

General EM

- $X = \text{observed}$
- $Z = \text{latent}$

• Goal: Maximize $p(x|\theta) = \sum_z p(x, z|\theta)$

- Decompose $\log(p(x|\theta))$:

$$\log p(x|\theta) = \mathcal{L}(q, \theta) + \text{KL}(q||p) \quad \text{where } q \text{ is any p.d.f.}$$

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log \left(\frac{p(x, z|\theta)}{q(z)} \right) \quad \text{ELBO}$$

$$\text{KL}(q||p) = - \sum_z q(z) \log \left(\frac{p(z|x, \theta)}{q(z)} \right) \quad \text{KL-Divergence}$$

↑ this holds for any distribution over z ($q(z)$)
(not just the true distribution of z)

proof:

$$\mathcal{L}(q, \theta) = \sum_z q(z) \left(\log \frac{p(x, z|\theta)}{q(z)} \right) = \sum_z q(z) \log \left(\frac{p(z|x, \theta) p(x|\theta)}{q(z)} \right)$$

$$= \sum_z q(z) \log \left(\frac{p(z|x, \theta)}{q(z)} \right) + \sum_z q(z) \underbrace{\log p(x|\theta)}_{=1}$$

$$= -\text{KL}(q||p) + \log p(x|\theta) //$$

Note: $KL(q||p)$ is always ≥ 0

proof:

$$KL(q||p) = E \left[-\log \left(\frac{p(z|x,\theta)}{q(z)} \right) \right] \quad -\log \text{ is convex}$$

$$\Rightarrow \text{jensen's inequality } \geq -\log E \left[\frac{p(z|x,\theta)}{q(z)} \right]$$

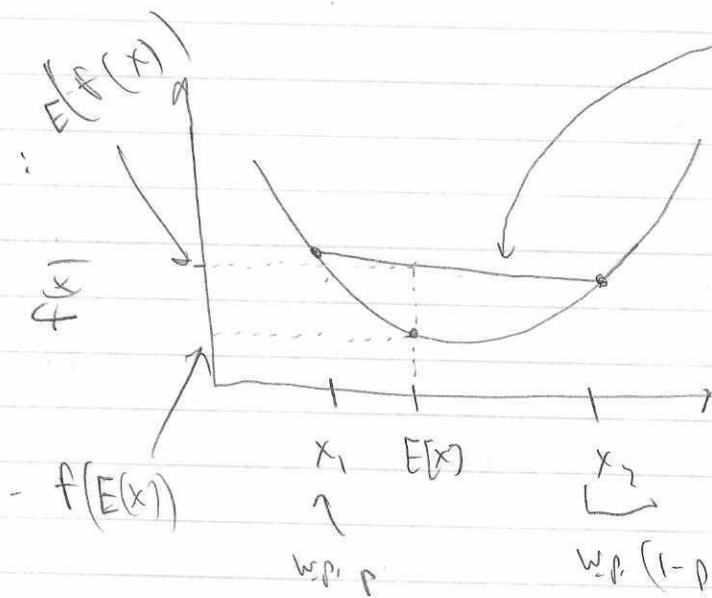
$$= -\log \left(\sum q(z) \frac{p(z|x,\theta)}{q(z)} \right)$$

$$= -\log \left(\underbrace{\sum p(z|x,\theta)}_{=1} \right)$$

$$= -\log(1)$$

$$= 0$$

Jensen's:



$$p f(x_1) + (1-p) f(x_2) \geq p f(x)$$

so $E(f(x)) \geq f(E(x))$ for $f()$ convex

$f() = -\log()$ here

• So $\mathcal{L}(q, \theta)$ is a lower bound on $\log p(x|\theta)$

$$\mathcal{L}(q, \theta) \leq \log p(x|\theta)$$

• EM Algorithm:

$$\log p(x|\theta) = \mathcal{L}(q, \theta) + \overbrace{\text{KL}(q||p)}^{\geq 0}$$

Step 1: Maximize $\mathcal{L}(q, \theta)$ over q , holding θ fixed

↑ during this step, $p(x|\theta)$ is constant, so $\mathcal{L}(q, \theta)$ is maximized when $\text{KL}(q||p)$ is 0

$$\text{KL}(q||p) = -\sum_z q(z) \log \left[\frac{p(z|x, \theta)}{q(z)} \right]$$

↖ if this equals $p(z|x, \theta)$, then

$$= -\sum_z q(z) \log(1) = 0$$

So E-step: set $q(z) = p(z|x, \theta)$

↑ posterior of $z|x$

Step 2: Maximize $\mathcal{L}(q, \theta)$ over θ holding q fixed.

