# PROBABILISTIC GRAPHICAL MODELS.

## Lecture 3: LEARNING PARAMETERS IN GRAPHICAL MODELS.

There are two aspects in the subject of learning graphical models: parameter learning and structure (of graph) learning. We concentrate on parameter learning (from samples) which is easier. In this lecture we suppose that all variables in the samples are "observed" or "visible". In lect 9 we will treat the case where some of the variables are "not accessible" or "hidden".

## I. The Kullback-Leibler divergence.

Let $p(\underline{x})$ and $q(\underline{x})$ two probability distributions over a discrete alphabet $\underline{x} = (x_1, \dots x_K) \in \mathcal{A}^K$ where $x_i \in \mathcal{A}$. By definition:

$$KL(p \| q) = \sum_{\underline{x}} p(\underline{x}) \log p(\underline{x}) - \sum_{\underline{x}} p(\underline{x}) \log q(\underline{x})$$

$$= \sum_{\underline{x}} p(\underline{x}) \log \left\{ \frac{p(\underline{x})}{q(\underline{x})} \right\}.$$

We also use the notation:

$$KL(p \| q) = \mathbb{E}_p \left[ \log p(x) \right] - \mathbb{E}_p \left[ \log q(x) \right]$$

$$= \mathbb{E}_p \left[ \log \left\{ \frac{p(x)}{q(x)} \right\} \right]$$

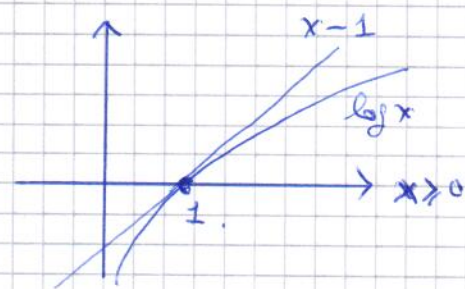or sometimes $\left\langle \log \frac{p(x)}{q(x)} \right\rangle_p$ .

## Main properties:

$$KL(p \| q) \geq 0 \quad \& \quad = 0 \quad \text{iff} \quad p(x) = q(x) \quad \forall x.$$

$$KL(p \| q) \neq KL(q \| p) \quad \text{Not symmetric.}$$

## Proof of positivity:

$$\log x \leq x - 1 \quad \text{for} \quad x \geq 0$$



$$\Rightarrow \log \frac{q(x)}{p(x)} \leq \frac{q(x)}{p(x)} - 1$$

$$\Rightarrow 1 - \frac{q(x)}{p(x)} \leq -\log \frac{q(x)}{p(x)} = \log \frac{p(x)}{q(x)}$$

$$\Rightarrow p(x) - q(x) \leq p(x) \log \frac{p(x)}{q(x)}$$

$$\Rightarrow \underbrace{\sum_x p(x)}_{1} - \underbrace{\sum_x q(x)}_{1} \leq \sum_x p(x) \log \frac{p(x)}{q(x)}$$

$$\Rightarrow 0 \leq \sum_x p(x) \log \frac{p(x)}{q(x)} = KL(p \| q)$$

Alternative proof by Jensen's inequality.

$$KL(p\|q) = \sum_{\underline{x}} p(\underline{x}) \log \frac{p(\underline{x})}{q(\underline{x})} = - \sum_{x} p(\underline{x}) \log \frac{q(\underline{x})}{p(\underline{x})}.$$

Note that $(-\log z)$ is a convex fct of $z > 0$ :



By convexity we have that the mean of $(-\log z)$ is bigger than the $(-\log)$ of the mean (see picture).

