

# Expectation Maximization (contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Expectation Maximization

- EM is a method to optimize  $\log p(\mathcal{D}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathcal{D}, \mathbf{Z}|\Theta)$  for point estimation of  $\Theta$
- EM optimizes  $\mathcal{L}(q, \Theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathcal{D}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \right\}$ , which is a lower bound on  $\log p(\mathbf{X}|\Theta)$

Data

Latent variables

1. Initialize  $\Theta$  as  $\Theta^{(0)}$  somehow (e.g., randomly), set  $t = 1$
2. Set  $q^{(t)} = p(\mathbf{Z}|\mathcal{D}, \Theta^{(t-1)}) \propto p(\mathcal{D}|\mathbf{Z}, \Theta^{(t-1)})p(\mathbf{Z}|\Theta^{(t-1)})$
3. Set  $\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathbb{E}_{q^{(t)}} [\log p(\mathcal{D}, \mathbf{Z}|\Theta)] = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{(t-1)})$
4. If not converged, set  $t = t + 1$  and go to step 2

Computing the CP of latent variables

Maximizing the expected CLL

- CP  $q^{(t)}$  in step 2 and expectation in step 3 may not be tractable. May need approximations



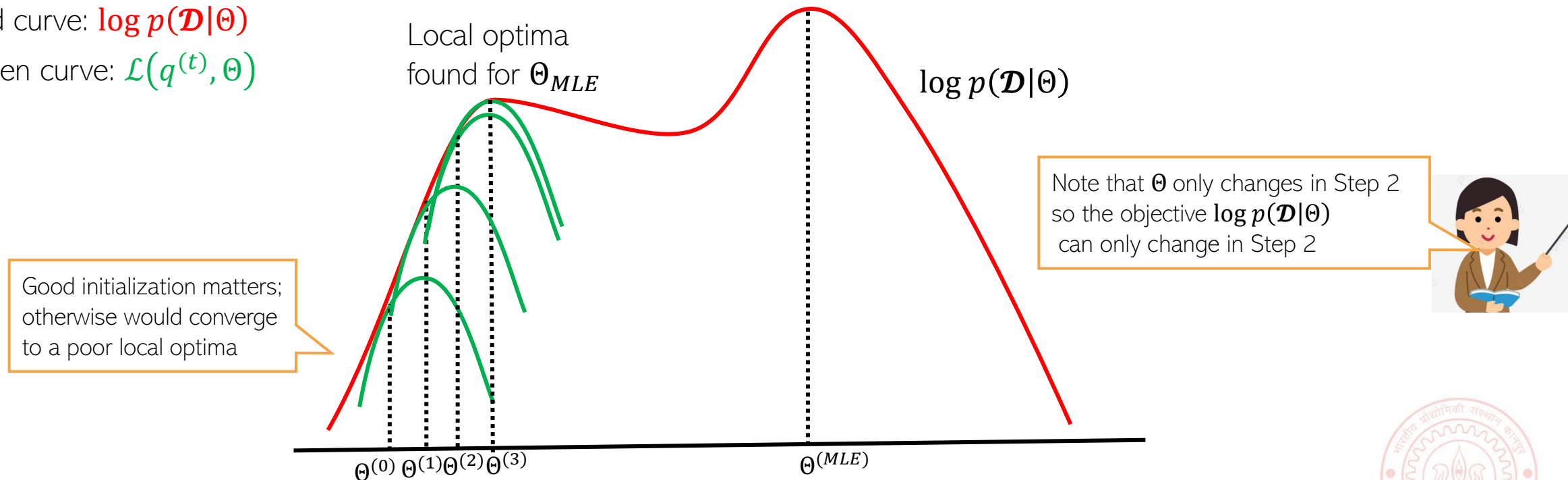
# EM is guaranteed to converge!

$$\log p(\mathcal{D}|\Theta) = \mathcal{L}(q, \Theta) + KL(q||p_z)$$

- Maximization of lower bound  $\mathcal{L}(q, \Theta)$  alternates between these two steps
  - Step 1 sets  $q^{(t)} = p(\mathbf{Z}|\mathcal{D}, \Theta^{(t-1)})$  so  $KL(q||p_z)$  becomes zero, and red and green curves touch at  $\Theta^{(t-1)}$
  - Step 2 sets  $\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathcal{L}(q^{(t)}, \Theta) = \operatorname{argmax}_{\Theta} \mathbb{E}_{q^{(t)}}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)] = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{(t-1)})$

Red curve:  $\log p(\mathcal{D}|\Theta)$

Green curve:  $\mathcal{L}(q^{(t)}, \Theta)$



- EM is guaranteed to converge (possibly to a local optima)

