

# CS 3750 Machine Learning

## Lecture 19

### Latent variable models

### Variational approximations.

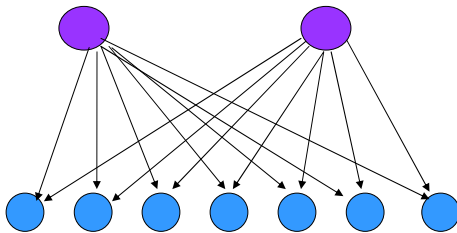
Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

CS 2750 Machine Learning

### Cooperative vector quantizer

**Latent variables (s):** binary vars  
**Dimensionality k**



**Observed variables x:** real valued vars  
**Dimensionality d**

---

CS 2750 Machine Learning

## Cooperative vector quantizer

### Model:

#### Latent var $s_i$ :

~ Bernoulli distribution  
parameter:  $\pi_i$

$$P(s_i | \pi_i) = \pi_i^{s_i} (1 - \pi_i)^{1-s_i}$$

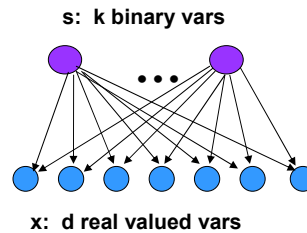
#### Observable variables $\mathbf{x}$ :

~ Normal distribution  
parameters:  $\mathbf{W}, \Sigma$

$$P(\mathbf{x} | \mathbf{s}) = N(\mathbf{W}\mathbf{s}, \Sigma)$$

We assume  $\Sigma = \sigma I$

$$\mathbf{W} = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1k} \\ w_{21} & & & \\ & \dots & & \\ w_{d1} & \dots & \dots & w_{dk} \end{pmatrix}$$



#### Joint for one instance of $\mathbf{x}$ and $\mathbf{s}$ :

$$P(\mathbf{x}, \mathbf{s} | \Theta) = (2\pi)^{d/2} \sigma^{d/2} \exp\left\{-\frac{1}{2\sigma^2} (\mathbf{x} - \mathbf{W}\mathbf{s})^T (\mathbf{x} - \mathbf{W}\mathbf{s})\right\} \prod_{i=1}^k \pi_i^{s_i} (1 - \pi_i)^{(1-s_i)}$$

CS 2750 Machine Learning

## Cooperative vector quantizer

### Our objective:

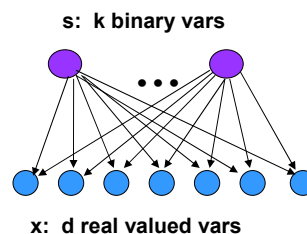
- Learn the parameters of the model  
 $\mathbf{W}, \pi, \sigma$
- One can use the data likelihood or loglikelihood and optimize ..

### Learning if $\mathbf{x}$ and $\mathbf{s}$ are observable

#### Log likelihood:

$$\sum_{n=1}^N \log P(\mathbf{x}^{(n)}, \mathbf{s}^{(n)} | \Theta) = \sum_{n=1}^N \left[ -d \log \sigma - \frac{1}{2\sigma^2} (\mathbf{x}^{(n)} - \mathbf{W}\mathbf{s}^{(n)})^T (\mathbf{x}^{(n)} - \mathbf{W}\mathbf{s}^{(n)}) + \sum_{i=1}^k s_i^{(n)} \log \pi_i + (1 - s_i^{(n)}) \log(1 - \pi_i) \right]$$

#### Solution: nice and easy



CS 2750 Machine Learning

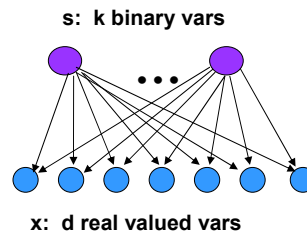
## Cooperative vector quantizer

Our objective:

- Learn the parameters of the model

$W, \pi, \sigma$

- One can use the data likelihood or loglikelihood and optimize ..



