

# Expectation- Maximization & Baum-Welch

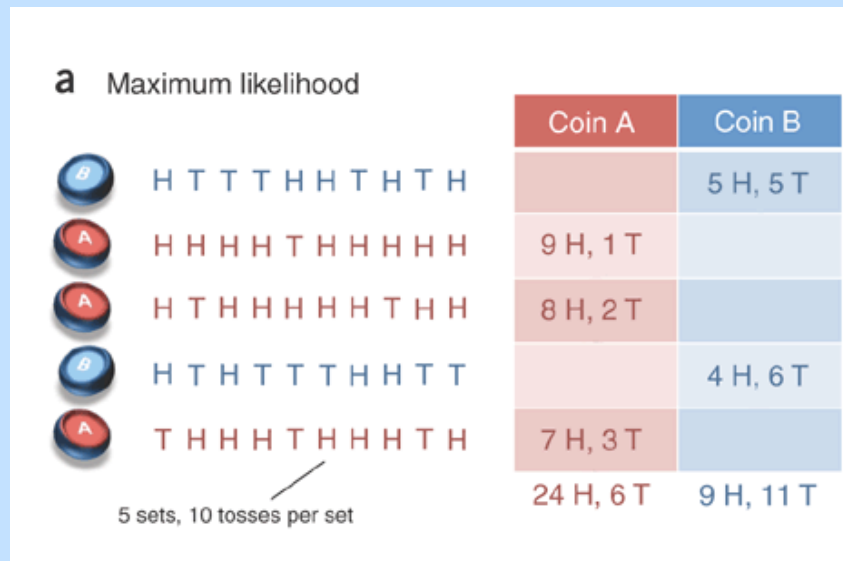
Slides: Roded Sharan, Jan 15; revised by Ron Shamir, Nov 15

# The goal

- Input: **incomplete** data originating from a probability distribution with some unknown parameters
- Want to find the parameter values that maximize the likelihood
- EM - approach that helps when maximum likelihood solution cannot be directly computed.
- Seeks a local maximum by iteratively solving two easier subproblems

# Coin flipping: complete data

- Coins A, B with unknown heads probs.  $\theta_A, \theta_B$
- Goal: estimate  $\theta_A, \theta_B$
- Experiment: Repeat x5: choose A or B with prob. 1/2, flip 10 times, record results.

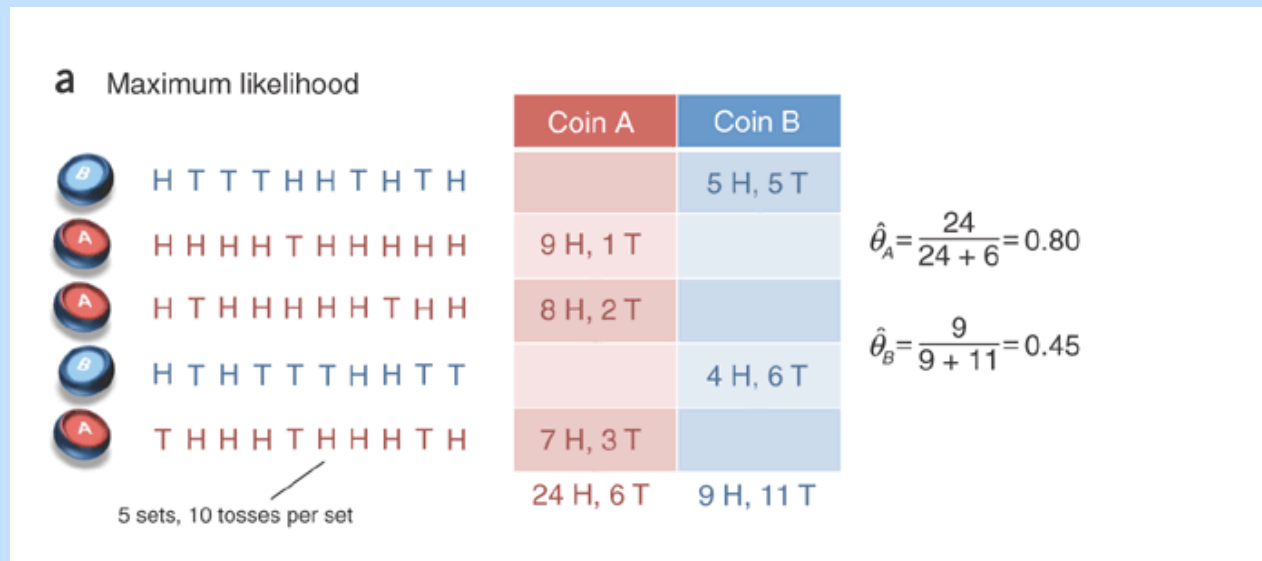


$x = (x_1, \dots, x_5)$  : no of H in set 1, ... 5

$Y = (y_1, \dots, y_5)$  : coin used in set 1, ... 5

# Coin flipping: complete data

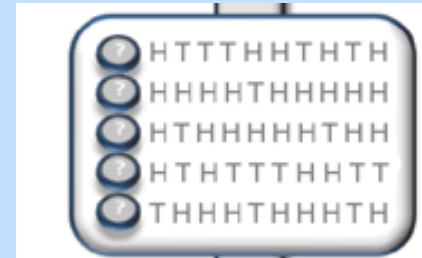
- Natural guess:  $\theta_i$  = fraction of H in flips of coin  $i$
- This is actually the ML solution: maximizes  $P(x,y|\theta)$  (ex.)



