

EM algorithm (Expectation-Maximization)

- arises in many apps. involving parameter estimation involving prob. models with incomplete obs. data
ex.: RNA-seq., (comp. bio lab), clustering & classification (ML), speaker ID, ...

Ex.: Recall our toy casino prob. with a "Fair" die & a "Loaded" die. If we know the "F" & "L" parameters (e.g. $P(6) = \frac{1}{2}$ for "L"), then finding MLSE using the Viterbi Alg. is "easy."

But what if you don't?

Need to est. or learn the parameters:
the EM alg. is popularly used.

Toy Ex.: Coin-flipping Experiment

- Given a pair of coins **A** and **B** having unknown biases θ_A and θ_B
- Goal: estimate (θ_A, θ_B) .

- Experiment (1) Choose one of the two coins at random with equal prob.
- (2) Perform 10 indep. coin flips with that coin.
- (3) Repeat (1) & (2) 5 times

- ML: when you get to observe which coin was used in Step (1): "Easy"

B	H T T T H H T H T H
A	H H H H T H H H H H
A	...
B	...
A	...

Coin	Coin
	5H, 5T
9H, 1T	
8H, 2T	
	4H, 4T
7H, 3T	
24H, 6T	9H, 11T

MLE: $\hat{\theta}_A = \frac{24}{30} = 0.8$ $\hat{\theta}_B = \frac{9}{20} = 0.45$

Aside:

$$\text{MLE}(\theta_A | H_A \text{ heads} \ \& \ T_A \text{ tails}) = ?$$

$$\theta_A^* = \underset{\theta}{\text{argmax}} \text{IP}(H_A \text{ heads}, T_A \text{ tails} | \theta)$$

$$= \underset{\theta}{\text{argmax}} \theta^{H_A} \cdot (1-\theta)^{T_A}$$

$$f(\theta) = \theta^{H_A} \cdot (1-\theta)^{T_A}$$

$$\log f(\theta) = H_A \log \theta + T_A \log(1-\theta)$$

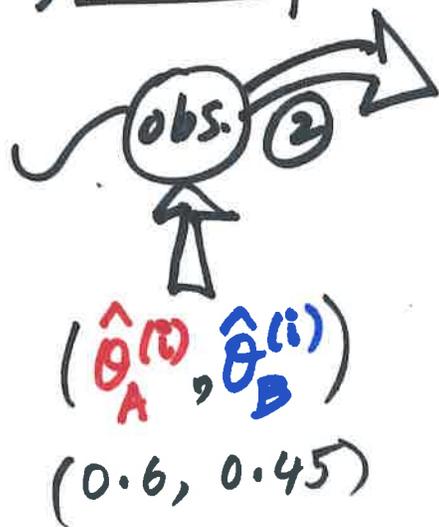
$$\frac{\partial \log f(\theta)}{\partial \theta} = \frac{H_A}{\theta} - \frac{T_A}{1-\theta} = 0$$

$$\Rightarrow \boxed{\theta_A^* = \frac{H_A}{H_A + T_A}}$$

b) "Hard-decision" EM (HEM)

1) Start with initial est. $(\hat{\theta}_A^{(0)}, \hat{\theta}_B^{(0)}) \stackrel{\text{e.g.}}{=} (0.6, 0.45)$

2) E-step:



"Hard" decision on which coin

0x(A)	1x(B)
1x(A)	0x(B)
1x(A)	0x(B)
0x(A)	1x(B)
1x(A)	0x(B)

Coin A	Coin B
	5H, 5T
9H, 1T	
8H, 2T	
	4H, 6T
7H, 3T	
24H, 6T	9H, 11T

3) $(\hat{\theta}_A^{(1)} = 0.8, \hat{\theta}_B^{(1)} = 0.45)$

3) M-step: Re-estimate $(\hat{\theta}_A^{(i+1)}, \hat{\theta}_B^{(i+1)})$ based on latest E-step.

4) Go back to E-step & iterate until convergence.