Learning if only  $x$  are observable

Log likelihood of data:

$$\log P(D|\Theta) = \sum_{n=1}^N \log P(\mathbf{x}^{(n)}|\Theta) = \sum_{n=1}^N \log \sum_{\{\mathbf{s}^n\}} P(\mathbf{x}^{(n)}, \mathbf{s}^{(n)}|\Theta)$$

Solution: does not let us benefit from the decomposition

EM: used to work in such cases ...

## EM

Let  $H$  – be a set of all variables with hidden or missing values

Derivation

$$P(H, D | \Theta, \xi) = P(H | D, \Theta, \xi) P(D | \Theta, \xi)$$

$$\log P(H, D | \Theta, \xi) = \log P(H | D, \Theta, \xi) + \log P(D | \Theta, \xi)$$

$$\log P(D | \Theta, \xi) = \log P(H, D | \Theta, \xi) - \log P(H | D, \Theta, \xi)$$



Log-likelihood of data

Average both sides with  $P(H | D, \Theta', \xi)$  for  $\Theta'$

$$E_{H|D, \Theta'} \log P(D | \Theta, \xi) = E_{H|D, \Theta'} \log P(H, D | \Theta, \xi) - E_{H|D, \Theta'} \log P(H | D, \Theta, \xi)$$

$$\underbrace{\log P(D | \Theta, \xi)} = F(\Theta | \Theta') = E(\Theta | \Theta') + H(\Theta | \Theta')$$

Log-likelihood of data

## EM algorithm

**Algorithm** (general formulation)

Initialize parameters  $\Theta$

Repeat

Set  $\Theta' = \Theta$

**1. Expectation step**

$$E(\Theta | \Theta') = \langle \log P(H, D | \Theta, \xi) \rangle_{P(H|D, \Theta')}$$

**2. Maximization step**

$$\Theta = \arg \max_{\Theta} E(\Theta | \Theta')$$

until no or small improvement in  $\Theta$  ( $\Theta = \Theta'$ )

**Problem:** posterior  $P(H | D, \Theta', \xi)$  is defined over  $2^k$  probabilities

## EM algorithm

Posterior  $P(H | D, \Theta', \xi)$  for our model

$$P(H | D, \Theta') = \prod_{n=1}^N P(s^{(n)} | x^{(n)}, \Theta')$$

- Each data point  $n=1, \dots, N$  requires us to calculate  $2^k$  probabilities
- If  $k$  is larger then this is a bottleneck!!!

## Variational approximation

Let  $H$  – be a set of all variables with hidden or missing values

### Derivation

$$\log P(D | \Theta, \xi) = \log P(H, D | \Theta, \xi) - \log P(H | D, \Theta, \xi)$$

 **Log-likelihood of data**

Average both sides with  $Q(H | \lambda)$

$$E_{H|\lambda} \log P(D | \Theta, \xi) = E_{H|\lambda} \log P(H, D | \Theta, \xi) - E_{H|\lambda} \log P(H | \Theta, \xi) \\ + E_{H|\lambda} \log Q(H | \lambda) - E_{H|\lambda} \log Q(H | \lambda)$$

$$\log P(D | \Theta, \xi) = F(P, Q) + KL(Q, P)$$

**Log-likelihood of data**

## Variational approximation

Let  $H$  – be a set of all variables with hidden or missing values

$$\log P(D | \Theta, \xi) = E_{H|\lambda} \log P(H, D | \Theta, \xi) - E_{H|\lambda} \log P(H | \Theta, \xi) \\ + E_{H|\lambda} \log Q(H | \lambda) - E_{H|\lambda} \log Q(H | \lambda)$$

$$\log P(D | \Theta, \xi) = F(Q, \Theta) + KL(Q, P)$$

$$F(Q, \Theta) = \sum_{\{H\}} Q(H | \lambda) \log P(H, D | \Theta, \xi) - \sum_{\{H\}} Q(H | \lambda) \log Q(H | \lambda)$$

$$KL(Q, P) = \sum_{\{H\}} Q(H | \lambda) [\log Q(H | \lambda) - \log P(H | D, \Theta)]$$

