

EM (Expectation Maximization) algorithm (& its diverse applications in bioinformatics)

Optional !Saturday Track!
CS5691 PRML Jul-Nov 2025

“Algorithmic Thinking in Bioinformatics” Course
Manikandan Narayanan
Associate Professor, Computer Science and
Engineering, IIT Madras

Context: Seen so far...

- Session on Motivation and Background for “Algorithmic Thinking in Bioinformatics” course
- Sessions on several biological questions & their algorithmic solutions
- This session: Take a ***computation-first*** approach:
 - computational problem (“*learn the parameters of a mixture model aka latent variable model from data*”),
 - its algorithmic solution (*EM algorithm*), and
 - several biological applications (from motif-finding to soft k-means clustering, genes/isoforms’ expression quantitation, haplotype estimation, etc.).

Lecture style and sources

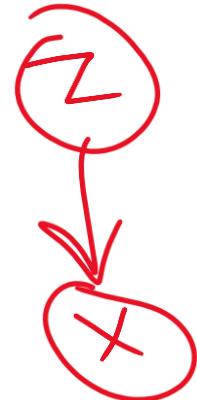
- Probabilistic (mixture) modeling and EM algo. is front-and-center:
 - Specifically, generative modeling + inference + parameter-learning using a mixture/latent-variable model.
 - Propose EM algorithm as a generic tool to learn the (MLE) parameters of a mixture model from data.
 - Apply this generic EM algo. to analyze various types of biological datasets with missing/incomplete data, and hidden structure/patterns.
- Sources for this lecture:
 - Main sources:
 - Do & Batzoglou. What is the expectation maximization algorithm?. *Nat Biotechnol* **26**, 897–899 (2008). - **cited as [DB08]**
 - Fan, Yuan & Liu. The EM Algorithm and the Rise of Computational Biology. *Statistical Science*. 25(4): 476-491 (2010). – **cited as [FYL10]**
 - Other sources:
 - David J. C. MacKay. Information Theory, Inference and Learning Algorithms (Chapters 20-22). 2003. – **cited as [DJM]** for content/figures taken as is from this book.
 - Stanford CS228 Probabilistic Graphical Models Course Lecture Notes. <https://ermongroup.github.io/cs228-notes/> - **cited as [ECL]**

Outline of EM Algo. & Appns.

- **Expectation-Maximization (EM) Algorithm (Algo.)**
 - **Background: Mixture Model (MM) or Latent Variable Model (LVM)**
 - Generic algorithm for parameter-learning of a MM
 - Algorithm analysis (correctness)
- Bioinformatics Applications (Appns.) of EM Algo.
 - Application 0: Bernoulli mixture of two coins (toy/warmup appn.)
 - Application 1: EM meets motif-finding in DNA sequences
 - Application 2: Soft k-means clustering of gene expression data (revisit appn.)
 - Application 3: Genes/isoforms' expression quantitation (deconvolution)
 - Application 4: Haplotype estimation from genotype data
 -
 -
 -

Recall: Probabilistic Mixture Model (MM) or Latent Variable Model (LVM)

- What? A complex yet elegant model built from simpler (tractable) components.
 - Any model with joint distbn. $P(X, Z) = P(Z)P(X|Z)$, where X is observed and Z is latent.
 - Marginal $P(X) = \sum_z P(Z)P(X|Z)$
- Why? i.e., why include variables that are always unobserved/hidden/latent in a model?
 - To make model more expressive, yet interpretable; and
 - To capture prior knowledge about hidden patterns/structure in our data (e.g., on clusters in a dataset)!



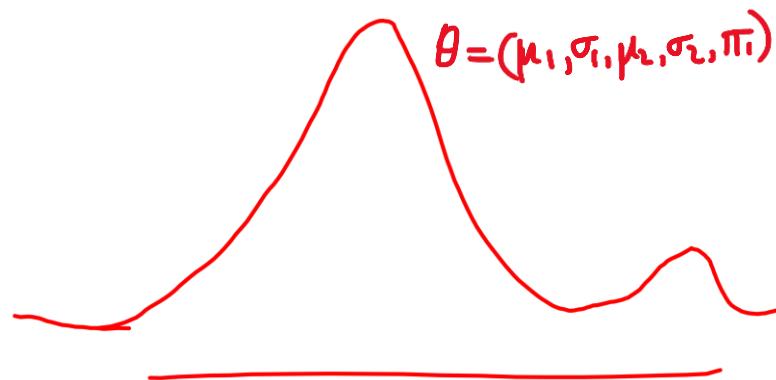
Recall: MM or LVM Tasks

- Generate *iid* data from model parameters (generative model): $\theta \rightarrow \mathbf{x} := \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$
 $\{\pi_k\}_{k=1}^K \rightarrow z^{(n)} (= k')$
 $\text{Model}_{k'}(\theta_{k'}) \rightarrow x^{(n)}$
- Inference:
$$r_k^{(n)} := P(Z = k \mid X = x^{(n)}; \theta) = \frac{\pi_k P(X = x^{(n)} \mid Z = k; \theta)}{\sum_{k'=1}^K \pi_{k'} P(X = x^{(n)} \mid Z = k'; \theta)}$$

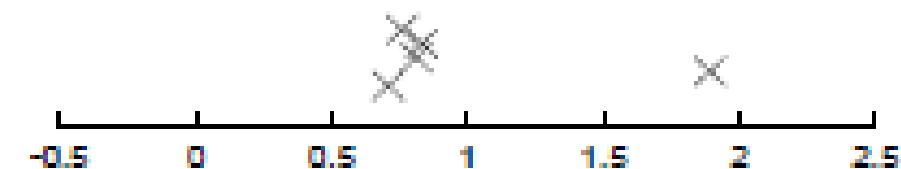
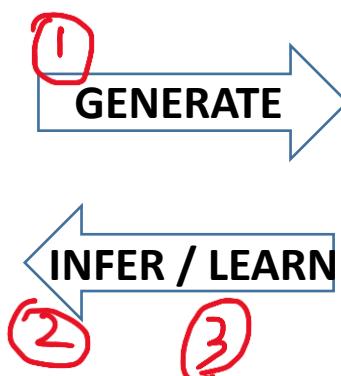
 $P(X = x^{(n)}; \theta) = \text{Dr. of above expn.}$
- Learning model parameters from data (MLE approach; EM algorithm)
$$\hat{\theta}_{ML} = \arg \max_{\theta} P(\mathbf{x}; \theta)$$

$$= \arg \max_{\theta} \sum_{\mathbf{z}} P(\mathbf{x}, \mathbf{z}; \theta)$$

Recall: MM Tasks – Example (Mixture of two 1D Gaussians - GMM)



$$\text{Cov: } \frac{\pi_1}{\pi_2} \frac{\mu_1 - \mu_2}{\sigma_1^2 + \sigma_2^2}$$
$$z = z_1$$

