# PROBABILISTIC GRAPHICAL MODELS

## Lecture 4 : LEARNING PARAMETERS WITH HIDDEN VARIABLES.

In this chapter we consider learning parameters from iid samples of the form

$$\left\{ \underline{v}^{(1)} \underline{h}^{(1)} ; \underline{v}^{(2)} \underline{h}^{(2)} ; \dots ; \underline{v}^{(N)} \underline{h}^{(N)} \right\}$$

$$\left\{ \underline{x}^{(1)} ; \underline{x}^{(2)} ; \dots ; \underline{x}^{(N)} \right\} .$$

where $v^{(1)} \dots v^{(N)}$ are accessible or <u>visible</u> and $\underline{h}^{(1)} \dots \underline{h}^{(N)}$ are not accessible or <u>hidden</u> (or erased).

It turns out that directly maximizing the log-likelihood is intractable because the marginalization over hidden variables is difficult. We will use so-called <u>variational methods</u>.

<u>The setting</u> is as follows : we suppose the data comes from a probabilistic model with visible and hidden variables

$$P(\underline{V}, \underline{h}, \vartheta) = p(\underline{v}, \underline{h} \mid \vartheta) \, p_0(\vartheta)$$

visible ↗ ↑ prior.
hidden

Here we assume a prior $p_0(\theta)$ over parameters $\theta$. The case of ML training corresponds a flat-prior $p_0(\theta) \sim$ constant.

Examples of models: RBM, HMM, Gaussian Mixture Model,

The variational method replaces the distribution

$$p(\underline{h}^{(1)} \cdots \underline{h}^{(N)} \theta \mid \underline{v}^{(1)} \cdots \underline{v}^{(N)})$$ which is not

available because the $\underline{h}$'s are not observed (i.e empirical frequencies are not available) by a "variational class" of distributions of the form

$$\left\{ \prod_{m=1}^{N} q_m(\underline{h}^{(m)}) \right\} q(\theta)$$

and optimize $L(\theta)$ over $q_m$'s and $q(\theta)$.

This will lead us to the "variational Bayes Expectation-Maximization algorithm". For the class $q(\theta) = \delta_{\theta \theta_*}$ and flat $p_0(\theta)$ this reduces to the standard Expectation-Maximization (EM) algorithm.

Remark: $\underline{h}^{(m)}, m=1,\ldots N$ are dummy variables that are NOT observed - Only $\underline{v}^{(m)}$ are known.

# I. VARIATIONAL BOUND.

Assume iid data samples $\underline{v}^{(1)}, \underline{h}^{(1)} \dots \underline{v}^{(N)}, \underline{h}^{(N)}$ where $\underline{v}^{(1)} \dots \underline{v}^{(N)}$ are visible and $\underline{h}^{(1)} \dots \underline{h}^{(N)}$ are hidden. We have

$$
\log \mathbb{P}(\underline{v}^{(1)} \dots \underline{v}^{(N)}) \geqslant \sum_{m=1}^{N} \mathbb{E}_{q_m(\underline{h}^m) q(\vartheta)} \left[ \log p(\underline{h}^m, \underline{v}^m | \vartheta) \right]
$$

$$
- \sum_{m=1}^{N} \mathbb{E}_{q_m(\underline{h}^m)} \left[ \log q_m(\underline{h}^m) \right] - \mathbb{E}_{q(\vartheta)} \left[ \log \frac{q(\vartheta)}{p_0(\vartheta)} \right]
$$

- first term is called energy term.
- second " " entropy term.
- optimization of the lower bound on $q_m(\underline{h}^m)$ & $q(\vartheta)$ will give the variational Bayes EM algorithm.

## Proof of the variational bound:

The starting point is to write

$$
KL \left( \underbrace{\prod_{m=1}^{N} q_m(\underline{h}^m) \, q(\vartheta)}_{\text{variational distr}} \,\middle\|\, p(\underline{h}^{(1)} \dots \underline{h}^{(N)}, \vartheta | \underline{v}^{(1)} \dots \underline{v}^{(N)}) \right) \geqslant 0
$$

$$KL = \mathbb{E}_{\left\{\prod_{m=1}^{N} q_m(h^m)\right\} q(\vartheta)} \left[ \log \left\{ \prod_{m=1}^{N} q_m(h^m) \right\} q(\vartheta) \right]$$

$$- \mathbb{E}_{\left\{\prod_{m=1}^{N} q_m(h^m)\right\} q(\vartheta)} \left[ \log p(h^{(1)} \dots h^{(N)}, \vartheta \mid v^{(1)} \dots v^{(N)}) \right]$$

$$= \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \left[ \log q_m(h^m) \right] + \mathbb{E}_{q(\vartheta)} \left[ \log q(\vartheta) \right]$$

$$- \mathbb{E}_{\left\{\prod_{m=1}^{N} q_m(h^m)\right\} q(\vartheta)} \left[ \log \frac{p(h^{(1)} \dots h^{(N)}, \vartheta, v^{(1)} \dots v^{(N)})}{p(v^{(1)} \dots v^{(N)})} \right]$$

