# Partial Recovery

Consider the constant degree case

$$p = \frac{a}{n}, \quad q = \frac{b}{n}$$

and $a, b = O(1)$

<u>Claim</u>: Exact recovery is impossible

<u>Proof</u>: The degree of a node

$$\sim Bin\left(p, \frac{n}{2}\right) + Bin\left(q, \frac{n}{2}\right)$$

$$\sim Poi\left(\frac{a+b}{2}\right)$$

where $X \sim Poi(\lambda)$ has $\mathbb{P}[x = k] = \frac{\lambda^k e^{-\lambda}}{k!}$

Thus $\Omega(1)$ fraction of nodes are isolated. ☒

But can we still solve <u>partial</u>* recovery?

   * i.e. our partition is $\frac{1}{2} + \varepsilon$ correlated with
      the true bisection

                                   Kesten-Stigum
                                      ↓ Bound

<u>Conjecture</u> [Decelle et al] If $(a-b)^2 > 2(a+b)$
then there is a polynomial time algorithm
to solve partial recovery. Else if $(a-b)^2 < 2(a+b)$
it's information theoretically impossible

Where does this conjecture come from?

## Belief Propagation

Let $\phi^i_{u \leftarrow v} \overset{\Delta}{=}$ node u's belief that it is in community $i$, if v were not there

Assumption: neighbors of u's communities are independent conditioned on community of u

Hope: It's true on trees, and sparse random graphs are locally tree-like

We get the following update rules

$$\phi^i_{u \leftarrow v} \propto \prod_{\substack{w \neq v \\ \text{s.t } (w,u) \in E}} \sum_{j=1}^{2} \phi^j_{w \leftarrow u} P_{i,j}$$

↑ probability of edge btwn community i & j

Iterate until convergence, and compute marginals

$$\phi^i_u \propto \prod_{\substack{w \text{ s.t.} \\ (w,u) \in E}} \sum_{j=1}^{2} \phi^j_{w \leftarrow u} P_{i,j}$$

Some effect of missing edges — global interaction

But there is a trivial fixed point

$$\phi^i_{u \leftarrow v} = \frac{1}{2} \quad \forall i, u, v$$

i.e. no one knows anything

Decelle et al: The trivial fixed point is unstable iff $(a-b)^2 > 2(a+b)$

① If BP doesn't get stuck here, maybe it solves partial recovery?

② If BP does get stuck, maybe the problem is impossible?

Thm [Mossel, Neeman, Sly] [Massoulie]
Both parts of the conjecture are true

However spectral partitioning does not work — the maximum degree is $\Theta\left(\frac{\log n}{\log \log n}\right)$ and top eigenvectors are localized

## Non Backtracking walks

Simpler approach following [Hopkins, Steurer]

Let $d = \frac{a+b}{2}$ = avg. degree

Then $p = \frac{(1+\epsilon)d}{n}$, $q = \frac{(1-\epsilon)d}{n}$

Now the Kesten-Stigum Bound becomes $\epsilon^2 d > 1$

<u>Goal</u>: Find a polynomial in $A_{ij}$'s
that can be used to estimate $x_i x_j$

$\uparrow$ $\uparrow$

$\pm 1$ community membership

<u>Consider</u> $P_{ij}(A) = \frac{n}{\varepsilon d}\left(A_{ij} - \frac{d}{n}\right)$

we have $\mathbb{E}[P_{ij}(A)] = x_i x_j$ but its variance

$\sim \mathbb{E}[P_{ij}(A)^2] \approx \left(\frac{n}{\varepsilon d}\right)^2 \left(\frac{d}{n}\right) = \frac{n}{\varepsilon^2 d}$

is too large

<u>Main Idea</u>: Average over many walks

$$P_\alpha(A) = \prod_{(a,b) \in \alpha} P_{ab}(A)$$

It is still an unbiased estimator since

$$\mathbb{E}[P_\alpha(A)] = \prod_{(a,b) \in \alpha} \mathbb{E}[P_{ab}(A)]$$

$$= \prod_{(a,b) \in \alpha} x_a x_b = x_i x_j$$

It's variance is even larger

$$\mathbb{E}[P_\alpha(A)^2] \approx \left(\frac{n}{\varepsilon^2 d}\right)^\ell$$

