

①

Lecture - 6a

Predicting MBSM Links

Feb 20, 2024.

Reminder: Project proposal start.
PSA

actual $\Rightarrow G = (V, E)$

observed $\Rightarrow G^* = (V, E^*)$, where $E^* \subseteq E$ ($E^* = f(E)$)

unconnected pairs $Y = V \times V - E^*$ $\sim \Theta(n^2)$ haystack
MBSM links $X = E - E^*$ $\sim O(n)$ needles.baseline accuracy by guessing $O(\frac{1}{n})$

[so hard!]

MBSMness
function.

Reasons we care:

- most ~~G~~ are incomplete!
- use correlations to marshall resources.
- compare methods or models. (cross-validation)

ex:
 - assembly ecosystem
 - gene-gene interaction
 - NSA / CIA

1. want to predict
 2. need to think of trade.
 3. so

Idea of a predictor is to optimize a score function.

score(i, j) : Y $\rightarrow \mathbb{R}$

Key message:

We can define a "spurious link prob", but it will
(by changing the score function)

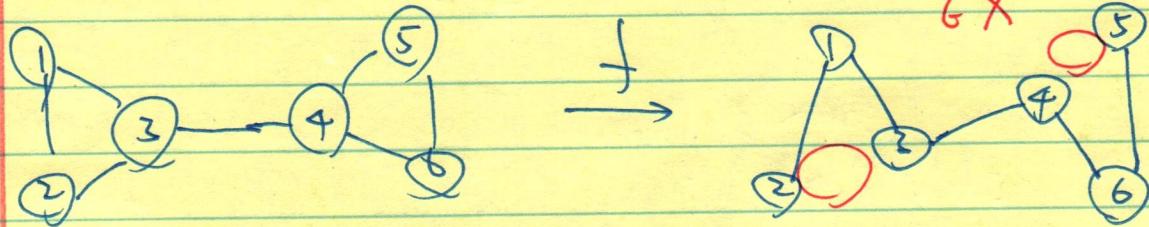
be easier than link prediction. (why?)

 $\Theta(n^2)$ to explore v.s. $\Theta(n)$ to explore.

(2) A baseline algm.

if $i, j \in \gamma$, $\text{score}(i, j) = \text{Uniform}(0, 1)$

↓
independent of G^* !



Score table

i	j	$\text{score}(i, j)$
1	5	$r = \text{Uniform}(0, 1)$
1	4	r
1	6	r

$\epsilon \gamma$

:	:	r
2	3	r
:	:	r
4	5	:

3

Topological predictors and local smoothness.

the idea

predict a missing link to occur where it would cluster
with other (\leftarrow served) edges

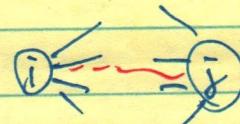
define this as
a predictor

Two ways to operationalize this idea :

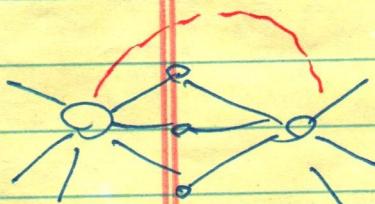
1) Jaccard coefficient.

2) degree product = $k_i k_j$

$$\text{Jaccard}(i, j) = \frac{|nb(i) \cap nb(j)|}{|nb(i) \cup nb(j)|}$$



$$\text{score}(i, j) = k_i k_j + \text{Uniform}(0, \epsilon) \quad \epsilon \ll 1$$



$$\text{scores}(i, j) = \text{Jaccard}(i, j) + \text{Uniform}(0, \epsilon)$$

↑
break ties randomly

$$\epsilon \ll \frac{1}{n-2}$$

④

an Example Jaccard =
Dp :

another example:

(5)

Measuring performance : the AUC
 predicting missing links is a binary classifier. (reg. only TP / FP).