- What if we do not know which coin was used in each round?

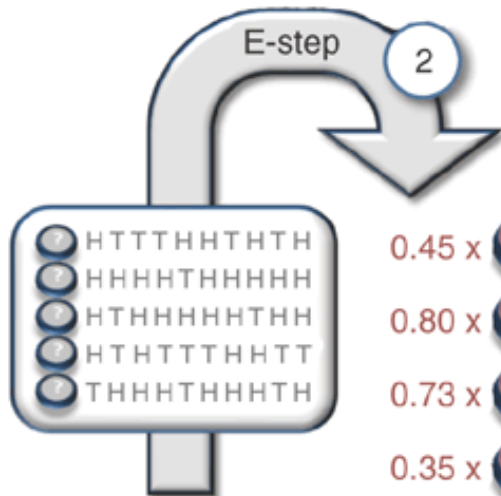
# Coin flipping: incomplete data

- Now  $(y_1, \dots, y_5)$  are **hidden / latent variables**.
- Cannot compute H prob for each coin
- If we guessed  $Y$  correctly - we could.
- Idea: guess initial  $\theta^0_A, \theta^0_B$ 
  - Use  $\theta^t_A, \theta^t_B$  to compute the most likely coin for each set, get new  $y$
  - Use the resulting  $y$  to recompute  $\theta_A, \theta_B$  using ML, get  $\theta^{t+1}_A, \theta^{t+1}_B$
  - Repeat till convergence
- EM: use probabilities rather than the single most likely completion  $y$



# Coin flipping: incomplete data

## b Expectation maximization

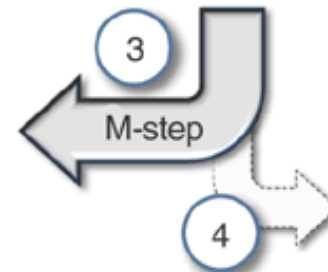


	Coin A	Coin B
	≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
	≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
	≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
	≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
	≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
	≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T



$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$



$$\hat{\theta}_A^{(10)} \approx 0.80$$

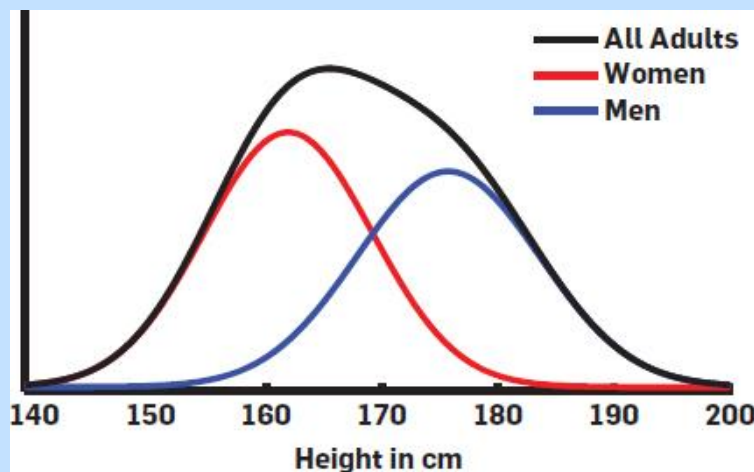
$$\hat{\theta}_B^{(10)} \approx 0.52$$

# The probabilistic setting

Input: data  $X$  coming from a probabilistic model with hidden information  $y$

Goal: Learn the model's parameters  $\theta$  so that the likelihood is maximized.

# Mixture of two Gaussians



Kalai et al. Disentangling Gaussians CACM 2012

Our input generates the black distribution. We want to color each sample red/blue and find the parameters of the two distributions to maximize the data probability. (assume  $\sigma$  known)

$$P(y_i = 1) = p_1 ; P(y_i = 2) = p_2 = 1 - p_1$$

$$P(x_i | y_i = j) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(x_i - \mu_j)^2}{2\sigma^2}\right)$$


$$\theta = (p_1, \mu_1, \mu_2)$$

# The likelihood function

$$P(y_i = 1) = p_1 ; P(y_i = 2) = p_2 = 1 - p_1$$

$$P(x_i | y_i = j) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(x_i - \mu_j)^2}{2\sigma^2}\right)$$

$$L(\theta) = \prod_i P(x_i | \theta) = \prod_i \sum_j P(x_i, y_i = j | \theta)$$

$$\log L(\theta) = \sum_i \log \left( \sum_j \frac{p_j}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(x_i - \mu_j)^2}{2\sigma^2}\right) \right)$$


To be continued...

