

## Lecture 13

### Bound optimization

General idea:

Goal: maximize some function

$\ell(\vec{\theta})$  wrt  $\vec{\theta}$ .

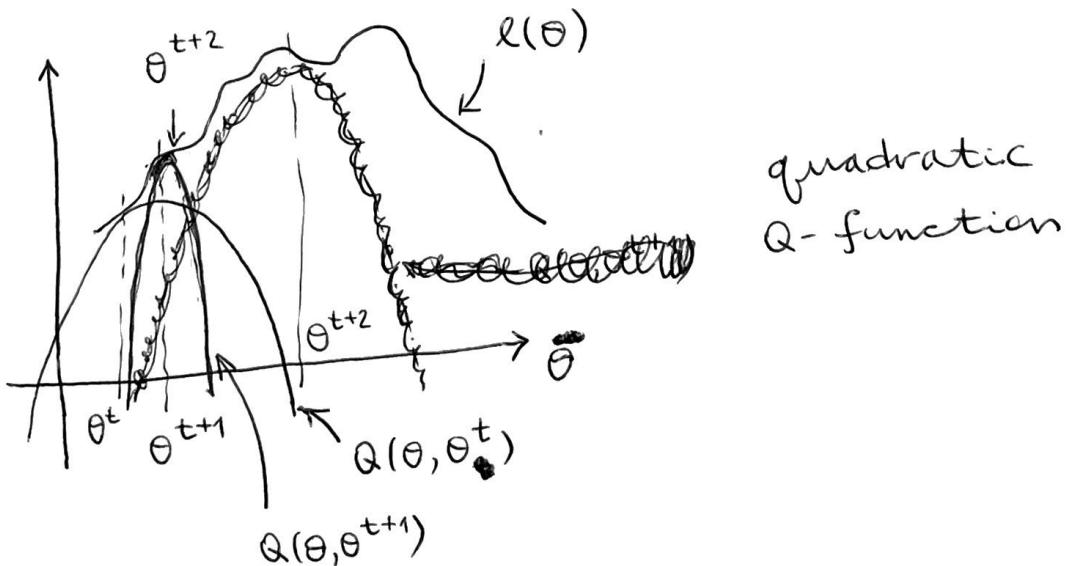
→ Construct a surrogate function  $Q(\vec{\theta}, \vec{\theta}^t)$   
s.t.  $\begin{cases} \ell(\vec{\theta}) \geq Q(\vec{\theta}, \vec{\theta}^t) \text{ and} \\ \ell(\vec{\theta}^t) = Q(\vec{\theta}^t, \vec{\theta}^t). \end{cases}$   
current  
prm values

→ Find  $\vec{\theta}^{t+1} = \arg \max_{\vec{\theta}} Q(\vec{\theta}, \vec{\theta}^t)$

Then  $\ell(\vec{\theta}^{t+1}) \geq Q(\vec{\theta}^{t+1}, \vec{\theta}^t) \geq Q(\vec{\theta}^t, \vec{\theta}^t) = \ell(\vec{\theta}^t)$

$\ell(\vec{\theta}^{t+1}) \geq \ell(\vec{\theta}^t)$  monotonically at each step

Ex.



[Expectation maximization (EM)]  
algorithm

$\begin{cases} \vec{y}_n = \text{observable data } (n=1, \dots, N) \\ \vec{z}_n = \text{hidden variables} \end{cases}$

alternation between

→ E-step (expectation step)  $\Rightarrow$  estimate hidden vars

→ M-step (maximization step)  $\Rightarrow$  compute MLE or MAP parameter estimates given the hidden vars & the observed data

Maximize log-likelihood:

$$\ell(\vec{\theta}) = \sum_{n=1}^N \log p(\vec{y}_n | \vec{\theta}) = \sum_{n=1}^N \log \left( \sum_{\vec{z}_n} p(\vec{y}_n, \vec{z}_n | \vec{\theta}) \right)$$