where $\alpha$ is a length $\ell$ path

But if they were pairwise independent
we could take

$$\underbrace{\frac{1}{|P_{ij}^\ell|}}_{\wedge} \sum_{\alpha \in P_{ij}^\ell} P_\alpha(A)$$

paths of length $\ell$ from $i$ to $j$

which would have variance

$$\sim \left(\frac{1}{n^{\ell-1}}\right)\left(\frac{n}{\varepsilon^2 d}\right)^\ell$$

which is $o(1)$ when $\varepsilon^2 d > 1$

Comment: This does <u>not</u> work because
up to scaling, we are computing the $(i,j)$
entry of

$$(A - \frac{d}{n}\vec{1}\vec{1}^\top)^\ell$$

and we know spectral methods fail b/c of
high degree nodes!

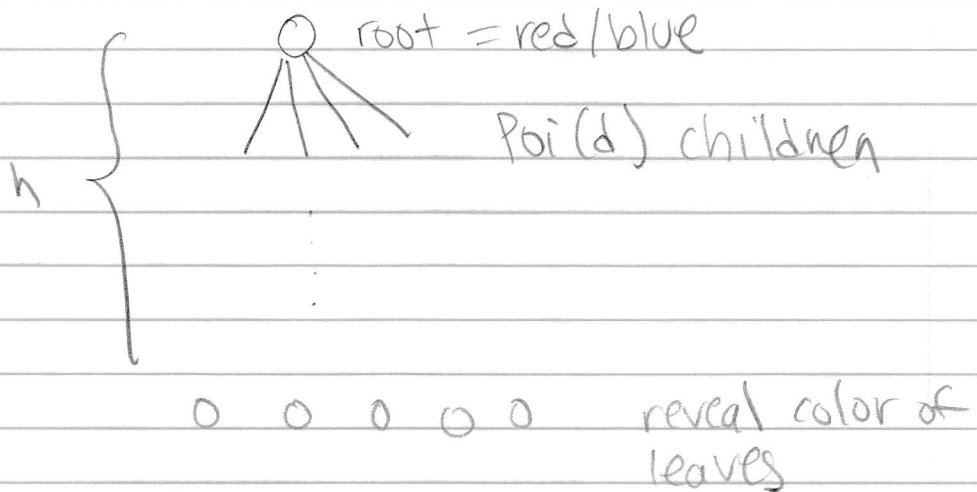But if you use self avoiding walks, they're
close enough to pairwise independent

Essentially, want to bound

$$\sum_{\alpha, \beta \in P_{ij}^\ell} \mathbb{E}[P_\alpha(A) P_\beta(A)]$$

and for SAW, if $\alpha$ and $\beta$ share r edges they must share at least r vertices, which reduces the number of possibilities by $n^r$

## Broadcast Tree

Another view



root = red/blue

Poi(d) children

reveal color of leaves

Each child has probability $\frac{a}{a+b}$ of being same color, o.w. different

Main Question: when can you guess the label of the root $\geq \frac{1}{2} + \varepsilon$ probability independent of h?

Thm [Kesten-Stigum] Can solve partial recovery if $(a-b)^2 > 2(a+b)$

Thm [Evans et al]: If $(a-b)^2 \leq 2(a+b)$ it's impossible

For upperbound, can take majority vote

[Kesten-Stigum] proved CLT for multi-type branching processes

Aside: Majority vote does not achieve optimal accuracy but BP/dynamic programming does

