

Lec 19: Learning Graphical Models

- Parameters learning:
Graph fixed, learning params of BN, MN

- Structure learning:
find the graph i.e. how variables depend on each other

In PAC learning theory, examples labelled by unknown function.
Goal was to find that function

We have some samples, we want to learn the joint distribution

Maximum likelihood

(c) iid
dataset D sampled from p^*
Goal: find dist p which is "close" to p^*
clustering of two distributions:

KL divergence

$$D_{KL}(p^* || p) = \sum_x p^*(x) \log(p^*(x)/p(x))$$

$$= \sum_x p^*(x) \log p^*(x) - \sum_x p(x) \log p(x)$$

$$= -H(p^*) - \mathbb{E}_{x \sim p^*} \text{entropy}$$

$- p^*$ is unknown, we cannot compute $H(p^*)$

- Same as maximize $\mathbb{E}_{x \sim p^*} \log p(x)$

- $\mathbb{E}_{x \sim p^*} \log p(x)$ is "log likelihood"
find dist p that maximizes the probability of the dataset

- Compute by taking average over samples

$$\mathbb{E}_{x \sim p^*} \log p(x) \approx \frac{1}{|D|} \sum_{x \in D} \log p(x)$$

D is iids from p^*

Max. Lik. learning

$$\max_w \frac{1}{|D|} \sum_{x \in D} \log p(x)$$

Empirical Distribution

$$\tilde{p}(x) = \frac{1}{|D|} \sum_{x \in D} \delta(x, x)$$

Let w be parameters of the dist. to be learnt.

Theorem: $\arg\max_w \frac{1}{N} \sum_{i=1}^N \log \tilde{p}(x_i | w)$

$$= \arg\min_w D_{KL}(\tilde{p}(x) || p(w))$$

Proof: $RHS = \sum_x \tilde{p}(x) \log \tilde{p}(x)$ $\stackrel{\text{does not depend on } w}{\sim}$

$$= \max_w \sum_x \tilde{p}(x) \log p(x | w)$$

$$= \max_w \sum_x \frac{1}{|D|} \sum_{y \in D} \delta(x, y) \log p(y | w)$$

$$= \max_w \frac{1}{|D|} \sum_{y \in D} \log p(y | w)$$

General loss function

$$\mathbb{E}_{x \sim p^*} [L(x, p)]$$

$$L(x, p) = -\log p(x)$$

Conditional Random Fields (CRFs)

x, y

Interested in finding $p(y|x)$

$$L(x, y, p) = -\log p(y|x)$$

Max. Lik. learning for BN

$$p_\theta(x) = \prod_{j=1}^n \phi_j(x_j | p_{\theta(j)})$$

$$D = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$$

$$L(\theta, D) = -\sum_{i=1}^m \log p_\theta(x^{(i)})$$

$$= -\sum_{i=1}^m \sum_{j=1}^n \log \phi_j(x_j^{(i)} | p_{\theta(j)})$$

$$= -\sum_{j=1}^n \left(\sum_{i=1}^m \sum_{x_j \in p_{\theta(j)}} \log \phi_j(x_j^{(i)} | p_{\theta(j)}) \right)$$

$$\text{minimized when } \sum_j \#(x_j, x_{p_{\theta(j)}}) \log \phi_j(x_j | p_{\theta(j)})$$

minimized at

$$\theta^*_{x_j | p_{\theta(j)}} = \#(x_j, x_{p_{\theta(j)}}) / \#x_{p_{\theta(j)}}$$

empirical value
for $p_{x_j | p_{\theta(j)}}$

Learning in Markov Networks

$$p(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \exp(\theta^\top f(x))$$

Suppose there are m edges,
 $\Rightarrow 4m$ parameters

$$\phi_{12}(0,0) = \theta_{12,00}$$

$$\phi_{12}(0,1) = \theta_{12,01}$$

$$\vdots$$

$$\phi_{23}(0,0) = \theta_{23,00}$$

$$\vdots$$

$$\phi_{23}(1,1) = \theta_{23,11}$$

$$\Rightarrow \frac{1}{Z(\theta)} \exp \left(\sum_{e \in E} \log \phi_e(x_e; \theta) \right)$$

$$= \frac{1}{Z(\theta)} \exp(\theta^\top f(x))$$

4m dim vector

$$x = (0 \ 1 \ 0)$$

0	1	0
0	0	1
0	0	0

$$y$$

0	0	0
0	0	1
0	1	0

$$\begin{aligned} f(0 \ 1 \ 0) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1 + \theta_2 + \theta_3 \\ &= \theta_0 + \theta_1 + \theta_2 + \theta_3 \end{aligned}$$

$$\begin{aligned} \hat{p}(0 \ 1 \ 0) &= \frac{1}{Z(\theta)} e^{(\theta_0 + \theta_1 + \theta_2 + \theta_3)} \\ &\text{f(x) : sufficient statistic} \\ &\text{"exponential families"} \\ \text{Log Likelihood} &= \log p(0, \theta) \\ (\text{LL}) &= \frac{1}{|D|} \sum_{x \in D} \log p(x; \theta) \end{aligned}$$

$$= \left(\frac{1}{|D|} \sum_{x \in D} \theta^\top f(x) \right) - \log Z(\theta)$$

linear in θ

Learning by gradient descent

$$\nabla_\theta (\text{LL}) = \frac{1}{|D|} \sum_{x \in D} f(x) - \nabla_\theta \log Z(\theta)$$

θ^* sample mean of sufficient statistics

$$\log Z(\theta) = \log \sum_x \exp(\theta^\top f(x)) = \frac{(\theta^\top f(x))}{e^{\theta^\top f(x)}}$$

$$\nabla_\theta \log Z(\theta) = \frac{1}{Z(\theta)} \sum_x \exp(\theta^\top f(x)) f(x)$$

$$= \frac{1}{Z(\theta)} \sum_x \frac{\exp(\theta^\top f(x))}{\exp(\theta^\top f(x))} f(x) = \frac{1}{Z(\theta)} \sum_x f(x) = \mathbb{E}_{x \sim p(\theta)} f(x)$$

expected value of $f(x)$ according to $p(\theta)$

When $\nabla_\theta (\text{LL}) = 0$

sample mean of sufficient statistics = mean according MN of sufficient statistic

"Moment matching of sufficient statistic."

$$\mathbb{E}_\theta f(x) \approx \frac{1}{|D|} \sum_{x \in D} f(x)$$

Can generate samples from $p(\theta)$
by MCMC method