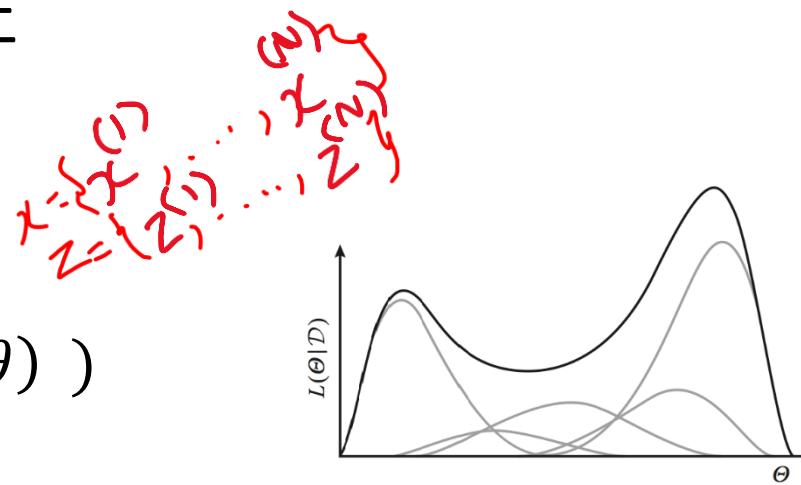


Outline of EM Algo. & Appns.

- **Expectation-Maximization (EM) Algorithm (Algo.)**
 - Background: Mixture Model (MM) or Latent Variable Model (LVM)
 - **Generic algorithm for parameter-learning of a MM**
 - Algorithm analysis (correctness)
- Bioinformatics Applications (Appns.) of EM Algo.
 - Application 0: Bernoulli mixture of two coins (toy/warmup appn.)
 - Application 1: EM meets motif-finding in DNA sequences
 - Application 2: Soft k-means clustering of gene expression data (revisit appn.)
 - Application 3: Genes/isoforms' expression quantitation (deconvolution)
 - Application 4: Haplotype estimation from genotype data
 -
 -
 -

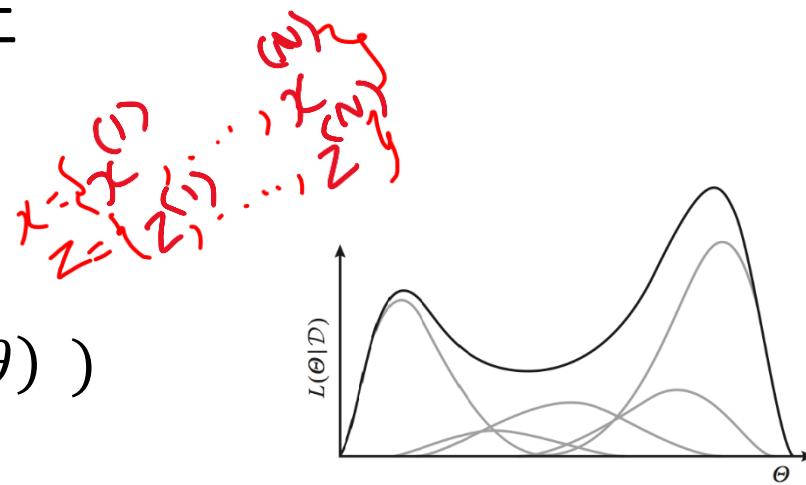
Goal: Parameter Learning of a Mixture model (aka MM or LVM such as GMM) via MLE

$$\begin{aligned}\hat{\theta}_{ML} &= \arg \max_{\theta} \log (P(x; \theta)) \\ &= \arg \max_{\theta} \log (\sum_z P(x, z; \theta))\end{aligned}$$



Goal: Parameter Learning of a Mixture Model (aka MM or LVM such as GMM) via MLE

$$\begin{aligned}\hat{\theta}_{ML} &= \arg \max_{\theta} \log (P(x; \theta)) \\ &= \arg \max_{\theta} \log (\sum_z P(x, z; \theta))\end{aligned}$$

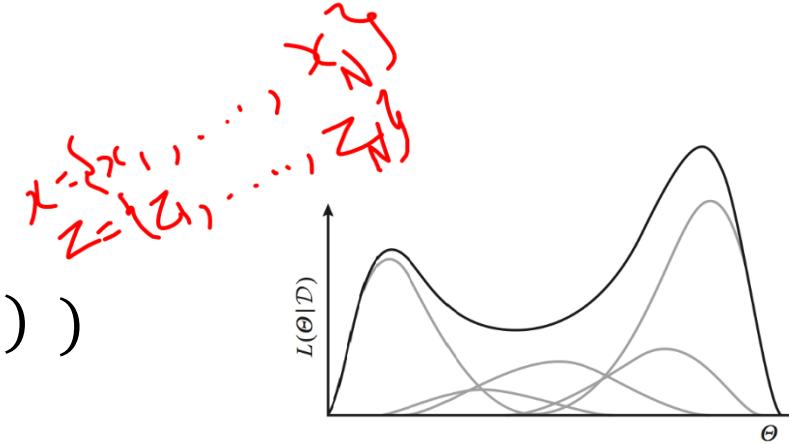


If iid:

$$\begin{aligned}P(x; \theta) &= \prod_{n=1}^N P(x^{(n)}; \theta) \\ &= \frac{1}{N} \sum_{z=1}^N P(x^{(n)}, z^{(n)}; \theta)\end{aligned}$$

Parameter-learning of a mixture model: Why is it hard?

$$\begin{aligned}\hat{\theta}_{ML} &= \arg \max_{\theta} \log (P(x; \theta)) \\ &= \arg \max_{\theta} \log (\sum_z P(x, z; \theta))\end{aligned}$$



- “log of sum”
 - calls for an approximate learning approach
 - EM algo. is one such approach that pushes the log across the sum!

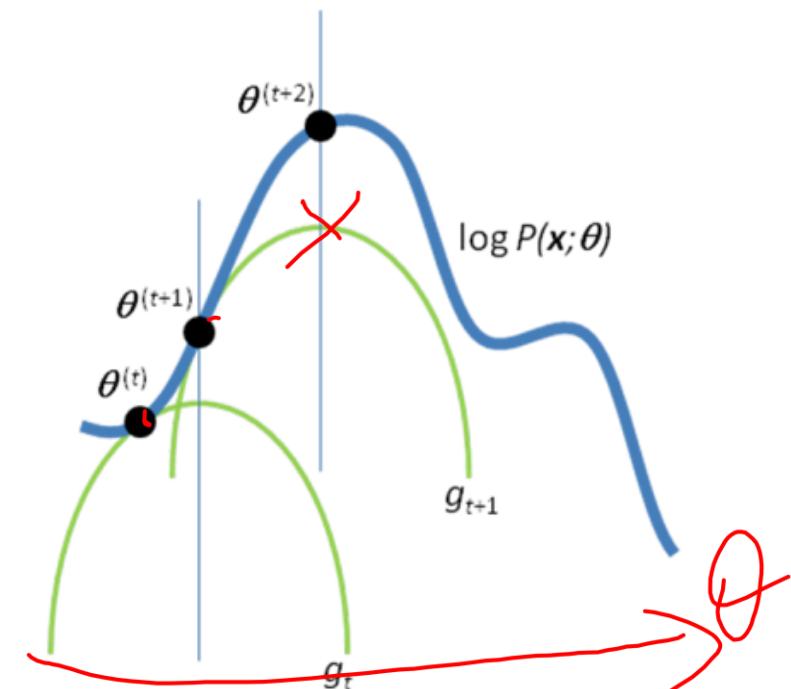
How to solve the MLE optimization problem of a mixture model in general?

