

Machine Learning

4771

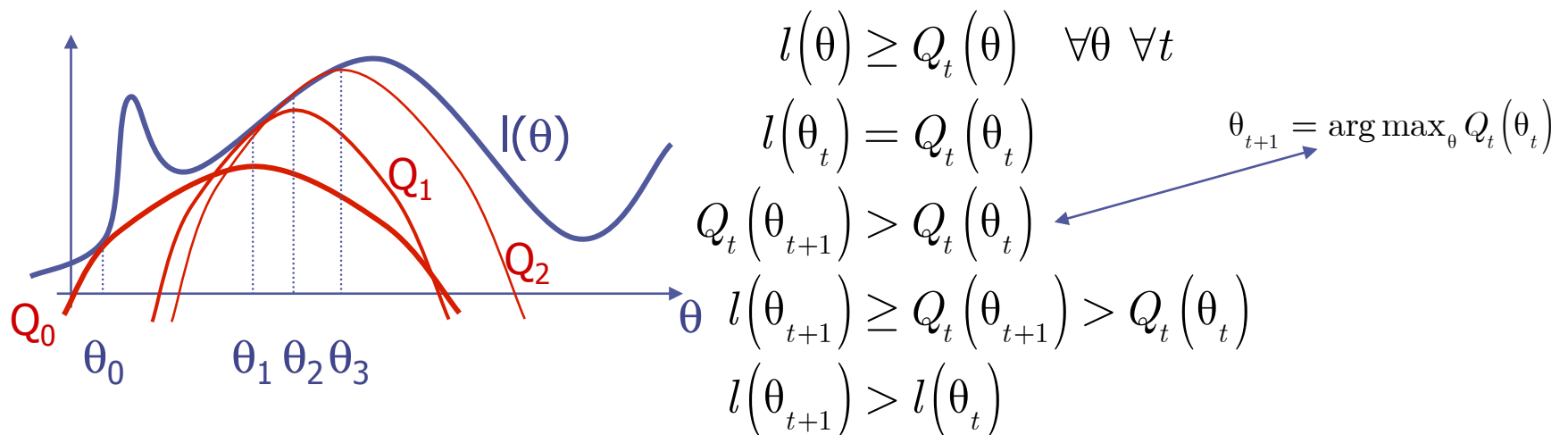
Instructor: Tony Jebara

Topic 13

- Expectation Maximization as Bound Maximization
- EM for Maximum A Posteriori

EM as Bound Maximization

- Let's now show that EM indeed maximizes likelihood
- **Bound Maximization:** optimize a lower bound on $l(\theta)$
- Since log-likelihood $l(\theta)$ not concave, can't max it directly
- Consider an auxiliary function $Q(\theta)$ which is concave
- $Q(\theta)$ kisses $l(\theta)$ at a point and is less than it elsewhere



- Monotonically increases log-likelihood
- But how to find a bound and guarantee we max it?



Jensen's Inequality

- An important general bound from Jensen (1906)

- For convex f : $f\left(E\{x\}\right) \leq E\{f(x)\}$

- For concave f : $f\left(E\{x\}\right) \geq E\{f(x)\}$

- Expectation in discrete case is sum weight by probability

- For convex f : $f\left(\sum_{i=1}^M p_i x_i\right) \leq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1$, $p_i \geq 0$

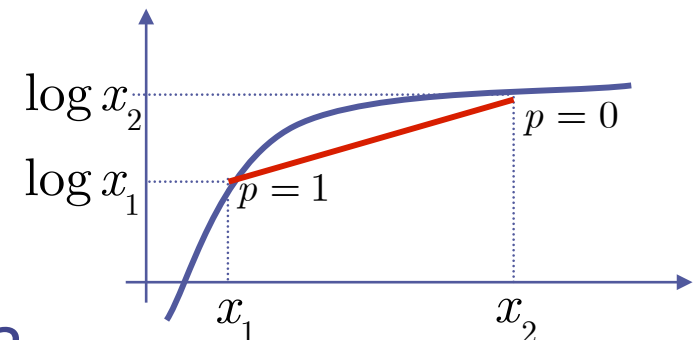
- For concave f : $f\left(\sum_{i=1}^M p_i x_i\right) \geq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1$, $p_i \geq 0$

- Example: $f(x) = \log(x)$ = concave and $M=2$

$$\log\left(px_1 + (1-p)x_2\right) \geq p \log x_1 + (1-p) \log x_2$$

- Bound $\log(\text{sum})$ with $\text{sum}(\log)$

- How to apply this to mixture models?



