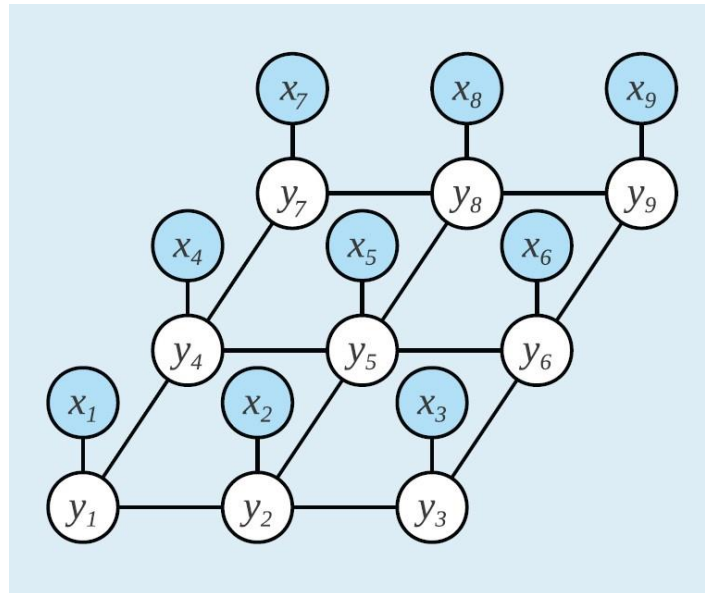


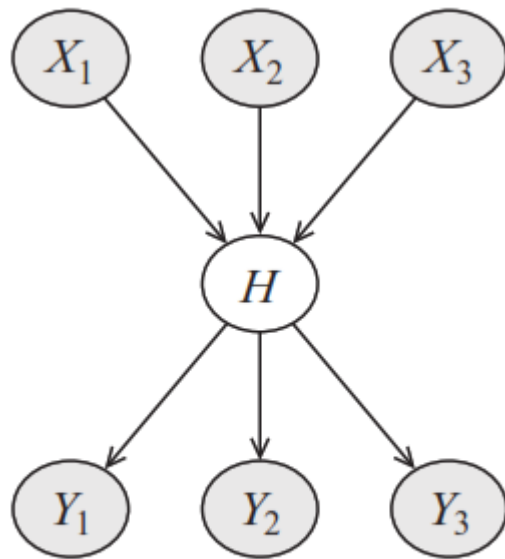
Probabilistic Graphical Models in Bioinformatics

Lecture 10: Expectation-maximization; regulatory motif finding

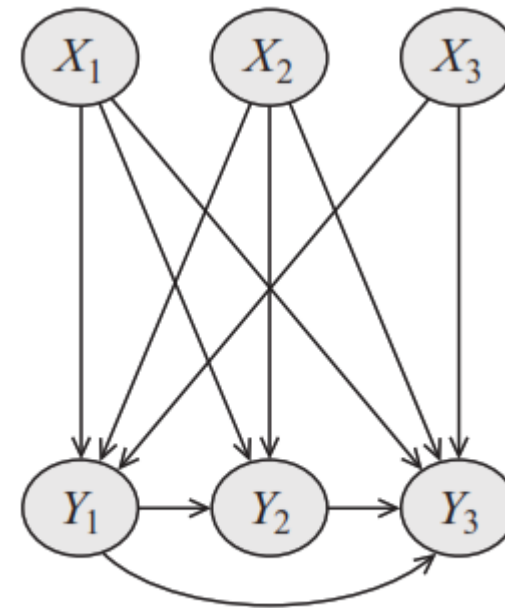


	Known structure (Parameter estimation)	Unknown structure (Structure learning)
Fully observable	<i>MLE</i> <i>Bayesian methods</i>	<i>Constraint-based methods</i> <i>Score-based methods e.g. hill climbing</i>
Partially observable		