Then we get for the log-likelihood of the visible data:

$$\log p(v^{(1)} \dots v^{(N)}) = KL\left( \prod_{m=1}^{N} q_m(h^m) q(\vartheta) \,\|\, p(h^{(1)} \dots h^{(N)} \vartheta \mid v^{(1)} \dots v^{(N)}) \right)$$

$$+ \mathbb{E}_{\prod_{m=1}^{N} q_m(h^m) q(\vartheta)} \left( \log p(h^{(1)} \dots h^{(N)}, \vartheta, v^{(1)} \dots v^{(N)}) \right)$$

$$- \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \left( \log q_m(h^m) \right) - \mathbb{E}_{q(\vartheta)} \left[ \log q(\vartheta) \right]$$

Assume now that the Data is iid :

$$p(h^{(1)} \cdots h^{(N)}, \vartheta, v^{(1)} \cdots v^{(N)}) = p(h^{(1)} \cdots h^{(N)}, v^{(1)} \cdots v^{(N)} | \vartheta) p_0(\vartheta)$$

$$\overset{iid}{=} \left\{ \prod_{m=1}^{N} p(h^{(m)}, v^{(m)} | \vartheta) \right\} p_0(\vartheta).$$

We then get :

$$\log p(v^{(1)} \cdots v^{(N)}) = KL\left( \prod_{m=1}^{N} q_m(h^m) \, q(\vartheta) \,\Big\|\, p(h^{(1)} \cdots h^{(N)}, \vartheta | v^{(1)} \cdots v^{(N)}) \right)$$

$$+ \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m) q(\vartheta)} \left[ \log p(h^m, v^m | \vartheta) \right] \qquad \leftarrow \text{energy term}$$

$$- \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \left[ \log q_m(h^m) \right] \qquad \leftarrow \text{entropy term.}$$

$$- \mathbb{E}_{q(\vartheta)} \left[ \log \frac{q(\vartheta)}{p_0(\vartheta)} \right]$$

The variational bound simply follows from $KL \geq 0$.

## II. Optimization of lower bound and the variational Bayes EM algorithm.

We do a "coordinate-wise" optimization over $\{q_m(h), m=1\cdots N\}$ and $q(\vartheta)$ of the lower bound (on page 55).

### E-step: Fix $q(\vartheta)$ and optimize over $\{q_m(h), m=1\cdots N\}$

We look at terms which depend on $q_m$:

$$(*) = \sum_{m=1}^{N} \underbrace{\mathbb{E}_{q_m(h^m) q(\vartheta)}\left[\log p(h^m, v^m, \vartheta)\right]}_{} - \mathbb{E}_{q_m(h^m)}\left[\log q_m(h^m)\right]$$

$$\mathbb{E}_{q(\vartheta)}\left[\log p(h^m, v^m, \vartheta)\right] = \log\left\{\frac{1}{Z}\exp \mathbb{E}_{q(\vartheta)}\left[\log p(h^m, v^m, \vartheta)\right]\right\}$$

$$+ \log Z$$

Thus:

$$(*) = \sum_{m=1}^{N}\left\{\mathbb{E}_{q_m(h^m)}\left[\log \frac{1}{Z}\exp \mathbb{E}_{q(\vartheta)}\left[\log p(h^m, v^m|\vartheta)\right]\right]\right.$$

$$\left. - \mathbb{E}_{q_m(h^m)}\left[\log q_m(h^m)\right]\right\}$$

$$= - \sum_{m=1}^{N} KL\left(q_m(h^m) \;\|\; \frac{1}{Z}\exp \mathbb{E}_{q(\vartheta)}\left(\log p(h^m, v^m|\vartheta)\right)\right).$$

Therefore to maximize (*) over $q_m$ we have
to minimize the KL, i.e, we set:

$$q_m(h^m) = \frac{1}{Z} \exp \mathbb{E}_{q(\vartheta)}\left[\log p(h^m v^m | \vartheta)\right]$$

where $Z =$ Normalization factor

$$= \sum_{h^m} \exp \mathbb{E}_{q(\vartheta)}\left[\log p(h^m v^m | \vartheta)\right].$$