For lower bound, intuition is

$$I(\sigma(\rho); \sigma(x)) = \left(\frac{a-b}{a+b}\right)^{2h}$$

↑ mutual information   ↑ ↑ color of root/leaf

There are $d^h$ leaves, so we think

$$I(\sigma(x); \underbrace{\sigma(x_1), \ldots \sigma(x_{d^n})}_{\text{all leaves}}) \approx d^h \left(\frac{a-b}{a+b}\right)^{2h}$$
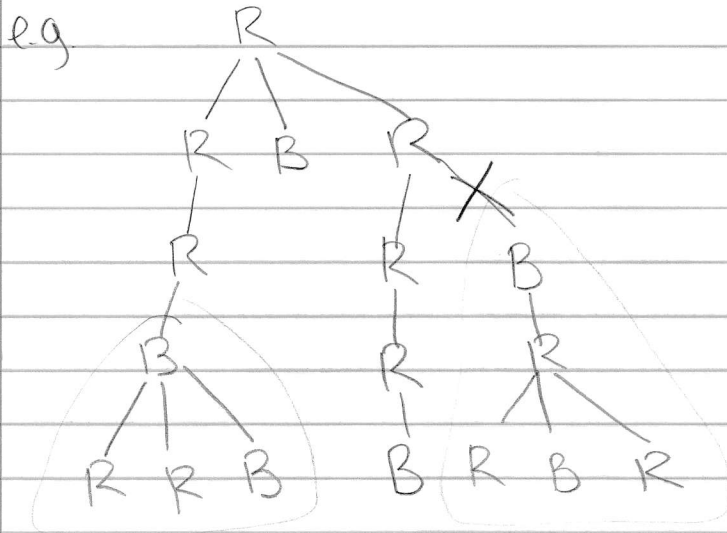
$$= \left(\frac{a+b}{2}\right)^h \left(\frac{a-b}{a+b}\right)^{2h}$$

$$= \left[\frac{(a-b)^2}{2(a+b)}\right]^h$$

Mutual information doesn't satisfy subadditivity, but [Evans et al] give a coupling argument
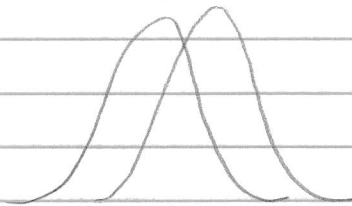
How robust is the Kesten-Stigum bound?

def: A monotone adversary can cut edges btwn nodes of opposite colors, remove subtree
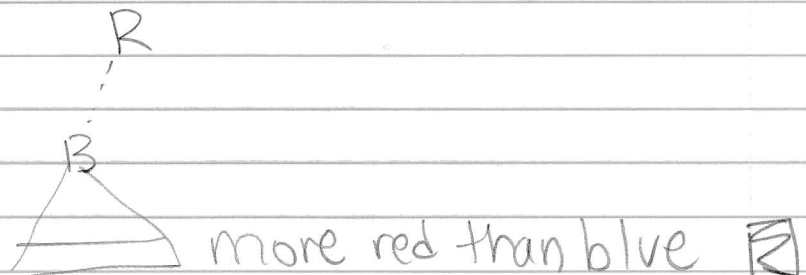
e.g



Claim: The adversary can whp flip the majority vote near the Kesten-Stigum bound
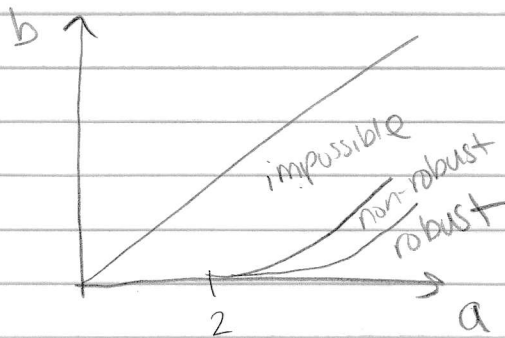
Proof: [sketch] # red - blue distributed as



Likely to have subtrees where



more red than blue

Theorem [Moitra, Perry, Wein] Reconstruction in semirandom broadcast free model is impossible for $(a-b)^2 \leq C_{a,b}(a+b)$ for some $C_{a,b} \geq 2$



These are the first random vs. semirandom separations

Main Complication: Need carefully designed adversary that maps into nice distribution

Theorem [Moitra, Perry, Wein] Recursive majority succeeds in semirandom broadcast tree model if $(a-b)^2 > (2+o(1))(a+b) \log \frac{a+b}{2}$