Hidden variables can greatly simplify the structure



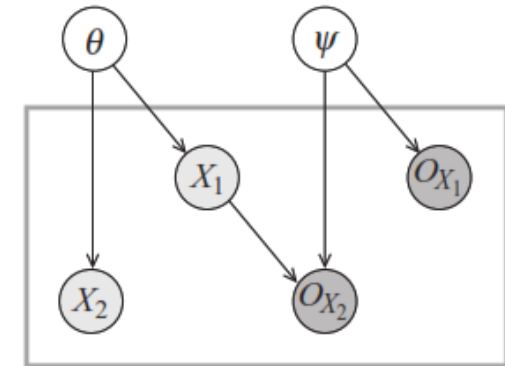
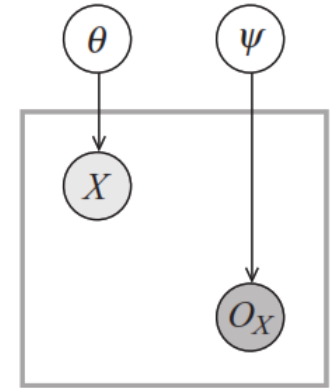
17 parameters



59 parameters

Observation mechanism

- **Missing completely at random (MCAR):** the observation mechanism is completely independent of observed and hidden variables.
 - The missing data are just a random subset of the data.
- **Missing at random (MAR):** the observation mechanism depends on some of the observed data but not directly on the hidden variables
 - Male participants of a study are more likely to tell their weight (weight is MAR).
- **Missing not at random (MNAR):** the observation mechanism depends on the hidden variables
 - People with low IQ values have missing observations for this variable (IQ is MNAR).



Expectation maximization (EM)

Expectation maximization (EM)

- EM is not a general-purpose algorithm for nonlinear function optimization.
- Tailored specifically to optimizing likelihood functions.
- Intuition
 - Parameter estimation from complete data is easy.
 - Guessing unobserved variables given parameters is possible (inference)
 - EM solves this “chicken and egg” problem in an iterative approach.

EM overview

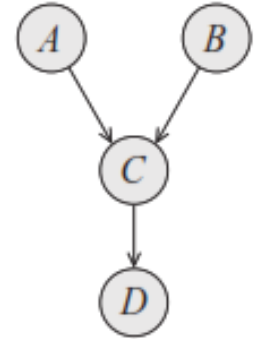
- Pick an initial values for parameters
- Iterate
 - **E-step (expectation)**: complete data using current parameters
 - **M-step (maximization)**: estimate parameters from the current complete data
- Each iteration is guaranteed to improve the log-likelihood function

Expectation Maximization (EM)

- In the fully observed case

$$\hat{\theta}_{d^1|c^0} = \frac{M[d^1, c^0]}{M[c^0]} = \frac{\sum_{m=1}^M \mathbf{I}\{\xi[m] \langle D, C \rangle = \langle d^1, c^0 \rangle\}}{\sum_{m=1}^M \mathbf{I}\{\xi[m] \langle C \rangle = c^0\}}$$

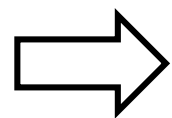
m^{th} training example



- Now consider $o = \langle a^1, ?, ?, d^0 \rangle$ is given.
 - There are four possible realizations to B and C ($\langle b^1, c^1 \rangle, \langle b^1, c^0 \rangle, \langle b^0, c^1 \rangle, \langle b^0, c^0 \rangle$)
 - We do not know which of them is true. Question: how do we know which one is more likely?
 - We can compute distribution of hidden variables given observations (inference task), $Q(B, C) = P(B, C \mid o, \theta)$

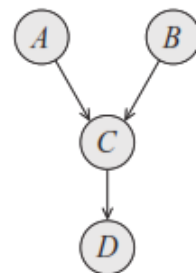
First observation: $o = \langle a^1, ?, ?, d^0 \rangle$

$$\begin{aligned}\theta_{a^1} &= 0.3 & \theta_{b^1} &= 0.9 \\ \theta_{d^1|c^0} &= 0.1 & \theta_{d^1|c^1} &= 0.8 \\ \theta_{c^1|a^0, b^0} &= 0.83 & \theta_{c^1|a^1, b^0} &= 0.6 \\ \theta_{c^1|a^0, b^1} &= 0.09 & \theta_{c^1|a^1, b^1} &= 0.2,\end{aligned}$$



