

**CS 3750 Machine Learning  
Lecture 19**

**Latent variable models  
Variational approximations.**

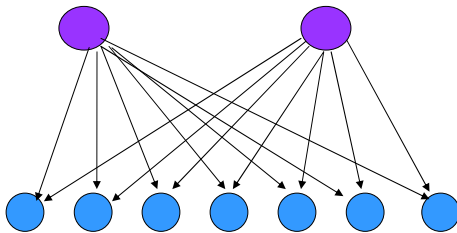
Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

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:  $\pi_i$

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

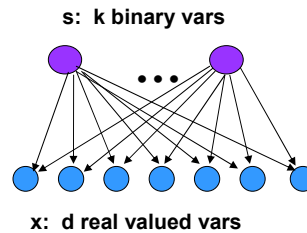
#### Observable variables $\mathbf{x}$ :

~ 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:

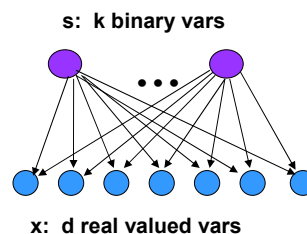
- Learn the parameters of the model  
 $\mathbf{W}, \pi, \sigma$
- One can use the data likelihood or loglikelihood and optimize ..

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

#### Log likelihood:

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

#### Solution: nice and easy



CS 2750 Machine Learning

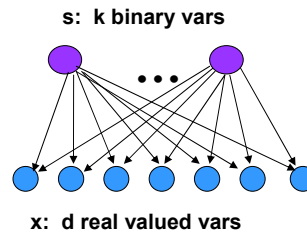
## Cooperative vector quantizer

**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

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

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

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

**Log-likelihood of data**

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

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

$$\underbrace{\log P(D | \Theta, \xi)}_{\text{Log-likelihood of data}} = F(\Theta | \Theta') = E(\Theta | \Theta') + H(\Theta | \Theta')$$

**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

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

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

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

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

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

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

**Why?**  $\log P(D | \Theta, \xi) \geq F(Q, \Theta)$

Maximization of F pushes up the lower bound on the log-likelihood

## Variational approximation

- **Comparison:**

- **EM uses true posterior**  $P(H | D, \Theta', \xi)$
- **Variational EM uses a surrogate posterior**  $Q(H | \lambda)$

**EM:**

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

**Variational EM:**

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

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

## Variational EM

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

- **E step:**

- Optimize

$$F(Q, \Theta) \text{ with respect to } \lambda \text{ while keeping } \Theta \text{ fixed}$$

- **M step**

- Optimize

$$F(Q, \Theta) \text{ with respect to } \Theta \text{ while keeping } \lambda \text{ s}$$

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  $Q(H | \lambda)$ 
  - the 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

$$Q(\mathbf{H} | \lambda) = \prod_{n=1, \dots, N} Q(\mathbf{s}^{(n)} | \lambda^{(n)})$$

$$Q(\mathbf{s}^{(n)} | \lambda^{(n)}) = \prod_{i=1, \dots, d} Q(s_i^{(n)} | \lambda_i^{(n)})$$

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

## Mean Field Approximation

**Functional F for the mean field:**

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

Assume just one data point  $\mathbf{x}$  and corresponding  $\mathbf{s}$  :

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

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

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

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

## Mean Field Approximation

**Functional F. Part 1:**

$$\left\langle -d \log \sigma - \frac{1}{2\sigma^2} (\mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i)^T (\mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i) \right\rangle_{Q(s|\lambda)} =$$

$$= \left\langle -d \log \sigma - \frac{1}{2\sigma^2} (\mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i)^T (\mathbf{x} - \sum_{i=1}^k s_i \mathbf{w}_i) \right\rangle_{Q(s|\lambda)}$$

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

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

$$\langle s_i \rangle_{Q(s|\lambda_i)} = \lambda_i \quad \langle s_i s_j \rangle_{Q(s|\lambda)} = \lambda_i \lambda_j + \delta_{ij} (\lambda_i - \lambda_i^2)$$

## 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:

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

## Mean Field Approximation

### Functional F:

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

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

Mean field parameters:  $\lambda$

## Mean Field Approximation

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

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

**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  $u=1..k$  and for all  $n$

CS 2750 Machine Learning

## Variational EM: M step

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

Start with  $\pi$ :

For  $N$  data points

$$\frac{\partial}{\partial \pi_u} F = \sum_{n=1}^N \lambda_u^{(n)} \log \frac{1}{\pi_u} - (1 - \lambda_u^{(n)}) \log \frac{1}{(1 - \pi_u)}$$