# KL divergence

Def: The **Kullback-Liebler divergence (aka relative entropy)** of discrete probability distributions  $P$  and  $Q$ :

$$D_{KL}(P \parallel Q) = \sum_{i \in \mathcal{I}} P(x_i) \cdot \log \frac{P(x_i)}{Q(x_i)} \quad \begin{array}{l} i \in \mathcal{I}: \text{sum over } x \text{ s.t } P(x) > 0 \\ Q(x) = 0 \rightarrow P(x) = 0, 0 \log 0 = 0 \end{array}$$

Lemma: KL divergence is nonnegative

$\log(x) \leq x - 1$  for all  $x > 0$  with equality iff  $x = 1$

$$\begin{aligned} -D_{KL}(P \parallel Q) &= \sum_{i \in \mathcal{I}} P(x_i) \cdot \log \frac{Q(x_i)}{P(x_i)} \leq \sum_{i \in \mathcal{I}} P(x_i) \cdot \left( \frac{Q(x_i)}{P(x_i)} - 1 \right) \\ &= \sum_{i \in \mathcal{I}} Q(x_i) - \sum_{i \in \mathcal{I}} P(x_i) = \sum_{i \in \mathcal{I}} Q(x_i) - 1 \leq 0 \end{aligned}$$

With equality iff  $P \equiv Q$

# The EM algorithm (i)

**Goal:**  $\max \log P(x|\theta) = \log (\sum P(x,y|\theta))$

Strategy: guess an initial  $\theta$  and iteratively adjust it, making sure that the likelihood always improves.

Assume we have a model  $\theta^t$  that we wish to improve to a new value.

Bayes rule:  $P(x|\theta) = P(x,y|\theta) / P(y|x,\theta)$

Take log and multiply both sides by  $P(y|x,\theta^t)$

$$P(y|x,\theta^t) \cdot \log P(x|\theta) = P(y|x,\theta^t) \cdot \log P(x,y|\theta) - P(y|x,\theta^t) \cdot \log P(y|x,\theta)$$

$$\sum_y P(y|x,\theta^t) \cdot \log P(x|\theta) = \sum_y P(y|x,\theta^t) \cdot \log P(x,y|\theta) - \sum_y P(y|x,\theta^t) \cdot \log P(y|x,\theta)$$

$$\log P(x|\theta) = \sum_y P(y|x,\theta^t) \cdot \log P(x,y|\theta) - \sum_y P(y|x,\theta^t) \cdot \log P(y|x,\theta)$$



# The EM algorithm (ii)

$$\log P(x | \theta) = \sum_y P(y | x, \theta^t) \cdot \log P(x, y | \theta) - \sum_y P(y | x, \theta^t) \cdot \log P(y | x, \theta)$$

$$\log P(x | \theta^t) = \sum_y P(y | x, \theta^t) \cdot \log P(x, y | \theta^t) - \sum_y P(y | x, \theta^t) \cdot \log P(y | x, \theta^t)$$

Want  $P(x | \theta) > P(x | \theta^t)$

Define  $\Delta = \log P(x | \theta) - \log P(x | \theta^t)$

Define  $Q(\theta | \theta^t) = \sum_y P(y | x, \theta^t) \cdot \log P(x, y | \theta)$

$$\Delta = Q(\theta | \theta^t) - Q(\theta^t | \theta^t) + \sum_y P(y | x, \theta^t) \cdot \log \frac{P(y | x, \theta^t)}{P(y | x, \theta)}$$

↑  
KL Divergence  $\geq 0$

$$\Rightarrow \Delta \geq Q(\theta | \theta^t) - Q(\theta^t | \theta^t)$$

# The EM algorithm (iii)

Main component:

$$Q(\theta | \theta^t) = \sum_y P(y | x, \theta^t) \cdot \log P(x, y | \theta)$$

$\log P(x, y | \theta)$  is called the **complete log likelihood** function  
→ Q is the expectation of the complete log likelihood over the distribution of y given the current parameters  $\theta^t$

The algorithm:

repeat

- **E-step:** Calculate the Q function
- **M-step:** Maximize  $Q(\theta | \theta^t)$  with respect to  $\theta$
- Stopping criterion: improvement in log likelihood  $\leq \epsilon$

Note: local optimum guaranteed to be reached, not global.

Starting point matters! Try many..



# Back to the Gaussian mixture model

$$Q(\theta | \theta^t) = \sum_y P(y | x, \theta^t) \cdot \log P(x, y | \theta)$$

$$P(x, y | \theta) = \prod_i P(x_i, y_i | \theta) = \prod_i \prod_j P(x_i, y_i = j | \theta)^{y_{ij}}$$

$$y_{ij} = \begin{cases} 1 & y_i = j \\ 0 & y_i \neq j \end{cases}$$

$$\log P(x, y | \theta) = \sum_i \sum_j y_{ij} \log P(x_i, y_i = j | \theta)$$

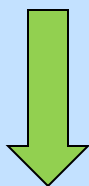
$$Q(\theta | \theta^t) = \sum_y P(y | x, \theta^t) \sum_i \sum_j y_{ij} \log P(x_i, y_i = j | \theta) =$$

$$\sum_i \sum_j \left( \sum_y P(y | x, \theta^t) y_{ij} \right) \log P(x_i, y_i = j | \theta)$$