Consider  $\{q_{fn}(\vec{z}_n)\}_{n=1}^N$  - a set of some prob. distributions over hidden vars  $\vec{z}_n: \sum_{\vec{z}_n} q_{fn}(\vec{z}_n) = 1$

$$\begin{aligned} \text{Then } \ell(\vec{\theta}) &= \sum_n \log \left( \sum_{\vec{z}_n} q_{fn}(\vec{z}_n) \frac{p(\vec{y}_n, \vec{z}_n | \vec{\theta})}{q_{fn}(\vec{z}_n)} \right) \geq \\ &\geq \sum_{n, \vec{z}_n} q_{fn}(\vec{z}_n) \log \frac{p(\vec{y}_n, \vec{z}_n | \vec{\theta})}{q_{fn}(\vec{z}_n)} \equiv \underbrace{\sum_n \tilde{I}(\vec{\theta}, q_{fn})}_{\text{evidence lower bound (ELBO)}} \end{aligned}$$

Jensen's  
inequality

Now, consider

E-step

$$\begin{aligned}\tilde{\mathcal{I}}(\vec{\theta}, q_{\theta n}) &= \sum_{\vec{z}_n} q_{\theta n}(\vec{z}_n) \log \frac{p(\vec{z}_n | \vec{y}_n, \vec{\theta}) p(\vec{y}_n | \vec{\theta})}{q_{\theta n}(\vec{z}_n)} = \\ &= \sum_{\vec{z}_n} q_{\theta n}(\vec{z}_n) \log p(\vec{y}_n | \vec{\theta}) + \sum_{\vec{z}_n} q_{\theta n}(\vec{z}_n) \log \frac{p(\vec{z}_n | \vec{y}_n, \vec{\theta})}{q_{\theta n}(\vec{z}_n)} \quad \text{(ii)} \\ \textcircled{=} \quad \log p(\vec{y}_n | \vec{\theta}) - \underbrace{D_{KL}(q_{\theta n}(\vec{z}_n) || p(\vec{z}_n | \vec{y}_n, \vec{\theta}))}_{\geq 0}.\end{aligned}$$

↑  
 $\sum_{\vec{z}_n} q_{\theta n}(\vec{z}_n) = 1$       ↙ tight lower bound

or  $q_{\theta n}(\vec{z}_n) = p(\vec{z}_n | \vec{y}_n, \vec{\theta}),$

$$\underline{\tilde{\mathcal{I}}(\vec{\theta}, q_{\theta n})} = \log p(\vec{y}_n | \vec{\theta}). \quad (*)$$

Finally, define  $Q(\vec{\theta}, \vec{\theta}^t) = \sum_n \tilde{\mathcal{I}}(\vec{\theta}, \underbrace{p(\vec{z}_n | \vec{y}_n, \vec{\theta}^t)}_{q_{\theta n}(\vec{z}_n)})$

Indeed, by Eq. (\*)  $Q(\vec{\theta}^t, \vec{\theta}^t) = \sum_n \log p(\vec{y}_n | \vec{\theta}^t) = \ell(\vec{\theta}^t).$

Moreover,  $Q(\vec{\theta}, \vec{\theta}^t) \leq \ell(\vec{\theta})$  due to the  $\geq 0 D_{KL}$  term.

## M-step

Need to maximize  $\sum_n \tilde{L}(\vec{\theta}, p(\vec{z}_n | \vec{y}_n, \vec{\theta}^t))$

wrt  $\vec{\theta}$ . Since  $-\sum_{\vec{z}_n} q_{\theta_n}(\vec{z}_n) \log q_{\theta_n}(\vec{z}_n) = H(q_{\theta_n})$  does not depend on  $\vec{\theta}$ , it suffices to maximize

$$\sum_{n, \vec{z}_n} q_{\theta_n}(\vec{z}_n) \log p(\vec{y}_n, \vec{z}_n | \vec{\theta}) \text{ to get } \vec{\theta}^{t+1}.$$

↑  
expectation value wrt  $q_{\theta_n}(\vec{z}_n)$

Ex. Gaussian mixture model (GMM)

$$p(\vec{y} | \vec{\theta}) = \sum_{k=1}^K \pi_k N(\vec{y} | \vec{\mu}_k, \Sigma_k)$$

application:  
clustering of  $\vec{y}_n \in \mathbb{R}^D$   
associate each data point  $\vec{y}_n$  with a hidden var  $z_n \in \{1, \dots, K\}$  which labels clusters:

$$p(z_n = k | \vec{y}_n, \vec{\theta}) = \frac{p(z_n = k | \vec{\theta}) p(\vec{y}_n | z_n = k, \vec{\theta})}{\sum_{k'=1}^K p(z_n = k' | \vec{\theta}) p(\vec{y}_n | z_n = k', \vec{\theta})}$$

responsibility  
of cluster  $k$  for  
datapoint  $n$

E-step

$$r_{nk}^t = \frac{\pi_k^t p(\tilde{y}_n | \tilde{\theta}_k^{(t)})}{\sum_{k'} \pi_{k'}^t p(\tilde{y}_n, \tilde{\theta}_{k'}^{(t)})}, \text{ where}$$

$\tilde{\theta}_k$  = prms describing mixture component  $k$ .

