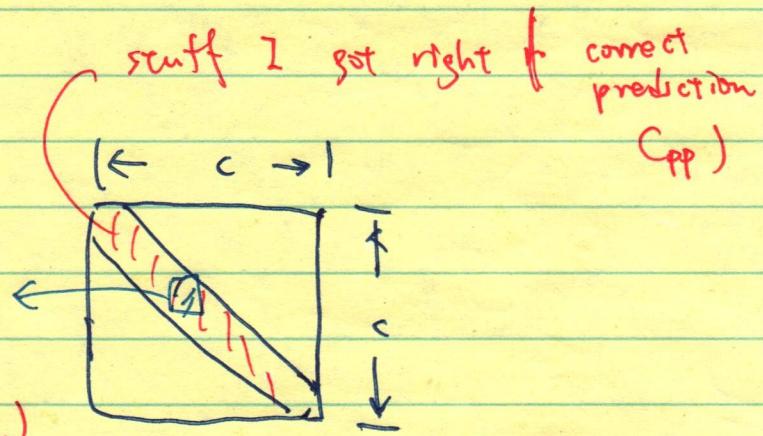
Feb 15

- only dive in neural #'; this semester.
- How do we measure the performance of what we came up?

measuring performance

confusion matrix

$$C_{P\neq q} = \# \text{ of inputs w/ (predicted label } p) \text{ and (actual label } q)$$



complete summary of the performance. $\underline{\text{H}}(\text{c}^2)$ in size

$C_{P\neq q}$ = misclass = things I thought should be p but actually q .

$$\text{accuracy (Acc)} = \frac{\# \text{ current}}{\# \text{ predictions}} = \frac{1}{N} \sum_p C_{Pp}.$$

"Unbalanced"

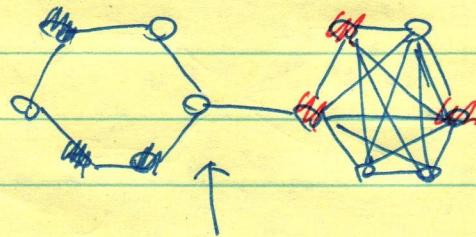
(what does $C_{P\neq q}$ count?)

4:03
F22

{ Imbalance

(3)

ex:



this one go last! (when ~~teaching~~^{teaching})

— take red b/c "summarizing all possible (label prop. paths)" (which we usually don't do)

If assume actual :

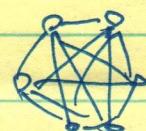
		TP actual	
		r	b
prod	r	3	1
	b	0	2
		FN	TN

← FP



blue

N



red.

P

$$\text{accuracy} = \frac{5}{6} = 0.83$$

If "binary classification" (and only if)