$$\mathbb{E}_z(-\log z) \geq - \log \mathbb{E}_z(z).$$

i.e
$$\sum_{x} p(\underline{x}) \left(- \log \frac{q(\underline{x})}{p(\underline{x})}\right) \geq - \log \sum_{x} p(\underline{x}) \frac{q(\underline{x})}{p(\underline{x})}$$

$$= - \log \sum_{x} q(\underline{x})$$

$$= - \log 1 = 0.$$

We have Thus found $\qquad KL(p\|q) \geq 0$

# II. Maximum Likelihood Method and KL minimization

Given a set of samples $\underline{x}^{(1)} \ldots \underline{x}^{(N)}$ from distribution $p(\underline{x} \mid \vartheta)$ where $\vartheta$ denotes the set of parameters of $p$ (say weights, biases in a Boltzman machine ...) The log-likelihood of the data is by definition:

$$L(\vartheta) = \log \text{Prob}(\underline{x}^{(1)} \ldots \underline{x}^{(N)})$$

$$= \log \left\{ \prod_{m=1}^{N} p(\underline{x}^{(m)} \mid \vartheta) \right\}$$

under iid
assumption for
data samples

$$= \sum_{m=1}^{N} \log p(\underline{x}^{(m)} \mid \vartheta)$$

### Maximum likelihood principle:

Set estimate $\hat{\vartheta} = \underset{\vartheta}{\arg\max} \, L(\vartheta)$.

### KL minimization:

Set $\hat{\vartheta} = \underset{\vartheta}{\arg\min} \, KL(q_{emp} \parallel p)$

where $q_{emp}(\underline{x}) = \frac{1}{N} \sum_{m=1}^{N} \delta_{\underline{x}, \underline{x}^{(m)}}$

is the empirical distr of the data.

Claim : ML maximization $\iff$ KL minimization.

Proof : $$KL(q_{emp} \| p) = \mathbb{E}_{q_{emp}} \log\left\{\frac{q_{emp}(\underline{x})}{p(\underline{x} | \vartheta)}\right\}$$

$$= \mathbb{E}_{q_{emp}}\left(\log q_{emp}(\underline{x})\right) - \underbrace{\mathbb{E}_{q_{emp}}\left(\log p(\underline{x} | \vartheta)\right)}$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\text{independent of } \vartheta}$$

Note that $\mathbb{E}_{q_{emp}}\left(\log p(\underline{x})\right) = \frac{1}{N} \sum_{m=1}^{N} \log p(\underline{x}^{(m)} | \vartheta)$.

$$= \frac{1}{N} L(\vartheta).$$

So :

$$KL(q_{emp} \| p) = \mathbb{E}_{q_{emp}}\left(\log q_{emp}(\underline{x})\right) - \frac{1}{N} L(\vartheta).$$

$$\Rightarrow \quad \max_{\vartheta} KL \quad \text{is equiv to} \quad \min_{\vartheta} L(\vartheta) ,$$

## III. ML Training of Belief Networks.

Recall for a BN we have a prob distr of the form

$$p(\underline{x}) = \prod_{i=1}^{K} p(x_i \mid (pa)_i)$$

where $(pa)_i$ are the variables which are parents of $x_i$. We assume that each $p(x_i \mid (pa)_i)$ depends on a set of parameters that we call $\theta_{i \mid (pa)_i}$.

We apply the ML principle or equivalently we want to minimize $KL(q_{emp} \| p)$ where $q_{emp}$ is the empirical distribution of data samples $\underline{x}^{(1)} \ldots \underline{x}^{(N)}$.

---

**Lemma:** $KL(q_{emp} \| p)$ is minimized

for $p(x_i \mid (pa)_i) = q_{emp}(x_i \mid (pa)_i)$.

---

In practice we have

$$q_{emp}(x_i = s \mid (pa)_i = t) = \frac{\sum_{m=1}^{N} \mathbb{1}(x_i^{(m)} = s ; pa_i^{(m)} = t)}{\sum_{m=1}^{N} \mathbb{1}((pa)_i^{(m)} = t)}$$

which is therefore known. We solve for $\theta_{i|\text{pa}_i}$ from the equation of the Lemma.

