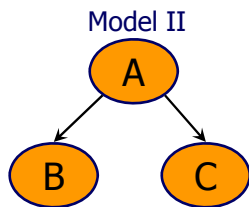


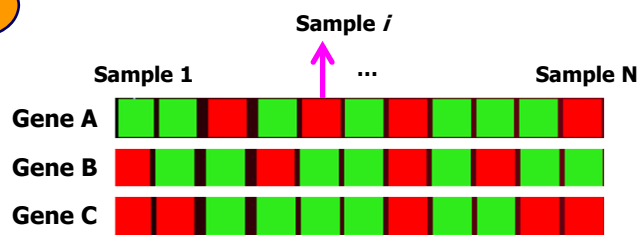


Johnson Hall (JHN) 022

## Computing $P(\text{Data} \mid \text{model II is true})$



- $P(A, B, C \mid \text{model II is true}) = ?$ 
  - $P(A)P(B|A)P(C|A)$
  - $P(A=\text{high})P(B=\text{low} \mid A=\text{high})P(C=\text{low} \mid A=\text{high})$



3

## Outline

- Probabilistic models in biology
  - Model selection problem
- Mathematical foundations
- Bayesian networks
- Learning from data
  - Maximum likelihood estimation
  - Maximum a posteriori (MAP)
  - Expectation and maximization



4

# Parameter Estimation

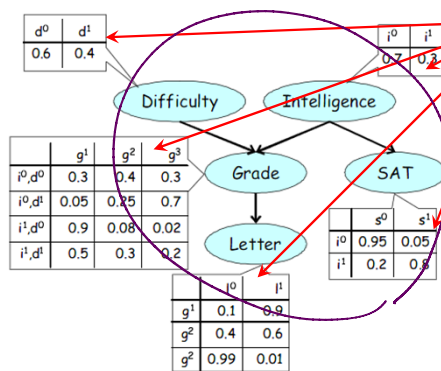
## Assumptions

- Fixed network structure
- Fully observed instances of the network variables:  $D = \{d[1], \dots, d[M]\}$
- Maximum likelihood estimation (MLE)!

For example,  
 $\{i^0, d^1, g^1, l^0, s^0\}$

strong data

$M=100$   
 $M=1000$



"Parameters" of the Bayesian network

from Koller & Friedman

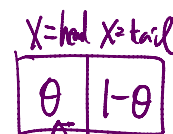
5

## The *Thumbtack* example

### Parameter learning for a single variable.

#### Variable

- $X$ : an outcome of a thumbtack toss
- $\text{Val}(X) = \{\text{head}, \text{tail}\}$



#### Data

- A set of thumbtack tosses:  $x[1], \dots, x[M]$

$x[i] \in \{\text{head}, \text{tail}\}$



6

## Maximum likelihood estimation

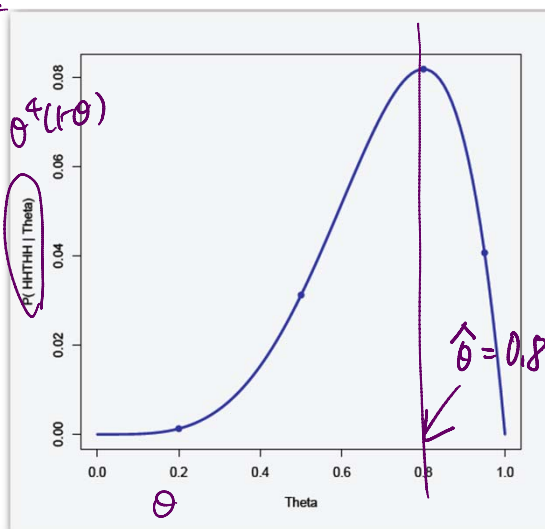
- Say that  $P(x=\text{head}) = \theta$ ,  $P(x=\text{tail}) = 1-\theta$ 
  - $P(\text{HHTTHHH} \dots \langle M_h \text{ heads}, M_t \text{ tails} \rangle; \theta) = \theta \theta (1-\theta) \dots = \theta^{M_h} (1-\theta)^{M_t}$
- Definition:** The likelihood function
  - $L(\theta : D) = P(D : \theta) = \theta^{M_h} (1-\theta)^{M_t}$   
 $P(D:\theta)$ ,  $\theta$  is a var in  $P(D:\theta)$   $\frac{P(D:\theta)}{P(\theta)}$
- Maximum likelihood estimation (MLE)
  - Given data  $D = \text{HHTTHHH} \dots \langle M_h \text{ heads}, M_t \text{ tails} \rangle$ , find  $\theta$  that maximizes the likelihood function  $L(\theta : D)$ .  
 $\theta^{M_h} (1-\theta)^{M_t}$   $P(D:\theta)$   $\theta \in [0, 1]$