**Approximation: maximize**  $F(Q, \Theta)$

**Parameters:**  $\Theta, \lambda$

## Variational EM

Let  $H$  be a set of all variables with hidden or missing values

- **E step:**
  - Optimize  $F(Q, \Theta)$  with respect to  $\lambda$  while keeping  $\Theta$  fixed
- **M step**
  - Optimize  $F(Q, \Theta)$  with respect to  $\Theta$  while keeping  $\lambda$ s

Note: if  $Q(H)$  is the posterior then the variational EM reduces to the standard EM

## Variational EM

- So what is the deal?
  - Why should we use the variational EM?
- Hope:
  - If we choose  $Q(H | \lambda)$  well the optimization of both  $\lambda$  and  $\Theta$  will become easy
- A well behaved choice for  $Q(H | \lambda)$ 
  - Use mean field approximation

## Mean Field Approximation

### Assumption:

- $Q(H|\lambda)$  is the mean field approximation.
- Variables in the  $Q(H)$  distribution are independent variables  $H_i$ .
- $Q$  is completely factorized:

$$Q(H | \lambda) = \prod_i Q_i(H_i | \lambda_i)$$

- For our CVQ model
  - Hidden variables are binary sources

$$Q(\mathbf{s} | \lambda) = \prod_i Q_i(s_i | \lambda_i)$$

$$Q(s_i | \lambda_i) = \lambda_i^{s_i} (1 - \lambda_i)^{1-s_i}$$

## Mean Field Approximation

### Functional F for the mean field:

$$F(Q, \Theta) = \sum_{\{H\}} Q(H | \lambda) \log P(H, D | \Theta, \xi) - \sum_{\{H\}} Q(H | \lambda) \log Q(H | \lambda)$$

$$F(Q, \Theta) = \langle \log P(\mathbf{x}, \mathbf{s} | \Theta) \rangle_{Q(s|\lambda)} - \langle \log Q(\mathbf{s} | \lambda) \rangle_{Q(s|\lambda)}$$

$$= \left\langle -d \log \sigma - \frac{1}{2\sigma^2} (\mathbf{x} - \mathbf{W}\mathbf{s})^T (\mathbf{x} - \mathbf{W}\mathbf{s}) \right\rangle_{Q(s|\lambda)} \quad (1)$$

$$+ \left\langle \sum_{i=1}^k s_i \log \pi_i + (1 - s_i) \log(1 - \pi_i) \right\rangle_{Q(s|\lambda)} \quad (2)$$

$$- \left\langle \sum_{i=1}^k s_i \log \lambda_i + (1 - s_i) \log(1 - \lambda_i) \right\rangle_{Q(s|\lambda)} \quad (3)$$

## Mean Field Approximation

### Functional F. Part 1:

$$\begin{aligned}
 & \left\langle -d \log \sigma - \frac{1}{2\sigma^2} \left( \mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i \right)^T \left( \mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i \right) \right\rangle_{Q(s|\lambda)} = \\
 & = \left\langle -d \log \sigma - \frac{1}{2\sigma^2} \left( \mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i \right)^T \left( \mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i \right) \right\rangle_{Q(s|\lambda)} \\
 & = \left\langle -d \log \sigma - \frac{1}{2\sigma^2} \left[ \mathbf{x}^T \mathbf{x} - 2 \sum_{i=1}^k (s_i \mathbf{w}_i)^T \mathbf{x} + \sum_{i=1}^k \sum_{j=1}^k s_i s_j \mathbf{w}_i^T \mathbf{w}_j \right] \right\rangle_{Q(s|\lambda)} \\
 & = -d \log \sigma - \frac{1}{2\sigma^2} \left[ \mathbf{x}^T \mathbf{x} - 2 \sum_{i=1}^k \langle s_i \rangle_{Q(s_i|\lambda_i)} \mathbf{w}_i^T \mathbf{x} + \sum_{i=1}^k \sum_{j=1}^k \langle s_i s_j \rangle_{Q(s_i|\lambda_i)} \mathbf{w}_i^T \mathbf{w}_j \right] \\
 & \quad \langle s_i \rangle_{Q(s_i|\lambda_i)} = \lambda_i \quad \langle s_i s_j \rangle_{Q(s_i|\lambda_i)} = \lambda_i \lambda_j + \delta_{ij} (\lambda_i - \lambda_i^2)
 \end{aligned}$$