Why is loglikl. **hard to optimize**?: $\hat{\theta}_{ML} = \arg \max_{\theta} \log (\sum_z P(x, z; \theta))$

EM algorithm finds a stationary point (mostly a local maxima) of the log likelihood function.

Pushes log across the sum

Finds a series of “nice” lower bounds (iterative E/M steps)



[from Do & Batzoglou. What is the expectation maximization algorithm?. *Nat Biotechnol* 2008]

How to solve the MLE optimization problem of a mixture model in general?

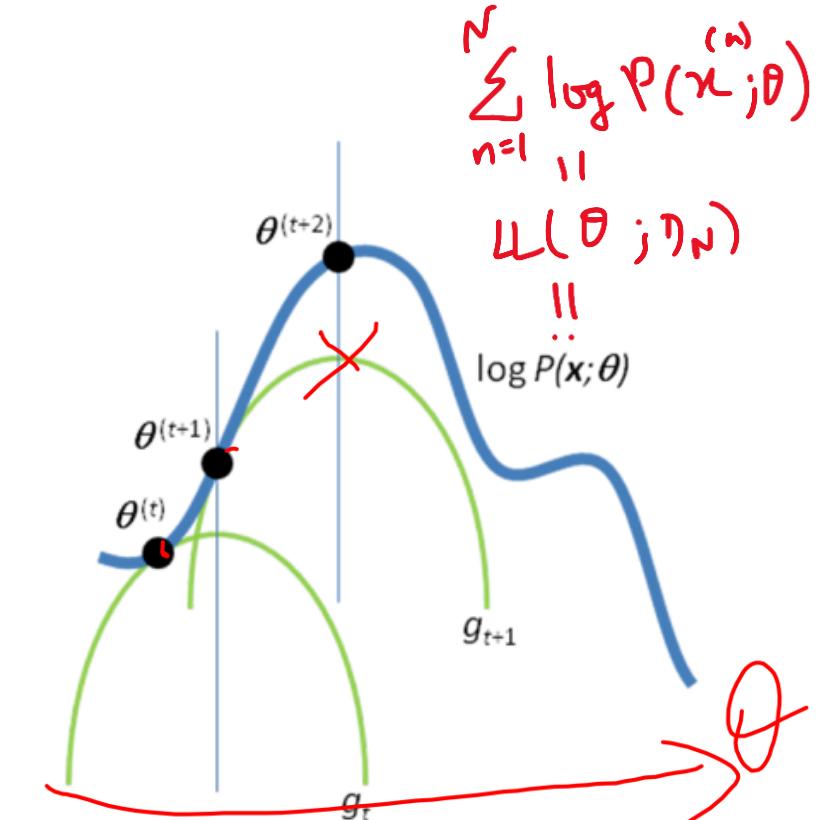
Why is loglikl. hard to optimize?: $\hat{\theta}_{ML} = \arg \max_{\theta} \log (\sum_z P(x, z; \theta))$

EM algorithm finds a stationary point (mostly a local maxima) of the log likelihood function.

Pushes log across the sum

Finds a series of “nice” lower bounds (iterative E/M steps)

E-step: $g_t \rightarrow g_{t+1} \dots$
M-step: $\theta_t \rightarrow \theta_{t+1}$



[from Do & Batzoglou. What is the expectation maximization algorithm?. *Nat Biotechnol* 2008]

EM algo. – generic version (*iid* data)

- Versatile and wide-spread approach for
 - (approximate) learning of parameters θ in not just Gaussian Mixture Models (GMM), but also any other MM or LVM; and so
 - enjoys many applications in bioinformatics & beyond
- EM algo. – generic version pseudocode:
 - Starting at an initial θ_0 , repeat until convergence for $t = 1, 2, \dots$:
 - *E-Step*: For each $x \in D$, compute the posterior $p(z | x; \theta_t)$.
 - *M-Step*: Compute new weights via

$$\theta_{t+1} = \arg \max_{\theta} \sum_{x \in D} \mathbb{E}_{z \sim p(z|x; \theta_t)} \log p(x, z; \theta).$$

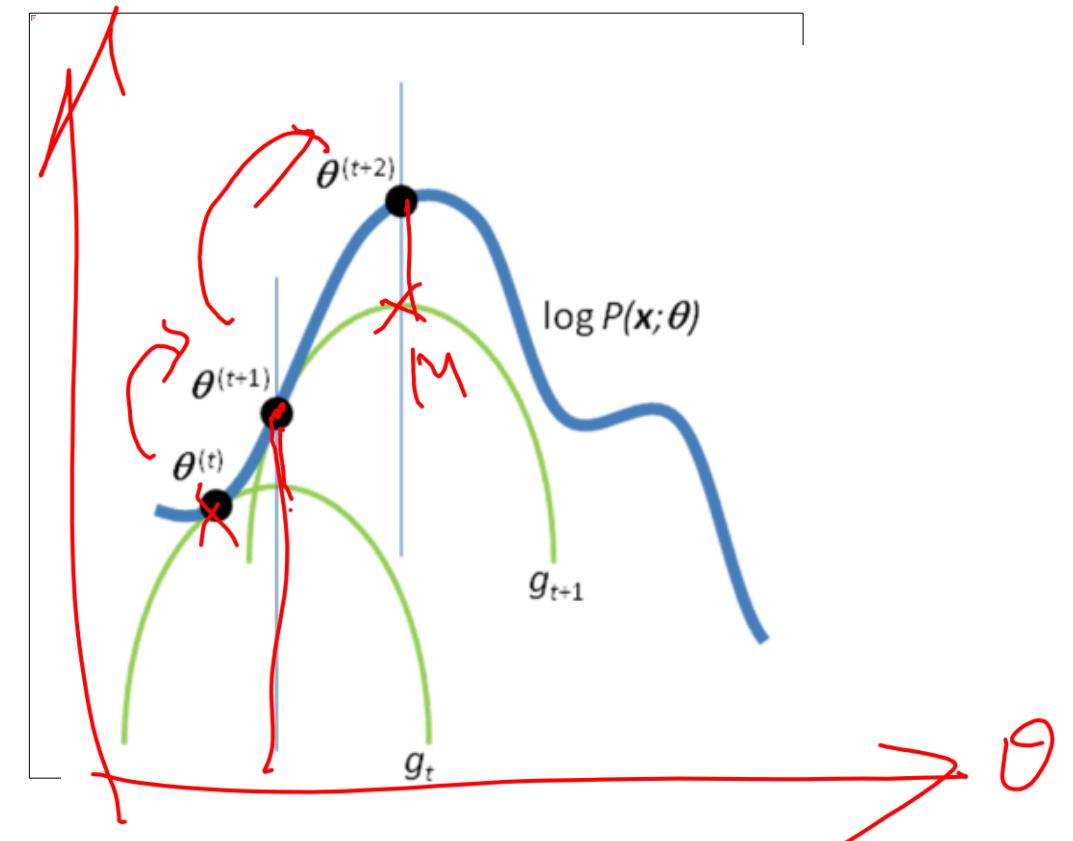
Outline of EM Algo. & Appns.