$$\begin{aligned}Q(\langle b^1, c^1 \rangle) &= 0.3 \cdot 0.9 \cdot 0.2 \cdot 0.2 / 0.2196 = 0.0492 \\ Q(\langle b^1, c^0 \rangle) &= 0.3 \cdot 0.9 \cdot 0.8 \cdot 0.9 / 0.2196 = 0.8852 \\ Q(\langle b^0, c^1 \rangle) &= 0.3 \cdot 0.1 \cdot 0.6 \cdot 0.2 / 0.2196 = 0.0164 \\ Q(\langle b^0, c^0 \rangle) &= 0.3 \cdot 0.1 \cdot 0.4 \cdot 0.9 / 0.2196 = 0.0492,\end{aligned}$$

$$P(a^1, d^0 | \theta)$$



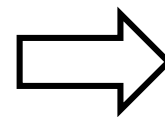
Second observation: $o = \langle ?, b^1, ?, d^1 \rangle$

$$Q'(A, C) = P(A, C | b^1, d^1, \theta)$$

$$\begin{aligned}Q'(\langle a^1, c^1 \rangle) &= 0.3 \cdot 0.9 \cdot 0.2 \cdot 0.8 / 0.1675 = 0.2579 \\ Q'(\langle a^1, c^0 \rangle) &= 0.3 \cdot 0.9 \cdot 0.8 \cdot 0.1 / 0.1675 = 0.1290 \\ Q'(\langle a^0, c^1 \rangle) &= 0.7 \cdot 0.9 \cdot 0.09 \cdot 0.8 / 0.1675 = 0.2708 \\ Q'(\langle a^0, c^0 \rangle) &= 0.7 \cdot 0.9 \cdot 0.91 \cdot 0.1 / 0.1675 = 0.3423.\end{aligned}$$

$$\tilde{\theta}_{d^1|c^0} = \frac{\bar{M}_{\theta}[d^1, c^0]}{\bar{M}_{\theta}[c^0]}$$

$$\begin{aligned}\bar{M}_{\theta}[d^1, c^0] &= Q'(\langle a^1, c^0 \rangle) + Q'(\langle a^0, c^0 \rangle) \\ &= 0.1290 + 0.3423 = 0.4713\end{aligned}$$



$$\tilde{\theta}_{d^1|c^0} = \frac{0.4713}{1.4057} = 0.3353.$$

$$\begin{aligned}\bar{M}_{\theta}[c^0] &= Q(\langle b^1, c^0 \rangle) + Q(\langle b^0, c^0 \rangle) + Q'(\langle a^1, c^0 \rangle) + Q'(\langle a^0, c^0 \rangle) \\ &= 0.8852 + 0.0492 + 0.1290 + 0.3423 = 1.4057.\end{aligned}$$

Sufficient statistics

- A function $s(\xi)$ from instances to a vector in \mathcal{R}^k if for any two datasets D and D' and any $\theta \in \Theta$ we have that

$$\sum_{x[i] \in D} s(x[i]) = \sum_{x[i] \in D'} s(x[i]) \Rightarrow L(\theta: D) = L(\theta: D')$$

- $\sum_{x[i] \in D} s(x[i])$ is called the sufficient statistics of data set D .
- Multinomial example:
 - $\langle M[1], \dots, M[K] \rangle$ is the sufficient statistics from the data
 - It is computed by instance level statistics

$$s(x^k) = (0, \dots, 0, \overset{\substack{\nearrow \\ k^{th} \text{ position}}}{1}, 0, \dots, 0)$$

$$L(\theta: D) = \prod_k \theta_k^{M[k]}$$

Fully vs partially observed case

- For each data case $\langle o[m], h[m] \rangle$
 - we define weight as $Q(h[m]) = P(h[m] \mid o[m], \theta)$
- *Expected sufficient statistics*

$$\bar{M}_{\theta}[\mathbf{y}] = \sum_{m=1}^M \sum_{\mathbf{h}[m] \in \text{Val}(\mathbf{H}[m])} Q(\mathbf{h}[m]) \mathbf{I}\{\xi[m] \langle \mathbf{Y} \rangle = \mathbf{y}\}.$$

$$M[c^1] = \sum_{m=1}^M \mathbf{I}\{\xi[m] \langle C \rangle = c^1\}.$$

Fully observed case

$$\bar{M}_{\theta}[c^1] = \sum_{m=1}^M P(c^1 \mid \mathbf{o}[m], \theta).$$


Partially observed case

The EM algorithm for Bayesian networks

- Begin with some initial parameter assignment θ^0
- The algorithm then iterates for $t = 0, 1, \dots$
 - Expectation (E-step): the algorithm uses the current parameters θ^t to compute the *expected sufficient statistics*.