CS 2750 Machine Learning

## Mean Field Approximation

### Functional F. Part 2:

$$\begin{aligned}
 \left\langle \sum_{i=1}^k s_i \log \pi_i + (1-s_i) \log(1-\pi_i) \right\rangle_{Q(s|\lambda)} &= \sum_{i=1}^k \langle s_i \rangle_{Q(s_i|\lambda_i)} \log \pi_i + (1 - \langle s_i \rangle_{Q(s_i|\lambda_i)}) \log(1-\pi_i) \\
 &= \sum_{i=1}^k \lambda_i \log \pi_i + (1-\lambda_i) \log(1-\pi_i)
 \end{aligned}$$

### Functional F. Part 3:

$$\left\langle \sum_{i=1}^k s_i \log \lambda_i + (1-s_i) \log(1-\lambda_i) \right\rangle_{Q(s|\lambda)} = \sum_{i=1}^k \lambda_i \log \lambda_i + (1-\lambda_i) \log(1-\lambda_i)$$

CS 2750 Machine Learning

## Mean Field Approximation

**Functional F:**

$$\begin{aligned} F(Q, \Theta) &= \langle \log P(\mathbf{x}, \mathbf{s} | \Theta) \rangle_{Q(\mathbf{s}|\lambda)} - \langle \log Q(\mathbf{s} | \lambda) \rangle_{Q(\mathbf{s}|\lambda)} \\ &= -d \log \sigma - \frac{1}{2\sigma^2} \left[ \mathbf{x}^T \mathbf{x} - 2 \sum_{i=1}^k \lambda_i \mathbf{w}_i \mathbf{x} + \sum_{i=1}^k \sum_{j=1}^k [\lambda_i \lambda_j + \delta_{ij} (\lambda_i - \lambda_i^2)] \mathbf{w}_i^T \mathbf{w}_j \right] \\ &\quad + \sum_{i=1}^k \lambda_i \log \pi_i + (1 - \lambda_i) \log(1 - \pi_i) \\ &\quad + \sum_{i=1}^k \lambda_i \log \lambda_i + (1 - \lambda_i) \log(1 - \lambda_i) \end{aligned}$$

**Parameters:  $\mathbf{W}, \pi, \sigma$**

**Mean field parameters:  $\lambda$**

---

CS 2750 Machine Learning

## Mean Field Approximation

**Functional F (for all data points):**

$$\begin{aligned} F(Q, \Theta) &= \sum_{n=1}^N \langle \log P(\mathbf{x}, \mathbf{s} | \Theta) \rangle_{Q(\mathbf{s}|\lambda)} - \langle \log Q(\mathbf{s} | \lambda) \rangle_{Q(\mathbf{s}|\lambda)} \\ &= -d \log \sigma - \frac{1}{2\sigma^2} \left[ \mathbf{x}^{(n)T} \mathbf{x}^{(n)} - 2 \sum_{i=1}^k \lambda_i^{(n)} \mathbf{w}_i \mathbf{x}^{(n)} + \sum_{i=1}^k \sum_{j=1}^k [\lambda_i^{(n)} \lambda_j^{(n)} + \delta_{ij} (\lambda_i^{(n)} - \lambda_i^{(n)2})] \mathbf{w}_i^T \mathbf{w}_j \right] \\ &\quad + \sum_{i=1}^k \lambda_i^{(n)} \log \pi_i + (1 - \lambda_i^{(n)}) \log(1 - \pi_i) \\ &\quad + \sum_{i=1}^k \lambda_i^{(n)} \log \lambda_i^{(n)} + (1 - \lambda_i^{(n)}) \log(1 - \lambda_i^{(n)}) \end{aligned}$$