- **Expectation-Maximization (EM) Algorithm (Algo.)**
 - Background: Mixture Model (MM) or Latent Variable Model (LVM)
 - Generic algorithm description
 - **Algorithm analysis (correctness)**
- Bioinformatics Applications (Appns.) of EM Algo.
 - Application 0: Bernoulli mixture of two coins (toy/warmup appn.)
 - Application 1: EM meets motif-finding in DNA sequences
 - Application 2: Soft k-means clustering of gene expression data (revisit appn.)
 - Application 3: Genes/isoforms' expression quantitation (deconvolution)
 - Application 4: Haplotype estimation from genotype data
 - .
 - .
 - .

EM algorithm Correctness (Idea in pictures)

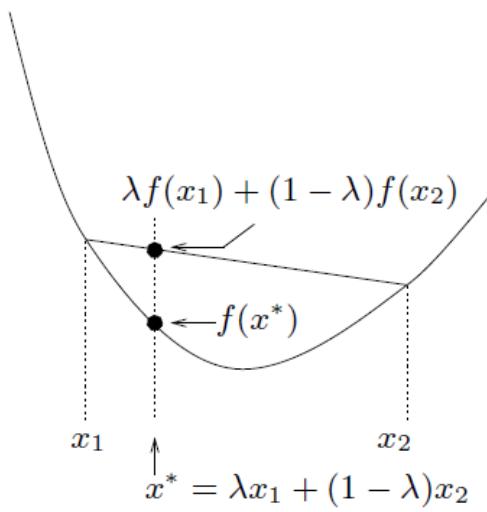
$$\hat{\theta}_{ML} = \arg \max_{\theta} \log \left(\sum_z P(x, z; \theta) \right)$$

Find a series of “nice” lower bounds
(based on Jensen’s Inequality)!

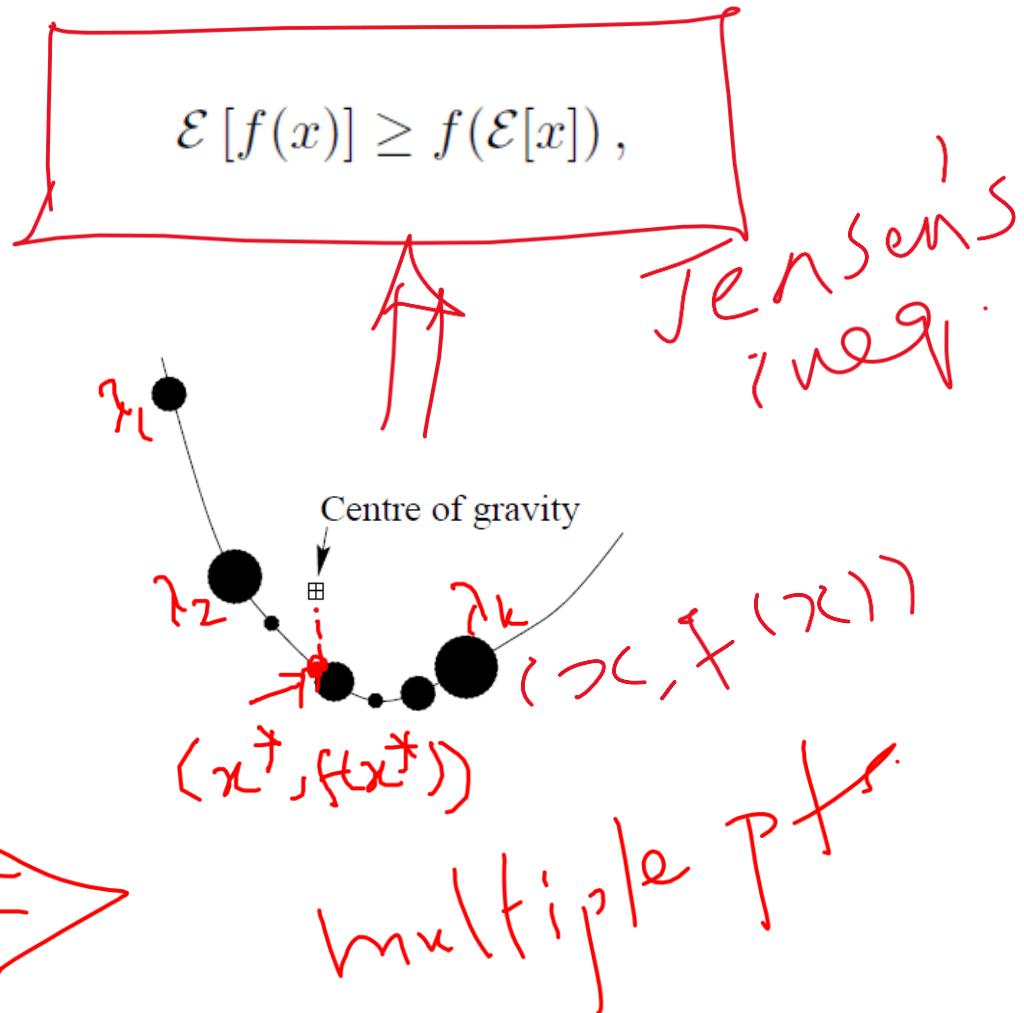


Jensen's inequality background

Jensen's inequality for convex functions f – in pictures

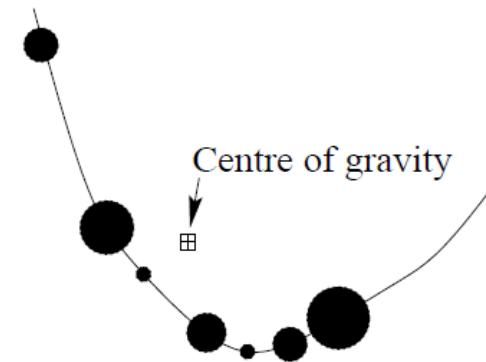


of convex fn.



When is Jensen's inequality an equality??

$$\mathcal{E}[f(x)] \geq f(\mathcal{E}[x]),$$



Jensen's inequality - formally

Convex ↘ functions. A function $f(x)$ is convex ↘ over (a, b) if every chord of the function lies above the function, as shown in figure 2.10; that is, for all $x_1, x_2 \in (a, b)$ and $0 \leq \lambda \leq 1$,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2). \quad (2.47)$$

A function f is strictly convex ↘ if, for all $x_1, x_2 \in (a, b)$, the equality holds only for $\lambda = 0$ and $\lambda = 1$.