## Proof of Lemma.

$$KL(q_{emp} \| p) = \mathbb{E}_{q_{emp}}(\log q_{emp}) - \mathbb{E}_{q_{emp}}(\log p(\mathbf{x}))$$

$$= \mathbb{E}_{q_{emp}}(\log q_{emp}) - \mathbb{E}_{q_{emp}}\left(\sum_{i=1}^{K} \log p(x_i | \text{pa}_i)\right)$$

$$= \mathbb{E}_{q_{emp}}(\log q_{emp}) - \sum_{i=1}^{K} \underbrace{\mathbb{E}_{q_{emp}}\left(\log p(x_i | \text{pa}_i)\right)}$$

$$\mathbb{E}_{q_{emp}(x_i, \text{pa}_i)}\left(\log p(x_i | \text{pa}_i)\right)$$

because $p(x_i | \text{pa}_i)$ depends only on $(x_i, \text{pa}_i)$.

Now we add and subtract an appropriate term:

$$KL(q_{emp} \| p) = \mathbb{E}_{q_{emp}}(\log q_{emp}) - \sum_{i=1}^{K} \mathbb{E}_{q_{emp}(x_i, \text{pa}_i)}\left[\log q_{emp}(x_i | \text{pa}_i)\right]$$

$$- \left\{ \sum_{i=1}^{K} \mathbb{E}_{q_{emp}(x_i, \text{pa}_i)}\left[\log p(x_i | \text{pa}_i)\right] - \mathbb{E}_{q_{emp}(x_i, \text{pa}_i)}\left[\log q_{emp}(x_i | \text{pa}_i)\right] \right\}$$

The first two terms are independent of the parameters $\theta_{i|pa_i}$ and play the role of a constant when we minimize.

For the last two terms we remark that by Bayes law:

$$\mathbb{E}_{q_{emp}(x_i, (pa)_i)} = \mathbb{E}_{q_{emp}(pa_i)} \mathbb{E}_{q_{emp}(x_i|pa_i)}$$

Thus

$$KL(q_{emp} \| p) = \text{constant} +$$

$$\sum_{i=1}^{k} \mathbb{E}_{q_{emp}(pa_i)} \left\{ \mathbb{E}_{q_{emp}(x_i|pa_i)} \log q_{emp}(x_i|pa_i) \right.$$

$$\left. - \mathbb{E}_{q_{emp}(x_i|pa_i)} \log p(x_i|pa_i) \right\}$$

$$= \text{constant} + \sum_{i=1}^{K} \mathbb{E}_{q_{emp}(pa_i)} \left[ KL\left( q_{emp}(x_i|pa_i) \| p(x_i|pa_i) \right) \right].$$

The $KL \geq 0$ is minimized (vanishes) for a set of parameters $\theta_{i|pa_i}$ such that

$$p(x_i|pa_i) = q_{emp}(x_i|pa_i)$$

# IX ML Training of MRF or Factor graph models.

Here $\quad p(\underline{x}|\vartheta) = \frac{1}{Z(\vartheta)} \prod_C \psi_C(x_C|\vartheta_c)$.

For iid samples $(\underline{x}^{(1)} \cdots \underline{x}^{(N)})$ we have

$$L(\vartheta) = \sum_{m=1}^{N} \log p(\underline{x}^{(m)})$$

$$= \sum_C \sum_{m=1}^{N} \log \psi_C(x_C^{(m)}|\vartheta_c) - N \log Z(\vartheta).$$

where

$$Z(\vartheta) = \sum_{\underline{x} \in \mathcal{A}^K} \prod_C \psi_C(x_C|\vartheta_c).$$

Now $\log Z(\vartheta)$ is intractable and all parameters are coupled. One can use <u>gradient ascent</u> in order to <u>maximize</u> $L(\vartheta)$. This involves computing $\nabla_\vartheta L(\vartheta)$.