7

## Likelihood function

Probability of  $\text{HHTTHH}$  given  $P(H) = \theta$ :

$\theta$	$\theta^4(1-\theta)$
0.2	0.0013
0.5	0.0313
0.8	0.0819
0.95	0.0407



8

## MLE for the *Thumbtack* problem

- Given data  $D = \text{HHTTTHH}\dots$  ( $M_h$  heads,  $M_t$  tails),
  - MLE solution  $\hat{\theta} = M_h / (M_h + M_t)$ .

- Proof:

$$\begin{aligned}
 L(\theta; D) &= p(D; \theta) \\
 &= \theta^{M_h} (1-\theta)^{M_t} \\
 \rightarrow \log(\cdot) \quad \log L(\theta; D) &= \log p(D; \theta) \\
 &= M_h \log \theta + M_t \log(1-\theta) \\
 \frac{\partial}{\partial \theta} \log L(\theta; D) &= 0 \\
 \frac{M_h}{\theta} - \frac{M_t}{1-\theta} &= 0
 \end{aligned}$$

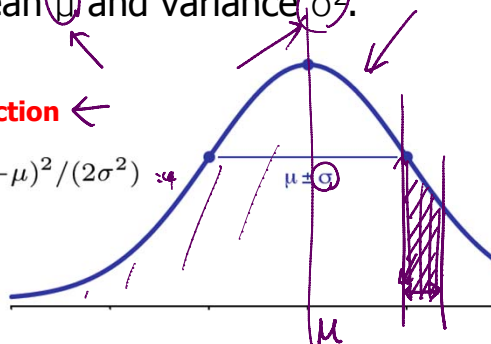
9

## Continuous Space

- Assuming sample  $x_1, x_2, \dots, x_n$  is from a parametric distribution  $f(x|\theta)$ , estimate  $\theta$ .
- Say that the  $n$  samples are from a normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

Probability density function

$$\begin{aligned}
 f(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)} \\
 \theta &= (\mu, \sigma^2)
 \end{aligned}$$



## Continuous Space (cont.)

- Let  $\Theta_1 = \mu$ ,  $\Theta_2 = \sigma^2$

$$L(\theta_1, \theta_2 : x_1, x_2, \dots, x_n) = p(x_1, \dots, x_n : \theta_1, \theta_2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\log L(\theta_1, \theta_2 : x_1, x_2, \dots, x_n) = \sum_{i=1}^n \left[ \log \frac{1}{\sqrt{2\pi\sigma^2}} + \log e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right]$$

$$\frac{\partial}{\partial \theta_1} \log L(\theta_1, \theta_2 : x_1, x_2, \dots, x_n) = \frac{\partial}{\partial \mu} \left[ \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2} \right] = 0 \quad \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial}{\partial \theta_2} \log L(\theta_1, \theta_2 : x_1, x_2, \dots, x_n) = \frac{\partial}{\partial \sigma^2} \left[ -\frac{n}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 \right] = 0$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

11

## Any Drawback?

- Is it biased?

- Is it? Yes. As an extreme, when  $n = 1$ ,  $\hat{\theta}_2 = 0$ .