Expectation-Maximization

$$\begin{aligned}
 l(\theta) &= \sum_{n=1}^N \log p(x_n | \theta) && \leftarrow \text{Original Log-Likelihood} \\
 &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) && \leftarrow \text{Has Hidden Variables (messy)} \\
 &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) \frac{p(z | x_n, \theta_t)}{p(z | x_n, \theta_t)} && \leftarrow \text{Multiply by 1} \\
 & && \leftarrow \text{Ratio of hidden posterior density} \\
 &= \sum_{n=1}^N \log \sum_z p(z | x_n, \theta_t) \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \leftarrow \text{Rearrange} \\
 &\geq \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \leftarrow \text{Jensen log}(\sum_i p_i x_i) \\
 &= \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(x_n, z | \theta) \\
 & \quad - \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(z | x_n, \theta_t) \\
 &= Q(\theta | \theta_t) - \text{const} && \leftarrow \text{New auxiliary function called Q (not messy)}
 \end{aligned}$$

EM as Bound Maximization

- Now have the following bound and maximize it:

$$l(\theta) \geq Q(\theta | \theta_t) - \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(z | x_n, \theta_t)$$

$$\begin{aligned} \theta^{t+1} &= \arg \max_{\theta} Q(\theta | \theta_t) = \arg \max_{\theta} \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(x_n, z | \theta) \\ &= \arg \max_{\theta} \sum_{n=1}^N \sum_z \tau_{n,z} \log p(x_n, z | \theta) \end{aligned}$$

- $Q(\theta | \theta_t)$ is called **Auxiliary Function**... take derivatives of it
- This is easy for e-families... just weighted max likelihood!
- For example, Gaussian mixture:

$$\frac{\partial Q(\theta)}{\partial \vec{\mu}_k} = \frac{\partial}{\partial \vec{\mu}_k} \sum_{n=1}^N \sum_k \tau_{n,k} \log \pi_k N(\vec{x}_n | \vec{\mu}_k, \Sigma_k)$$

$$0 = \sum_{n=1}^N \tau_{n,k} \frac{\partial}{\partial \vec{\mu}_k} \left(-\frac{1}{2} (\vec{x}_n - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x}_n - \vec{\mu}_k) \right)$$

$$\vec{\mu}_k = \frac{\sum_{n=1}^N \tau_{n,k} \vec{x}_n}{\sum_{n=1}^N \tau_{n,k}}$$

... similarly get π_k and Σ_k

EM as Expected Log-Likelihood

- **Incomplete Log-Likelihood** $l(\theta) = \log p(\text{observed} | \theta) = \sum_{n=1}^N \log \sum_z p(x_n, z | \theta)$
- **Complete Log-Likelihood** $l^C(\theta) = \log p(\text{observed, hidden} | \theta) = \sum_{n=1}^N \log p(x_n, z_n | \theta)$
- We don't know the hidden variables z
- EM computes expected values of hidden z under current θ_t
- EM chooses Q to be the **Expected Complete Log-Likelihood**

$$\begin{aligned}
 E \{l^C(\theta)\} &= \sum_{\text{hidden}} p(\text{hidden} | \text{observed}, \theta_t) l^C(\theta) \\
 &= \sum_{z_1} \cdots \sum_{z_N} p(z_1, \dots, z_N | x_1, \dots, x_N, \theta_t) l^C(\theta) \\
 &= \sum_{z_1} \cdots \sum_{z_N} \prod_n p(z_n | x_n, \theta_t) l^C(\theta) \\
 &= \sum_{z_1} \cdots \sum_{z_N} \prod_n p(z_n | x_n, \theta_t) \sum_n \log p(x_n, z_n | \theta) \\
 &= \sum_n \sum_{z_n} p(z_n | x_n, \theta_t) \log p(x_n, z_n | \theta) \sum_{z_1} \cdots \sum_{z_{i \neq n}} \cdots \sum_{z_N} \prod_{i \neq n} p(z_i | x_i, \theta_t) \\
 &= \sum_n \sum_{z_n} p(z_n | x_n, \theta_t) \log p(x_n, z_n | \theta) = Q(\theta | \theta_t)
 \end{aligned}$$

EM for Max A Posteriori

- We can also do MAP instead of ML with EM (stabilizes sol'n)

$$\log \text{posterior}(\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) + \log p(\theta)$$

- Prior doesn't have log-sum
- The E-step remains the same: lower bound log-sum

$$\log \text{posterior}(\theta) = l(\theta) + \log p(\theta) \geq E\{l^C(\theta)\} + \text{const} + \log p(\theta)$$

- The M-step becomes slightly different for each model

- For example, mixture of Gaussians with prior on covariance

$$\log \text{posterior}(\theta) = \sum_{n=1}^N \log \sum_k \pi_k N(\vec{x}_n | \vec{\mu}_k, \Sigma_k) + \log \prod_k p(\Sigma_k | S, \eta)$$

$$\log \text{posterior}(\theta) \geq \sum_{n=1}^N \sum_k \tau_{n,k} \log \pi_k N(\vec{x}_n | \vec{\mu}_k, \Sigma_k) + \sum_k \log p(\Sigma_k | S, \eta) + \text{const}$$

- Updates on π and μ stay the same, only Σ is:

$$\Sigma_k \leftarrow \frac{1}{\sum_{n=1}^N \tau_{n,k} + \eta} \left(\sum_{n=1}^N \tau_{n,k} (\vec{x}_n - \vec{\mu}_k)(\vec{x}_n - \vec{\mu}_k)^T + \eta S \right)$$

- Typically, we use the identity matrix I for S and a small eta.