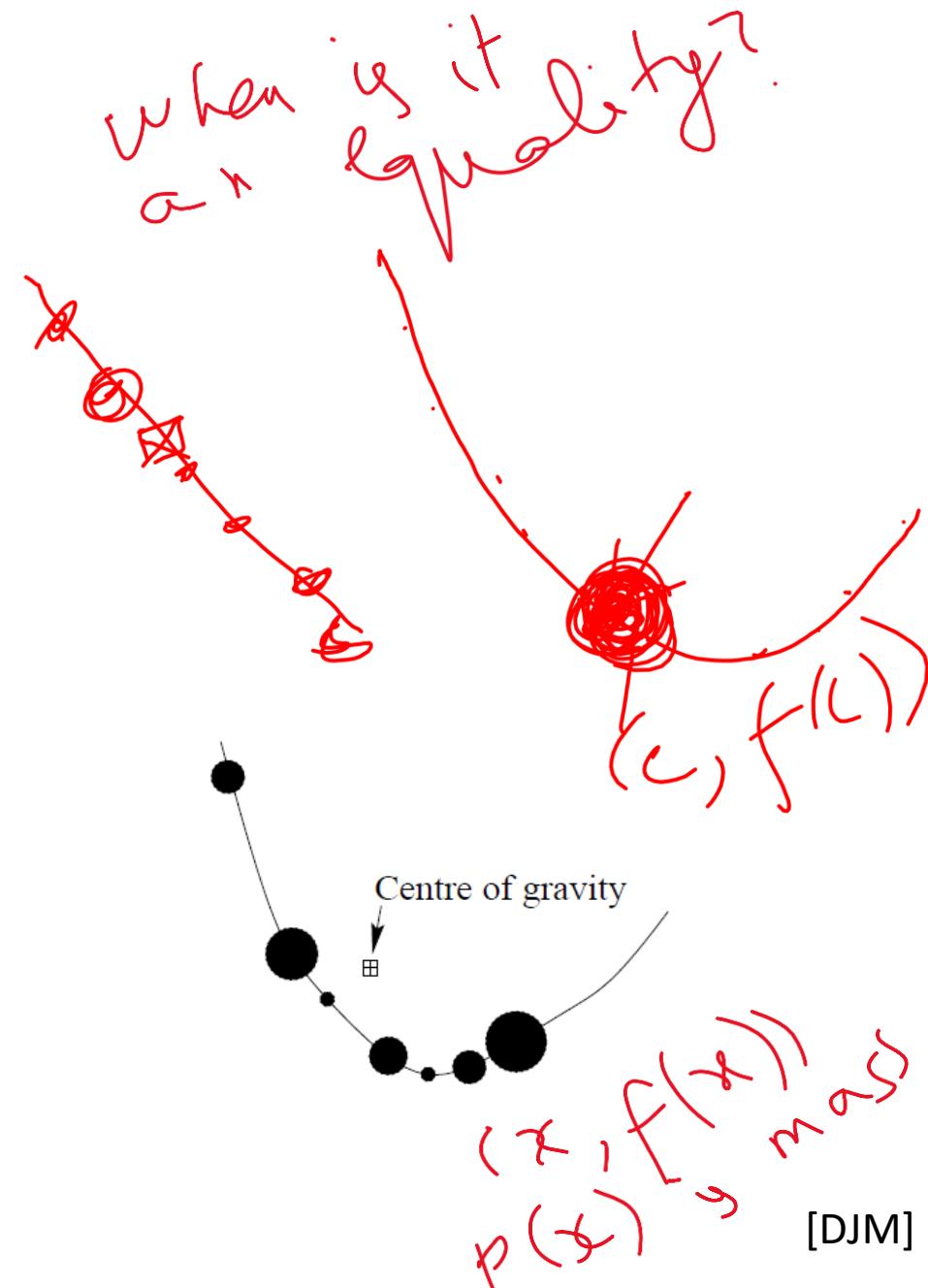
Similar definitions apply to concave ↗ and strictly concave ↗ functions.

Jensen's inequality. If f is a convex ↘ function and x is a random variable then:

$$\mathcal{E}[f(x)] \geq f(\mathcal{E}[x]), \quad (2.48)$$

where \mathcal{E} denotes expectation. If f is strictly convex ↘ and $\mathcal{E}[f(x)] = f(\mathcal{E}[x])$, then the random variable x is a constant.

Jensen's inequality can also be rewritten for a concave ↗ function, with the direction of the inequality reversed.



EM algorithm Correctness (Idea in eqns):

- Jensen's inequality (RHS is called Evidence Lower BOund or ELBO):

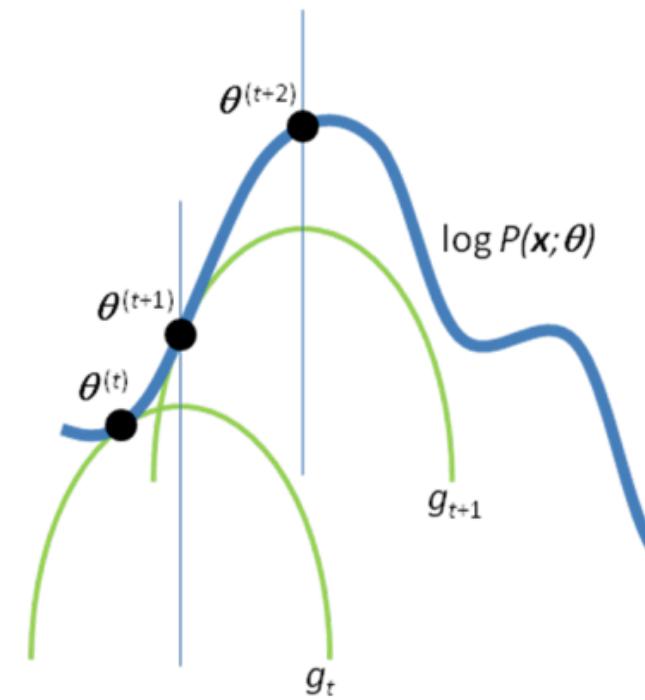
$$\log \left(\sum_z P(x, z; \theta) \right) = \log \left(\sum_z Q(z) \cdot \frac{P(x, z; \theta)}{Q(z)} \right) \geq \sum_z Q(z) \log \left(\frac{P(x, z; \theta)}{Q(z)} \right).$$

EM algo. Correctness – view as alternate max.