The E-step is an iterative step; at time $t$
we do

$$\begin{cases} q_m^{t+1}(h^m) = \dfrac{1}{Z_t} \exp \mathbb{E}_{q_t(\vartheta)}\left[\log p(h^m, v^m | \vartheta)\right] \\[4mm] Z_t = \displaystyle\sum_{h^m} \exp \mathbb{E}_{q_t(\vartheta)}\left[\log p(h^m, v^m | \vartheta)\right] \end{cases}$$

<u>M-step :</u>  Fix $\{q_m(h^m), m=1 \cdots N\}$ and

optimize over $q(\vartheta)$.

We must maximize terms in lower bound (on page 55)
which depend on $q(\vartheta)$ :

$$(\times\times) = \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)\, q(\vartheta)} \left[ \log p(h^m, v^m \,|\, \vartheta) \right] - \mathbb{E}_{q(\vartheta)} \left[ \log \frac{q(\vartheta)}{p_0(\vartheta)} \right]$$

$$= -\mathbb{E}_{q(\vartheta)}\{\log q(\vartheta)\} + \mathbb{E}_{q(\vartheta)} \left\{ \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \left[ \log p(h^m v^m \,|\, \vartheta) \right] \right.$$

$$\left. + \log p_0(\vartheta) \right\}$$

$$+ \log \widetilde{Z} + \log \left[ \frac{1}{\widetilde{Z}} \exp \left\{ \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \left[ \log p(h^m v^m \,|\, \vartheta) \right] + \log p_0(\vartheta) \right\} \right]$$

$$= + \log \widetilde{Z} + \log \left[ \frac{p_0(\vartheta)}{\widetilde{Z}} \exp \left\{ \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \left[ \log p(h^m v^m \,|\, \vartheta) \right] \right\} \right]$$

with the normalization factor

$$\widetilde{Z} = \int d\vartheta\, p_0(\vartheta)\, \exp \left[ \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)} \log p(h^m v^m \,|\, \vartheta) \right].$$

Note that $\tilde{Z}$ is independent of $q(\vartheta)$. We see that maximizing $(**)$ over $q(\vartheta)$ is equivalent to maximizing

$$- \mathbb{E}_{q(\vartheta)}\left[\log q(\vartheta)\right] + \mathbb{E}_{q(\vartheta)}\left[\log \frac{P_0(\vartheta)}{\tilde{Z}} \prod_{m=1}^{N} \exp\left(\mathbb{E}_{q_m(h^m)} \log p(h^m v^m | \vartheta)\right)\right]$$

$$= -KL\left(q(\vartheta) \,\|\, \frac{P_0(\vartheta)}{\tilde{Z}} \prod_{m=1}^{N} \exp\left(\mathbb{E}_{q_m(h^m)} \log p(h^m v^m | \vartheta)\right)\right)$$

This is achieved by setting

$$q(\vartheta) = \frac{P_0(\vartheta)}{\tilde{Z}} \prod_{m=1}^{N} \exp\left(\mathbb{E}_{q_m(h^m)} \log p(h^m v^m | \vartheta)\right).$$

$$\boxed{\begin{array}{l} \underline{\text{M - step iteration}} : \quad \text{at time } t \quad \text{we do} : \\[2em] \left\{ \begin{array}{l} q_{t+1}(\vartheta) = \dfrac{P_0(\vartheta)}{\tilde{Z}_{t+1}} \prod_{m=1}^{N} \exp\left(\mathbb{E}_{q_m^{t+1}(h^m)} \log p(h^m v^m | \vartheta)\right) \\[3em] \tilde{Z}_{t+1} = \displaystyle\int d\vartheta \, P_0(\vartheta) \exp\left\{\sum_{m=1}^{N} \mathbb{E}_{q_m^{t+1}(h^m)} \log p(h^m v^m | \vartheta)\right\} . \end{array} \right. \end{array}}$$

<u>Remarks:</u>   The iterative E & M steps can only increase the lower bound and will converge to some local maximum. There is no guarantee that the global maximum of the lower bound is found. Also, there is no garantee that the log-likelihood $\log p(V^{(1)} \cdots V^{(N)})$ increases along E & M steps.