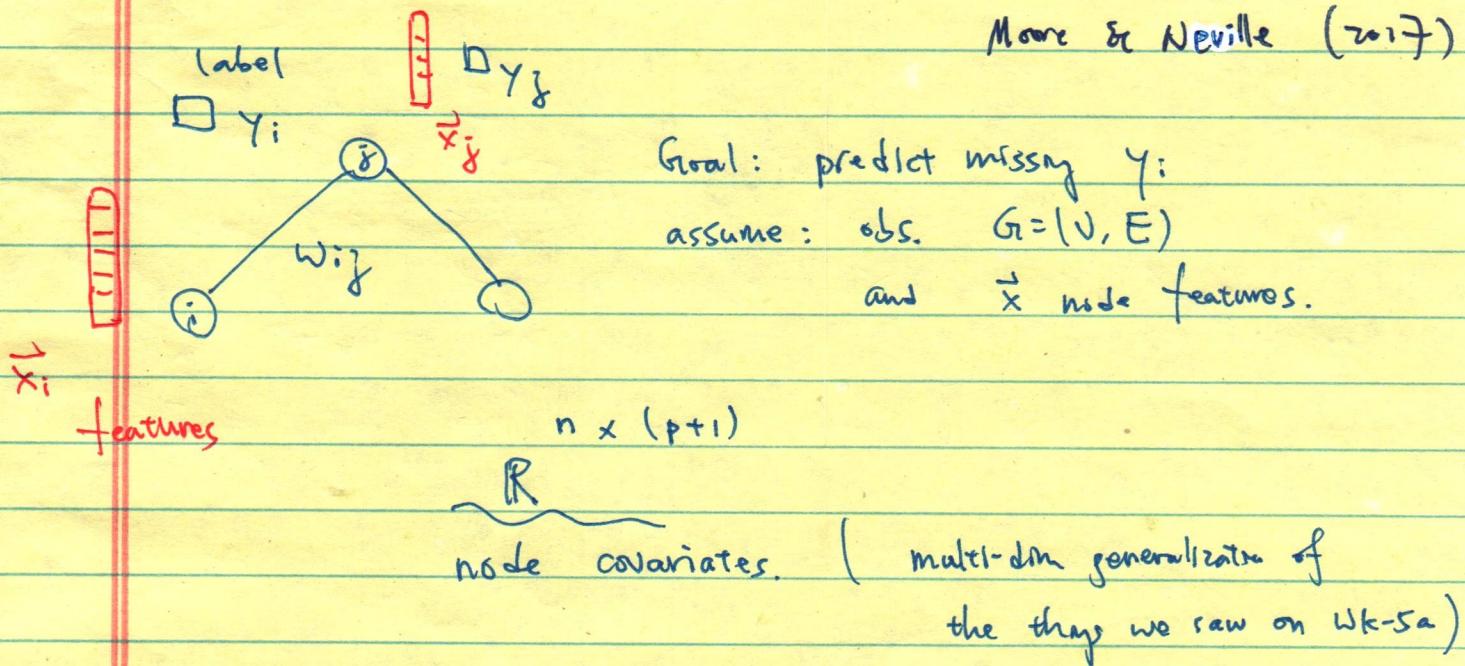
3

{ Penguin example. (see Notes) }

- 2020 H. Sayama & Michelle Girvan / Fang, Jint
- label propagation fails (I know penguins in this)
- lesson: at least look at the data once.
(Exploratory data analysis)
plot / compute assortativity coefficient, etc.

{ Modern techniques (2019) }

- predicting missing node attri } node classification
(Deep) } collective preference



→ ① $\text{cov}(\vec{x}_i)$ within features of a node

→ ② $\text{cov}(\vec{x}_i, \vec{x}_j)$, conditioned on $(i, j) \in E$

4

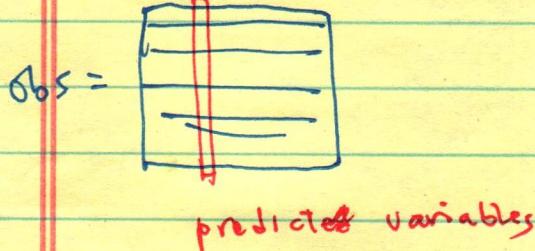
— At least ⁽⁵⁾ yr history.

"GNN" 2009

{ Deep learning came for graph (2014)

complicated.
not scalable.

"Deepwalk" graphs are hard, vectors are easy.
(spaces)



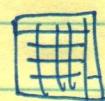
$$f(\boxed{\text{graph}} | \{\theta\}) = \boxed{\text{v}}$$

↑
not graphs. / not adj matrix.

"point cloud machine learning"

- SUM,
- random forest,
- regression
- table prop.

Deepwalk ($G, \{\theta\}$) \rightarrow output



"embedding"

(penguin)

is also an embedding

Solve optimization prob. $\left(\text{low-dim structure in high-dim space} \right) \text{by LR - SVD}$

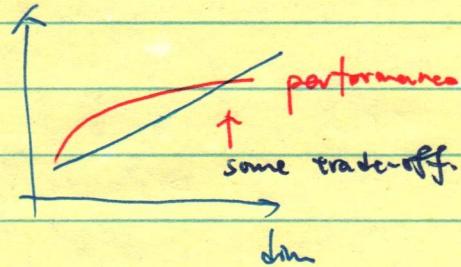
find $\mathbb{R}^{|V| \times d}$ matrix,

s.t. $f(\text{embedding})$ is minimized!

(how good is the embedding) \rightarrow variety of methods !

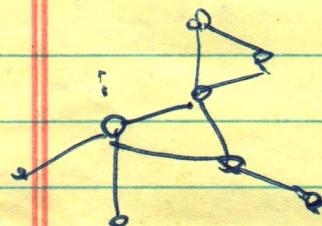
5

also:

— $\text{node2vec}(G, \theta)$ computation
cost

— How do we generate ?

"Random Walk"



deepwalk uses

Word2Vec (\circ)(random walks for NN embeddings,
but works on sentences)

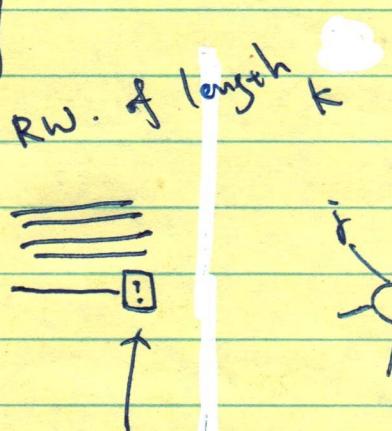
- NLP

a quick brown fox jump on

words : nodes

sentence : seq. of nodes. from a RW. of length k

corpus : a walk of sentences

height \times width

INPUT $\begin{cases} (n \times r) \times k \\ (\text{node}) \quad (\text{sentences for every node}) \\ \quad \quad \quad \{ \text{each of length } k \} \end{cases}$

select uniformly at random
from $n(i)$

OUTPUT embeddy

{ Deepwalk is just a wrapper for

word2vec

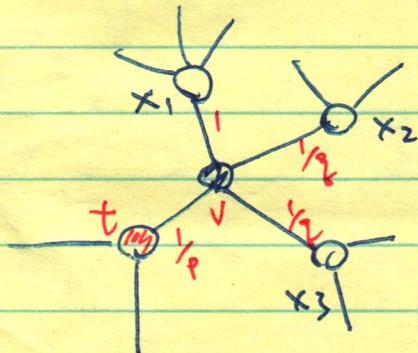
(6)

seg. encodes the structure of graph.

