

Graphical Models

powerful model for describing high-dimensional distributions...

... via their conditional independence structure

def: The Ising model is a distribution on $\{\pm 1\}^n$ described as

$$P[\sigma] = \frac{e^{-\beta H(\sigma)}}{Z_{\beta}}$$

where (1) $H(\sigma) = - \sum_{ij} J_{ij} \sigma_i \sigma_j - \sum_j h_j \sigma_j$

is the Hamiltonian

$$(2) Z_{\beta} = \sum_{\sigma} e^{-\beta H(\sigma)}$$

is the partition function

(3) and β is the inverse temperature

Properties

The case where $J_{ij} \geq 0$ is called ferromagnetic

i.e. the σ_i 's want to point in the same dir.

Similarly $J_{ij} \leq 0$ is called antiferromagnetic

one of the original motivations was to give a microscopic explanation of spontaneous magnetization:

At what temperature is the average magnetization $\frac{\sum \sigma_i}{n}$ typically far from zero?

Claim: As $B \rightarrow \infty$ the distribution is concentrated on the ground states i.e. minimizers of $H(\sigma)$

key property: Let $G = (V, E)$ where $V = [n]$

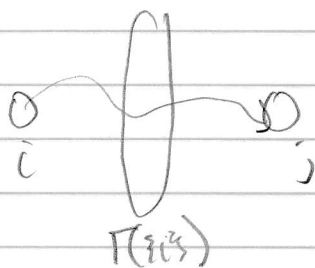
$$E = \{(i, j) \mid J_{ij} \neq 0\}$$

Proposition: Let $T = \Gamma(\xi_i's) \cup \xi_i's$. Then neighbors

$$\sigma_i \perp \sigma_{V \setminus T} \mid \sigma_{T(\xi_i's)}$$

"once you know the spins of the neighbors of i , its spin is independent of everything else"

Thus independence relations follow from graph separation



all paths are blocked

Proof: If we fix $\sigma_{T \setminus \{i\}}$ the new Hamiltonian is

$$- [\underbrace{\sigma_i, \dots}_{\sigma_j \text{'s for } j \in V \setminus i}] \begin{bmatrix} J_{ii} & 0 \\ 0 & \tilde{J} \end{bmatrix} \begin{bmatrix} \sigma_i \\ \vdots \end{bmatrix} - [\sigma_i, \dots] \begin{bmatrix} h_i \\ \vdots \end{bmatrix} + C$$

$$= H_i(\sigma_i) + H_{V \setminus i}(\sigma_{V \setminus i})$$

Hence the distribution factorizes \square

Takeaway: If you know the independence structure, you can ^{often} write the distribution more concisely than specifying $2^n - 1$ values

Algorithmic questions

① Given an Ising model, can you sample from it?

② Can you compute the posterior, given a partial assignment?

i.e. inference

③ Can you learn it from samples?

In many situations ① \approx ②

Hardcore model

Given a bounded degree graph $G=(V,E)$
 $\Delta = \text{max degree}$

define

$$P[\mathcal{I}] = \frac{\lambda^{|\mathcal{I}|}}{Z_\lambda}$$

↑
independent sets

Thm [Jerrum, Sinclair] [Weitz] If $\Delta \geq 3$ then
for any $\lambda < \lambda_c(\Delta) \triangleq \frac{(\Delta-1)^{\Delta-1}}{(\Delta-2)^\Delta}$

there is an efficient algorithm for sampling

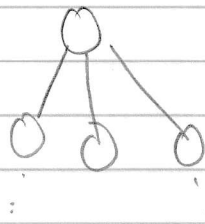
Thm [Sly, Sun] If $\Delta \geq 3$ then for any
 $\lambda > \lambda_c(\Delta)$

It is NP-hard to approximately sample

what is happening at $\lambda_c(\Delta)$?

Tree Uniqueness

Consider a $(\Delta-1)$ -ary tree



degree Ising model

Note: The constants depend on upper/lower bounds on non zero J_{ij} 's and upper bounds on $|J_{ij}|$'s, which is necessary

Takeaway: You can learn even when you can't sample from it

Warm-Up: Hardcore Model

def: A distribution \mathcal{D} on S^n is δ -unbiased

if for any i , and assignment $x \in S^{n-1}$ we have

$$\mathbb{P}[X_i = \alpha \mid X_{-i} = x] \geq \delta \quad \forall \alpha \in S$$

i.e. all variables have some randomness

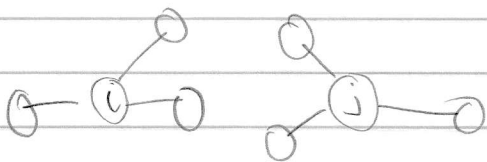
closed
under
taking
marginals

Proposition: If the hardcore model is δ -unbiased, can learn its edges

Proof: (sketch) Consider any pair i, j

① If ever $X_i = X_j = 1 \Rightarrow (i, j) \notin E$

② Conversely if $(i, j) \notin E$ consider



Fix a partial assignment where none of the neighbors of i or j are in I . Then

$$P[X_i = X_j = 1 \mid X_{-ij} = x] \geq \delta^2 \quad \square$$

Finally note we can prove a lower bound on δ in terms of Δ and an upperbound on λ