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

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

## Variational EM: M step

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

Parameters  $\mathbf{w}$ :

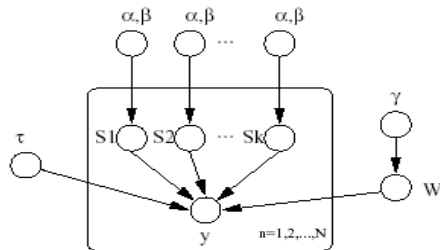
$$\frac{\partial}{\partial w_{uv}} F = \sum_{n=1}^N -\frac{1}{2\sigma^2} \left[ \lambda_v^{(n)} x_u^{(n)} + 2 \sum_{j \neq v} \lambda_v^{(n)} \lambda_j^{(n)} w_{uj} + 2 \lambda_v^{(n)} w_{uv} \right] = 0$$

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

For each variable  $v$ :

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

## Bayesian CVQ Model



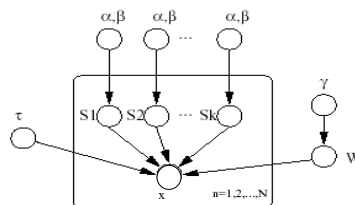
$$y = \sum_{k=1}^K s_k w_k + \varepsilon$$

Bayesian model:  
Distributions over parameters

$$P(y | S, \theta) \sim N\left(\sum_{k=1}^K s_k w_k, \tau^{-1} I\right)$$

CS 2750 Machine Learning

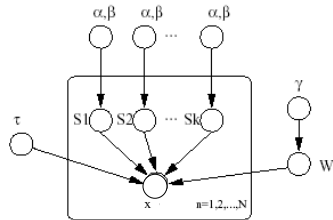
## Model Specification



$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$	observed data
$\mathbf{S} = \{s_1, \dots, s_k\}$	latent sources
$\boldsymbol{\pi} = \{\pi_1, \pi_2, \dots, \pi_k\}$	probability of $s_k = 1$
$\mathbf{W} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k\}$	$D \times K$ weight matrix
$\boldsymbol{\gamma} = \{\gamma_1, \gamma_2, \dots, \gamma_k\}$	Variance of $\mathbf{W}$
$\tau$	Precision of noise

CS 2750 Machine Learning

## Priors



$$P(\boldsymbol{\pi}) = \prod_{k=1}^K \text{Beta}(\pi_k | \alpha, \beta)$$

$$P(\mathbf{W}) = \prod_{k=1}^K N(w_k | 0, \gamma_k)$$

$$P(\boldsymbol{\gamma}) = \prod_{k=1}^K \text{Gamma}(\gamma_k | a_\gamma, b_\gamma)$$

$$P(\tau) = \text{Gamma}(\tau | c_\tau, d_\tau)$$

CS 2750 Machine Learning

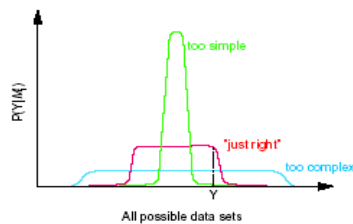
## Why the Bayesian model ?

Very useful for Bayesian Model Selection

- Assume we do not know the number of sources
- Bayesian score tells us how good the structure is

Benefits of the Bayesian score:

- Embodies Occam's Razor
- Prevents overfit



$$P(M_i | D) = \frac{P(D | M_i)P(M_i)}{P(D)}$$

$$P(D | M_i) = \int_{\theta} P(D | \theta, M_i)P(\theta | M_i) \leftarrow \text{Marginal likelihood}$$

CS 2750 Machine Learning

## Variational approximation

- Approximation: loglikelihood of data

$$\begin{aligned}\log P(X) &= \log \int_{\theta} P(X, \theta) d\theta \\ &= \log \int_{\theta} \sum_H P(X, H, \theta) d\theta \\ &= \log \int_{\theta} \sum_H P(X, H | \theta) P(\theta) d\theta \\ &\geq \int_{\theta} \sum_H Q(H, \theta) \log \frac{P(X, H | \theta) P(\theta)}{Q(H, \theta)} d\theta = F(Q)\end{aligned}$$

Where Q is a distribution with different parameterization

## Variational approximation

- Approximation: loglikelihood of observable data

$$\log P(X) = F(Q) + KL(Q(H, \theta), P(H, \theta))$$

- Optimization of F(Q) is pushing up the lower bound on the loglikelihood of observable data
- How to choose Q ?