- E-step maximizes ELBO over $Q(z)$ (setting it to $P(z|x; \theta_t)$), and
- M-step maximizes ELBO over θ

$$g_t(\theta) = \sum_z P(z|x; \hat{\theta}^{(t)}) \log \left(\frac{P(x, z; \theta)}{P(z|x; \hat{\theta}^{(t)})} \right).$$

*X, Q, X
ELBO*



[DB08]

ELBO

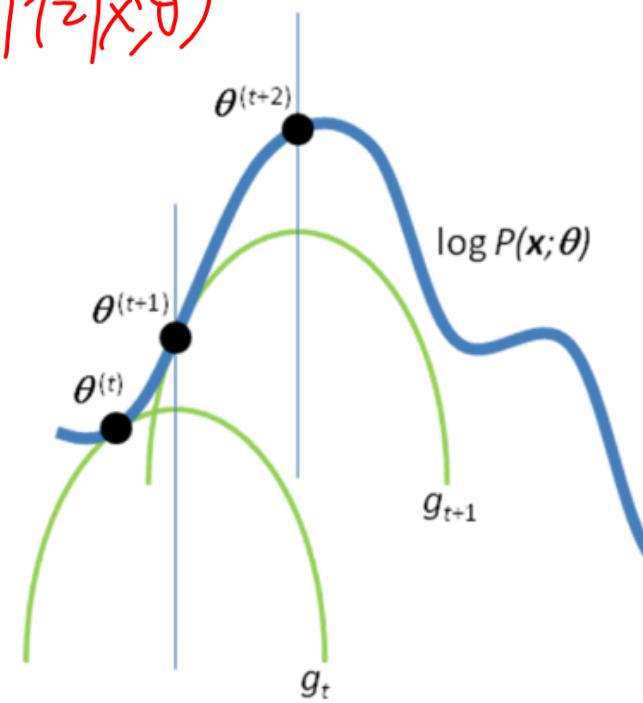
$$\log P(x; \theta) \geq \mathbb{E}_Q[z] \log \frac{P(x, z; \theta)}{Q(z)}$$

Blank space for illustration

- To prove:

EM algo. convergence to a stationary point

$$\begin{aligned}
 & \rightarrow \boxed{\log P(x; \theta)} \geq \boxed{\sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)}} \quad \forall x, Q, \theta \leftarrow g_t : \text{tight} \\
 & \text{To prove: } \log P(x; \theta^{t+1}) \geq \log P(x; \theta^t) \quad Q(z) = P(z|x; \theta^t) \\
 & \log P(x; \theta^{t+1}) \geq g_t(\theta^{t+1}) \quad \text{ELBO} \\
 & \geq g_t(\theta^t) \quad \max \cdot g_t \\
 & = \log P(x; \theta^t) \quad \text{tight, LB}
 \end{aligned}$$



Two final details in the algorithm's correctness

- M-step ($\operatorname{argmax}_{\theta} g_t(\theta)$) can be simplified
- EM algo. also works with non-iid data (e.g., Hidden Markov Models with structured data)

Recap: EM algo. – generic version (*iid* data)

- Versatile and wide-spread approach for
 - (approximate) learning of parameters θ in not just GMM, but also other LVMs; and so
 - enjoys many appns. in clustering, bioinformatics & beyond
- EM algo. – generic version pseudocode:
 - Starting at an initial θ_0 , repeat until convergence for $t = 1, 2, \dots$.
 - *E-Step*: For each $x \in D$, compute the posterior $p(z | x; \theta_t)$.
 - *M-Step*: Compute new weights via

$$\theta_{t+1} = \arg \max_{\theta} \sum_{x \in D} \mathbb{E}_{z \sim p(z|x; \theta_t)} \log p(x, z; \theta).$$

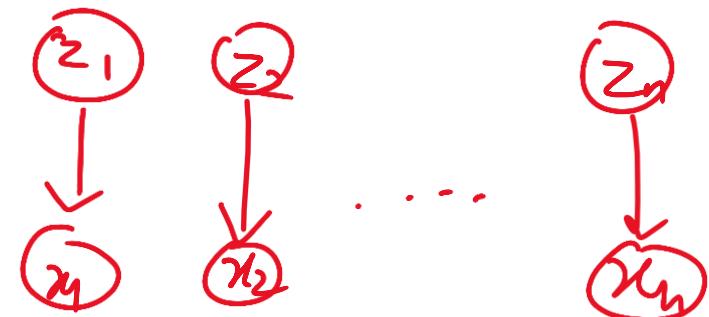
$\tilde{g}_x(\theta)$

$\{\tilde{g}_i : i \in \text{if well}\}$

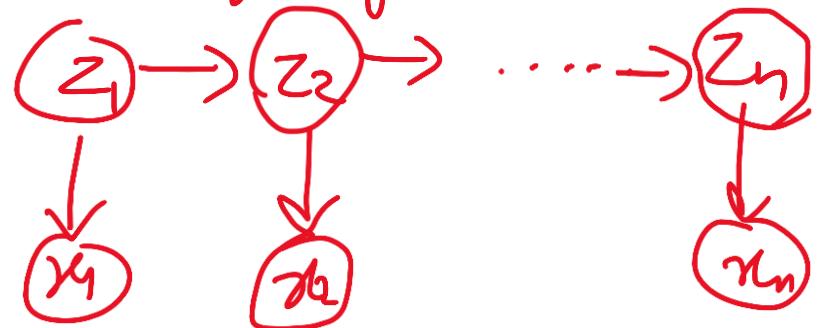
Brief aside: iid vs. structured data

If iid:

$$\begin{aligned} P(\mathbf{x}; \theta) &= \prod_{n=1}^N P(x_n; \theta) \\ &= \prod_n \sum_{z_n} P(x_n, z_n; \theta) \end{aligned}$$



MLE (EM algo.) also works if data is not iid, e.g. if data is structured:



EM algo. – really generic version (non-*iid* or structured data)

- Versatile and wide-spread approach for
 - (approximate) learning of parameters θ in not just GMM, but also other LVMs; and so
 - enjoys many appns. in clustering, bioinformatics & beyond

- EM algo. – generic version pseudocode:

- Starting at an initial θ_0 , repeat until convergence for $t = 1, 2, \dots$:

- *E-Step*: For each $x \in D$, compute the posterior $p(z | x; \theta_t)$.
 - *M-Step*: Compute new weights via

$$\theta_{t+1} = \arg \max_{\theta} \sum_{x \in D} \mathbb{E}_{z \sim p(z|x; \theta_t)} \log p(x, z; \theta).$$

repeat

$$\vec{x} : p(\vec{z} / \vec{x}; \theta_t)$$

$$\theta_{t+1} = \arg \max_{\theta} \sum_{\vec{z} \sim P(\vec{z} | \vec{x}; \theta_t)} \mathbb{E}_{\vec{x}} [\log p(\vec{x}, \vec{z}; \theta)]$$

Fine print: For certain models like HMM where data points $x = \{x\}_{x \in D}$ are **not iid**, and $z = \{z\}$ are also **not indep.**, we deal with the whole dataset like so: $\theta_{t+1} = \arg \max_{\theta} E_{z \sim p(z|x; \theta_t)} \log p(x, z; \theta)$.

[ECL]

Outline of EM Algo. & Appns.

- Expectation-Maximization (EM) Algorithm (Algo.)
 - Background: Mixture Model (MM) or Latent Variable Model (LVM)
 - Generic algorithm description
 - Algorithm analysis (correctness)
- **Bioinformatics Applications (Appns.) of EM Algo.**
 - **Application 0: Bernoulli mixture of two coins (toy/warmup appn.)**
 - Application 1: EM meets motif-finding in DNA sequences
 - Application 2: Soft k-means clustering of gene expression data (revisit appn.)
 - Application 3: Genes/isoforms' expression quantitation (deconvolution)
 - Application 4: Haplotype estimation from genotype data
 - .
 - .
 - .

