

CSE 527

Lecture 7

Relative entropy
Convergence of EM
Weight matrix motif models

Talk this week

- COMBI/GS Seminar
Thomas R. Gingeras, Ph.D.
"Empirical Analysis of Sites of RNA Transcription for 30% of the Human Genome: The Changing Landscape of the Human Genome Annotations"

Wednesday, October, 20, 2004
3:30 pm, Hitchcock 132
- Refreshments in lobby at 3:20

Relative Entropy

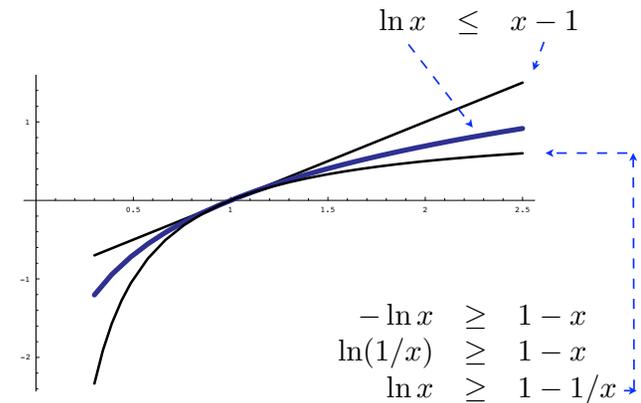
- AKA Kullback-Liebler Distance/Divergence, AKA Information Content
- Given distributions P, Q

$$H(P||Q) = \sum_{x \in \Omega} P(x) \log \frac{P(x)}{Q(x)}$$

Notes:

Let $P(x) \log \frac{P(x)}{Q(x)} = 0$ if $P(x) = 0$ [since $\lim_{y \rightarrow 0} y \log y = 0$]

Undefined if $0 = Q(x) < P(x)$



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

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

Furthermore: $H(P||Q) = 0$ if and only if $P = Q$

EM Convergence

Visible x
 hidden y
 Parameters θ

Goal Maximum likelihood estimate of θ
 i.e. Find θ maximizing $\Pr(x|\theta)$ (or $\log P(x|\theta)$)

$\forall y:$ $P(y|x) = P(x,y)/P(x)$ so $P(x) = P(x,y)/P(y|x)$

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

$$\log P(x|\theta) =$$

$$\underbrace{\sum_y P(y|x,\theta^*) \cdot \log P(x,y|\theta)}_{Q(\theta|\theta^*)} - \sum_y P(y|x,\theta^*) \cdot \log P(y|x,\theta)$$

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

Key trick: Q is easier to optimize than whole thing

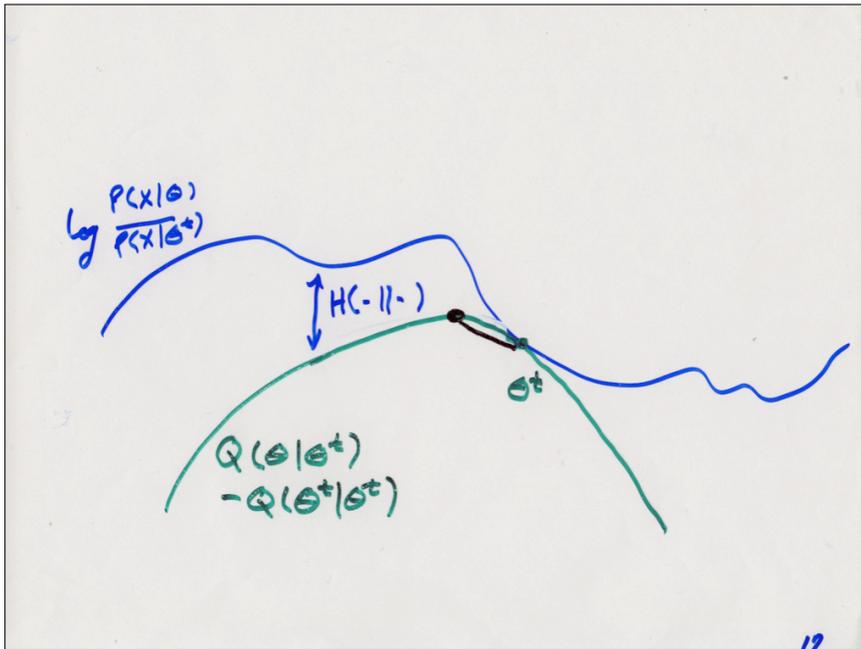
$$\textcircled{1} \log P(x|\theta) - \log P(x|\theta^*) =$$

$$\textcircled{2} Q(\theta|\theta^*) - Q(\theta^*|\theta^*)$$

$$+ \sum_y P(y|x,\theta^*) \log \frac{P(y|x,\theta^*)}{P(y|x,\theta)}$$

$$H(P(y|x,\theta^*) || P(y|x,\theta)) \geq 0$$

$\therefore \textcircled{1} \geq 0$ if $\textcircled{2} \geq 0$



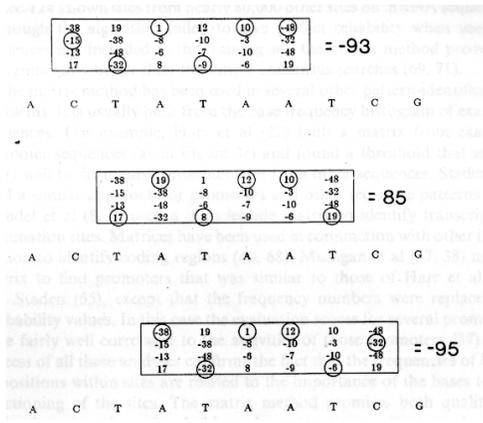
Sequence Motifs

- ## *E. coli* Promoters
- “**TATA Box**” - consensus TATAAT ~ 10bp upstream of transcription start
 - Not exact: of 168 studied
 - nearly all had 2/3 of TAxyzT
 - 80-90% had all 3
 - 50% agreed in each of x,y,z
 - **no** perfect match
 - Other common features at -35, etc.

TATA Box Frequencies

pos base	1	2	3	4	5	6
A	2	95	26	59	51	1
C	9	2	14	13	20	3
G	10	1	16	15	13	0
T	79	3	44	13	17	96

Scanning for TATA



Stormo, Ann. Rev. Biophys. Biophys Chem, 17, 1988, 241-263

Weight Matrices: Statistics

- Assume:

$f_{b,i}$ = frequency of base b in position i

f_b = frequency of base b in all sequences

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

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

Weight Matrices: Chemistry

- Experiments show ~80% correlation of log likelihood weight matrix scores to measured binding energy of RNA polymerase to variations on TATAAT consensus