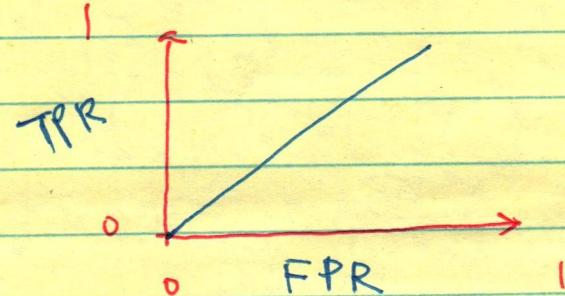
ask ??

we can use Area Under the Curve (AUC)

nice properties

→ Receiver operating characteristics (ROC) curve.

- 1) scale invariant (only cares about the order)
- 2) threshold invariant (general measure relative of distinguishability)



Mathematically

$$\text{AUC} = \Pr \left(\text{score(TP)} > \text{score(TN)} \right)$$

$\text{AUC} = 1 \Rightarrow$ perfect distinguishability

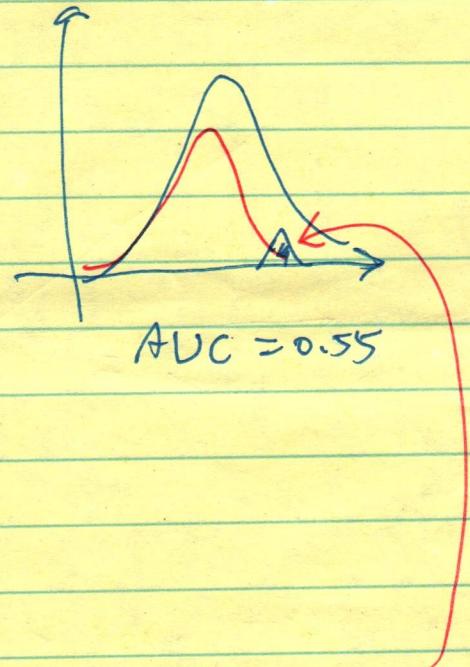
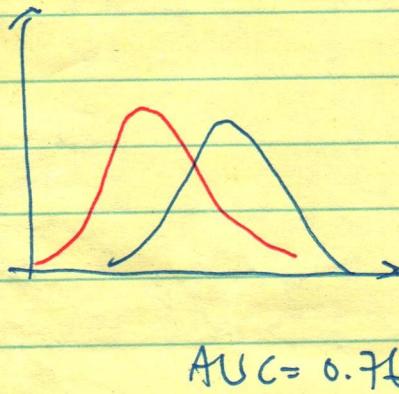
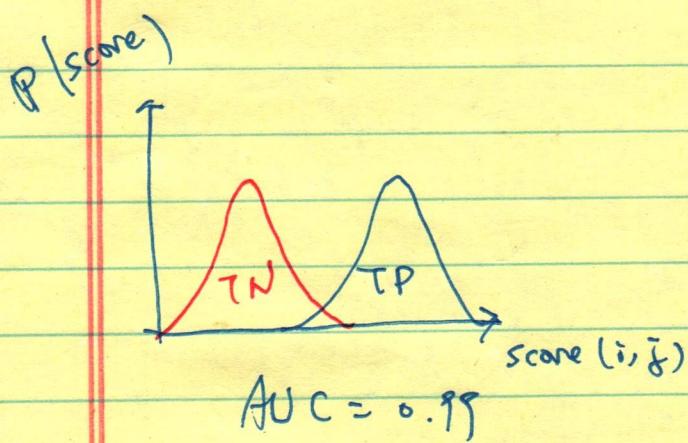
$= \frac{1}{2} \Rightarrow$ if I cannot tell a TP from TN.

In link prediction

$$|y-x|$$

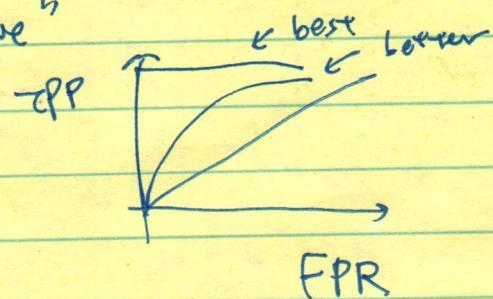
$$\begin{aligned} \text{TN} : & \quad \text{H}(x^2) \\ & |x| \end{aligned}$$

$$\text{TP} : \quad O(n)$$



$AUC \Rightarrow$ useful for company predictors
(agnostically) actually ...

if we want to see how much better, we need to look at
"the curve"



①

Wk- 6b.