Bernoulli mixture (of two coins): generative model

a Maximum likelihood

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

5 sets, 10 tosses per set

Bernoulli mixture (of two coins): inference

a Maximum likelihood

x_1 H T T T H H T H T H
 x_2 H H H H T H H H H H
⋮
 x_N T H H H T H H H T H

N m

5 sets, 10 tosses per set

	Col A	Col B
		5 H, 5 T
	9 H, 1 T	
	8 H, 2 T	
		4 H, 6 T
	7 H, 3 T	

$x_h = \# \text{ of heads in } n^{\text{th}} \text{ set}$
 $m = \# \text{ of tosses per set}$

Bernoulli mixture (of two coins): parameter-learning

a Maximum likelihood

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

5 sets, 10 tosses per set

MLE from (i) complete data (warmup), and
 (ii) observed data (using EM algorithm)

Bernoulli mixture (of two coins): MLE from complete data (\mathbf{x} , and labels \mathbf{z})

a Maximum likelihood

x_1 H T T T H H T H T H
 x_2 H H H H T H H H H H
 \vdots H T H H H H H T H H
 \vdots H T H T T T H H T T
 x_N T H H H T H H H T H

5 sets, 10 tosses per set

$N \quad m$

	Coin A	Coin B
x_1		5 H, 5 T
x_2	9 H, 1 T	
\vdots	8 H, 2 T	
\vdots		4 H, 6 T
x_N	7 H, 3 T	
	24 H, 6 T	9 H, 11 T

$$\hat{\theta}_A = \frac{24}{24+6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9+11} = 0.45$$

$$L(\theta_1, \theta_2 | \vec{x}, \vec{z}) = \prod_{n=1}^N \theta_{z_n}^{x_n} (1-\theta_{z_n})^{m-x_n} \cdot \left(\frac{1}{2}\right)$$

$$\hat{\theta}_k = \frac{\sum_{n=1}^N r_k(n) x_n}{r_k \cdot m}$$

Bernoulli mixture (of two coins): MLE from observed data (obs. \mathbf{x} , but labels \mathbf{z} hidden) – calls for EM algo.

a Maximum likelihood

x_1 H T T T H H T H T H
 x_2 H H H H T H H H H H
 \vdots H T H H H H H T H H
 \vdots H T H T T T H H T T
 x_N T H H H T H H H T H

5 sets, 10 tosses per set

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

$$\mathcal{L}(\theta_1, \theta_2 | \vec{x}, \vec{z}) = \prod_{n=1}^N \sum_{z_n} \theta_{z_n}^{x_n} (1-\theta_{z_n})^{m-x_n} \left(\frac{1}{2}\right)$$

Recap: EM algo. – generic version (*iid* data)

- Versatile and wide-spread approach for
 - (approximate) learning of parameters θ in not just GMM, but also other LVMs; and so
 - enjoys many appns. in clustering, bioinformatics & beyond
- EM algo. – generic version pseudocode:
 - Starting at an initial θ_0 , repeat until convergence for $t = 1, 2, \dots$:
 - *E-Step*: For each $x \in D$, compute the posterior $p(z | x; \theta_t)$.
 - *M-Step*: Compute new weights via

$$\theta_{t+1} = \arg \max_{\theta} \sum_{x \in D} \mathbb{E}_{z \sim p(z|x; \theta_t)} \log p(x, z; \theta).$$

Bernoulli mixture (of two coins): M-step optimization: MLE from hallucinated/expected complete data

a Maximum likelihood

x_1  H T T T H H T H T H

x_2  H H H H T H H H H H

:  H T H H H H H T H H

:  H T H T T T H H T T

x_N  T H H H T H H H T H

	Coin A	Coin B
x_1		5 H, 5 T
x_2	9 H, 1 T	
:	8 H, 2 T	
:		
x_N	4 H, 6 T	
	7 H, 3 T	
	24 H, 6 T	9 H, 11 T

5 sets, 10 tosses per set
 $N \quad m$

Hallucinated/expected complete data := obs. data x + hidden labels z “drawn/hallucinated”

from $P(z | x; \theta^t)$

$$\hat{\theta}_k = \frac{\sum_{n=1}^N r_k(n) x_n}{R_k \cdot m}$$

Bernoulli mixture (of two coins): EM Algo.

a Maximum likelihood

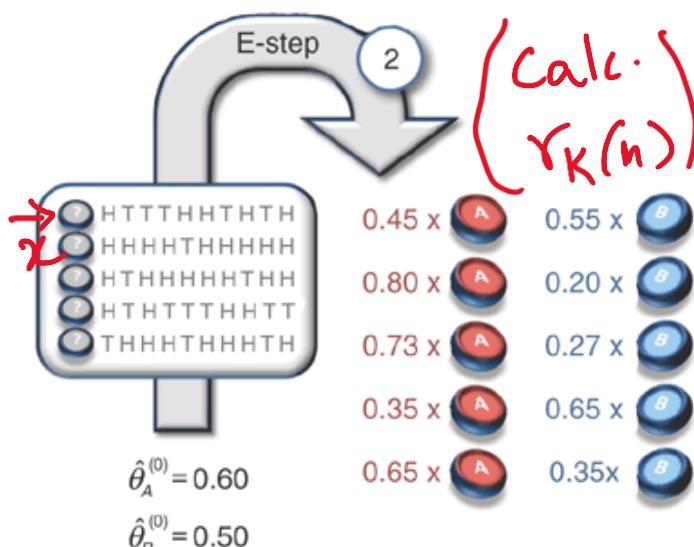
x_1 H T T T H H T H T H
 x_2 H H H H T H H H H H
 \vdots H T H H H H H T H H
 \vdots H T H T T T H H T T
 x_N T H H H T H H H T H
 5 sets, 10 tosses per set $N M$

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

$$\hat{\theta}_A = \frac{24}{24+6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9+11} = 0.45$$

b Expectation maximization



Coin A	Coin B
≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
≈ 21.3 H, 8.6 T ≈ 11.7 H, 8.4 T	

$$\hat{\theta}_k = \frac{\sum_{n=1}^N r_{k(n)} x_n}{r_k \cdot M} \quad (\text{Update } \hat{\theta}_k)$$

Outline of EM Algo. & Appns.

- Expectation-Maximization (EM) Algorithm (Algo.)
 - Background: Mixture Model (MM) or Latent Variable Model (LVM)
 - Generic algorithm description
 - Algorithm analysis (correctness)
- **Bioinformatics Applications (Appns.) of EM Algo.**
 - Application 0: Bernoulli mixture of two coins (toy/warmup appn.)
 - Application 1: EM meets motif-finding in DNA sequences
 - **Application 2: Soft k-means clustering of gene expression data (revisit appn.)**
 - Application 3: Genes/isoforms' expression quantitation (deconvolution)
 - Application 4: Haplotype estimation from genotype data
 -
 -
 -

Cancer example: Can we group cancer samples into different (sub)types?



Golub et al. Science 1999

Recall: MLE for one 1D Gaussian (closed-form soln.) –
How to change it to learn mixture of Gaussians?

