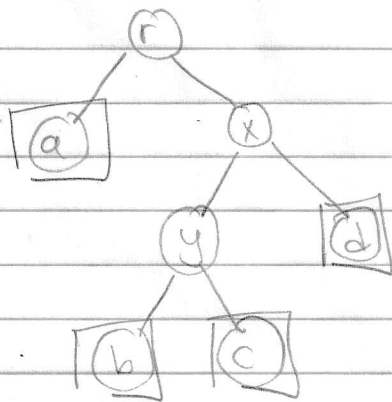# Applications to Learning

Next we'll study applications of tensor decomp.
to learning latent variable models

## Phylogenetic Trees

Setup: There is a rooted binary tree with
root $r$



$0 = $ extinct

$\boxed{0} = $ extant

An alphabet $\Sigma$ of states, e.g. $\Sigma = \{A, C, G, T\}$,
and a markov model consisting of

    (1) An initial distribution $\pi : \Sigma \to \mathbb{R}^{\geq 0}$
        for the symbol at the root

    (2) Along each edge, a conditional distribution
$$P_{ij}^{uv} = \mathbb{P}\left[s(u) = j \mid s(u) = i\right]$$

Can we recover the model from the joint
distribution on the extant nodes?

Our algorithm will work in two phases

(1) recover the topology of the tree

The Sometimes called "tree of life", identifies
speciation events and ancestral relationships

(2) estimate the markov parameters

## Evolutionary Distance

Steel introduced a distance function between
pairs of nodes on the tree with the properties

(a) it is nonnegative

(b) it can be evaluated for any pair,
just given their joint distribution

definition: steel's evolutionary distance on
an edge $(u,v)$ is

$$\psi_{uv} = -\ln|\det(P^{uv})| - \frac{1}{2}\ln\left(\prod_{i=1}^{k}\pi_u(i)\right)_{0 \text{ if leaf}}$$

$$+ \frac{1}{2}\ln\left(\prod_{i=1}^{k}\pi_v(j)\right)$$

conditional distribution $u \to v$

marginal on $v$

Assuming all $P^{uv}$'s are full rank,

Lemma: Steel's evolutionary distance satisfies

(1) $\psi_{uv}$ is nonnegative

(2) for any pair $(a,b)$ we have

$$\psi_{ab} \overset{\Delta}{=} -\ln|\det(F^{ab})| = \sum_{(u,v) \in P_{ab}} \psi_{uv}$$
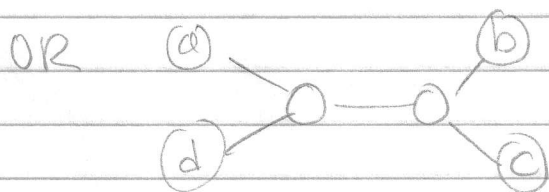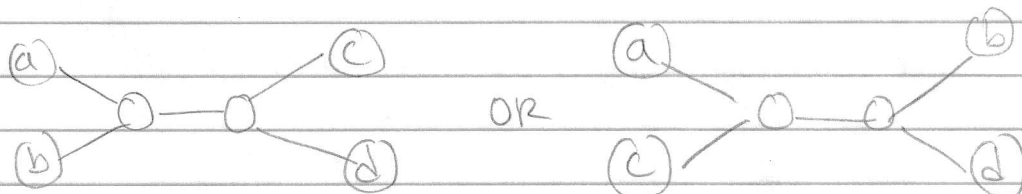
joint distribution on $(a,b)$

where $P_{ab}$ is the path connecting $a$ and $b$

[Erdos, Steel, Szekely, Warnow] used Steel's distance and quartet tests to reconstruct the topology

Lemma: If all distances are strictly positive, it is possible to determine the induced topology on any four nodes given an oracle for computing pairwise distances

Proof: there are three possible induced topologies



OR



OR

i.e. if we delete edges not on shortest path, and contract paths to a single edge

Easy to check we're in first case iff.