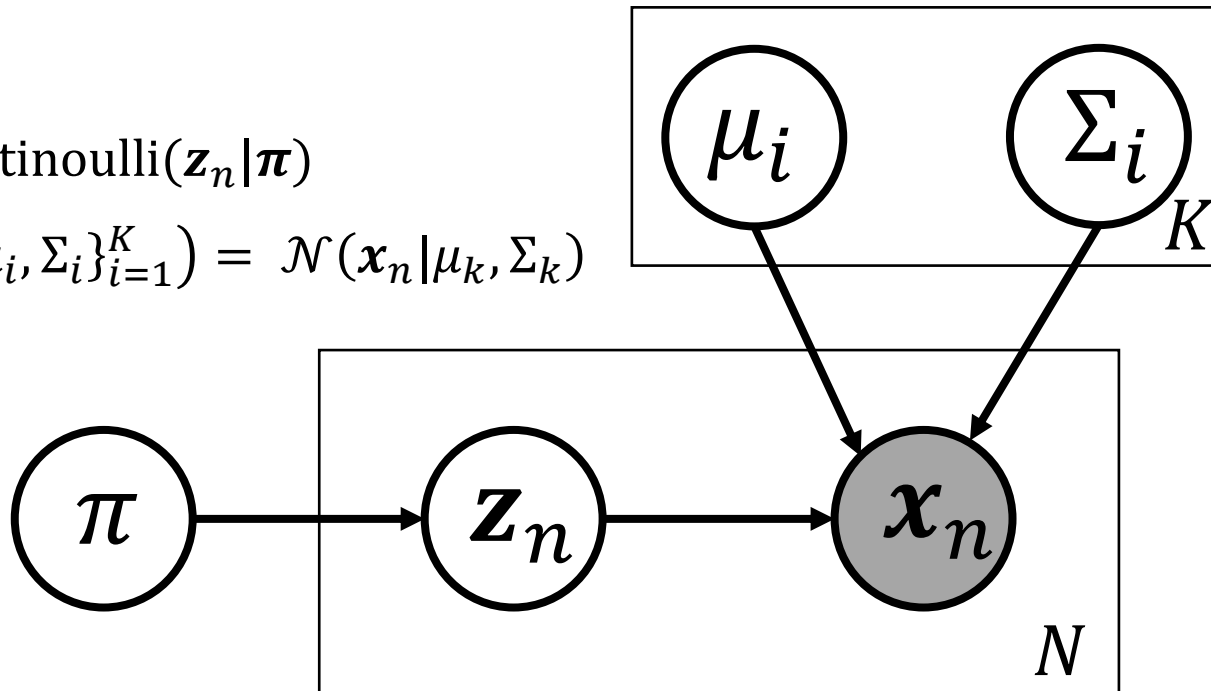


# Gaussian Mixture Model (GMM)

- $N$  observations  $\{\mathbf{x}_n\}_{n=1}^N$  each from one of the  $K$  Gaussians  $\{\mathcal{N}(\mu_i, \Sigma_i)\}_{i=1}^K$
- We don't know which Gaussian each observation  $\mathbf{x}_n$  comes from
- Assume  $\mathbf{z}_n \in \{1, 2, \dots, K\}$  denotes which Gaussian generated  $\mathbf{x}_n$
- Suppose we want to do point estimation for the parameters  $\{\mu_i, \Sigma_i\}_{i=1}^K$

$$p(\mathbf{z}_n | \boldsymbol{\pi}) = \text{multinoulli}(\mathbf{z}_n | \boldsymbol{\pi})$$

$$p(\mathbf{x}_n | \mathbf{z}_n = k, \{\mu_i, \Sigma_i\}_{i=1}^K) = \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$



$$p(\mathbf{x}_n | \{\pi_i, \mu_i, \Sigma_i\}_{i=1}^K)$$

$$= \sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{x}_n | \mu_i, \Sigma_i)$$

$$\log p(\mathbf{x}_n | \boldsymbol{\Theta}) = \log \sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{x}_n | \mu_i, \Sigma_i)$$

Can use gradient based optimization for MLE of  $\boldsymbol{\Theta}$  but the update equations are a bit complicated