Feb 22, 2024.

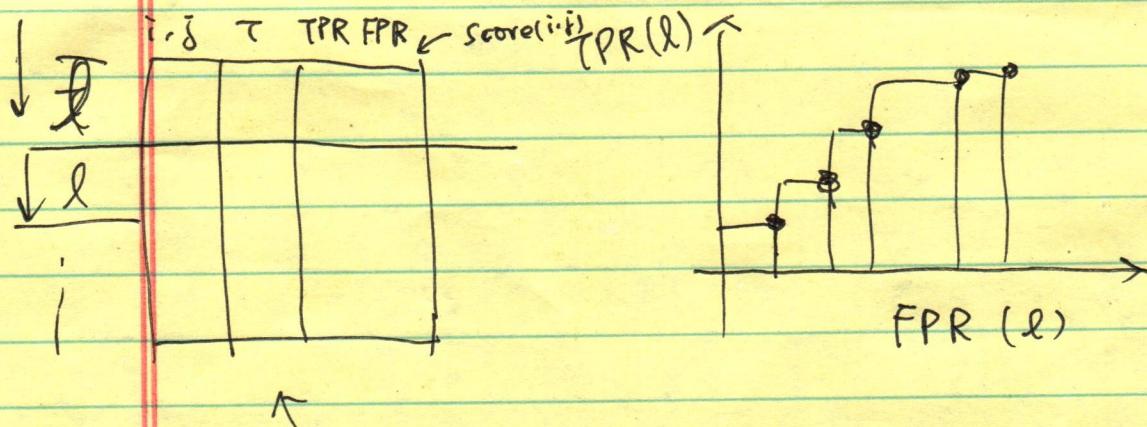
$$\text{TPR}(\ell) = \frac{1}{T} \sum_{k=1}^{\ell} \tau_k \quad ; \quad \tau_k = \begin{cases} 1 & \text{if it's missing a edge} \\ 0 & \text{otherwise} \end{cases}$$

TP

$$\text{FPR}(\ell) = \frac{1}{F} \sum_{k=1}^{\ell} (1 - \tau_k) ; \quad \text{FPR}(0) \equiv 0.$$

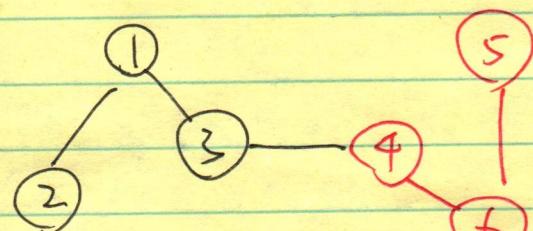
FP

$$\text{AUC} = \sum_{\ell=1}^{|\mathcal{T}|} \text{TPR}(\ell) \times \underbrace{\left(\text{FPR}(\ell) - \text{FPR}(\ell-1) \right)}_{\gamma \Delta x}$$



ROC coordinates

Example (See Clauset's notes)

Observed $G^+ = (V, E^+)$

(2)

Wk 6 - 6.

Feb 22, 2024.

{ The world of missing link predictors.

Three main classes.

1) topological : score (i, j) use only G^* (local measure of structure around i, j)

Jaccard, deg prod, $\frac{1}{2} + \frac{\deg(i) \deg(j)}{\deg(i) + \deg(j)}$

node-level statistic

(implicitly assumes an edge generating mechanism)

2) model-based

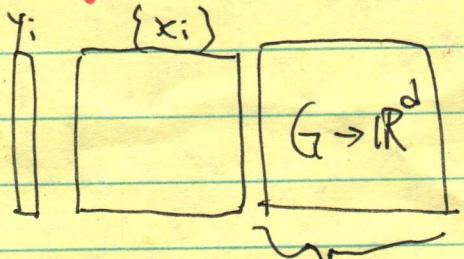
$$\Pr(G|\theta) = \prod_{i < j} \Pr(i \rightarrow j | \theta) \quad \text{edges are}$$

elaborate random graph models \rightarrow conditionally indep
(assume some explicit edge generating mechanism)

Stochastic block model (SBM)

3) embeddings

$$d(i, j)^{-1} \propto \text{score}(i, j)$$



assumes a spatial
homophily mechanism
for edge
generating.