graph \rightarrow seg \rightarrow vectors.

perhaps graph \rightarrow vectors.

{ node2vec (generalization of DeepWalk)
"parametrizes the walk"



return
p task parameter
q in-out

Deepwalk : 1st-order RW.

node2vec : 2nd-order RW.

... + v [?]

Assign weights:
according to

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } dtx = 0 \\ 1 & \text{if } dtx = 1 \\ \frac{1}{q} & \text{if } dtx = 2 \end{cases}$$

v to decide which vertex to visit next,
given that it came from t.

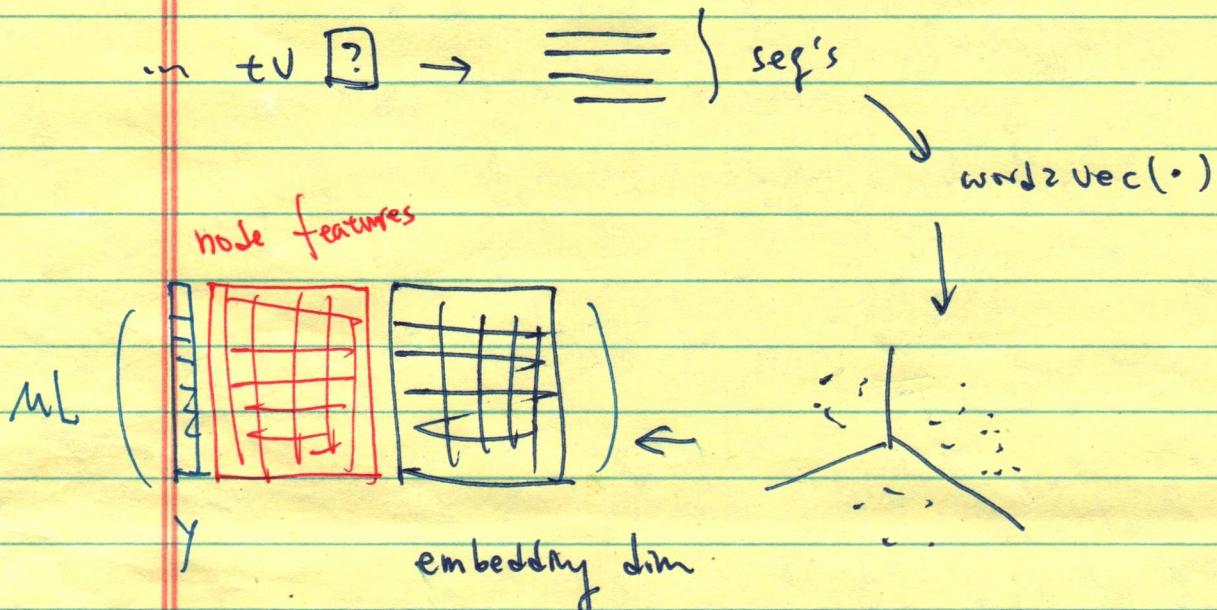
- feed to word2vec to move on!

blow your mind!

as A
(2014)
DeepWalk (2014) 10517 citations
node2vec (2016) 11306

7

- One argument of Dw/NV \Rightarrow they are scalable !
- hierarchical softmax ... but there are libraries !
stop gram



not many theory about why it works ... until ...

{ Third paper ! Jia-Benson (2022)

Great paper, see figure 1.

(and Table 4)

(that LR better

MRF \rightarrow NN than LP)

parametrized by
smoothness and
noise levels by
cross-validation

see Fig 3 of node2vec for different embedding

^{with}
X Prop. V
V X V
linear
regression NN or MRF

different structures

maybe existence of other embedding to explore structure.

trade-off general, ok performance method v.s. specialized method

If you want to explore in projects, that's good.

ex: they focus on prediction, but not how the parameters scale !

~~Can take-home messages.~~

- ML papers brag that their method is the best.
 - (but, truth is, some are good in one set of networks, some good at another.)
 - "no-free lunch" theorem.
 - cross-validation in graphs.
 - math V.S. NN-representation (slow/fast v. Bayesian / assumptions)
structure
 - A lot of variations for RW methods,
see Huang-Silva-Singh (2021)