- Log likelihood:

$$\ln P\left(\left\{x^{(n)}\right\}_{n=1}^N \mid \mu, \sigma\right) = -N \ln(\sqrt{2\pi}\sigma) - \sum_n (x^{(n)} - \mu)^2 / (2\sigma^2)$$

- MLE estimates:

$$\hat{\mu} = \sum_{n=1}^N \frac{x^{(n)}}{N}, \quad \hat{\sigma}_N^2 = \frac{\sum_{n=1}^N (x^{(n)} - \hat{\mu})^2}{N}$$

Recall: MLE for one 1D Gaussian (closed-form soln.) –
How to change it to learn mixture of Gaussians?

- Log likelihood:

$$\ln P\left(\left\{x^{(n)}\right\}_{n=1}^N \mid \mu, \sigma\right) = -N \ln(\sqrt{2\pi}\sigma) - \sum_n (x^{(n)} - \mu)^2 / (2\sigma^2)$$

θ^t = current params. $r_k^{(n)} = p(z^{(n)}=k \mid x=x^{(n)}; \theta^t)$

- ↓
• MLE estimates:

$$\hat{\mu}_N^{(k)} = \frac{\sum_{n=1}^N r_k^{(n)} x^{(n)}}{\sum_{n=1}^N r_k^{(n)}}, \quad \hat{\sigma}_N^{(k)} = \frac{\sum_{n=1}^N r_k^{(n)} (x^{(n)} - \hat{\mu})^2}{\sum_{n=1}^N r_k^{(n)}}$$

(wted. sample mean & variance)

EM-algorithm in Action: Parameter-learning of a Mixture of 1D Gaussians (aka univariate GMM)

Calculate:

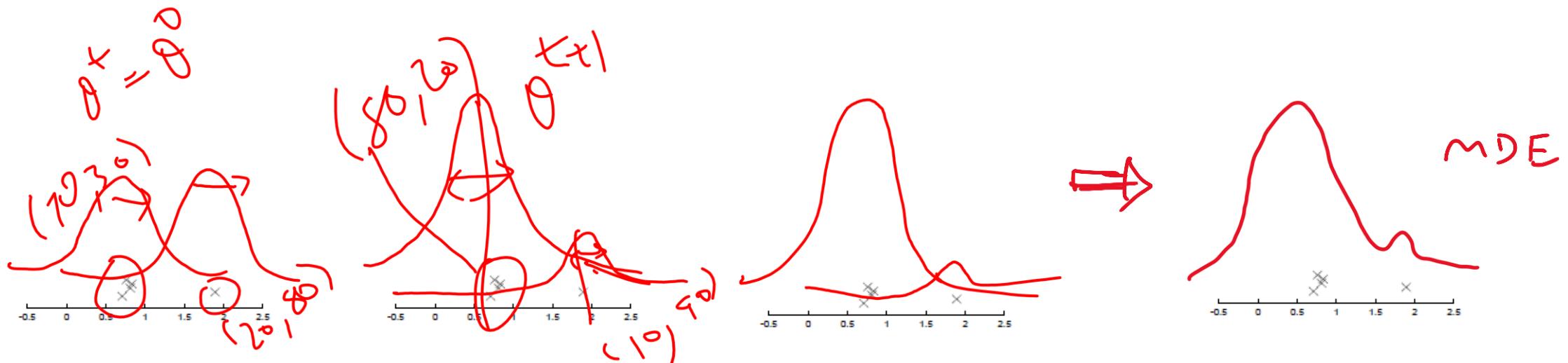
$$r_k(n) = p(z=k \mid \theta_t) = \frac{\pi_k N(x \mid \mu_k, \sigma_k^2)}{\sum_{k=1}^K \pi_k N(x \mid \mu_k, \sigma_k^2)}$$

E-step (Assign)

$$\hat{\mu}_k = \frac{\sum_{n=1}^N r_k(n) x_n}{R_k}; \hat{\sigma}_k^2 = \frac{\sum_{n=1}^N r_k(n) (x_n - \hat{\mu}_k)^2}{R_k}$$

$$\pi_k = \frac{R_k}{\sum_{k=1}^K R_k} = \frac{R_k}{N} \quad (\theta_{t+1})$$

M-step (Update)



EM-algorithm in Action: Parameter-learning of a Mixture of 1D Gaussians (aka univariate GMM)

Calculate:

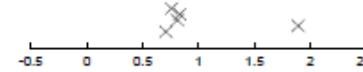
$$r_k(n) = p(z=k \mid \theta_t) = \frac{N(x \mid \mu_k, \sigma_k^2)}{\sum_{k=1}^K N(x \mid \mu_k, \sigma_k^2)}$$

E-step (Assign)

$$\hat{\mu}_k = \frac{\sum_{n=1}^N r_k(n) x_n}{R_k};$$

(θ_{t+1})

M-step (Update)



EM algo. – GMM version (aka) Vanilla Soft K-means

Initialization: Set K means $\{m^{(k)}\}$ to random values.

Assignment:
(E-step)

$$r_k^{(n)s} = \frac{\exp(-\beta d(\mathbf{m}^{(k)}, \mathbf{x}^{(n)}))}{\sum_{k'} \exp(-\beta d(\mathbf{m}^{(k')}, \mathbf{x}^{(n)}))}.$$

(here $\beta = 1/(2\sigma^2)$)

Prove: Let $r_k^{(n)s} = r_k^{(n)h}$
 $\beta \rightarrow \infty$ [hard k-mean]

Update:
(M-step)

$$\mathbf{m}^{(k)} = \frac{\sum_n r_k^{(n)} \mathbf{x}^{(n)}}{R^{(k)}}$$

where $R^{(k)}$ is the total responsibility of mean k ,

$$R^{(k)} = \sum_n r_k^{(n)}.$$

Thank you!

Backup slides

EM algorithm Correctness (Idea in eqns)

- Jensen's inequality: RHS is called ELBO (Evidence Lower Bound).

$$\log P(x; \theta) = \log \left(\sum_z P(x, z; \theta) \right) = \log \left(\sum_z Q(z) \cdot \frac{P(x, z; \theta)}{Q(z)} \right) \geq \sum_z Q(z) \log \left(\frac{P(x, z; \theta)}{Q(z)} \right).$$

$\geq \bar{z} = g(z)$

$$\begin{aligned} & \log (E_{z \sim Q} [z]) \\ &= \log (E_{z'} [z']) \\ &\geq E_{z'} (\log [z']) \\ &= E_z (\log z') \end{aligned}$$

Note: By law of the unconsc. statistician, $E_{g(z)} [g(z)] = E_z [g(z)]$.

When is Jensen's ineq. an equality for us?

ELBO

$$\log P(x; \theta) \geq \sum_z Q(z) \underbrace{\log \frac{P(x, z; \theta)}{Q(z)}}_{\text{Jensen's Ineq.}} \quad \forall x, Q, \theta$$

$$\frac{P(x, z; \theta)}{Q(z)} = \text{const.} <$$

$$\Rightarrow Q(z) \propto P(x, z; \theta)$$

$$\sum_z Q(z) = 1$$

$$Q(z) = \frac{P(x, z; \theta)}{\sum_z P(x, z; \theta)}$$

$$Q(z) = P(z|x; \theta)$$