EM algorithm used when you do not get to observe which coin was used in step ① (i.e. coin label info is "hidden")

c) "Soft-Decision" EM (SEM)

"Soft-decision" on coin label

	Coin A	Coin B
(5H, 5T)	0.45 x (A)	0.55 x (B)
(9H, 1T)	0.8 x (A)	0.2 x (B)
(8H, 2T)	0.73 x (A)	0.27 x (B)
(4H, 6T)	0.35 x (A)	0.65 x (B)
(7H, 3T)	0.65 x (A)	0.35 x (B)
	21.3H, 8.6T	11.7H, 8.4T

$$P_{(0)}(A | 5H, 5T) = \frac{P(5H, 5T | A) P(A)}{P(5H, 5T | A) P(A) + P(5H, 5T | B) P(B)}$$

$$= \frac{(0.6)^5 (0.4)^5 \frac{1}{2}}{(0.6)^5 (0.4)^5 \frac{1}{2} + (0.45)^5 (0.55)^5 \frac{1}{2}}$$

$$= 0.45$$

(Bayes')

(3) M-step: Coin A: $\hat{\theta}_A^{(1)} = \frac{21.3}{21.3 + 8.6} \cong 0.71$
 Coin B: $\hat{\theta}_B^{(1)} = \frac{11.7}{11.7 + 8.4} \cong 0.58$

Continue ... $\hat{\theta}_A^{(10)} = 0.8$; $\hat{\theta}_B^{(10)} = 0.52$ CONVERGED!

Hard-EM

$$* P(A | \underline{\theta}^{(i)}, \underline{Y}) \underset{B}{\overset{A}{\geq}} P(B | \underline{\theta}^{(i)}, \underline{Y})$$

Soft-EM

$$P(A | \underline{\theta}^{(i)}, \underline{Y}) = \frac{P(\underline{Y} | A; \underline{\theta}^{(i)}) P(A)}{P(\underline{Y}; \underline{\theta}^{(i)})}$$

$$P(B | \underline{\theta}^{(i)}, \underline{Y}) = \frac{P(\underline{Y} | B; \underline{\theta}^{(i)}) P(B)}{P(\underline{Y}; \underline{\theta}^{(i)})}$$

(BAYES' THM.)

(* Ex. $\underline{Y} = (5H, 5T)$ (first expt.)

$$P(A | \hat{\theta}_A^{(0)} = 0.6, \hat{\theta}_B^{(0)} = 0.45, 5 \text{ Heads}, 5 \text{ Tails})$$

$$< P(B | \hat{\theta}_A^{(0)} = 0.6, \hat{\theta}_B^{(0)} = 0.45, 5H, 5T)$$

$$\begin{aligned} (+) \text{ Ex. } P_{\theta^{(0)}}(A | 5H, 5T) &= \frac{P(5H, 5T | A) P(A)}{P(5H, 5T | A) P(A) + P(5H, 5T | B) P(B)} \\ &= \frac{(0.6)^5 (0.4)^5 \frac{1}{2}}{(0.6)^5 (0.4)^5 \frac{1}{2} + (0.45)^5 (0.55)^5} = \boxed{0.45} \end{aligned}$$

Goal: $\max_{\theta} P(\underbrace{\text{observed data}}_{=X} | \theta)$



In our coin-flipping ex.

Complete data $\equiv (X, Z)$

X is labeled "COIN TOSSES" and Z is labeled "COIN LABELS"

e.g. $((15H, 5T), \dots)$ and e.g. $(BABBA)$

$\operatorname{argmax}_{\theta} P(X | \theta)$ is "hard"

but $\operatorname{argmax}_{\theta} P(X, Z | \theta)$ is "easy"

e.g. in our coin-flipping ex., $\operatorname{MLE}_{\theta} \left(\begin{array}{c} \text{coin} \\ \text{tosses} \\ \text{AND} \\ \text{coin labels} \end{array} \right)$ is "easy" ($\theta^* = \frac{\#H}{\#H + \#T}$)

$$\theta^* = \operatorname{argmax}_{\theta} \log \left[\sum_z f(x|z; \theta) p(z; \theta) \right]$$

← GOAL
(COMPUTATIONALLY INFEASIBLE)

H-EM: Replace \sum_z by a single-term z^*

where

$$z^* = \operatorname{argmax}_z p(z|x; \theta^{(m)})$$

current est. of θ

"E-step"

then improve guess

$$\theta^{(m+1)} = \operatorname{argmax}_{\theta} f(x|z^*; \theta) p(z^*; \theta)$$

"M-step"

where z^* is from the previous step

Iterate between (1) & (2) until convergence.