EM would give simpler updates



# Detour: MLE for GMM when $\mathbf{Z}$ is known

GMM then is just like generative classification with Gaussian class conditionals

- Derivation of the MLE solution for  $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  when  $\mathbf{Z}$  is known

$$\begin{aligned}\hat{\Theta} &= \operatorname{argmax}_{\Theta} p(\mathbf{X}, \mathbf{Z} | \Theta) = \operatorname{argmax}_{\Theta} \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \Theta) \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \underbrace{p(\mathbf{z}_n | \Theta)}_{\text{multinoulli}} \underbrace{p(\mathbf{x}_n | \mathbf{z}_n, \Theta)}_{\text{Gaussian}} \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \prod_{k=1}^K p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta)^{z_{nk}} \\ &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta)]^{z_{nk}} \\ &= \operatorname{argmax}_{\Theta} \log \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta)]^{z_{nk}} \\ &= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K z_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]\end{aligned}$$

In general, in models with probability distributions from the **exponential family**, the MLE problem will usually have a simple analytic form

Also, due to the form of the likelihood (Gaussian) and prior (multinoulli), the MLE problem had a nice separable structure after taking the log

Can see that, when estimating the parameters of the  $k^{\text{th}}$  Gaussian  $(\pi_k, \mu_k, \Sigma_k)$ , we only will only need training examples from the  $k^{\text{th}}$  class, i.e., examples for which  $z_{nk} = 1$



# EM for Gaussian Mixture Model (GMM)

1. Initialize  $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  as  $\Theta^{(0)}$ . Set  $t = 1$
2. Set CP  $q^{(t)} = p(\mathbf{Z}|\mathbf{X}, \Theta^{(t-1)})$ . Assuming i.i.d. data, this means computing  $\forall n, k$

Probability of data point  $n$   
belonging to the  $k$ -th Gaussian

"Soft-clustering"

Same as writing  $z_n = k$

$$p(\mathbf{z}_{nk} = 1 | \mathbf{x}_n, \Theta^{(t-1)}) \propto p(\mathbf{z}_{nk} = 1 | \Theta^{(t-1)}) p(\mathbf{x}_n | \mathbf{z}_{nk} = 1, \Theta^{(t-1)})$$

$$= \pi_k^{(t-1)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})$$

3. Set  $\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathbb{E}_{q^{(t)}} [\log p(\mathbf{X}, \mathbf{Z} | \Theta)] = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{(t-1)})$

This only required expectation for EM  
for GMM is  $\mathbb{E}[\mathbf{z}_{nk}]$  which can be  
computed easily using the CP of  $\mathbf{z}_n$

EM for GMM does **two**  
**things**: soft-clustering  
and estimating the  
density  $p(\mathbf{X} | \Theta)$



$$\Theta^{(t)} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \Theta^{(t-1)})} [\log p(\mathbf{x}_n, \mathbf{z}_n | \Theta)]$$

$$= \operatorname{argmax}_{\Theta} \mathbb{E} \left[ \sum_{n=1}^N \sum_{k=1}^K \mathbf{z}_{nk} \left[ \log \pi_k^{(t-1)} + \log \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)}) \right] \right]$$

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[\mathbf{z}_{nk}] [\log \pi_k^{(t-1)} + \log \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})]$$

$$\pi_k^{(t)} = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}]$$

$$\mu_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}] \mathbf{x}_n$$

$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}] (\mathbf{x}_n - \mu_k^{(t)}) (\mathbf{x}_n - \mu_k^{(t)})^\top$$

$N_k = \sum_{n=1}^N \mathbb{E}[\mathbf{z}_{nk}]$   
denotes the effective  
number of points  
from  $k$ -th Gaussian

4. Go to step 2 if not converged



# Bayesian Linear Regression (Revisited)

$N \times D$  input matrix

$N \times 1$  responses

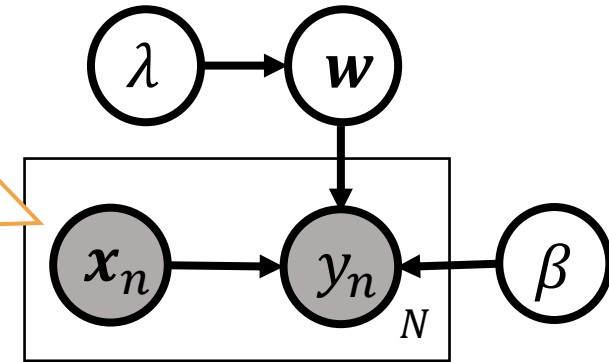
- $N$  observations  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$  from a lin-reg model with weights  $\mathbf{w}$
- Suppose the hyperparameters are also unknown, so need to estimate  $\mathbf{w}, \beta, \lambda$

$$p(y_n | \mathbf{x}_n, \mathbf{w}, \beta) = \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad p(\mathbf{w} | \lambda) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I})$$

CP of  $\mathbf{w}$ :  $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\boldsymbol{\Sigma} = (\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \quad \boldsymbol{\mu} = \beta \mathbf{A}^{-1} \mathbf{X}^\top \mathbf{y}$$