# Application (cont.)

$$Q(\theta | \theta^t) = \sum_i \sum_j P(y_{ij} = 1 | x_i, \theta^t) \log P(x_i, y_i = j | \theta)$$

$$w_{ij}^t := P(y_{ij} = 1 | x_i, \theta^t) = \frac{P(x_i, y_i = j | \theta^t)}{\sum_j P(x_i, y_i = j | \theta^t)}$$

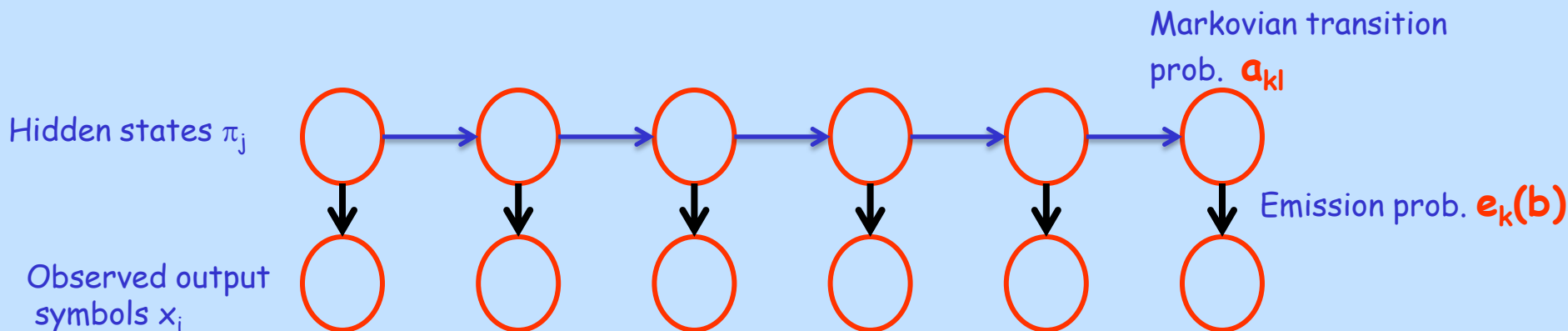


$$Q(\theta | \theta^t) = \sum_i \sum_j w_{ij}^t \left( \log \frac{1}{\sqrt{2\pi}} - \log \sigma + \log p_j - \frac{(x_i - \mu_j)^2}{2\sigma^2} \right)$$

Now write the derivatives and equate to zero to get the optimal parameters  $\theta^{t+1} = (\mu_1^{t+1}, \mu_2^{t+1}, p_1^{t+1})$

# EM for HMM: The Baum-Welch algorithm

# Reminder: HMM



$$\text{Model} = (\Sigma, Q, \Theta)$$

path  $\Pi = \pi_1, \dots, \pi_M$

Given sequence  $X = (x_1, \dots, x_M)$ :

- $a_{kl} = P(\pi_i = l \mid \pi_{i-1} = k)$ ,
- $e_k(b) = P(x_i = b \mid \pi_i = k)$

$$P(X, \Pi) = a_{\pi_0, \pi_1} \cdot \prod_{i=1}^{L-1} e_{\pi_i}(x_i) \cdot a_{\pi_i, \pi_{i+1}}$$

Goal: Finding path  $\Pi^*$  maximizing  $P(X, \Pi)$



# Max likelihood in HMM

- $\gamma = \pi, \theta = (a_{kl}, e_k(b))$
- the log likelihood is

$$\log P(x | \theta) = \log \sum_{\pi} P(x, \pi | \theta)$$

And the Q function is:

$$Q(\theta | \theta^t) = \sum_{\pi} P(\pi | x, \theta^t) \cdot \log P(x, \pi | \theta)$$

# Computing Q

$$P(x, \pi | \theta) = \prod_{k=1}^M \prod_b [e_k(b)]^{E_k(b, \pi)} \cdot \prod_{k=1}^M \prod_{l=1}^M a_{kl}^{A_{kl}(\pi)}$$

Emission  
probability, state k  
character b

Number of times we  
saw b from k in path  $\pi$

Transition  
probability, state k  
to state l

Number of transitions  
from k to l in path  $\pi$

# Computing Q (ii)

$$Q(\theta | \theta^t) = \sum_{\pi} P(\pi | x, \theta^t) \cdot \left[ \sum_{k=1}^M \sum_b E_k(b, \pi) \cdot \log(e_k(b)) + \sum_{k=1}^M \sum_{l=1}^M A_{kl}(\pi) \cdot \log a_{kl} \right] =$$

$$= \sum_{k=1}^M \sum_b \sum_{\pi} \underbrace{P(\pi | x, \theta^t) \cdot E_k(b, \pi)}_{\downarrow} \cdot \log(e_k(b)) + \sum_{k=1}^M \sum_{l=1}^M \sum_{\pi} \underbrace{P(\pi | x, \theta^t) \cdot A_{kl}(\pi)}_{\downarrow} \cdot \log a_{kl}$$

$$\sum_{\pi} \underbrace{P(\pi | x, \theta^t)}_{\substack{\uparrow \\ \text{probability}}} \cdot \underbrace{E_k(b, \pi)}_{\substack{\uparrow \\ \text{value}}} = \underbrace{E_k(b)}_{\substack{\uparrow \\ \text{expectation}}}$$

$$\sum_{\pi} \underbrace{P(\pi | x, \theta^t)}_{\substack{\uparrow \\ \text{probability}}} \cdot \underbrace{A_{kl}(\pi)}_{\substack{\uparrow \\ \text{value}}} = \underbrace{A_{kl}}_{\substack{\uparrow \\ \text{expectation}}}$$