**Parameters:  $\mathbf{W}, \pi, \sigma$**

**Mean field parameters:  $\lambda = \lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(N)}$**

---

CS 2750 Machine Learning

## Variational EM: E step

Optimization of the functional F with respect to  $\lambda$ :

$$\frac{\partial}{\partial \lambda_u} F = \frac{1}{\sigma^2} (\mathbf{x} - \sum_{j \neq u} \lambda_j \mathbf{w}_j)^T \mathbf{w}_u - \frac{1}{2\sigma^2} \mathbf{w}_u^T \mathbf{w}_u + \log \frac{\pi_u}{1 - \pi_u} - \log \frac{\lambda_u}{1 - \lambda_u}$$

$$\text{set } \frac{\partial}{\partial \lambda_u} F = 0$$

$$\lambda_u = g \left( \frac{1}{\sigma^2} (\mathbf{x} - \sum_{j \neq u} \lambda_j \mathbf{w}_j)^T \mathbf{w}_u - \frac{1}{2\sigma^2} \mathbf{w}_u^T \mathbf{w}_u + \log \frac{\pi_u}{1 - \pi_u} \right)$$

$$g(x) = \frac{1}{1 + e^{-x}}$$

Defines a fixed point equation

Iterate a set fixed point equations for all indexes  $u=1..k$

## Variational EM: M step

Optimization of the functional F with respect to  $\Theta$ .

Start with  $\pi$ :

$$\frac{\partial}{\partial \pi_u} F = \lambda_u \log \frac{1}{\pi_u} - (1 - \lambda_u) \log \frac{1}{(1 - \pi_u)} \quad \text{For 1 data point}$$

$$\frac{\partial}{\partial \pi_u} F = \sum_{n=1}^N \lambda_u^n \log \frac{1}{\pi_u} - (1 - \lambda_u^n) \log \frac{1}{(1 - \pi_u)} \quad \text{For N data points}$$

$$\text{set } \frac{\partial}{\partial \pi_u} F = 0$$

$$\pi_u = \frac{\sum_{n=1}^N \lambda_u^{(n)}}{N} \quad \text{Closed form solution}$$

## Variational EM: M step

Optimization of the functional  $F$  with respect to  $\Theta$ .

Parameters  $\mathbf{w}$ :

$$\frac{\partial}{\partial w_{uv}} F = \sum_{n=1}^N -\frac{1}{2\sigma^2} \left[ \lambda_v^{(n)} x_u^{(n)} + 2 \sum_{j \neq v} \lambda_v^{(n)} \lambda_j^{(n)} w_{uj} + 2 \lambda_v w_{uv} \right] = 0$$

$$\mathbf{W} = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1k} \\ w_{21} & & & \\ & \dots & & \\ w_{d1} & \dots & \dots & w_{dk} \end{pmatrix} \quad \mathbf{W} = (\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_k)$$

For each variable  $v$ :

The equations define a set of  $k$  linear equations that can be solved

## Mean Field approximation

Let  $H$  – be a set of all variables with hidden or missing values

• E step:

$$\text{– Optimize } KL(Q,P) = \sum_{\{H\}} Q(H | \lambda) [\log Q(H | \lambda) - \log P(H | D, \Theta)]$$

$$\log P(D | \Theta, \xi) = F(P, Q) + KL(Q, P)$$

$$F(P, Q) = \sum_{\{H\}} Q(H | \lambda) \log P(H, D | \Theta, \xi) - \sum_{\{H\}} Q(H | \lambda) \log Q(H | \lambda)$$

$$KL(Q,P) = \sum_{\{H\}} Q(H | \lambda) [\log Q(H | \lambda) - \log P(H | D, \Theta)]$$

**Approximation: maximize**  $F(P, Q)$

**Parameters:**  $\Theta, \lambda$