$$Q(H, \theta) = Q_{\theta}(\theta) Q_H(H)$$

- Then:
$$F(Q) = \int_{\theta} Q_{\theta}(\theta) \left[ \sum_H Q_H(H) \log \frac{P(X, H | \theta)}{Q_H(H)} \right] d\theta$$
$$+ \int_{\theta} Q_{\theta}(\theta) \log \frac{Q_{\theta}(\theta)}{P(\theta)} d\theta \quad \leftarrow \text{KL distance}$$

## Variational Bayes approximation

- Evaluation of  $Q(H, \theta)$  is intractable
- Meanfield approximation

$$Q(H, \theta) = \prod_{k=1}^K Q(H_k) \prod_{i=1}^P Q(\theta_i)$$

- Allows analytical evaluation of  $F(Q)$

## VB learning

Learn Model with an EM like algorithm

(1) VBE – Optimize  $Q(H)$

Estimate state of latent variables

$$Q^*_H(H) \propto \exp\langle \log P(D, H | \theta) \rangle_{Q_\theta(\theta)}$$

(2) VBM – Optimize  $Q(\Theta)$

Estimate parameters

$$Q^*_\theta(\theta) \propto P(\theta) \exp\langle \log P(D, H | \theta) \rangle_{Q_H(H)}$$

## VBE

$$\begin{aligned}\log \frac{\lambda_k}{1-\lambda_k} &= \left\langle \log \frac{\pi_k}{1-\pi_k} \right\rangle_{Q_\theta(\pi)} + y^T \langle \tau \rangle_{Q_\theta(\tau)} \langle \mathbf{w}_k \rangle_{Q_\theta(W)} \\ &\quad - \sum_{j \neq k} \lambda_j \langle \tau \rangle_{Q_\theta(\tau)} \text{tr} \left( \langle \mathbf{w}_j \mathbf{w}_k^T \rangle_{Q_\theta(W)} \right) \\ &\quad - \frac{1}{2} \langle \tau \rangle_{Q_\theta(\tau)} \text{tr} \left( \langle \mathbf{w}_k \mathbf{w}_k^T \rangle_{Q_\theta(W)} \right)\end{aligned}$$

## VBM

$$Q_\pi(\boldsymbol{\pi}) = \prod_{k=1}^K \text{Beta}(\pi_k \mid \tilde{\alpha}_k, \tilde{\beta}_k)$$

$$Q_W(\mathbf{W}) = \prod_{d=1}^D N(\mathbf{w}_k \mid \tilde{\mathbf{m}}_w^{(d)}, \tilde{\boldsymbol{\Sigma}}_w^{(d)})$$

$$Q_\gamma(\boldsymbol{\gamma}) = \prod_{k=1}^K \text{Gamma}(\gamma_k \mid \tilde{a}_{\gamma k}, \tilde{b}_{\gamma k})$$

$$Q_\tau(\tau) = \text{Gamma}(\tau \mid \tilde{c}_\tau, \tilde{d}_\tau)$$

## VBM (cont')

$$\tilde{\alpha}_k = \alpha_k + \sum_{n=1}^N \langle s_k^{(n)} \rangle$$

$$\tilde{\beta}_k = \beta_k + N - \sum_{n=1}^N \langle s_k^{(n)} \rangle$$

$$\tilde{\Sigma}_w^{(d)} = \left( \text{diag}(\langle \gamma \rangle) + \langle \tau \rangle \sum_{n=1}^N \langle \mathbf{s}^n \mathbf{s}^{nT} \rangle \right)^{-1}$$

$$\tilde{\mathbf{m}}_w^{(d)} = \tilde{\Sigma}_w^{(d)} \langle \tau \rangle \sum_{n=1}^N \langle \mathbf{s}^n \rangle \mathbf{x}_d^n$$

## VBM (cont')

$$\tilde{a}_{\gamma k} = a_{\gamma k} + \frac{D}{2}$$

$$\tilde{b}_{\gamma k} = b_{\gamma k} + \frac{\langle \|\mathbf{w}_k\|^2 \rangle}{2}$$

$$\tilde{c}_\tau = c_\tau + \frac{ND}{2}$$

$$\begin{aligned} \tilde{d}_\tau = d_\tau + \frac{1}{2} \sum_{n=1}^N \left\{ \|\mathbf{x}\|^2 - 2\mathbf{y}^{nT} \langle \mathbf{W} \rangle \langle \mathbf{s}^n \rangle \right. \\ \left. + \text{tr} \left( \langle \mathbf{W}^T \mathbf{W} \rangle \langle \mathbf{s}^n \mathbf{s}^{nT} \rangle \right) \right\} \end{aligned}$$