$$\bar{M}_{\theta^t}[x, \mathbf{u}] = \sum_m P(x, \mathbf{u} \mid \mathbf{o}[m], \theta^t).$$

Current set of
parameters



- Maximization (M-step): treat the current *expected sufficient statistics* as observed, and perform *maximum likelihood estimation*

$$\theta_{x|\mathbf{u}}^{t+1} = \frac{\bar{M}_{\theta^t}[x, \mathbf{u}]}{\bar{M}_{\theta^t}[\mathbf{u}]}$$

The EM algorithm for Bayesian networks

Procedure Expectation-Maximization (

\mathcal{G} , // Bayesian network structure over X_1, \dots, X_n

θ^0 , // Initial set of parameters for \mathcal{G}

\mathcal{D} // Partially observed data set

)

1 **for** each $t = 0, 1, \dots$, until convergence

2 // E-step

3 $\{\bar{M}_t[x_i, \mathbf{u}_i]\} \leftarrow \text{Compute-ESS}(\mathcal{G}, \theta^t, \mathcal{D})$

4 // M-step

5 **for** each $i = 1, \dots, n$

6 **for** each $x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})$

7 $\theta_{x_i|\mathbf{u}_i}^{t+1} \leftarrow \frac{\bar{M}_t[x_i, \mathbf{u}_i]}{\bar{M}_t[\mathbf{u}_i]}$

8 **return** θ^t

Procedure Compute-ESS (

\mathcal{G} , // Bayesian network structure over X_1, \dots, X_n

θ , // Set of parameters for \mathcal{G}

\mathcal{D} // Partially observed data set

)

1 // Initialize data structures

2 **for** each $i = 1, \dots, n$

3 **for** each $x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})$

4 $\bar{M}[x_i, \mathbf{u}_i] \leftarrow 0$

5 // Collect probabilities from all instances

6 **for** each $m = 1 \dots M$

7 Run inference on $\langle \mathcal{G}, \theta \rangle$ using evidence $\mathbf{o}[m]$

8 **for** each $i = 1, \dots, n$

9 **for** each $x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})$

10 $\bar{M}[x_i, \mathbf{u}_i] \leftarrow \bar{M}[x_i, \mathbf{u}_i] + P(x_i, \mathbf{u}_i | \mathbf{o}[m])$

11 **return** $\{\bar{M}[x_i, \mathbf{u}_i] : \forall i = 1, \dots, n, \forall x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})\}$

EM – General case

- X : observed data, Z : hidden variables, θ : parameters
- Complete data: $\{X, Y\}$
- **E-step:** we compute the expectation of the complete-data likelihood

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta).$$

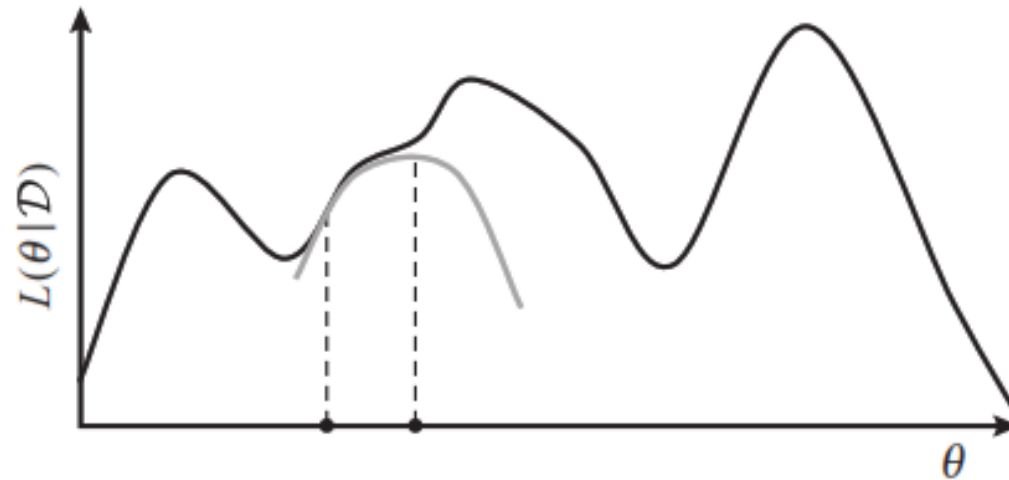
Posterior of hidden variables
given observations

