

Lecture 4+5 Classical Inference & Structure Learning on

Trees

Last time: undirected graphical models, conditional indep., testing Ising vs uniform.

What else do we want to with graphical models?

Separate 2 settings:

① Willing to assume the world is described by some **known** graphical model.

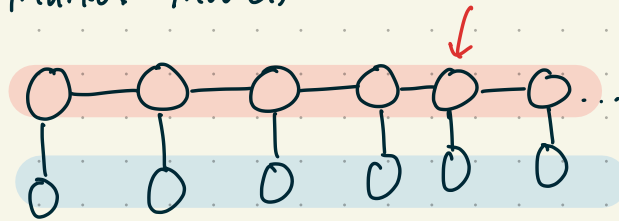
② Believe world is described by some **unknown** graphical model.

① Observe values of random variables corresponding to a subset of nodes. **Infer** something about observations / rest of graphical model.

Example models :

- Hidden Markov Models

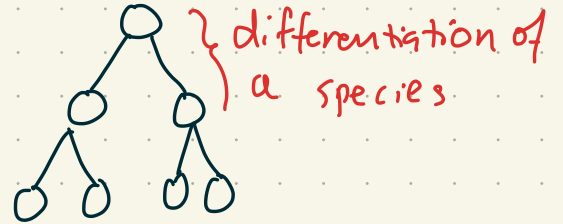
State of the world, evolving in time



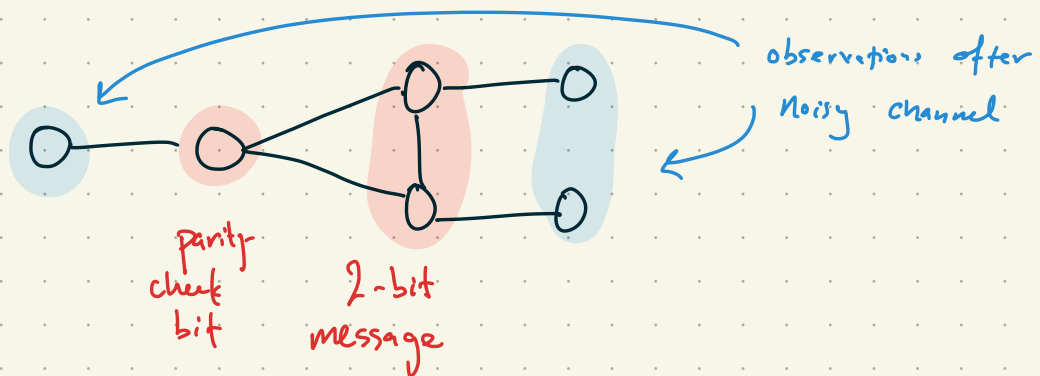
Observation / "emission"

Speech recognition

- Evolution of genomes



- Error-correcting codes



Example inference tasks:

- Compute $\Pr(\text{Observations}) = \sum_{x \text{ making observed values on observed nodes}} \Pr(x)$

- Compute marginal distribution $\Pr(x_A)$ for some $A \subseteq V$

- Compute conditional distribution $\Pr(x_A | x_B = y)$

"Posterior inference"

- Compute most likely x - "mode".

These are classical - a whole course at MIT, "Algorithms for

Inference". So we will only scratch the surface here, then move on to other topics.

Naive algorithms: involve sum or maximization over all possible values for x - intractable, b/c $(\# \text{ of vals per node})^{|V|}$ possible values

In general shouldn't hope to beat naive algs by much - NP, #P hardness

But, can do (much) better in special cases.

On **trees**, all of these can be solved by dynamic programming!

message passing, belief propagation, Viterbi, sum-product, max-product, junction tree, ...

Dynamic Programming for Marginal & conditional distributions

Computing marginal & conditional is \approx same - just question of

whether we fix values of some nodes by introducing potential

$$\delta(x) = \begin{cases} 1 & \text{if } x = \text{desired value} \\ 0 & \text{o.w.} \end{cases}$$

Problem: Given a tree T and factors $\{\gamma_c\}_{c \in \text{cliques of } T}$

compute $\{\Pr(x_v = y)\}_{v \in T, y}$ - ie all 1-wise marginals.

Assuming discrete distributions over universe Ω .

Observation: only cliques in T are edges + individual nodes.

Won't use this, but for intuition, therefore,

$$\Pr(x) \propto \prod_i \psi_i(x_i) \cdot \prod_{i,j} \psi_{i,j}(x_i, x_j)$$

i, j adjacent in tree

Root the tree at v . For any $u \in T$, let T_u be the graphical model we get by restricting to subtree rooted at u .