For the simplest possible variational class $q(\vartheta) = \delta_{\vartheta, \vartheta_*}$ where we assume that $\vartheta$ takes a single value $\vartheta_*$ and with a flat prior $p_o(\vartheta) \sim$ constant the algorithm reduces to the <u>standard EM algorithm</u>. We will show that in this case not only the lower bound will only increase (until we reach a local maximum) but remarquably <u>also the likelihood itself increases</u>. This is a very nice feature of the EM standard algo-

## III. Reduction to standard EM algorithm.

If we assume that $\vartheta$ is characterized by not a full distribution but by a single numbers $\vartheta_t$, we set $q_t(\vartheta) = \delta_{\vartheta, \vartheta_t}$ and the E-step becomes:

$$
\begin{cases}
q_m^{t+1}(h^m) = \frac{1}{Z_t} \cdot p(h^m, v^m \mid \vartheta_t) = p(h^m \mid v^m \vartheta_t), \\[3mm]
Z_t = \sum_{h^m} p(h^m, v^m \mid \vartheta_t).
\end{cases}
$$

To get the M-step we note that $\vartheta_{t+1} = \text{argmax}_\vartheta \, \delta_{\vartheta, \vartheta_{t+1}}$ so with $p_0(\vartheta) = \text{constant}$:

$$
\vartheta_{t+1} = \underset{\vartheta}{\text{argmax}} \sum_{m=1}^{N} \mathbb{E}_{q_m^{t+1}(h^m)} \left[ \log p(h^m, v^m \mid \vartheta) \right].
$$

---

**Standard EM algorithm:**

E-step : $\quad q_m^{t+1}(h^m) = p(h^m \mid v^m \vartheta_t)$

M-step : $\quad \vartheta_{t+1} = \underset{\vartheta}{\text{argmax}} \sum_{m=1}^{N} \mathbb{E}_{q_m^{t+1}(h^m)} \left[ \log p(h^m v^m \mid \vartheta) \right]$

## A closer look at EM lower bound:

In standard ML training we seek to maximize $\log p(v^{(1)}\cdots v^{(N)} \mid \vartheta)$. Since we do not have access to hidden variable and cannot estimate $p(h^{(1)}\cdots h^{(N)} \mid v^{(1)}\cdots v^{(N)}\vartheta)$ we replace it by the variational class $\prod_{m=1}^{N} q_m(h^m)$.

Starting from

$$KL\left(\prod_{m=1}^{N} q_m(h^m) \,\middle\|\, p(h^{(1)}\cdots h^{(N)} \mid v^{(1)}\cdots v^{(N)}\vartheta)\right) \geq 0$$

we get by identical (and simpler) calculations than before the variational lower bound (iid samples):

$$L(\vartheta) = \log p(v^{(1)}\cdots v^{(N)} \mid \vartheta) = KL\left(\prod_{m=1}^{N} q_m(h^m) \,\middle\|\, p(h^{(1)}\cdots h^{(N)} \mid v^{(1)}\cdots v^{(N)}\vartheta)\right)$$

$$+ \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)}\left[\log p(h^m, v^m \mid \vartheta)\right] \qquad \leftarrow \text{ energy term}$$

$$- \sum_{m=1}^{N} \mathbb{E}_{q_m(h^m)}\left[\log q_m(h^m)\right] \qquad \leftarrow \text{ entropy term}$$

$$\geq (\text{energy term}) \quad + \quad (\text{entropy term}).$$

We have the nice properties:

(i) During EM steps lower bound only increases until it reaches a local stat pt. This is obvious because EM are coordinate-wise optimization of lower bound

(ii) (less obviously, also $\log p(v^{(1)} \ldots v^{(N)} | \theta_t) =_, L(\theta_t)$ increases (→ popular algo.).

Proof of (ii).

Note that the formula for $L(\theta_t)$ is valid for any variational distr $q_m(h^m)$ → so it is valid when we use $q_m^t(h^m)$ as well as $q_m^{t+1}(h^m)$ !.

At time $t$ of algorithm:

$$L(\theta_t) = \text{expression}(\theta_t, q_m^t) = \text{expression}(\theta_t, q_m^{t+1})$$

but here E-step: $P_m(h^m | v^m \theta_t) = q_m^{t+1}$

Thus

$$KL\left( \prod_{m=1}^{N} q_m^{t+1} \| p(h^1 \ldots h^N | v^1 \ldots v^N \theta) \right) = 0$$

$\Rightarrow \text{expression}(\theta_t, q_m^{t+1}) = \text{expression of lower Bound}(\theta_t, q_m^{t+1})$

(M-step) $\leadsto \leq \text{expression of lower Bound}(\theta_{t+1}, q_m^{t+1})$

$\leq L(\theta_{t+1})$.