(3)

But even more excitingly are **meta-learning** solutions

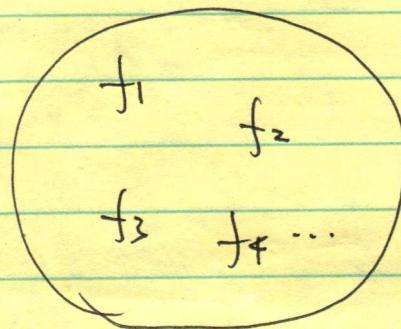
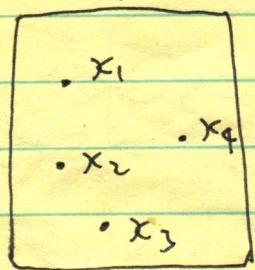
(aka. ensemble learning)

→ Key assumption: edges are the **observations** and the **predictors**
brief review

model selection via cross validation

observed data X

$\{x_i\}$



candidate models. F .

$\{f_n\}$

random forest,
neural network,
regression, etc.

* Which f_k provides the best **out-of-sample prediction accuracy** on not-yet-observed x_i ?

Solution: ① partition X in test and train data, randomly $X_{train} \cup X_{test}$

② learn apply each f_k to $X_{train} \mapsto \theta_k$

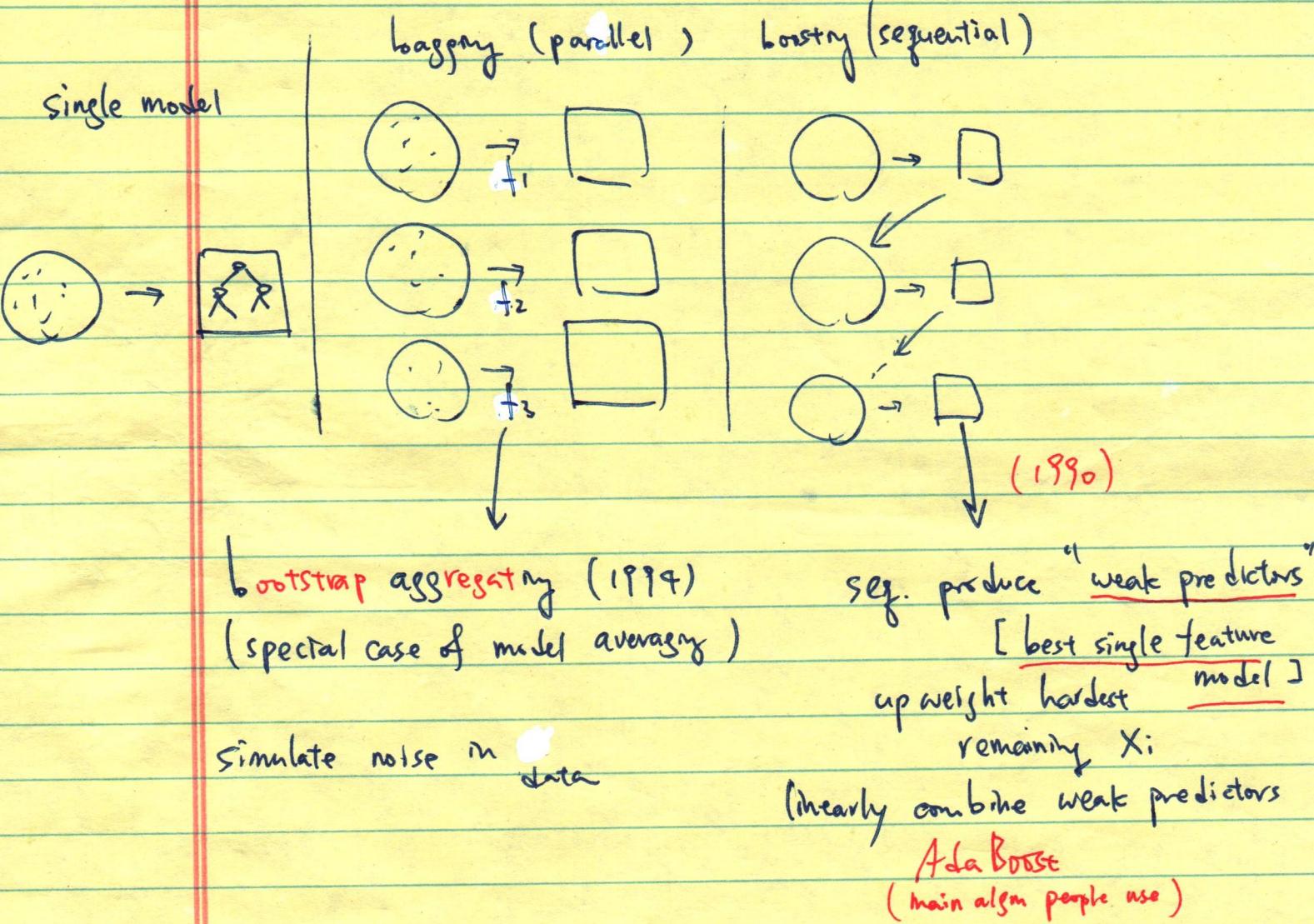
③ measure accuracy of each θ_k on $X_{test} \mapsto \varphi_k$