Let $x_v \in \Omega$.

$$\Pr(X_v = x_v) = \frac{1}{Z} \sum_{x \in \Omega^{T \setminus \{v\}}} \prod_i \psi_i(x_i) \prod_{i,j} \psi_{i,j}(x_i, x_j)$$

$$= \frac{1}{Z} \psi_v(x_v) \cdot \prod_{i \in \text{children}(v)} \left[\sum_{x_i \in \Omega} \psi_{i,v}(x_i, x_v) \cdot \underbrace{\psi_i(x_i) \cdot \sum_{x \in \Omega^{T_i}} \prod_{j \in T_i} \psi_j(x_j) \prod_{j \sim_{T_i} k} \psi_{j,k}(x_j, x_k)}_{\text{the things we would've multiplied to compute } \Pr_{T_i}(X_i = x_i)} \right]$$

the things we would've multiplied to compute $\Pr_{T_i}(X_i = x_i)$

If we knew \bullet , could compute \square

in $O(|\Omega|)$ time, and \square in $O(|\Omega| \cdot \text{degree})$ time.


Can use a dynamic program, computing \bullet for each choice of $x_i \in \Omega$ and each subtree T_i , where computation for T_i happens before its parent.

$$\text{Time: } O(n \cdot \text{degree} \cdot |\Omega|^2)$$

(compute Z by adding appropriate table entries.)

- Makes it look like would need $O(n^2)$ to compute all marginals, but there is a clever way to do all at once in same $O(n \cdot \text{degree} \cdot \log)$ time. (Can Google "sum product" or "Belief Propagation")

What happens on non-trees?

- Can view  as a "message" passed by x_i to its parent, constructed from similar "messages" it received from its children.
- Could use same formula for constructing messages and passing them around, but now on non-trees. "loopy BP".

Heuristic, sometimes seems ok in practice, maybe

expected to work if graph has no short cycles

("locally tree-like") and weak long-range

correlations.

- Can try other algs - MCMC, variational inference, ...

always heuristic, maybe w/ guarantees in special cases.

take "Algorithms for Inference" @ MIT.

Moving on to ②:

Learning Graphical Models

Assume getting samples $X_1 \dots X_n$ iid from some unknown graphical model.

What can we learn about it?

Fully-connected graph \Leftrightarrow represent any distribu. So need some assumptions.

2 learning tasks:

① TV learning

② Structure learning - find the underlying graph

- how to distinguish no edge, edge w/ $\psi_{ij} \approx 1$?

- need assumptions on ψ_{ij} 's.

Today: trees.

Chow-Liu (infinite sample version):

Instead of iid samples, let's pretend we get access to ^{marginal} \wedge distribution of every pair of variables X_i, X_j .

- Compute $I(X_i; X_j) = \mathbb{E}_{x_i, x_j} \log \frac{\Pr(x_i, x_j)}{\Pr(x_i) \Pr(x_j)} = KL(\{x_i, x_j\} \parallel \{x_i\} \otimes \{x_j\})$

- Let G be a graph where weight of edge ij is $I(X_i; X_j)$

- Output maximum spanning tree of G

Reminder: $= \arg \max_{\text{Tree } T \text{ on vertices of } G} \text{weight}(T)$

Theorem: Suppose T is a tree-structured graphical model. Then Chow-Liu, run on marginals of T , returns T . (Exception: if there is another tree T' which can represent same distn, can get T' - MST will not be unique.)

Proof: Follows from 2 key claims:

① If S is another tree-structured graphical model on the same set of variables, s.t. for every edge $i, j \in S$, $\{x_i, x_j\}_S = \{x_i, x_j\}_T$ and $\{x_i\}_S = \{x_i\}_T$ for all i , then joint dist'n of x_i, x_j under S

the distribution of S = distribution of T iff S is a maximum spanning tree in G .

② For every spanning tree S of G \exists a distribution which is Markov w.r.t. S satisfying hypotheses of ①.

So, let S be MST in T , define a dist'n Markov w.r.t. S as in ②.

Then that dist'n must = T .