Complete-data likelihood

- **M-step:** we maximize this function

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}).$$

Mathematical intuition



Theorem 19.5

During iterations of the EM procedure of algorithm 19.2, we have that

$$\ell(\boldsymbol{\theta}^{t+1} : \mathcal{D}) - \ell(\boldsymbol{\theta}^t : \mathcal{D}) \geq \mathbf{E}_{P(\mathcal{H}|\mathcal{D}, \boldsymbol{\theta}^t)} [\ell(\boldsymbol{\theta}^{t+1} : \mathcal{D}, \mathcal{H})] - \mathbf{E}_{P(\mathcal{H}|\mathcal{D}, \boldsymbol{\theta}^t)} [\ell(\boldsymbol{\theta}^t : \mathcal{D}, \mathcal{H})].$$

As a consequence, we obtain that:

$$\ell(\boldsymbol{\theta}^t : \mathcal{D}) \leq \ell(\boldsymbol{\theta}^{t+1} : \mathcal{D}).$$



Each iteration improves the likelihood

Hard-assignment EM

- It iterates over two steps
 - **Completing the data:** for each data instance $o[m]$, select the single assignment $h[m]$ that maximizes $P(h \mid o[m], \theta^t)$.
 - Estimating the new parameters θ^{t+1} using the complete data
- The objective is to maximize the likelihood of the complete data

$$\max_{\theta, \mathcal{H}} \ell(\theta : \mathcal{H}, \mathcal{D}).$$

- **The “Soft-assignment” EM objective:** attempts to maximize $l(\theta : D)$, averaging over all possible completions of the data

- Both the scoring function and the search procedure is considerably more complicated in the case of incomplete data.
- Scoring structures
 - The likelihood score does not penalize more complex models
 - Recall Bayesian score

$$\text{score}_{\mathcal{B}}(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{D} | \mathcal{G}) + \log P(\mathcal{G})$$

where

$$P(\mathcal{D} | \mathcal{G}) = \int_{\Theta_{\mathcal{G}}} P(\mathcal{D} | \boldsymbol{\theta}_{\mathcal{G}}, \mathcal{G}) P(\boldsymbol{\theta}_{\mathcal{G}} | \mathcal{G}) d\boldsymbol{\theta}_{\mathcal{G}}.$$

Likelihood function does not decompose for incomplete data

- Different approximations for dealing this issue such as Laplace approximation

$$\text{score}_{\text{Laplace}}(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{G}) + \log P(\mathcal{D} | \tilde{\boldsymbol{\theta}}_{\mathcal{G}}, \mathcal{G}) + \frac{\dim(\mathbf{C})}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{C}|,$$