S-EM

(A) Replace

$$\log \{ \mathbb{E}_{z|\theta}^{(A)} [f(x|z;\theta)] \}$$

by

$$\mathbb{E}_{z|\theta}^{(B)} \{ \log [f(x|z;\theta)] \}$$

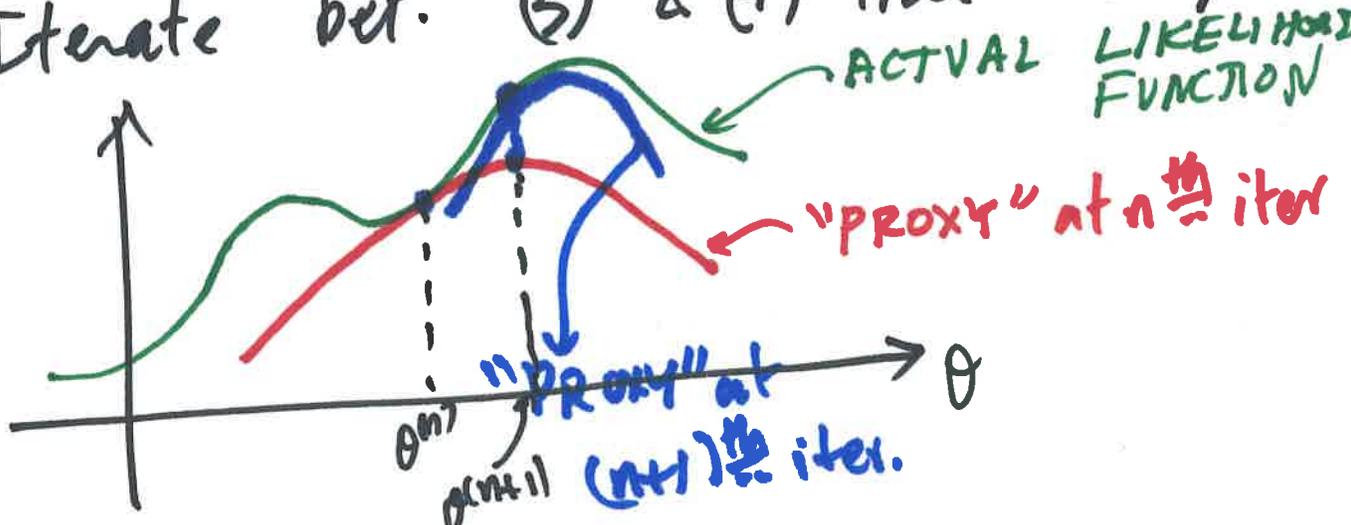
"actual" (A)
"proxy" (B)

(B) Replace $\mathbb{E}_{z|\theta} \{ \cdot \}$ by $\mathbb{E}_{z|x;\theta^{(m)}} \{ \cdot \}$

(3) E-step: $Q(\theta | \theta^{(m)}) = \mathbb{E}_{z|x;\theta^{(m)}} \left[\log f(x|z;\theta) \right]$

(4) M-step: $\theta^{(m+1)} = \operatorname{argmax}_{\theta} Q(\theta | \theta^{(m)})$

Iterate bet. (3) & (4) till convergence.



$$L(X; \theta) \triangleq p(X|\theta) = \sum_{\mathbf{z}} p(X, \mathbf{z}|\theta)$$

$$\text{E-step: } Q(\theta|\theta^{(t)}) = \mathbb{E}_{\mathbf{z}|X, \theta^{(t)}} \{ \log L(X, \mathbf{z}; \theta) \}$$

$$\text{M-step: } \theta^{(t+1)} = \arg \max_{\theta} Q(\theta|\theta^{(t)})$$

(E-M alg)

Pf. of correctness: EM iterates and improves on $Q(\theta|\theta^{(t)})$ rather than $\log p(X|\theta)$. Improvements to the former \Rightarrow improvements to the latter.

$$\log p(X|\theta) = \log p(X, \mathbf{z}|\theta) - \log p(\mathbf{z}|X, \theta)$$

Multiply by $p(\mathbf{z}|X, \theta^{(t)})$ and sum over \mathbf{z}

$$\sum_{\mathbf{z}} p(\mathbf{z}|X, \theta^{(t)}) \log p(X|\theta) = \underbrace{\sum_{\mathbf{z}} p(\mathbf{z}|X, \theta^{(t)}) \log p(X, \mathbf{z}|\theta)}_{Q(\theta|\theta^{(t)})} - \underbrace{\sum_{\mathbf{z}} p(\mathbf{z}|X, \theta^{(t)}) \log p(\mathbf{z}|X, \theta)}_{H(\theta|\theta^{(t)})}$$

$$\log p(X|\theta) = Q(\theta|\theta^{(t)}) + H(\theta|\theta^{(t)}) \quad (A)$$

(A) is true for all values of θ , incl. $\theta = \theta^{(t)}$

$$\log p(X|\theta^{(t)}) = Q(\theta^{(t)}|\theta^{(t)}) + H(\theta^{(t)}|\theta^{(t)}) \quad (B)$$