- The MLE  $\hat{\theta}_2$  systematically underestimates  $\theta_2$ .

Why? A bit harder to see, but think about  $n = 2$ . Then  $\theta_1$  is exactly between the two sample points, the position that exactly minimizes the expression for  $\hat{\theta}_2$ . Any other choices for  $(\theta_1, \theta_2)$  make the likelihood of the observed data slightly lower. But it's actually pretty unlikely that two sample points would be chosen exactly equidistant, and on opposite sides of the mean, so the MLE  $\hat{\theta}_2$  systematically underestimates  $\theta_2$ .

12

## Maximum A Posteriori

- Incorporating priors. How?

MLE:  $p(D|\theta)$   $\log p(D|\theta)$

MAP:  $p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$

$\log p(D|\theta) + \log p(\theta)$

- MLE vs MAP estimation

$\theta \rightarrow 0.5$   
HHHHH  $\hat{\theta} = 1$

13

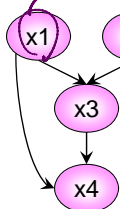
## MLE for General Problems

- Learning problem setting  $x_1 \dots x_n$ 
  - A set of random variables  $X$  from unknown distribution  $P^*$
  - Training data  $D = M$  instances of  $X$ :  $\{d[1], \dots, d[M]\}$   
 $\uparrow$   $\text{HvT} \rightarrow (0, 0, \dots)$
- A *parametric model*  $P(X; \theta)$  (a 'legal' distribution)
- Define the **likelihood function**:
  - $L(\theta : D) = p(D : \theta)$
  - $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D : \theta)$
- Maximum likelihood estimation
  - Choose parameters  $\hat{\theta}^*$  that satisfy:  $\frac{\partial}{\partial \theta} \log p(D : \theta) \Big|_{\theta = \hat{\theta}} = 0$

14

# MLE for Bayesian Networks

Structure G



$$P_G = P(x_1, x_2, x_3, x_4) = \prod_i P(x_i | \text{pa}_i)$$

$$= P(x_1) P(x_2) P(x_3 | x_1, x_2) P(x_4 | x_1, x_3)$$

More generally?

$$P_G = \prod_i P(x_i | \text{pa}_i)$$

Parameters  $\theta$

$$\theta_{x_1}, \theta_{x_2}, \theta_{x_3|x_1, x_2}, \theta_{x_4|x_1, x_3}$$

(more generally  $\theta_{x_i|\text{pa}_i}$ )

Given D:  $x[1], \dots, x[m], \dots, x[M]$ , estimate  $\theta$ .

$$L(\theta; D) = P(x[1], \dots, x[M] | \theta) = \prod_{m=1}^M P(x[m] | \theta_{x_1}, \dots, \theta_{x_4})$$

■ Likelihood decomposition:

$$P(x[1], \dots, x[M]) = \prod_{m=1}^M P(x_1[m], x_2[m], x_3[m], x_4[m])$$

■ The local likelihood function for  $X_i$  is:

$$L_i(\theta_{x_i} | D) = \prod_{m=1}^M P(x_i[m] | \theta_{x_i})$$

15

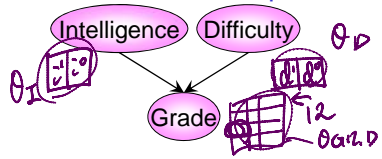
# Bayesian Network with Table CPDs

The Thumbtack example



vs

The Student example



Joint distribution

$$P(X)$$

$$P(I, D, G) = P(I) P(D) P(G | I, D)$$

Parameters

$$\theta$$

$$\theta_I, \theta_D, \theta_{G|I,D}$$

Data

$$D: \{H \dots x[m] \dots T\}$$

$$D: \{(i^1, d^1, g^1) \dots (i[m], d[m], g[m]) \dots\}$$

Likelihood function

$$L(\theta; D) = P(D; \theta)$$

$$\theta^{Mh} (1-\theta)^{Mt}$$

$$\theta_{I=i^1} \theta_{I=i^2} \dots \theta_{D=d^1} \theta_{D=d^2} \dots$$

MLE solution

$$\hat{\theta} = \frac{Mh}{Mh + Mt}$$

$$\theta_{G=A | I=i^1, D=d^1} = \frac{M_{G=A, I=i^1, D=d^1}}{M_{I=i^1, D=d^1}}$$

16

## Maximum Likelihood Estimation Review

- Find parameter estimates which make observed data most likely

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(D|\theta)$$

- General approach, as long as tractable likelihood function exists

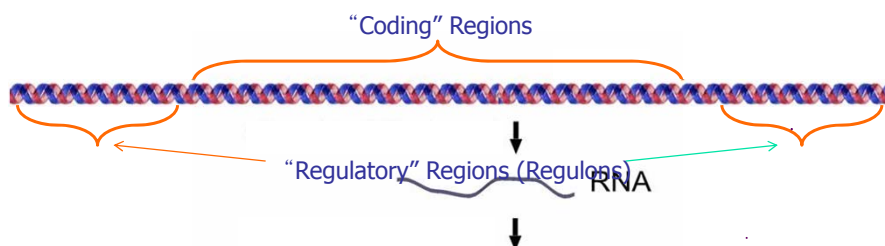
$$\frac{\partial}{\partial \theta} \log P(D|\theta)$$

- Can use all available information  $\leftarrow D$

17

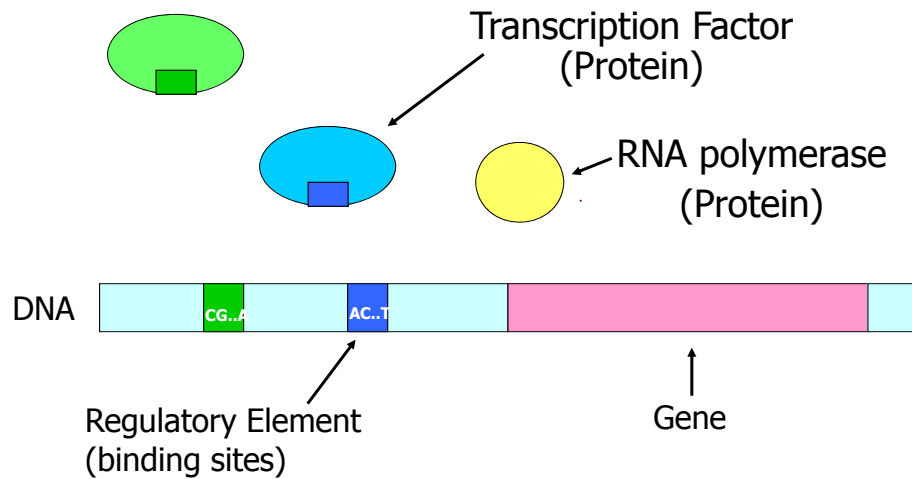
## Example – Gene Expression

- Instruction for making the proteins
- Instruction for when and where to make them



- Regulatory regions contain "binding sites" (6-20 bp).
- "Binding sites" attract a special class of proteins, known as "transcription factors".
- Bound transcription factors can initiate transcription (making RNA).
- Proteins that inhibit transcription can also be bound to their binding sites.