In practice expensive to compute

Structural EM

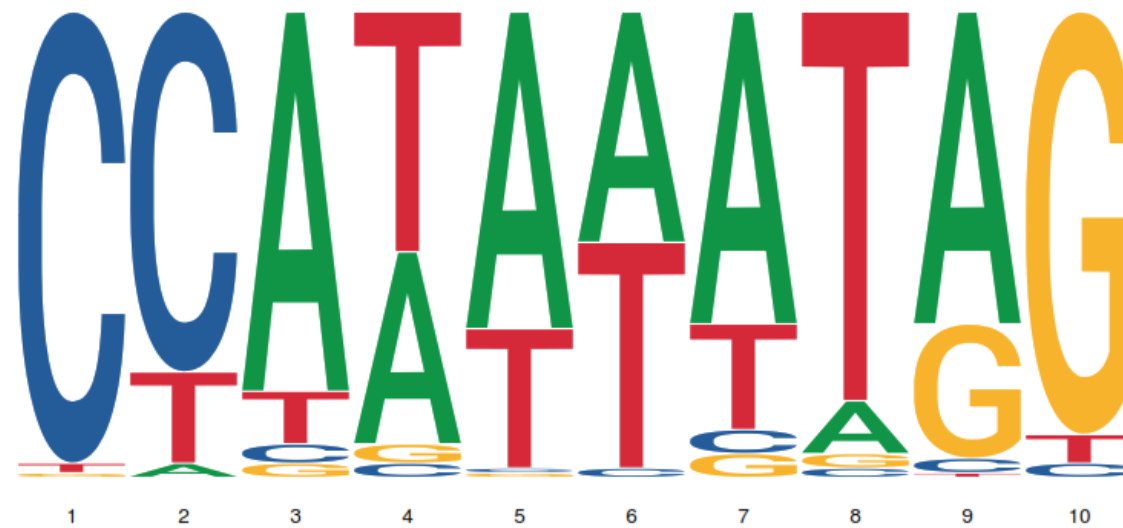
- E-step:
 - we use our current model to generate (perhaps implicitly) a completed data set
 - We may only need to compute expected sufficient statistics
 - M-step:
 - Structure learning (e.g. hill climbing)
 - Parameter estimation
- If we use the BIC score, the structural EM will improve the BIC score at each iteration.
 - Unlike the case EM, no proof the structure it finds is a local maximum.

Algorithm 19.3 The structural EM algorithm for structure learning

```
Procedure Structural-EM (  
     $\mathcal{G}^0$ , // Initial bayesian network structure over  $X_1, \dots, X_n$   
     $\theta^0$ , // Initial set of parameters for  $\mathcal{G}^0$   
     $\mathcal{D}$  // Partially observed data set  
)  
1  for each  $t = 0, 1 \dots$ , until convergence  
2      // Optional parameter learning step  
3       $\theta^{t'} \leftarrow \text{Expectation-Maximization}(\mathcal{G}^t, \theta^t, \mathcal{D})$   
4      // Run EM to generate expected sufficient statistics for  $\mathcal{D}_{\mathcal{G}^t, \theta^{t'}}$   
5       $\mathcal{G}^{t+1} \leftarrow \text{Structure-Learn}(\mathcal{D}_{\mathcal{G}^t, \theta^{t'}}^*)$   
6       $\theta^{t+1} \leftarrow \text{Estimate-Parameters}(\mathcal{D}_{\mathcal{G}^t, \theta^{t'}}^*, \mathcal{G}^{t+1})$   
7  return  $\mathcal{G}^t, \theta^t$ 
```

	Known structure (Parameter estimation)	Unknown structure (Structure learning)
Fully observable	<i>MLE</i> <i>Bayesian methods</i>	<i>Constraint-based methods</i> <i>Score-based methods e.g. hill climbing</i>
Partially observable	Gradient ascent Expectation maximization	Structural EM