↑ this increases the lower bound on $\log p(x|\theta)$

↑ As you change θ , $q(z)$ becomes no longer equal to posterior, so $KL \uparrow$

$$\log p(x|\theta) = \underbrace{\mathcal{L}(q, \theta)}_{\uparrow} + \underbrace{KL(q||p)}_{\uparrow}$$

$\Rightarrow \log p(x|\theta) \uparrow$

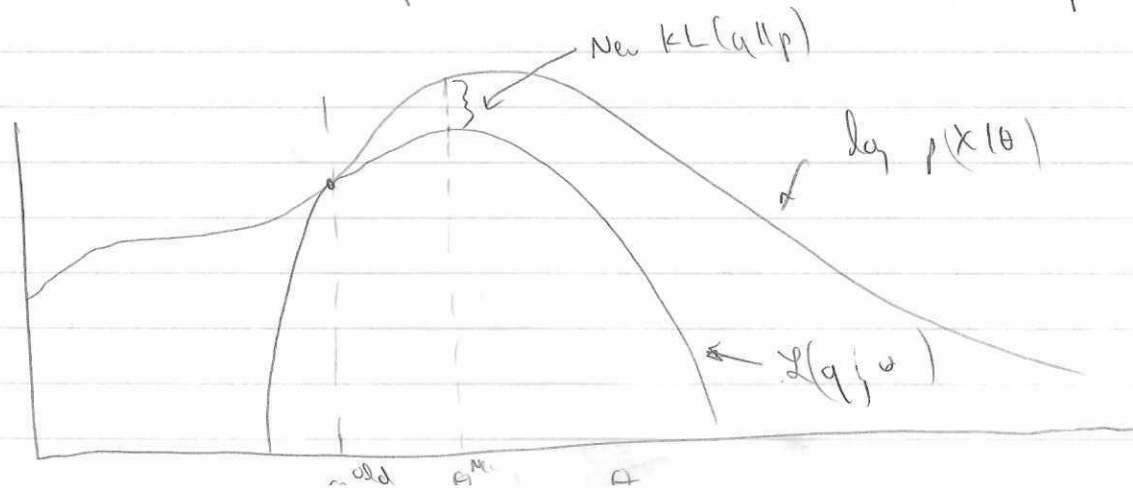
$$\bullet \max_{\theta} \mathcal{L}(q, \theta)$$

$$= \max_{\theta} \sum_z q(z) \log \left(\frac{p(x, z|\theta)}{q(z)} \right)$$

$$= \max_{\theta} \sum_z q(z) \log p(x, z|\theta) + \text{constant}$$

$$= \max_{\theta} E \left[\log p(x, z|\theta) \right]$$

↑ expectation taken over current posterior at θ



• HW: Implement EM algorithm from Li (2011)

Inputs: A , a matrix $A[i, k]$: genotype log-likelihood for individual i and dosage k at a SNP
 prior π , a k -vector containing initial values of genotype probabilities. Default $\text{rep}(\frac{1}{k+1}, k+1)$
 niter, an integer; # iterations of EM, default 100
 tol, stopping criterion. Default 0.001

Function

$$\pi_h^{(new)} = \frac{1}{n} \sum_{i=1}^n \frac{\pi_h^{old} \underbrace{a_{ih}}_{A[i, h]} \leftarrow e}{\sum_{k=0}^k \pi_h^{old} \underbrace{a_{ik}}_{A[i, k]} \leftarrow e}$$

Objective function: $\sum_{i=1}^n \log \left[\sum_{k=0}^k a_{ik} \pi_k \right]$

↑ stop when objective differs by less than tol in adjacent iterations.

Return

B , a matrix, $B[i, h]$ = log-posterior probability for individual i having dosage h at a SNP

$$b_{ih} = \frac{\pi_h a_{ih}}{\sum_{k=0}^k \pi_k a_{ik}}$$