

Last Time:

EM Algorithm: Based on Maximum Likelihood Estimation

$L(x_1 \dots x_n | \theta)$  function to be maximized

Mixture => complicated likelihood (hidden parameters):

- Not easy to solve analytically

Likelihood Surface (see example in slides):

- Potentially many bumps observed in distribution surface due to relabeling and more clusters than have been allowed for
- Local maxima are a problem for many optimization techniques
- Relabeling in particular is problematic for some, including Markov Chain Monte Carlo

EM: Estimates  $z$ 's (hidden parameters) then estimate the hidden parameters repeat (iteration)

How Convergence Works:

Goal: Maximum likelihood estimate of  $\theta$  i.e. find  $\theta$  maximizing  $\Pr(x|\theta)$  (or  $\log(\Pr(x|\theta))$ ).

$$P(X|Y) = P(X,Y)/P(Y) \quad P(X) = P(X,Y)/P(Y|X)$$

$$\log P(X|\theta) = \log P(X,Y|\theta) - \log P(Y|\theta)$$

$$\log P(X|\theta) = \sum_y P(y|X,\theta) \cdot \log P(y,X|\theta) - \sum_y P(y|\theta) \cdot \log P(y|\theta)$$

$$Q(\theta|\theta') = \sum_y P(y|X,\theta') \cdot \log P(y,X|\theta)$$

A key trick:  $Q$  is easier to optimize than the whole thing

1)  $\log P(X|\theta) \geq \log P(X|\theta')$

2) 
$$Q(\theta|\theta') \geq Q(\theta'|\theta') + \sum_y P(y|X,\theta') \cdot \log \frac{P(y|X,\theta')}{P(y|X,\theta)}$$

$$H(P(y|X,\theta') || Q(y|X,\theta)) = \sum_y P(y|X,\theta') \cdot \log \frac{P(y|X,\theta')}{P(y|X,\theta)} \geq 0 \quad \text{Relative entropy}$$

1)  $\geq 0$  if 2)  $\geq 0$

Find  $\theta$  that maximizes  $Q(\theta|\theta')$  by maximizing  $Q(\theta|\theta') - Q(\theta'|\theta')$

### Relative Entropy:

-Given  $P, Q$

$$H(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}$$

Let  $P(x) \log \{P(x)/Q(x)\} = 0$  if  $P(x) = 0$  (since  $\lim y \log y = 0$ )

And undef. If  $0 = Q(x) < P(x)$

Theorem:  $H(P \parallel Q) \geq 0$

Upper bound:  $\log x \leq x - 1$

Lower Bound:  $\log x \geq 1 - 1/x$

$$\begin{aligned} H(P \parallel Q) &= \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &\geq \sum_x P(x) (1 - Q(x)/P(x)) \\ &= \sum_x (P(x) - Q(x)) = \sum_x P(x) - \sum_x Q(x) = 1 - 1 = 0 \end{aligned}$$

Furthermore if  $H(P \parallel Q) = 0$  iff  $P=Q$ .

### **SEQUENCE ANALYSIS**

3 billion basepairs, 98% is not traditional “genes”

-“TATA” Box analysis shows consensus sequence TATAAT ~ 10 bp up stream from transcriptional start site

- NOT exact: of 168 studied

- Nearly all had 2/3 TaxyzT
- 80-90% had all 3
- 50% agreed in each of x,y,z
- no perfect match

-Other common feature at -35 ...

Turn frequencies into probabilistic model:

Assign probability to string based on frequencies

Weight matrices, I: Statistics

- Assume:

- $f_{b,i}$  = frequency of base in position
- $f_b$  = frequency of base  $b$  in all sequences

-Log likelihood ratio, given  $S = B_1, B_2, \dots, B_6$

$$\log \frac{P(s | promoter)}{P(S | nonpromoter)} = \log \frac{\prod_{i=1}^6 f_{B_i,i}}{\prod_{i=1}^6 f_{B_i}} = \sum_{i=1}^6 \log \frac{f_{B_i,i}}{f_{B_i}}$$

Weight Matrices, II: Chemistry

- Experiments Show ~ 80% correlation of log likelihood weight matrix scores to measure binding energy of RNA polymerase to variation on TATAAT consensus.

Why are Promoters so Fuzzy?

- Restriction enzymes, e.g., recognize very precise sequences?
- Why not polymerase/promoter?
- Likely: Variation gives ~1000 fold variation in expression level
- Possible consequence: rarely expressed genes will be hard to find