Now consider general Ising models, define

$$\text{width } \lambda = \max_{i,j} \left(\sum_j |A_{ij}| + |h_i| \right)$$

$$\text{and let } \eta = \min_{i,j \text{ st. } A_{ij} \neq 0} |A_{ij}|$$

We'll follow [Wu, Sanghavi, Dimakis]

def: The logistic function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Key Fact: The conditional distribution in an Ising model is logistic, i.e.

$$P[X_i = 1 \mid X_{-i} = x] = \sigma(\langle w, x' \rangle)$$

↑ padded
[x, 1]

Moreover $\|w\|_1 \leq 2\lambda$

Proof: By definition

$$P[X_i=1 | x_i=x] =$$

$$\frac{\exp(\sum_{j \neq i} A_{ij} x_j + h_i)}{\exp(\sum_{j \neq i} A_{ij} x_j + h_i) + \exp(-\sum_{j \neq i} A_{ij} x_j - h_i)}$$

$$\exp(\sum_{j \neq i} A_{ij} x_j + h_i) + \exp(-\sum_{j \neq i} A_{ij} x_j - h_i)$$

Thus we find

$$w = 2[A_{i1}, A_{i2}, \dots, A_{i(n-1)}, A_{i(n)}, \dots, A_{in}, h_i]$$

and so $\|w\|_1 \leq 2\lambda$. \square

Thus we can try to learn the parameters by logistic regression

$$\hat{w} = \underset{\hat{w}}{\operatorname{argmin}} \frac{1}{N} \sum_{s=1}^N \ell(y_s \langle \hat{w}, x'_s \rangle) \quad (*)$$

#samples label of i

where $\ell \equiv$ negative log-likelihood

$$= \ln(1 + e^{-y_s \langle \hat{w}, x'_s \rangle}) = \begin{cases} -\ln \sigma(\langle \hat{w}, x'_s \rangle) & \text{if } y_s = 1 \\ -\ln(1 - \sigma(\langle \hat{w}, x'_s \rangle)) & \text{else} \end{cases}$$

claim (*) is a convex optimization problem

Note: This would work, but would require

$m \approx n$ samples

and we can do much better, i.e. $m \approx C_{x, \Delta} \log n$

New approach, use l_1 -regularization

$$\hat{w} = \underset{\tilde{w}}{\operatorname{argmin}} \frac{1}{N} \sum_{s=1}^N \ell(y_s \langle \tilde{w}, x'_s \rangle)$$

$$\text{s.t. } \|\hat{w}\|_1 \leq 2\lambda$$

Main ingredients

(1) Generalization bounds for l_1 :

Lemma Suppose ℓ has Lipschitz constant L and \mathcal{D} is a distribution on $X \times Y$ with $X = \{x \mid \|x\|_\infty \leq R\}$. Then for any $w \in W$

$$L(w) \leq \hat{L}(w) + 2LRD \sqrt{\frac{2 \ln(2n)}{N}}$$

\uparrow true loss \uparrow empirical loss

$$+ LRD \sqrt{\frac{2 \ln(2/\delta)}{N}}$$

failure probability

and $W = \{w \mid \|w\|_1 \leq W\}$

Intuition: Instead of needing an ϵ -net over the unit ball in \mathbb{R}^n , just need a ϵ -net over the W -dilation of the l_1 ball, which has size $\sim \binom{n}{W}$

② Strong convexity: For any \hat{w} , want to show

$$L(\hat{w}) - L(w) \geq \frac{1}{n} \sum_{x'} \left[\left(\sigma(\langle \hat{w}, x \rangle) - \sigma(\langle w, x \rangle) \right)^2 \right]$$

Intuition: Consider two Bernoulli distributions with parameters a, b respectively. Then

$$d_{TV}(\text{Ber}(a), \text{Ber}(b)) = |a - b|$$

and the KL-divergence btwn general p, q is

$$d_{KL}(p||q) = \sum_x p(x) \ln \frac{p(x)}{q(x)}$$

which turns into

$$d_{KL}(\text{Ber}(a) || \text{Ber}(b)) = a \ln \left(\frac{a}{b} \right) + (1-a) \ln \left(\frac{1-a}{1-b} \right)$$

Fact [Pinsker] $d_{TV}^2 \leq \frac{1}{2} d_{KL}$

Finally since our loss function is the negative log-likelihood we have

$$L(\hat{w}) - L(w) = \frac{1}{n} \sum_{x'} \left[d_{KL}(\sigma(\langle \hat{w}, x \rangle) || \sigma(\langle w, x \rangle)) \right]$$