In this latent variable model, there are no local variables.  $\mathbf{w}, \beta, \lambda$  are all “global”



Many ways to optimize the marginal likelihood in MLE-II, e.g., gradient descent

MLE-II  $(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \log p(\mathbf{y} | \mathbf{X}, \beta, \lambda)$

EM solves the MLE-II problem by optimizing a lower bound on the log marginal likelihood, and gives simple update equations for  $\beta, \lambda$

EM

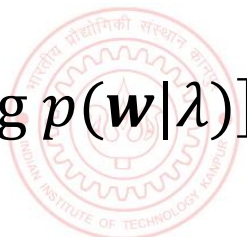
Expected CLL

$$(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \mathbb{E}_{p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)})} [\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \beta, \lambda)]$$

Data

$$= \operatorname{argmax}_{\beta, \lambda} \mathbb{E}_{p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)})} [\log p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \beta) + \log p(\mathbf{w} | \lambda)]$$

$\mathbf{w}$  treated as latent variable here



# EM for Bayesian Linear Regression

$$(\beta^{(t)}, \lambda^{(t)}) = \operatorname{argmax}_{\beta, \lambda} \mathbb{E}[\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \beta^{(t-1)}, \lambda^{(t-1)})]$$

1. Initialize  $\beta$  as  $\beta^{(0)}$  and  $\lambda$  as  $\lambda^{(0)}$ . Set  $t = 1$

2. Update the CP of  $\mathbf{w}$  as

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)}) = \mathcal{N}(\boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})$$

$$\boldsymbol{\Sigma}^{(t)} = (\beta^{(t-1)} \mathbf{X}^\top \mathbf{X} + \lambda^{(t-1)} \mathbf{I})^{-1} \quad \boldsymbol{\mu}^{(t)} = \beta^{(t-1)} \boldsymbol{\Sigma}^{(t)} \mathbf{X}^\top \mathbf{y}$$

3. Update  $\beta$  and  $\lambda$  as

$$\lambda^{(t)} = \frac{D}{\mathbb{E}[\mathbf{w}^\top \mathbf{w}]} = \frac{D}{\boldsymbol{\mu}^{(t)\top} \boldsymbol{\mu}^{(t)} + \operatorname{trace}(\boldsymbol{\Sigma}^{(t)})}$$

$$\beta^{(t)} = \frac{N}{\|\mathbf{y} - \mathbf{X} \boldsymbol{\mu}^{(t)}\|^2 + \operatorname{trace}(\mathbf{X}^\top \boldsymbol{\Sigma}^{(t)} \mathbf{X})}$$

4. If not converged, set  $t = t + 1$  and go to step 2

Note the dependence: CP of  $\mathbf{w}$  depends on current values of  $\beta, \lambda$  and their update depends on the CP on  $\mathbf{w}$



Less common but another alternative: Compute CP of  $\beta$  and  $\lambda$  in step 2, and compute MLE on  $\mathbf{w}$  in step 3. That would amount to doing MLE-II for  $\mathbf{w}$





# Extra: MLE-II for Bayesian Lin. Reg.

- The MLE-II problem for Bayesian linear regression

$$\begin{aligned}(\hat{\beta}, \hat{\lambda}) &= \operatorname{argmax}_{\beta, \lambda} \log p(\mathbf{y} | \mathbf{X}, \beta, \lambda) \\ &= \operatorname{argmax}_{\beta, \lambda} (2\pi)^{-\frac{N}{2}} |\beta^{-1} \mathbf{I} + \lambda^{-1} \mathbf{X}^T \mathbf{X}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} \mathbf{y}^T (\beta^{-1} \mathbf{I} + \lambda^{-1} \mathbf{X}^T \mathbf{X})^{-1} \mathbf{y} \right)\end{aligned}$$

- This objective doesn't have a closed form solution
- Solved using iterative/alternating optimization
  - Gradient descent for  $\lambda, \beta$
  - Alternating optimization ( $\lambda, \beta$  and the mean/covariance of the CP depend on each other) - similar to EM but with some differences – next slide
- EM is also a way to do MLE-II but EM doesn't optimize the marginal likelihood but a lower bound on the marginal likelihood



# An algorithm for MLE-II for Bayesian Lin. Reg.

1. Initialize  $\beta$  as  $\beta^{(0)}$  and  $\lambda$  as  $\lambda^{(0)}$ . Set  $t = 1$

$$(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \log p(\mathbf{y} | \mathbf{X}, \beta, \lambda)$$