Proof of ①: We have

$$KL(T \parallel S) = \mathbb{E}_{x \sim T} \log \frac{T(x)}{S(x)}$$

shorthand for $Pr_T(x)$

$$= \underbrace{\mathbb{E}_{x \sim T} \log T(x)}_{\text{indep of } S} - \mathbb{E}_{x \sim T} \log S(x)$$

$$= \mathbb{E}_{x \sim T} \log T(x) - \mathbb{E}_{x \sim T} \log \prod_i Pr_S(x_i | x_{\text{parent}_S(i)})$$

$$= \mathbb{E}_{x \sim T} \log T(x) - \mathbb{E}_{x \sim T} \log \prod_i Pr_T(x_i | x_{\text{parent}_S(i)})$$

$$= \mathbb{E}_{x \sim T} \log T(x) - \left(\mathbb{E}_{x \sim T} \sum_i \log \frac{Pr_T(x_i | x_{\text{parent}_S(i)})}{Pr_T(x_i) Pr_T(x_{\text{parent}_S(i)})} - \sum_i H(x_i) \right)$$

$$= \mathbb{E}_{x \sim T} \log T(x) + \sum_i H(x_i) - \sum_i I(x_i; x_{\text{parent}_S(i)})$$

If $\text{distn on } S = \text{distn on } T$, then this = 0. If not 0, then $\sum_i I(x_i; x_{\text{parent}_S(i)})$ must not be maximal. \square

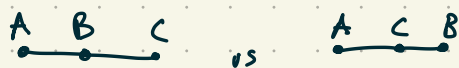
Proof of (2): define distribution via $\Pr_S(x) = \Pr_T(x_{\text{root}}) \cdot \prod_i \Pr_T(x_i | x_{\text{parent}_S(i)})$.

Marginals match by induction on depth.

What about finite samples?

- estimate $I(x_i; x_j)$ using empirical distns
- how accurately do we need them?

if interactions are really weak, might need a lot of samples to distinguish



2 options:

- add assumptions on interaction strength

- learn in TV - if $A \text{ --- } B \text{ --- } C$ vs $A \text{ --- } C \text{ --- } B$ hard to distinguish,

describe close-by distns conj.

also [Bhattacharyya - Gonen - Price - Vignodchandran]

Theorem [Daskalakis - Pan]: Assume alphabet $|\Omega| = 2$ (binary Ising model).

With $O(\frac{n \log n}{\epsilon^2})$ samples, Chow-Liu + empirical estimates for $I(\cdot; \cdot)$ learns a distribution which is ϵ -close to true one in TV dist.

(Aside: what "distribution" does Chow-Liu output? As described above it only gives the tree. Can get full distribution by estimating $\Pr(x_i | x_{\text{parent}(i)})$ from samples.)

This theorem is out of scope for this class. But we will discuss why

you can't beat $O(\frac{n \log n}{\epsilon^2})$ samples, returning to Le Cam's method from lecture 1.

[Koehler]

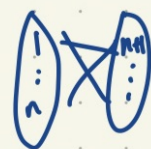
Theorem: To test between:

null: D is uniform on $\{\pm 1\}^n$

alternative: D is a free structured Ising model w/ $TV(D, \text{unif}) \geq 0.01$

requires $\Omega(n \log n)$ samples.

Idea:



always match
 $1 \dots n$ to $n+1 \dots 2n$

first, pick a random matching M on $\{1, 2, \dots, 2n\}$, from a set S of "allowed" matchings.

then, $X_1 \dots X_N \sim$ Ising model w/ $\Pr(X) \propto \exp\left(\frac{\beta}{n} \sum_{i \in n} x_i x_{M(i)}\right)$

Claim 1: If Hamming dist. between $M, M' \geq \Omega(n)$, then
 $TV(P_M, P_{M'}) \geq \Omega(1)$.

Proof: Consider the distribution of $\sum_i x_i x_{M(i)}$.

Under P_M , it is a sum of $\frac{n}{2}$ independent ± 1 bits,

$$\text{each w/ bias } \mathbb{E}_{x \sim P_M} x_i x_{M(i)} = \frac{\exp(\frac{\beta}{n}) - \exp(-\frac{\beta}{n})}{\exp(\frac{\beta}{n}) + \exp(-\frac{\beta}{n})} = \frac{\beta}{n} \pm O(\frac{1}{n})$$

Hence, $\frac{1}{n} \sum x_i x_{M(i)} \rightarrow N(\beta, O(1))$.

In particular, $\Pr\left(\frac{1}{n} \sum x_i x_{M(i)} \geq \beta\right) \rightarrow \frac{1}{2}$.