$$\nu_{a,b} + \nu_{c,d} < \min \{ \nu_{a,c} + \nu_{b,c}, \nu_{a,d} + \nu_{b,d} \}$$

Similar relation holds in other cases ☒

Lemma: If for any quadruple we can determine the induced topology, we can determine the overall topology

Proof: Strategy: determine which pairs of leaves are siblings, and so on

To do this, fix a pair $(a,b)$. They are siblings iff for every other pair $(c,d)$, the quartet test on $\{a,b,c,d\}$ determines we're in the first case

Now to continue, $(a,b)$ are non-siblings but share a grandparent iff for any pair $(c,d)$ neither of which is a sibling of $a$ or $b$, the quartet test on $\{a,b,c,d\}$ determines we're in first case.

And so on. ☒

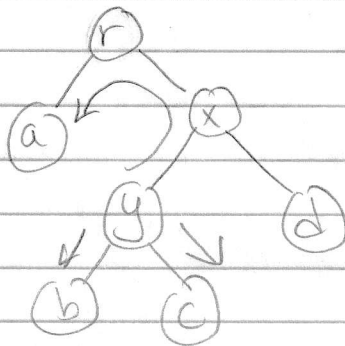Note: We can only estimate $\nu_{a,b}$ from samples, but [Erdos et al] show how to only use

quartet tests when distances are small

## Getting the Transition Matrices

If we know the topology, how can we estimate the $p^{uv}$'s?

Claim: If the $p^{uv}$'s are full rank, we can reroot the tree arbitrarily

Now consider any triple $(a,b,c)$ of leaves, e.g.



star test

reroot the tree at $y$ and consider the joint distribution on $(a,b,c)$

$$T^{abc}_{ijk} = \mathbb{P}\left[s(a)=i, s(b)=j, s(c)=k\right]$$

Then we have

$$T^{abc} = \sum_e \mathbb{P}\left[s(y)=\ell\right] \mathbb{P}\left[s(a)=\cdot \mid s(y)=\ell\right]$$
$$\mathbb{P}\left[s(b)=\cdot \mid s(y)=\ell\right] \mathbb{P}\left[s(c)=\cdot \mid s(y)=\ell\right]$$

But notice that

$$P[s(b) = \cdot \mid s(y) = \ell] = \ell^{th} \text{ row of } P^{yb}$$

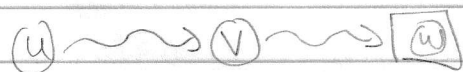$$P[s(c) = \cdot \mid s(y) = \ell] = \ell^{th} \text{ row of } P^{yc}$$

$$P[s(a) = \cdot \mid s(y) = \ell] = \ell^{th} \text{ row of } P^{yx} P^{xr} P^{ra}$$

Thus we have that

(1) $T^{abc}$ is a low rank tensor that can be computed from samples

(2) It's tensor decomposition determines properties of the parameters of the model

We can also determine the internal transitions up to equivalence. E.g. suppose

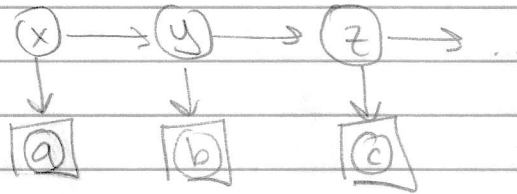$$\text{(u)} \rightsquigarrow \text{(v)} \rightsquigarrow \boxed{\text{w}}$$

We can find $P^{uw}$ and $P^{vw}$ using star tests, and then solve for

$$P^{uw} = P^{uv} P^{vw} \implies P^{uv} = P^{uw}(P^{vw})^{-1}$$

[Mosel, Roch] show how to recover transition matrices using only short paths, in Valiant's PAC model

# Hidden Markov Models

Setup: A hidden markov model is given by

$$x \rightarrow y \rightarrow z \rightarrow \cdots$$
$$\downarrow \quad\quad \downarrow \quad\quad \downarrow$$
$$\boxed{a} \quad\quad \boxed{b} \quad\quad \boxed{c}$$

where $\Sigma_s$ and $\Sigma_o$ are alphabets on hidden and observed states respectively.