2. Update the CP of  $\mathbf{w}$  as

$$p(\mathbf{w}^{(t)} | \mathbf{X}, \mathbf{y}, \beta^{(t-1)}, \lambda^{(t-1)}) = \mathcal{N}(\boldsymbol{\mu}^{(t)}, \mathbf{A}^{(t)^{-1}})$$

$$\mathbf{A}^{(t)} = \beta^{(t-1)} \mathbf{X}^\top \mathbf{X} + \lambda^{(t-1)} \mathbf{I} \quad \boldsymbol{\mu}^{(t)} = \beta^{(t-1)} \mathbf{A}^{(t)^{-1}} \mathbf{X}^\top \mathbf{y}$$

3. Update  $\beta, \lambda$  as

RHS depends on  $\beta$  and  $\lambda$ .  
Thus it is an implicit solution  
(though still in closed form)

$$\lambda^{(t)} = \frac{\gamma^{(t)}}{\boldsymbol{\mu}^{(t)\top} \boldsymbol{\mu}^{(t)}}$$

In practice, we can compute them in the beginning for  $\mathbf{X}^\top \mathbf{X}$  and multiply by  $\beta^{(t-1)}$  in this iteration to get  $\{\eta_d^{(t)}\}_{d=1}^D$

In each iteration, we need to compute the eigenvalues

$$\{\eta_d^{(t)}\}_{d=1}^D = \text{eigvals}(\beta^{(t-1)} \mathbf{X}^\top \mathbf{X})$$

where

$$\gamma^{(t)} = \sum_{d=1}^D \frac{\eta_d^{(t)}}{\lambda^{(t-1)} + \eta_d^{(t)}}$$

RHS depends on  $\beta$  and  $\lambda$ .  
Thus it is an implicit solution  
(though still in closed form)

$$\beta^{(t)} = \frac{N - \gamma^{(t)}}{\|\mathbf{y} - \mathbf{X} \boldsymbol{\mu}^{(t)}\|^2}$$

4. If not converged, set  $t = t + 1$  and go to step 2

Note that this MLE-II procedure for Bayesian linear regression looks very similar to the EM algo for BLR



# EM: Some other examples

- Problems with missing features (which are treated as latent variables)
  - Suppose each input  $\mathbf{x}_n$  has two parts - observed and missing:  $\mathbf{x}_n = [\mathbf{x}_n^{obs}, \mathbf{x}_n^{miss}]$
  - For such problems, MLE for a model  $p(\mathbf{X}|\Theta)$ , assuming i.i.d. data, would have the form

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n^{obs} | \Theta)$$

Suppose we are estimating the mean/covariance of a multivariate Gaussian given  $N$  input, with some inputs observations may have missing features

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log \int p([\mathbf{x}_n^{obs}, \mathbf{x}_n^{miss}] | \Theta) d\mathbf{x}_n^{miss}$$

- Here  $\mathbf{x}_n^{miss}$  can be treated as a latent variable
- The CP will be  $p(\mathbf{x}_n^{miss} | \mathbf{x}_n^{obs}, \Theta)$
- Using the CP, compute expected CLL and maximize it w.r.t.  $\Theta$
- Problems with missing labels (which are treated as latent variables)

An example of semi-supervised learning

This part is like GMM, thus EM can be used

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log p(x_n, y_n | \Theta) + \sum_{n=N+1}^{N+M} \log \sum_{c=1}^K p(x_n, y_n = c | \Theta)$$



# EM when CP and/or expectation is intractable

- EM solves the following step for estimating  $\Theta$

$$\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathbb{E}_{q^{(t)}}[\log p(\mathcal{D}, \mathbf{Z}|\Theta)] = \operatorname{argmax}_{\Theta} \int \log p(\mathcal{D}, \mathbf{Z}|\Theta) p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D}) d\mathbf{Z}$$

- The above problem may be difficult to solve if one/both of the following is true
  1. CP  $p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D})$  can't be computed exactly (Solution: Need to approximate the CP)
  2. Integral for the expectation is intractable (Solution: Use **Monte Carlo approximation**)
    - Draw  $M$  i.i.d. samples of  $\mathbf{Z}$  from the current (exact/approximate) CP  $p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D})$

$$\{\mathbf{Z}^{(i)}\}_{i=1}^M \sim p(\mathbf{Z}|\Theta^{(t-1)}, \mathcal{D})$$

- Use these samples to get a Monte-Carlo approximation of expected CLL and maximize

$$\Theta^{(t)} = \operatorname{argmax}_{\Theta} \frac{1}{M} \sum_{i=1}^M \log p(\mathcal{D}, \mathbf{Z}^{(i)}|\Theta)$$

- Monte-Carlo approximation is commonly used in such problems