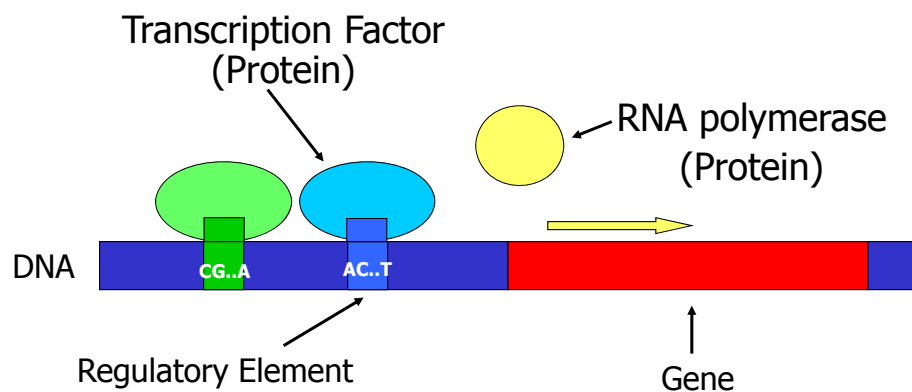
## Regulation of Genes



source: [M. Tompa](#), U. of Washington

19

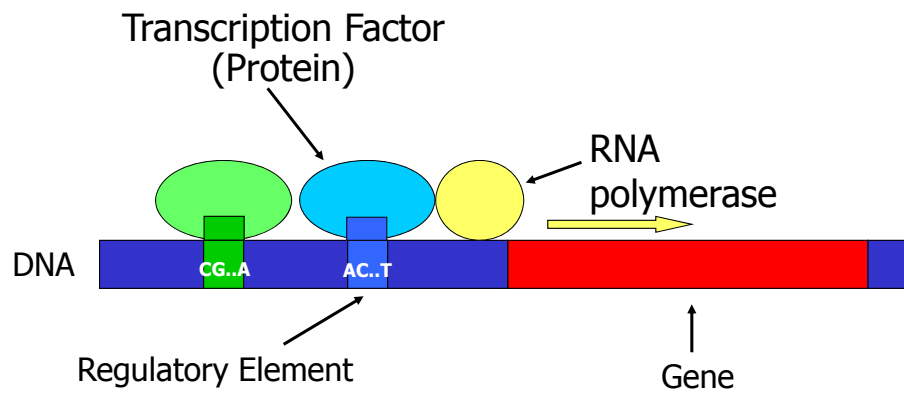
## Regulation of Genes



source: [M. Tompa](#), U. of Washington

20

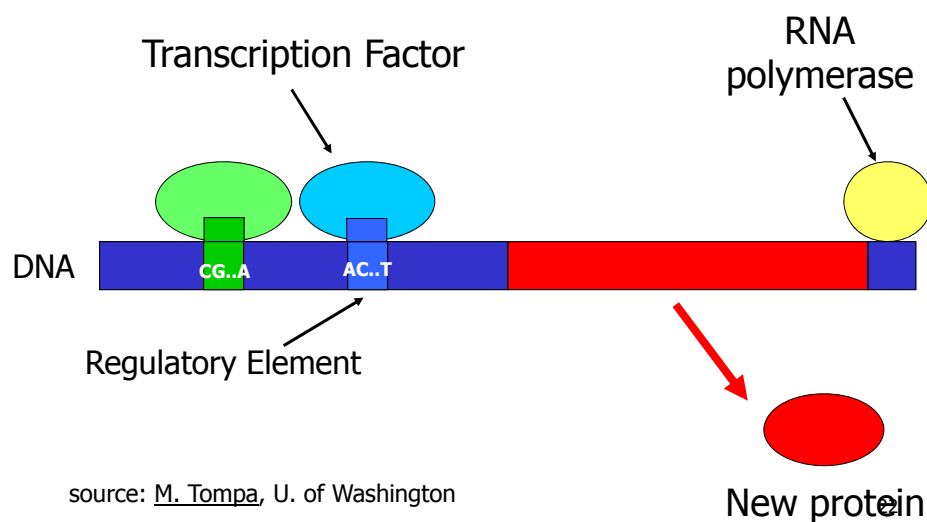
## Regulation of Genes



source: M. Tompa, U. of Washington

21

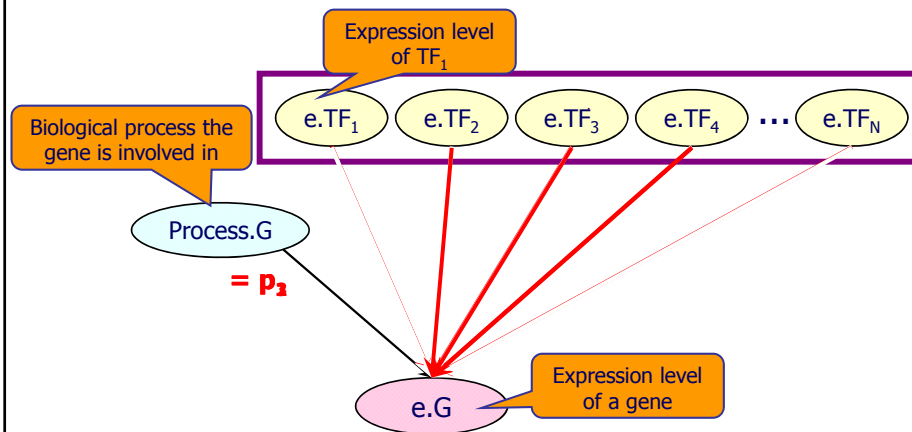
## Regulation of Genes



source: M. Tompa, U. of Washington

## The *Gene regulation* example

- What determines the expression level of a gene?
- What are observed and hidden variables?
  - e.G, e.TF's: observed; Process.G: hidden variables  $\Rightarrow$  want to infer!



23

## Not All Data Are Perfect

- Most MLE problems are simple to solve with complete data.
- Available data are "incomplete" in some way.