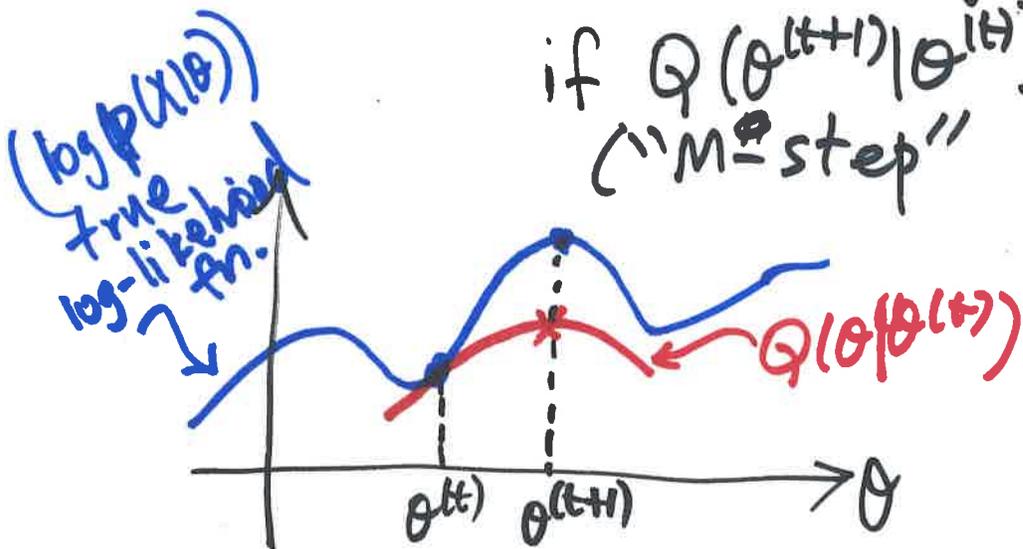
Subtract (B) from (A):

$$\log p(X|\theta) - \log p(X|\theta^{(t)}) = Q(\theta|\theta^{(t)}) - Q(\theta^{(t)}|\theta^{(t)}) + \underbrace{H(\theta|\theta^{(t)}) - H(\theta^{(t)}|\theta^{(t)})}_{\geq 0 \text{ (to be shown)}}$$

$$\Rightarrow \log p(X|\theta) - \log p(X|\theta^{(t)}) \geq \underbrace{Q(\theta|\theta^{(t)}) - Q(\theta^{(t)}|\theta^{(t)})}_{\geq 0 \text{ if } Q(\theta^{(t+1)}|\theta^{(t)}) \geq Q(\theta^{(t)}|\theta^{(t)})}$$

$$\Rightarrow \text{Thus, } \log p(X|\theta^{(t+1)}) \geq \log p(X|\theta^{(t)})$$

if $Q(\theta^{(t+1)}|\theta^{(t)}) \geq Q(\theta^{(t)}|\theta^{(t)})$
("M-step" in E-M alg.)



Proof of $H(\theta | \theta^{(t)}) - H(\theta^{(t)} | \theta^{(t)}) \geq 0$

defn. of $H(\cdot)$

$$\text{LHS} = \sum_{z} P(z | x, \theta^{(t)}) \log \left[\frac{P(z | x, \theta^{(t)})}{P(z | x, \theta)} \right]$$
$$= D(P(z | x, \theta^{(t)}) \| P(z | x, \theta)) \geq 0$$

Recall $D(p \| q) = - \sum_i p_i \log \frac{q_i}{p_i}$

True
Since $\log x \leq x - 1$



$$\geq \sum_i p_i \left[\frac{q_i}{p_i} - 1 \right]$$

$$= \sum_i q_i - \sum_i p_i = 0$$

Note: $H(\theta^{(t)} | \theta^{(t)}) = - \sum_z P(z | x, \theta^{(t)}) \log P(z | x, \theta^{(t)})$

= ENTROPY OF THE COND.
DIST. $\{z | x, \theta^{(t)}\}$