## Computation of $\nabla_\vartheta L(\vartheta)$: For $\vartheta_c$ we have

$$\nabla_{\vartheta_c} L(\vartheta) = \underbrace{\sum_{m=1}^{N} \nabla_{\vartheta_c} \log \psi_C(x_C^{(m)}|\vartheta_c)}_{\text{easy and explicit}} - \underbrace{N \nabla_{\vartheta_c} \log Z(\vartheta)}_{?}$$

$$\nabla_{\vartheta_c} \log Z(\vartheta) = \frac{1}{Z(\vartheta)} \nabla_{\vartheta_c} Z(\vartheta)$$

$$= \frac{1}{Z(\vartheta)} \sum_x \nabla_{\vartheta_c} \left\{ \prod_{c'} \psi_{c'}(x_{c'} | \vartheta_{c'}) \right\}$$

$$= \frac{1}{Z(\vartheta)} \sum_x \nabla_{\vartheta_c} \psi_c(x_c | \vartheta_c) \cdot \underbrace{\prod_{c' \neq c} \psi_{c'}(x_{c'} | \vartheta_{c'})}_{\dfrac{\prod_{c'} \psi_{c'}(x_{c'} | \vartheta_{c'})}{\psi_c(x_c | \vartheta_c)}}$$

$$= \frac{1}{Z(\vartheta)} \sum_x \left\{ \frac{\nabla_{\vartheta_c} \psi_c(x_c | \vartheta_c)}{\psi_c(x_c | \vartheta_c)} \right\} \left\{ \prod_c \psi_c(x_c | \vartheta_c) \right\}.$$

$$= \left\langle \nabla_{\vartheta_c} \log \psi_c(x_c | \vartheta_c) \right\rangle$$

where $\left\langle A(\underline{x}) \right\rangle \equiv \dfrac{1}{Z} \sum_x A(x) \prod_c \psi_c(x_c)$ is the

standard notation for Gibbs / MRF averages — note that

$$\left\langle \nabla_{\vartheta_c} \log \psi_c(x_c | \vartheta_c) \right\rangle = \underbrace{\mathbb{E}}_{P(x_c)} \left[ \nabla_{\vartheta_c} \psi_c(x_c | \vartheta_c) \right]$$

Marginal of $p(\underline{x} | \vartheta)$ over all
variables $(x_1 \cdots x_k) \setminus x_c$.

Summarizing we have for all cliques $C$ on

factor nodes $C$ :

$$\nabla_{\vartheta_C} L(\vartheta) = \sum_{m=1}^{N} \nabla_{\vartheta_C} \log \psi_C(x_C^{(m)}|\vartheta_C) - N \langle \nabla_{\vartheta_C} \log \psi_C(x_C|\vartheta_C) \rangle$$

$\underbrace{\phantom{xxxxxxxxxxxxxx}}$ easy

$\underbrace{\phantom{xxxxxxxxxxxxxx}}$ requires marginalisation

use Message passing,

or sampling, ...

(difficult in general).

## Example : Boltzman Machine or Ising Model.

$$p(\underline{x}) = \frac{1}{Z(W)} e^{+\frac{1}{2} \underline{x}^T W \underline{x}}$$

where $(W)_{ij}$ is a weight

matrix. (say $W_{ii} = 0$).

We have after application of above method (exercise) :

$$\frac{\partial L}{\partial W_{ij}} = \sum_{m=1}^{N} \left( x_i^{(m)} x_j^{(m)} - \langle x_i x_j \rangle \right) . \qquad \text{where}$$

$$\langle x_i x_j \rangle = \sum_{\underline{x}} x_i x_j \, p(\underline{x}) = \frac{1}{Z} \sum_{\underline{x}} x_i x_j \, e^{+\frac{1}{2} \underline{x}^T W \underline{x}}$$

$\underbrace{\phantom{xxxxxxxxxxxxxx}}$ difficult to compute in general.