③ unbiasedness

Proposition: Suppose \mathcal{D} is δ -unbiased. Then

$$\mathbb{E}_{x \sim \mathcal{D}} \left[\underbrace{(\sigma(\langle w, x \rangle + h))}_{\sigma} - \underbrace{(\sigma(\langle w', x \rangle + h'))}_{\sigma'} \right]^2 \leq \epsilon \quad (\square)$$

where $\epsilon \lesssim \delta e^{-2\|w\|_1 - 2|h'|}$ then

$$\|w - w'\|_0 \lesssim e^{\|w\|_1 + |h'|} \sqrt{\frac{\epsilon}{\delta}}$$

Proof: By subconditioning $(\square) \Rightarrow$

$$\epsilon \geq \mathbb{E}_{x_i} \left[\mathbb{E}_{x_c | x_i} \left[(\sigma - \sigma')^2 \right] \right]$$

$$= \mathbb{E}_{x_i} \left[\left(\sigma \Big|_{x_c=+1} - \sigma' \Big|_{x_c=+1} \right)^2 \mathbb{P}[x_c=+1 | x_i] \right]$$

$$+ \left(\sigma \Big|_{x_c=-1} - \sigma' \Big|_{x_c=-1} \right)^2 \mathbb{P}[x_c=-1 | x_i]$$

By δ -unbiasedness

$$\geq \delta \mathbb{E}_{x_i} \left[\left(\sigma \Big|_{x_c=+1} - \sigma' \Big|_{x_c=+1} \right)^2 + \left(\sigma \Big|_{x_c=-1} - \sigma' \Big|_{x_c=-1} \right)^2 \right]$$

Now by the Fact we have

Fact [Klivans, Meka]

$$|\sigma(a) - \sigma(b)| \geq e^{-|a|-3} \min(1, |a-b|)$$

$$\geq \delta e^{-2\|w\|, -2|h'|}$$

$$\mathbb{E}_{X_i} \left[\min(1, (\underbrace{\langle w, X_i^{i+} \rangle + h}_{\substack{\downarrow \\ X_i \text{ set to } +1}} - \underbrace{\langle w', X_i^{i+} \rangle - h'}_{\substack{\downarrow \\ \text{same, but } X_i \text{ set to } -1}})^2) \right]$$

$$+ \dots$$

$$\geq \delta e^{-2\|w\|, -2|h'|} \mathbb{E}_{X_i} \left[\min(1, 2(w_i - w'_i)^2) \right]$$

which follows b/c

$$\min(1, a^2) + \min(1, b^2) \geq \min(1, \frac{(a-b)^2}{2})$$

This bound holds for all i , and rearranging algebraically completes the proof. \square

Remark: Similar algorithms work for Markov random fields

$$\mathbb{P}[X=x] = \frac{e^{\beta P(x)}}{Z_\beta} \quad \leftarrow \text{degree } \leq d \text{ polynomial}$$

There are $C_{\Delta, d} n^{O(d)}$ time algorithms

Moreover there are $n^{\Omega(d)}$ lower bounds based on sparse parity with noise

Sufficient Statistics

Suppose we are given samples

$$X_1, \dots, X_N \sim P_\theta(x)$$

Is there a sufficient statistic

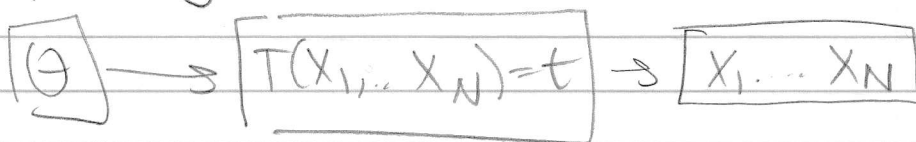
- i.e. we can compress to $T(X_1, \dots, X_N)$
w/o losing any information?

Factorization Theorem [Neyman, ...]

A statistic is sufficient iff

$$P_\theta(x_1, \dots, x_N) = u(x_1, \dots, x_N) v(T(x_1, \dots, x_N), \theta)$$

Graphically this means



i.e. $X_1, \dots, X_N \perp\!\!\!\perp \theta \mid T(X_1, \dots, X_N) = t$

There is a canonical way to satisfy this condition

def. An exponential family has the form

$$P_\theta(x) = \frac{h(x) e^{\langle \theta, T(x) \rangle}}{Z(\theta)}$$

e.g. for the Ising model, we can take

$$T(x) = [\text{vec}(xx^T), x]$$

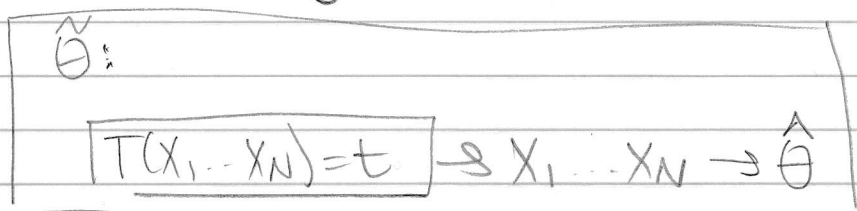
$$\Theta = [\text{vec}(A), h]$$

Corollary: For an exponential family, any estimator

$$\hat{\Theta}(x_1, \dots, x_N)$$

can be turned into another one $\hat{\Theta}(T(x_1, \dots, x_N) = t)$

Problem: Looking inside



But sampling can be hard

Thm [Montanari] [Bresler et al] There are graphical models that can be efficiently learned, but not if you reduce to sufficient statistics

Open: Are there computational-vs-statistical tradeoffs for learning exponential families?