# Computing Q (iii)

- So we want to find a set of parameters  $\theta^{t+1}$  that maximizes:

$$\sum_{k=1}^M \sum_b E_k(b) \cdot \log(e_k(b)) + \sum_{k=1}^M \sum_{l=1}^M A_{kl} \cdot \log a_{kl}$$

$$f_k(i) = P(x_0, \dots, x_i, \pi_i=k)$$

$$b_k(i) = P(x_{i+1}, \dots, x_L \mid \pi_i=k)$$

- $E_k(b)$ ,  $A_{kl}$  can be computed using forward/backward:

$$P(\pi_i=k, \pi_{i+1}=l \mid \mathbf{x}, \Theta^+) = [1/P(\mathbf{x})] \cdot f_k(i) \cdot a_{kl} \cdot e_l(x_{i+1}) \cdot b_l(i+1)$$

$$A_{kl} = [1/P(\mathbf{x})] \cdot \sum_i f_k(i) \cdot a_{kl} \cdot e_l(x_{i+1}) \cdot b_l(i+1)$$

$$\text{similarly, } E_k(b) = [1/P(\mathbf{x})] \cdot \sum_{\{i \mid x_i=b\}} f_k(i) \cdot b_k(i)$$

- For maximization, select:

$$a_{ij} = \frac{A_{ij}}{\sum_k A_{ik}}, \quad e_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')}$$

# Baum-Welch: EM for HMM

Maximize:  $\sum_{k=1}^M \sum_b E_k(b) \cdot \log(e_k(b)) + \sum_{k=1}^M \sum_{l=1}^M A_{kl} \cdot \log a_{kl}$

Multiply and divide by same factor

$$a_{ij} = \frac{A_{ij}}{\sum_k A_{ik}} \quad (\text{denote as } a_{ij}^{\text{chosen}}), \quad e_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')}$$

Difference between chosen set and some other:

$$\sum_{k=1}^M \sum_{l=1}^M A_{kl} \cdot \log \left( \frac{a_{kl}^{\text{chosen}}}{a_{kl}^{\text{other}}} \right) = \sum_{k=1}^M \sum_{k'} A_{kk'} \sum_{l=1}^M \frac{A_{kl}}{\sum_{k'} A_{kk'}} \log \left( \frac{a_{kl}^{\text{chosen}}}{a_{kl}^{\text{other}}} \right) =$$

$$= \sum_{k=1}^M \sum_{k'} A_{kk'} \sum_{l=1}^M a_{kl}^{\text{chosen}} \cdot \log \left( \frac{a_{kl}^{\text{chosen}}}{a_{kl}^{\text{other}}} \right)$$

→ always positive

# Summary: Parameter Estimation in HMM When States are Unknown

Input:  $X^1, \dots, X^n$  indep training sequences

Baum-Welch alg. (1972):

★ Expectation:

- compute expected no. of  $k \rightarrow l$  state transitions:  
$$P(\pi_i = k, \pi_{i+1} = l \mid X, \Theta) = [1/P(x)] \cdot f_k(i) \cdot a_{kl} \cdot e_l(x_{i+1}) \cdot b_l(i+1)$$
- $A_{kl} = \sum_j [1/P(X^j)] \cdot \sum_i f_k^j(i) \cdot a_{kl} \cdot e_l(x_{i+1}^j) \cdot b_l^j(i+1)$
- compute expected no. of symbol  $b$  appearances in state  $k$   
$$E_k(b) = \sum_j [1/P(X^j)] \cdot \sum_{\{i \mid x_{i+1}^j = b\}} f_k^j(i) \cdot b_k^j(i) \text{ (ex.)}$$

★ Maximization:

- re-compute new parameters from  $A, E$  using max. likelihood.

repeat (1)+(2) until improvement  $\leq \epsilon$



Leonard Baum, many years after the IDA



Lloyd Welch, USC Electrical Engineering

P value	HS binding sequence	MM binding sequence
1.0E-45		
1.0E-29		
1.0E-136		
1.0E-27		