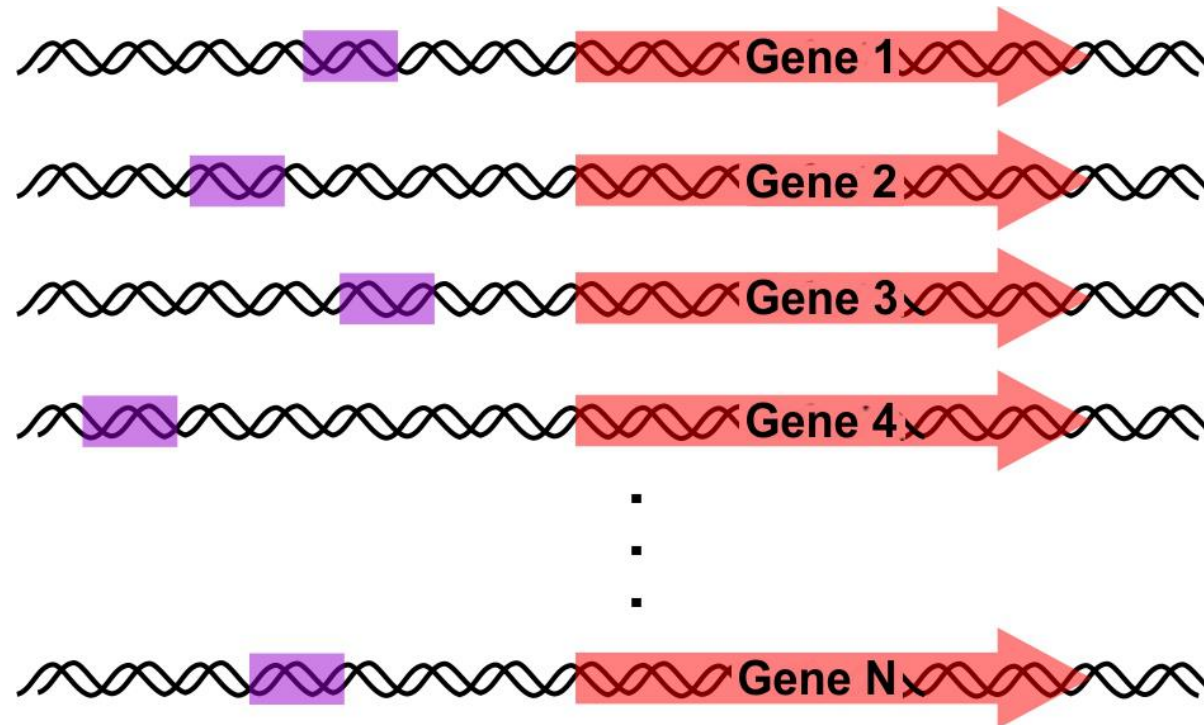
Motif finding



Regulatory motif

- A transcription factor regulates a gene by binding to a specific short DNA interval called a regulatory motif, or transcription factor binding site, in the gene's upstream region.
- For example, the transcription factor CCA1 binds to AAAATATCT in upstream of many genes

Finding regulatory motifs



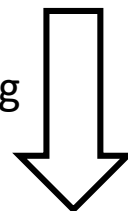
Given a collection of genes with common expression, **find** a common transcription factor binding motif!

```

1 atgaccgggatactgataaaaaaaagggggggcggtacacattagataaacgtatgaagtacgttagactcggcgccgccg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataaaaaaaagggggga
3 tgagtatccctgggatgacttaaaaaaaaggggggggtgctctccgattTTTgaatatgtaggatcattcgccagggtccga
4 gctgagaattggatgaaaaaaaggggggggtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTgcggaatgtgccgggaggtcggttacgtagggaagccctaacggacttaataaaaaaaagggggggcttatag
6 gtcaatcatgttcttTgtgaatggatttaaaaaaaaggggggggacgcttggcgacccaaattcagtggtggcgagcgcaa
7 cggTTTTggcccttTtagaggcccccgtaaaaaaaagggggggcaattatgagagagctaatactatcgcgTgcgtgttcat
8 aacttgagttaaaaaaaggggggggtggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
9 ttggccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataaaaaaaagggggggaccgaaagggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttaaaaaaaagggggga

```

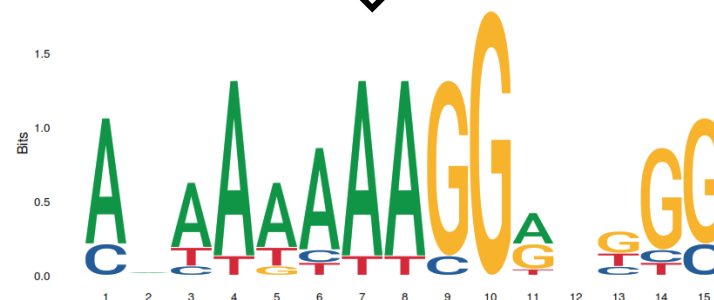
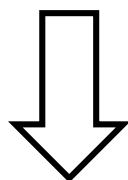
Motif finding



```

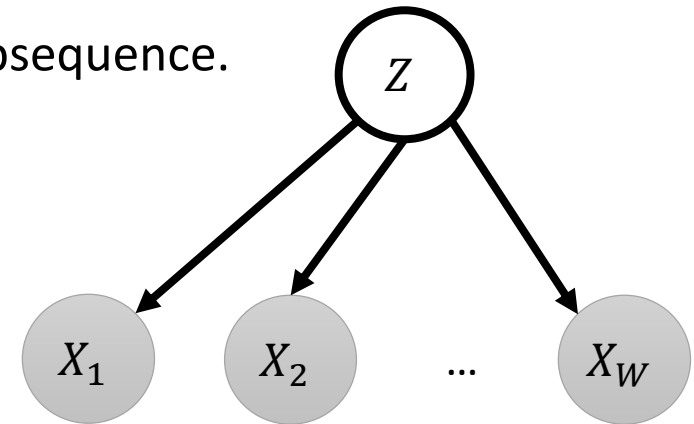
1 atgaccgggatactgatAgAAgAAAGGttGGGggcgtaacacattagataaacgtatgaagtacgttagactcggcgccgccg
2 acccctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaatacAAtAAAAGGcGGGa
3 tgagtatccctgggatgacttAAAAtAATGGaGtGGtgctctccgattTTTgaatatgtaggatcattcgccagggtccga
4 gctgagaattggatgcAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
5 tccctTTTgcggaatgtgccgggaggtcggttacgtagggaagccctaacggacttaatAtAAtAAAGGaaGGGcttatag
6 gtcaatcatgttcttTgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgacccaaattcagtggtggcgagcgcaa
7 cggTTTTggcccttTtagaggcccccgTAtAAAcAAGGaGGGccaattatgagagagctaatactatcgcgTgcgtgttcat
8 aacttgagttAAAAAAtAGGGAGccctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
9 ttggccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatActAAAAAGGaGcGGaccgaaagggaag
10 ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa

```

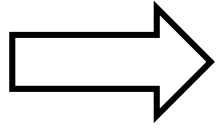


Expectation Maximization for motif finding

- Given:
 - A collection of sequences denoted as \mathcal{D} over the alphabet $\mathcal{A} = \{A, C, G, T\}$.
 - W : length of motif
- idea:
 - breaks up the sequences into N (overlapping) subsequences of length W . The data is $D = \{x[1], \dots, x[N]\}$ where $x[i] = (x_1[i], \dots, x_W[i])$.
 - Consider a two-component mixture model that assumes each subsequence is either an instance of the motif or background.
 - The hidden variable Z indicates which component generated the subsequence.
 - Perform EM for parameter estimation of the mixture model.



TT**ACCT**TAAC
G**ATGT**CTGTC
ACG**GCGT**TAG
CCCTA**ACGAG**



TTAC
TACC
ACCT
CCTT
CTTA
TTAA
TAAC

GATG
ATGT
TGTC
GTCT
TCTG
CTGT
TGTC

ACGG
CGGC
GGCG
GCGT
CGTT
GTTA
TTAG

CCCT
CCTA
CTAA
TAAC
AACG
ACGA
CGAG

each subsequence is either
an instance of the motif or
background.

TTAC
TACC
ACCT
CCTT
CTTA
TTAA
TAAC

GATG
ATGT
TGTC
GTCT
TCTG
CTGT
TGTC

ACGG
CGGC
GGCG
GCGT
CGTT
GTTA
TTAG

CCCT
CCTA
CTAA
TAAC
AACG
ACGA
CGAG



f_1	f_2	f_3	f_4
A	A	A	A
C	C	C	C
G	G	G	G
T	T	T	T

Motif model

$$\theta_m = (f_1, \dots, f_w)$$

f_0
A
C
G
T

Background model

$$\theta_{bg} = f_0$$

$$\Pr(TTAC \mid motif) = f_{1T}f_{2T}f_{3A}f_{4C}$$

$$\Pr(TTAC \mid background) = f_{0T}f_{0T}f_{0A}f_{0C}$$

Expectation Maximization for motif finding-2

- Conditional probability distributions:

$$P(z[i] = 1 \mid \lambda_m) = \lambda_m$$

$$P(x[i] \mid z[i] = 0) = \prod_j f_{0,x_j[i]}$$

$$P(x[i] \mid z[i] = 1) = \prod_j f_{1,x_j[i]}$$

- Questions:**

- Write the joint distribution of the following model.
 - Maximum likelihood estimation for λ_m, f_0, f_j .
 - Write down $P(z[i] = 1 \mid x[i])$.
- You will examine the EM algorithm for this problem in detail in project 3.