M-step  $\ell^t(\tilde{\theta}) = E \left[ \sum_n (\log p(z_n | \tilde{\mu}) + \log p(\tilde{y}_n | z_n, \tilde{\theta})) \right]$  wrt  $\tilde{\theta}$   $\begin{pmatrix} \pi_1 \\ \vdots \\ \pi_K \end{pmatrix}$  (1)

$$\therefore E \left[ \sum_n \log \left( \prod_k \pi_k^{z_{nk}} \right) + \sum_n \log \left( \prod_k N(\tilde{y}_n | \tilde{\mu}_k, \Sigma_k)^{z_{nk}} \right) \right]$$

$\uparrow z_{nk} = \mathbb{I}(z_n = k)$  'one-hot' encoding

$$= \sum_{n,k} \underbrace{E[z_{nk}]}_{r_{nk}^t} \log \pi_k + \sum_{n,k} E[z_{nk}] \log N(\tilde{y}_n | \tilde{\mu}_k, \Sigma_k) =$$

$$= \sum_{n,k} r_{nk}^t \log \pi_k - \frac{1}{2} \sum_{n,k} r_{nk}^t [(\tilde{y}_n - \tilde{\mu}_k)^T \Sigma_k^{-1} (\tilde{y}_n - \tilde{\mu}_k) + \log |\Sigma_k|] + \text{const}$$

=====

$\ell^t(\tilde{\theta})$  is maximized by

$$\left\{ \begin{array}{l} \tilde{\mu}_k^{t+1} = \frac{\sum_n r_{nk}^t \tilde{y}_n}{\sum_n r_{nk}^t}, \quad r_k^t = \sum_n r_{nk}^t \\ \Sigma_k^{t+1} = \frac{\sum_n r_{nk}^t (\tilde{y}_n - \tilde{\mu}_k^{t+1})(\tilde{y}_n - \tilde{\mu}_k^{t+1})^T}{\sum_n r_{nk}^t}, \\ \pi_k^{t+1} = \frac{1}{N} \sum_n r_{nk}^t = \frac{r_k^t}{N} \end{array} \right. \quad (+)$$

These priors can be used to estimate  $r_{nk}^{t+1}$  in the next E-step, etc.

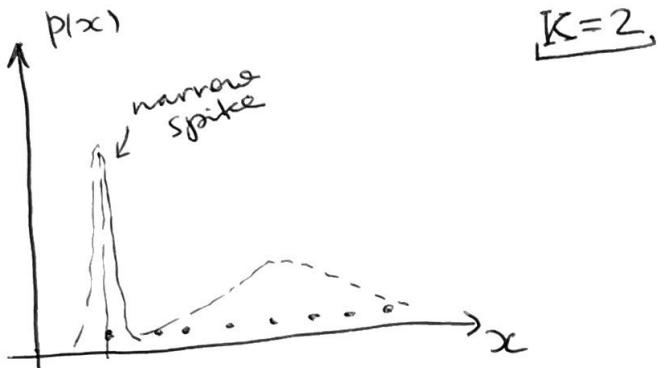
Note: MLE on the GMM might overfit.

Suppose that  $\Sigma_k = \sigma_k^2 \mathbb{I}$ ,  $\forall k$ .

It's possible to center one component on a single datapoint:  $\tilde{\mu}_{k'} = \tilde{y}_n$ , say.

$$\text{Then } N(\tilde{y}_n | \tilde{\mu}_{k'}, \sigma_{k'}^2 \mathbb{I}) = \frac{1}{\sqrt{2\pi\sigma_{k'}^2}}$$

This will  $\rightarrow \infty$  as  $\sigma_{k'} \rightarrow 0 \Rightarrow$  "collapsing variance problem".



Solution: switch from MLE to MAP estimation by introducing priors into Eq. (+); modified M-step.  $\downarrow$  pseudocounts

$$\text{For example, } \tilde{r}_{k'}^{t+1} = \frac{r_{k'}^t + d_{k'} - 1}{N + \sum_k d_k - K}.$$

Dirichlet prior