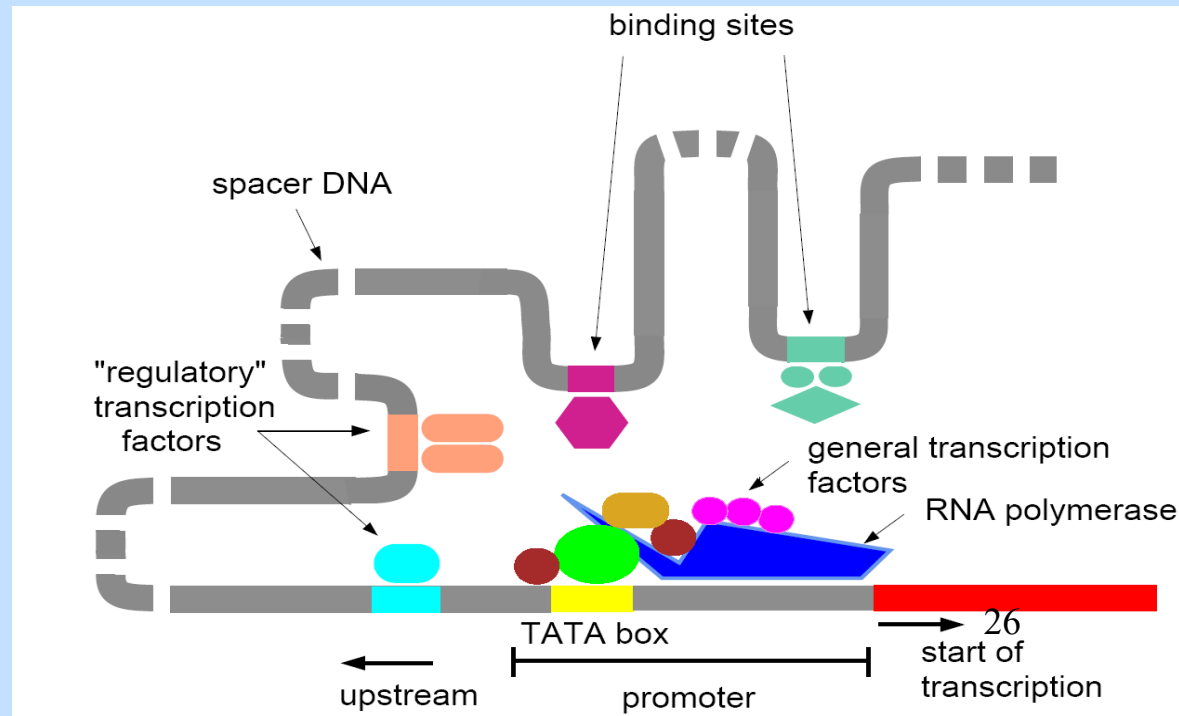
# de novo motif discovery using EM

Slides sources:  
Chaim Linhart, Danit Wider, Katherina Kechris

ARG81		PhyloCon
ARO80		PhyloCon
ARR1		Converge
ASH1		25Converge

# Transcription Factors

- A **transcription factor (TF)** is a protein that regulates a gene by binding to a **binding site (BS)** in its vicinity, specific to the TF.
- Binding sites vary in their sequences. Their sequence pattern is called a **motif**.



# Motif profile

Alignment

a	G	g	t	a	c	T	t
C	c	A	t	a	c	g	t
a	c	g	t	T	A	g	t
a	c	g	t	C	c	A	t
C	c	g	t	a	c	g	G

- Line up the patterns by their start indexes

$$\mathbf{s} = (s_1, s_2, \dots, s_t)$$

Profile

A	3	0	1	0	3	1	1	0
C	2	4	0	0	1	4	0	0
G	0	1	4	0	0	0	3	1
T	0	0	0	5	1	0	1	4

- Construct matrix profile with frequencies of each nucleotide in columns

**Motif finding:** Given a set of co-regulated genes, find a recurrent motif in their promoter regions.

# An example: Implanting Motif AAAAAAAAAGGGGGGG

atgaccgggatactgat**AAAAAAAAAGGGGGGG**ggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaata**AAAAAAAAAGGGGGGG**a  
tgagtatccctgggatgactt**AAAAAAAAAGGGGGGG**tgctctcccgatTTTTgaatatgtaggatcattcgccagggtccga  
gctgagaattggatg**AAAAAAAAAGGGGGGG**tccacgcaatcgcgaaccaacgcggacccaaaggcaagaccgataaaggaga  
tcctTTTgcggaatgtgccgggaggctggttacgtagggaagccctaacggacttaat**AAAAAAAAAGGGGGGG**cTTatag  
gtcaatcatgttcttTgtaatggattt**AAAAAAAAAGGGGGGG**gaccgcttggcgcacccaattcagtgtgggcgagcgcaa  
cgTTTTggcccttTtagaggccccgt**AAAAAAAAAGGGGGGG**caattatgagagagctaatttatcgctgctgttcat  
aacttgagtt**AAAAAAAAAGGGGGGG**ctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttTcat**AAAAAAAAAGGGGGGG**accgaaaggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagctt**AAAAAAAAAGGGGGGG**a

# Where is the Implanted Motif? (\*)

atgaccgggatactgataaaaaaagggggggggcgtagacattagataaacgtatgaagtacgtagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataaaaaaaaggggggga  
tgagtatccctgggatgacttaaaaaaaggggggggtgctctcccgatTTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgaaaaaaaggggggggtccacgcaatcgcgaaaccaacgcgacccaaaggcaagaccgataaaggaga  
tccTTTTgCGGtaatgtGCCgggaggctggttacgtaggaagccctaacggacttaataaaaaaaagggggggcttatag  
gtcaatcatgttcttTgtgaatggatttaaaaaaaaggggggggaccgcttggcgcacccaaattcagtgtgggCGagCGcaa  
CGTTTTgCCcttGttagaggccccgtaaaaaaaagggggggcaattatgagagagctaatttatCGcgtGcgtgttcat  
aacttgagttaaaaaaaagggggggctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataaaaaaaagggggggaccgaaaggaag  
ctggTgagcaacgacagattcttacgtgcattagctcGcttccggggatctaatagcacgaagcttaaaaaaaaggggggga