Under $P_{M'}$, bias is 0 for at least $\Omega(n)$ terms in the sum,

$$\text{so } \mathbb{E}_{P_{M'}} \frac{1}{n} \sum_i x_i x_{M'(i)} \leq (1 - \Omega(1)) \beta.$$

No longer a sum of independent terms, but variance is still $O(1)$.

$$\text{So, } \Pr_{P_{M'}}\left(\frac{1}{n} \sum x_i x_{M'(i)} \geq \beta\right) \leq \frac{O(1)}{\beta^2} \ll \frac{1}{2} \text{ if } \beta \gg 1. \quad \square$$

TV error ϵ w/ N samples, can identify underlying matching using N samples.

Now we need a new tool to show that identifying the underlying matching is not possible.

Fano's Inequality: Let M, X be joint random variables, M discrete taking values in finite set \mathcal{M} . Let $f(X) \in \mathcal{M}$ take values in \mathcal{M} . Then $\Pr(f(X) \neq M) \geq \frac{H(M) - I(M; X) - 1}{H(M)}$.

Intuition: if X doesn't contain much information about M , can't identify M using X .

How can we bound $I(M; X_1 \dots X_N)$?

Lemma $I(A; B) \leq \max_{a, a'} \text{KL}(\underbrace{\{B | A=a\}}_{\text{distribution of } B \text{ conditioned on } A=a} || \{B | A=a'\})$

Deferring proof of lemma for now, how do we use it?

$$KL(\{x_1, \dots, x_N | M\} \parallel \{x_1, \dots, x_N | M'\}) = N \cdot KL(\{x | M\} \parallel \{x | M'\})$$

by tensorization.

$$\text{Now, } KL(\{x | M\} \parallel \{x | M'\}) = \mathbb{E}_{x \sim P_M} \log \frac{\exp(\frac{\beta}{n} \sum_i x_i x_{M(i)})}{\exp(\frac{\beta}{n} \sum_i x_i x_{M'(i)})}$$

$$= \mathbb{E}_{x \sim P_M} \frac{1}{n} \sum_i x_i (x_{M(i)} - x_{M'(i)}) \leq O(1) \quad (\text{if } \beta = O(1))$$

Applying Fano, if we tried to use a function $f(x_1, \dots, x_N)$ to guess M , we would have $\Pr(f(x_1, \dots, x_N) \neq M) \geq 1 - \frac{O(N)}{\# \text{ possible matchings}}$

Fact: there is a set of $n^{2(n)}$ matchings all w/ Hamming dist $\geq \Omega(n)$

So, if $N = o(n \log n)$, can't identify M from x_1, \dots, x_N . \square

Loose ends: ① Proof of Fano's inequality

By data processing, $I(M; X) \geq I(M; f(X))$

$$= H(M) - H(M|f(X)) = (*)$$

Let E be a 0/1 r.v., $E = \begin{cases} 0 & \text{if } f(X) = m \\ 1 & \text{o.w.} \end{cases}$. Then $H(M|f(X)) = H(M, E|f(X))$,

$$\begin{aligned} \text{so } (*) &= H(M) - H(M, E|f(X)) \\ &= H(M) - H(E|f(X)) - H(M|E, f(X)) \\ &\geq H(M) - (1 - H(M|E, f(X))) \\ &= H(M) - (1 - \Pr(E=1) \cdot H(M|E=1, f(X)) - \underbrace{\Pr(E=0) \cdot H(M|E=0, f(X))}_{=0}) \\ &= H(M) - (1 - \Pr(E=1) \cdot H(M|E=1, f(X))). \end{aligned}$$

Rearranging, we get

$$\Pr(f(X) \neq m) = \Pr(E=1) \geq \frac{H(M) - 1 - I(M; X)}{H(M|E=1, f(X))} \geq \frac{H(M) - 1 - I(M; X)}{H(M)},$$

Since conditioning reduces information.

② Proof of Lemma:

$$I(A;B) = KL(\{A,B\} \parallel \{A\} \otimes \{B\})$$

$$= \mathbb{E}_{A,B} \log \frac{Pr(A,B)}{Pr(A)Pr(B)}$$

$$= \mathbb{E}_{a \sim A} \mathbb{E}_B \log \frac{Pr(B|a)}{Pr(B)}$$

$$= \mathbb{E}_{a \sim A} \mathbb{E}_B \log \frac{Pr(B|a)}{\mathbb{E}_{a' \sim A} Pr(B|a')}$$

$$= \mathbb{E}_{a \sim A} KL(\{B|a\} \parallel \mathbb{E}_{a' \sim A} \{B|a'\})$$

$$\leq \mathbb{E}_{a, a'} KL(\{B|a\} \parallel \{B|a'\})$$

$$\leq \max_{a, a'} KL(\{B|a\} \parallel \{B|a'\})$$

convexity of KL divergence!