Moreover let $P^{xy}$ be the transition matrices and $O^{xa}$ be the observation matrices

$$P^{xy}_{ij} = \mathbb{P}\left[ s(y)=j \mid s(x)=i \right]$$
$$O^{xa}_{ij} = \mathbb{P}\left[ s(a)=j \mid s(x)=i \right]$$

Theorem [Mossel, Roch] There is a polynomial time algorithm for learning HMMs when the observation / transition matrices are full rank and have lower bounded smallest singular value.

What about when the observation matrices are not full rank?

Consider the noisy parity problem

Setup. Let $S \subseteq [n]$ and for each sample, choose $X \in_{uar} \{0,1\}^n$ and set

$$y = \begin{cases} \chi_S(x) \triangleq \sum_{i \in S} X_i \mod 2 & \text{w/ prob } \frac{2}{3} \\ 1 - \chi_S(x) & \text{o.w} \end{cases}$$

**Theorem.** [Blum, Kalai, Wasserman] there is a $2^{n/\log n}$ time algorithm for finding $S$, i.e. the hidden noisy parity

This is the best known algorithm, and noisy parity is widely believed to be hard

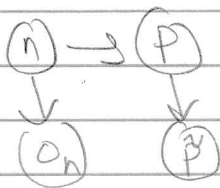Let's embed noisy parity into an HMM

$$\circled{i} = (X_i, P_i)$$

$i^{th}$ bit    running parity, i.e.

$$\sum_{\substack{j \in S \\ j \leq i}} X_j \pmod 2$$

Then the observation

$$\circled{i}$$
$$\downarrow$$
$$\circled{O_i} \quad O_i = X_i$$

we get the $i^{th}$ bit. This matrix is $4 \times 2$, and hence **not** full rank

Finally in the last step

$$\boxed{n} \to \boxed{p}$$
$$\downarrow \qquad \downarrow$$
$$\boxed{O_n} \qquad \boxed{\tilde{\beta}}$$

we set $p = P_n$, and get $\tilde{p} = \begin{cases} p & \text{w/ prob } 2/3 \\ 1-p & \text{o.w.} \end{cases}$

theorem [Mossel, Roch]: Learning general HMMs (i.e. without full rank observation matrices) is as hard as noisy parity

their proof uses the simulation above, and also the self-reducibility property of noisy parity, i.e. it's enough to distinguish noisy parity from the case when the label is random

Historical note: In phylogenetics, tensor methods are sometimes called Chang's lemma

## Mixtures of Spherical Gaussians

Another powerful application: Mixtures of spherical Gaussians

def: A Gaussian with mean $\mu$ and covariance $\Sigma$ has pdf

$$N(\mu, \Sigma, x) = \frac{e^{\frac{-(x-\mu)^T \Sigma^{-1}(x-\mu)}{2}}}{(2\pi)^{d/2} \det(\Sigma)^{1/2}}$$

In the special case where $\Sigma = I$, i.e. spherical, we get

$$N(\mu, I, x) = \frac{e^{-\frac{||x-\mu||^2}{2}}}{(2\pi)^{d/2}}$$

Setup: We get samples from a mixture model

$$M = \sum_{i=1}^{k} w_i \, N(\mu_i, \sigma^2 I, x)$$

↑ mixing weight

Can we learn the parameters/cluster?

Theorem [Hsu, Kakade] There is an algorithm with polynomial running time/sample complexity when the $M_i$'s have full rank

Note: The running time / sample complexity depend on $1/w_{min}$, $1/\sigma_{min}$, etc.

Lemma: If $\sigma^2$ is known, then the tensor

$$T = \sum_{i=1}^{k} w_i \, M_i^{\otimes 3}$$

can be expressed through moments of the mixture

Observe that this is the same recipe as before.