# Implanting Motif AAAAAAGGGGGGG with Four Mutations

atgaccgggatactgatAgAAgAAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataCAAtAAACGGcGGGga  
tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgattTTTTgaatatgtaggatcattcgccaggggccga  
gctgagaattggatgCAAAAAAAGGGattGttccacgcaatcgcgaaaccaacgcgaccxaaaggcaagaccgataaaggaga  
tcctTTTgCGgtaattgtgCCgggaggctggttacgtaggaagccctaacggacttaatAtAAtAAAGGaaGGGcttatag  
gtcaatcatgttcttTgtgaatggattAACAAtAAGGGctGGgaccgcttgCGcaccxaaattcagtgTggcgagCGcaa  
CGgtTTTgGCcttTgttagaggccccgtAtAAACaAGGaGGGccaattatgagagagctaatttatCGcgtGcgtgttcat  
aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta  
ttggccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataActAAAAAGGaGcGGgaccgaaaggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgttccggggatctaatagcagaagcttActAAAAAGGaGcGGga

# Where is the Motif???

atgaccgggatactgatagaagaaagggttggggggtacacattagataaacgtatgaagtacgtagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaatacaataaaaacggcgga  
tgagtatccctgggatgacttaaataatggagtggtgctctcccgatTTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgcaaaaaagggttgtccacgcaatcgcgaaccaacgcggaacccaaaggcaagaccgataaaggaga  
tccTTTTgcggaatgtgccgggaggctggttacgtaggaagccctaacggacttaataataaaaggaagggttatag  
gtcaatcatgttcttTggaatggatttaacaataagggtgggaccgcttggcgcacccaaattcagtgtgggagcgcaa  
cggTTTTggcccttTtagaggccccgtataaacaaggaggccaattatgagagagctaatttatcgcgTgcgtgttcat  
aacttgagttaaaaaataggagccctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatactaaaaaggagcggaccgaaagggaag  
ctggTgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttactaaaaaggagcgga

# MEME

## Multiple EM for Motif Elicitation

[Bailey, Elkan ISMB '94]

**Goal:** Given a set of sequences, find a motif (PWM) that maximizes the expected likelihood of the data

**Technique:** EM (Expectation Maximization)  
(based on [Lawrence, Reilly '90])

# The Mixture Model

Data:  $X = (X_1, \dots, X_n)$  :

all (overlapping)  $l$ -mers in the input sequences

Assume  $X_i$ 's were generated by a two-component mixture model -  $\theta = (\theta_1, \theta_2)$  :

Model #1:  $\theta_1 =$  motif model:

$f_{i,b}$  = prob. of base  $b$  at pos  $i$  in motif,  $1 \leq i \leq l$

Model #2:  $\theta_2 =$  background (BG) model:

$f_{0,b}$  = prob. of base  $b$

Mixing parameter:  $\lambda = (\lambda_1, \lambda_2)$

$\lambda_j$  = prob. that model # $j$  is used ( $\lambda_1 + \lambda_2 = 1$ )

Assume independence between  $l$ -mers



# Log Likelihood

Missing data:  $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)$  :

$\mathbf{Z}_i = (\mathbf{Z}_{i1}, \mathbf{Z}_{i2})$ ;  $\mathbf{Z}_{ij} = 1$  if  $X_i$  from model # $j$  ; 0 o/w

Complete Likelihood of model given data:

$$L(\boldsymbol{\theta}, \boldsymbol{\lambda} / \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} / \boldsymbol{\theta}, \boldsymbol{\lambda})$$

$$= \prod_{i=1 \dots n} p(\mathbf{X}_i, \mathbf{Z}_i / \boldsymbol{\theta}, \boldsymbol{\lambda})$$

$$p(\mathbf{X}_i, \mathbf{Z}_i / \boldsymbol{\theta}, \boldsymbol{\lambda}) = p(\mathbf{X}_i / \mathbf{Z}_i, \boldsymbol{\theta}, \boldsymbol{\lambda}) p(\mathbf{Z}_i) =$$

$$= \lambda_1 p(\mathbf{X}_i / \theta_1) \text{ if } \mathbf{Z}_{i1}=1; \lambda_2 p(\mathbf{X}_i / \theta_2) \text{ if } \mathbf{Z}_{i2}=1$$

$$\rightarrow \log L = \sum_{i=1 \dots n} \sum_{j=1,2} \mathbf{Z}_{ij} \log(\lambda_j p(\mathbf{X}_i / \theta_j))$$



# MEME: Algorithm

**Goal:** Maximize  $E[\log L]$

**Outline of EM algorithm:**

- Choose starting  $\theta, \lambda$
- Repeat until convergence of  $\theta$ :
  - E-step: Re-estimate  $Z$  from  $\theta, \lambda, X$
  - M-step: Re-estimate  $\theta, \lambda$  from  $X, Z$
- Repeat all of the above for various  $\theta, \lambda \dots$