24

## Outline

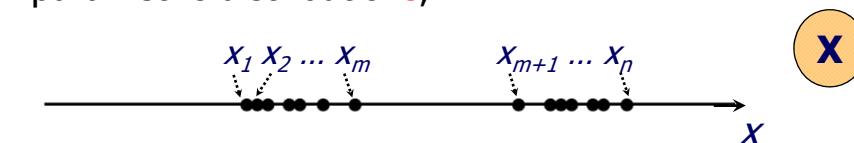
- Learning from data
  - Maximum likelihood estimation (MLE)
  - Maximum a posteriori (MAP)
  - Expectation-maximization (EM) algorithm



25

## Continuous Space Revisited...

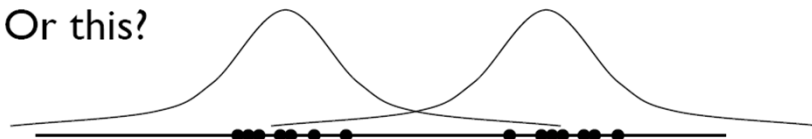
- Assuming sample  $x_1, x_2, \dots, x_n$  is from a mixture of parametric distributions,



This?



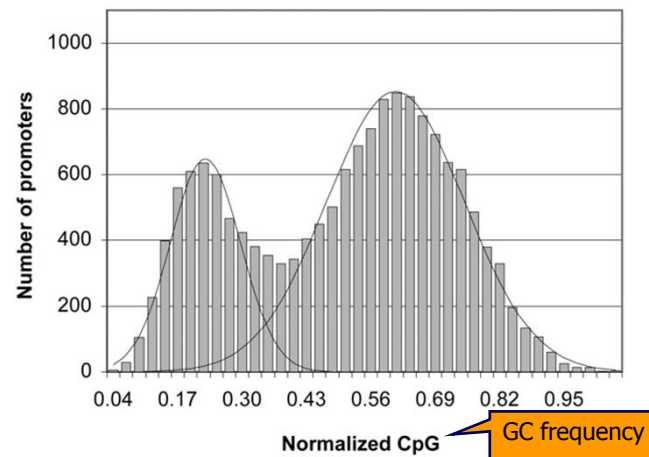
Or this?



26

## A Real Example

- CpG content of human gene promoters



"A genome-wide analysis of CpG dinucleotides in the human genome distinguishes two distinct classes of promoters" Saxonov, Berg, and Brutlag, PNAS 2006;103:1412-1417 27

## Acknowledgement

- Profs Daphne Koller & Nir Friedman, "Probabilistic Graphical Models"
- Prof Larry Ruzo, CSE 527, Autumn 2009