(1) there is a low rank tensor that can be estimated from samples

(2) It's tensor decomposition determines the parameters of the model

Proof: Consider $T_{abc}$

Case #1: $a, b, c$ are distinct

then $\mathbb{E}_{\mu}[X_a X_b X_c] = \sum_{i=1}^{k} w_i (\mu_i)_a (\mu_i)_b (\mu_i)_c$

this follows because the noise is independent across coordinates.

Written another way

$$\mathbb{E}_{\mu}[X_a X_b X_c] = \sum_{i=1}^{k} w_i (\mu_i^{\otimes 3})_{a,b,c}$$

Case #2 $a = b \neq c$

Then we have

$$\mathbb{E}_{\mu}[X_a X_b X_c] = \sum_{i=1}^{k} w_i ((\mu_i)_a^2 + \sigma^2)(\mu_i)_c$$

$$= \sum_{i=1}^{k} w_i (\mu_i^{\otimes 3})_{abc} + \sigma^2 \underbrace{\sum_{i=1}^{k} w_i (\mu_i)_c}_{\text{first moment}}$$

Case #3  $a = b = c$

then we have

$$\mathbb{E}_{\mu}[X_a X_b X_c] = \sum_{i=1}^{k} w_i (\mu_i^{\otimes 3})_{ax} + 3\sigma^2 \left( \underbrace{\sum_{i=1}^{k} w_i \mu_i}_{\text{first moment}} \right)_c$$

this follows because in one-dimension

$$\mathbb{E}[X^3] = \mu^3 + 3\mu\sigma^2$$
$$x \leftarrow N(\mu, \sigma^2)$$

Rearranging the inequalities, we get

$$T = \mathbb{E}_{\mu}[X^{\otimes 3}] - \sigma^2 \sum_{j=1}^{d} M_j$$

where $M_j = ( \mathbb{E}[X] \otimes e_j \otimes e_j + e_j \otimes \mathbb{E}[X] \otimes e_j$

$$+ e_j \otimes e_j \otimes \mathbb{E}[X] )$$

Now we can use tensor decomp.

## Further Applications

In a topic model we have

(1) there are $k$ topics $A_1, \ldots, A_k$ each associated with a distribution on words

(2) For each document, choose a single topic (pure topic model) or a mixture on topics (general topic models) and sample words from the associated distribution

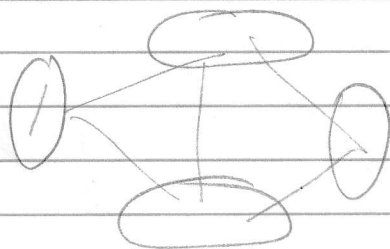Given many (short) documents, can we reconstruct the topics?

[Anandkumar, Hsu, Kakade] gave a polynomial time algorithm for pure topic models based on the identity

$$T_{ijk} \triangleq \mathbb{P}[\text{first three words are } (i,j,k) \text{ resp})$$

$$T = \sum_{j=1}^{k} \mathbb{P}[\text{topic} = j] \, A_j \otimes A_j \otimes A_j$$

Again, $T$ can be estimated from samples, and we can use tensor decomp.

In community detection, we have a hidden clustering



with a $k \times k$ matrix $R$ describing the connection probabilities.

Now we can partition the nodes into four sets
A, B, C, X and let

$$\Pi \in \{0,1\}^{n \times k} = \text{assignment of nodes to communities}$$

Then the $i^{th}$ column of $\Pi B$ denotes the probability a node in community $i$ has an edge to each of the rest of the nodes

Now let $(\Pi B)_i^A$ denote the restriction of the $i^{th}$ column to nodes in A, and consider

$$T = \sum_{i=1}^{k} p_i (\Pi R)_i^A \otimes (\Pi R)_i^B \otimes (\Pi R)_i^C$$

↑ fraction of nodes of community $i$ in X

Lemma: $T_{abc} = \mathbb{P}\left[ (x,a), (x,b), (x,c) \in E \right]$

"three star"

for a random $x \in X$, and over the randomness of G

Again, we have the same recipe

(1) Estimate a low rank tensor from counting three stars

(2) Use tensor decomp. to recover the community memberships