# E-step

Compute expectation of  $\log L$  over  $\mathbf{Z}$ :

$$E[\log L] = \sum_{i=1 \dots n} \sum_{j=1,2} \mathbf{Z}'_{ij} \log (\lambda_j p(\mathbf{X}_i / \theta_j))$$

where:

$$\begin{aligned} \mathbf{Z}'_{ij} &= p(\mathbf{Z}_{ij}=1 / \theta, \lambda, \mathbf{X}_i) = \\ &= p(\mathbf{Z}_{ij}=1, \mathbf{X}_i / \theta, \lambda) / p(\mathbf{X}_i / \theta, \lambda) = \\ &= p(\mathbf{Z}_{ij}=1, \mathbf{X}_i / \theta, \lambda) / \sum_{k=1,2} p(\mathbf{Z}_{ik}=1, \mathbf{X}_i / \theta, \lambda) = \\ &= \lambda_j p(\mathbf{X}_i / \theta_j) / \sum_{k=1,2} \lambda_k p(\mathbf{X}_i / \theta_k) \end{aligned}$$



# M-step

Find  $\theta, \lambda$  that maximize  $E[\log L] = Q(\theta, \lambda / \theta^t, \lambda^t)$ :

$$E[\log L] = \sum_{i=1 \dots n} \sum_{j=1,2} \mathbf{Z}'_{ij} \log (\lambda_j p(\mathbf{X}_i / \theta_j))$$

Finding  $\lambda$ :

Suffices to maximize  $L_1 = \sum_{i=1 \dots n} \sum_{j=1,2} \mathbf{Z}'_{ij} \log \lambda_j$

$$\lambda_1 + \lambda_2 = 1 \rightarrow L_1 = \sum_{i=1 \dots n} (\mathbf{Z}'_{i1} \log \lambda_1 + \mathbf{Z}'_{i2} \log (1 - \lambda_1))$$

$$dL_1 / d\lambda_1 = \sum_{i=1 \dots n} (\mathbf{Z}'_{i1} / \lambda_1 - \mathbf{Z}'_{i2} / (1 - \lambda_1))$$



# MEME: Algorithm

**M-step (cont.):**

$$dL_1/d\lambda_1 = \sum_{i=1\dots n} (\mathbf{Z}'_{i1} / \lambda_1 - \mathbf{Z}'_{i2} / (1-\lambda_1)) = 0$$

$$\rightarrow \lambda_1 \sum_{i=1\dots n} \mathbf{Z}'_{i2} = (1-\lambda_1) \sum_{i=1\dots n} \mathbf{Z}'_{i1}$$

$$\rightarrow \lambda_1 (\sum_{i=1\dots n} (\mathbf{Z}'_{i1} + \mathbf{Z}'_{i2})) = \sum_{i=1\dots n} \mathbf{Z}'_{i1}$$

$$\rightarrow \lambda_1 = (\sum_{i=1\dots n} \mathbf{Z}'_{i1}) / n$$

$$\lambda_2 = 1 - \lambda_1 = (\sum_{i=1\dots n} \mathbf{Z}'_{i2}) / n$$

**Finding  $\theta$ : ...**



## MEME Suite Menu

- Submit A Job
- Documentation
- Downloads
- User Support
- Alternate Servers
- Authors
- Citing
- New! Postdoc Available



# MEME

## Multiple Em for Motif Elicitation

Version 4.3.0

Use this form to submit DNA or protein sequences to MEME. MEME will analyze your sequences for similarities among them and produce a description (**motif**) for each pattern it discovers.

### Data Submission Form

#### Required

Your **e-mail address**:

Re-enter **e-mail address**:

Please enter the **sequences** which you believe share one or more motifs. The sequences may contain no more than **60000 characters** total in any of a large number of **formats**.

Enter the **name of a file** containing the sequences here:

or  
the **actual sequences** here (**Sample Protein Input Sequences**):

How do you think the occurrences of a single motif are **distributed** among the sequences?

- One per sequence  
 Zero or one per sequence  
 Any number of repetitions

MEME will find the optimum **width** of each motif within the limits you specify here:

**Minimum** width ( $\geq 2$ )

**Maximum** width ( $\leq 300$ )

Maximum **number of motifs** to find

#### Optional

**Description** of your sequences:

MEME will find the optimum **number of sites** for each motif within the limits you specify here:

**Minimum** sites ( $\geq 2$ )

**Maximum** sites ( $\leq 300$ )

**Shuffle** sequence letters

Enter the name of a file containing a **background Markov model**:

#### DNA-ONLY OPTIONS

(Ignored for protein searches)

- Search given **strand** only  
 Look for **palindromes** only

Version 4.3.0

Please send comments and questions to: [meme@nbcrc.net](mailto:meme@nbcrc.net)

Powered by Opal

# Tim Bailey, Charles Elkan



- **Senior Research Fellow  
Institute for Molecular  
Bioscience , University of  
Queensland, Brisbane,  
Australia**



- **Professor  
Department of Computer  
Science and Engineering  
University of California,  
San Diego**

FIN