④ Select $\underset{k}{\operatorname{argmax}} \varphi_k \rightarrow$ best model

"winner take all approach"

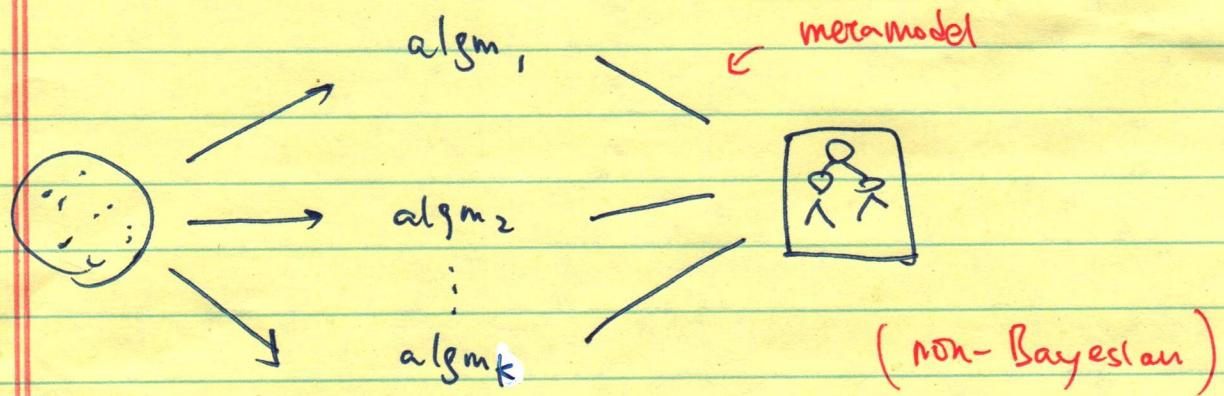
(4)

"Learning w/ ensembles"



(5)

Stacking (1992)



train many strong or weak predictors

learn a model of which predictor is best for a
given type of input.

Delphi method

→ constructing a predictive distribution

not model selection or model averaging, but model combination.

(6)

35 min - Start the slides.

		true label	
		0	1
predicted label	0	TN	FN
	1	FP	TP

② → positive.

$$\text{recall} = \frac{TP}{TP + FN}$$

$$\text{precision} = \frac{TP}{TP + FP}$$

» 203 features v. 25-30 features.

⑦

On real data :

- 1) inherently harder.
- 2) overfit - driven by algm.

80% social

remaining (b/s / info)
(easy others)

model class (U.S.) all data (No-Free lunch)