And also random vs. semirandom separators for community detection
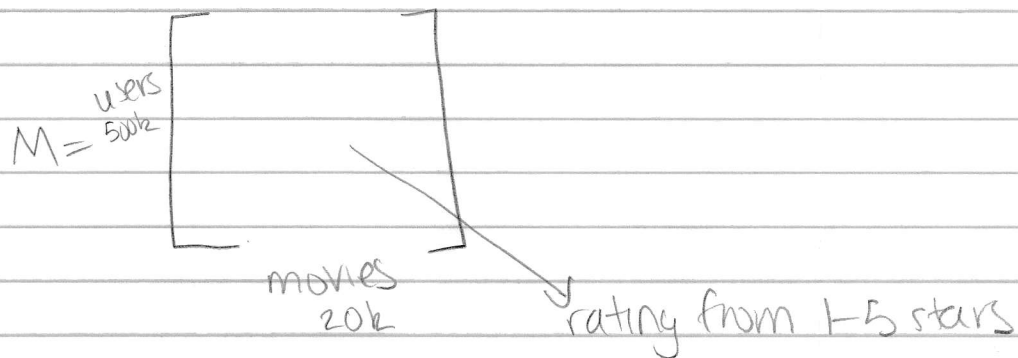
Aside: Above Average-Case Analysis

[Spielman, Teng '01].

"Explain why algorithms work well in practice, despite bad worst-case behavior"

usually called Beyond Worst-Case Analysis

Another example where semirandom models
are informative

## The Netflix Problem

$$M = \begin{matrix} \text{users} \\ 500k \end{matrix} \boxed{\phantom{xxxxxxxx}}$$

movies
20k          rating from 1-5 stars

Assume: M is low rank, e.g.

$$M = \left[\left[\begin{array}{c}[\quad]\end{array}\right]\right] + \left[\left[\begin{array}{c}[\quad]\end{array}\right]\right] + \cdots$$

$$\underbrace{\qquad}_{\text{drama}} \qquad \underbrace{\qquad}_{\text{comedy}}$$

and incoherent, ie its singular vectors far
from standard basis vectors

Theorem: [Candes, Tao]. Given $m \gtrsim nr\log^2 n$
u.a.r observations from M, there is
a polynomial time algorithm that whp
recovers M exactly

Their approach was based on convex programming

Let $\Omega$ = observed entries

$$\min \|X\|_*$$

$$\text{s.t. } X_{ij} = M_{ij} \quad \forall (i,j) \in \Omega$$

where $\|X\|_* = \sum_{i=1}^{n} \sigma_i(x)$ is the underline{nuclear norm}

Remark: This is a non-commutative generalization of $\ell_1$-norm minimization in compressed sensing

$$\text{sparsity: } \ell_1 \; :: \; \text{rank: } \| \; \|_*$$

Another powerful approach is alternating minimization

Repeat

$$U \leftarrow \underset{U}{\arg\min} \sum_{(i,j) \in \Omega} |(UV^T)_{ij} - M_{ij}|^2$$

$$V \leftarrow \underset{V}{\arg\min} \sum_{(i,j) \in \Omega} |(UV^T)_{ij} - M_{ij}|^2$$

Theorem [Keshavan et al] [Jain et al] [Hardt]
Alternating minimization with proper initialization succeeds whp given

$$C n r^2 \frac{\|M\|_F^2}{\sigma_r^2} \text{ observations}$$

What if a monotone adversary reveals more of M?

claim: Nuclear norm minimization still succeeds

Proof: It's just more constraints. ▨

Observation: Alternating minimization fails in the semirandom model

[Cheng, Ge] give a nearly linear time preprocessing step to fix the nonconvex approach

―――――

Another application: GMMs

def: A semirandom GMM proceeds as follows

(1) samples $X_i \sim F$

(2) adversary can move point in the direction of the center

$$y_i = (1-\lambda)X_i + \lambda M_j$$

Can we still find an accurate clustering?

[Awasthi, Vijayaraghavan] Yes if the reparation satisfies
$$\|M_i - M_j\| \geq \Delta \sigma$$

where $\|\xi_i\| \le \sigma$ and

$$\Delta \ge C\sqrt{\min(k,d)\log n}$$

moreover this is information-theoretically tight

Are there other stochastic models where semirandomness offers new insights?