

Inferring Models of cis-Regulatory Modules using Information Theory

BMI/CS 776

www.biostat.wisc.edu/bmi776/

Spring 2011

Mark Craven

craven@biostat.wisc.edu

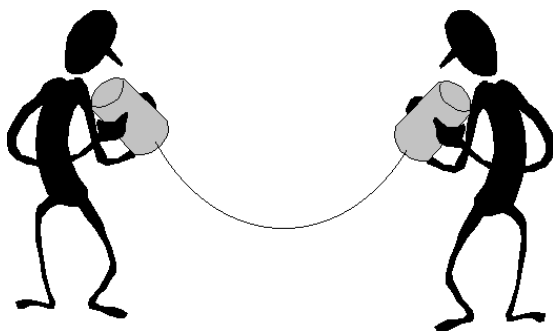
Goals for Lecture

the key concepts to understand are the following

- entropy
- mutual information
- motif logos
- using MI to identify CRM elements

Information Theory Background

- consider a problem in which you are using a code to communicate information to a receiver
- example: as bikes go past, you are communicating the manufacturer of each bike



Information Theory Background

- suppose there are only four types of bikes
- we could use the following code

<u>type</u>	<u>code</u>
Trek	11
Specialized	10
Cervelo	01
Serrota	00

- expected number of bits we have to communicate:
2 bits/bike

Information Theory Background

- we can do better if the bike types aren't equiprobable
- optimal code uses $-\log_2 P(c)$ bits for event with probability $P(c)$

Type/probability	# bits	code
$P(\text{Trek}) = 0.5$	1	1
$P(\text{Specialized}) = 0.25$	2	01
$P(\text{Cervelo}) = 0.125$	3	001
$P(\text{Serrota}) = 0.125$	3	000

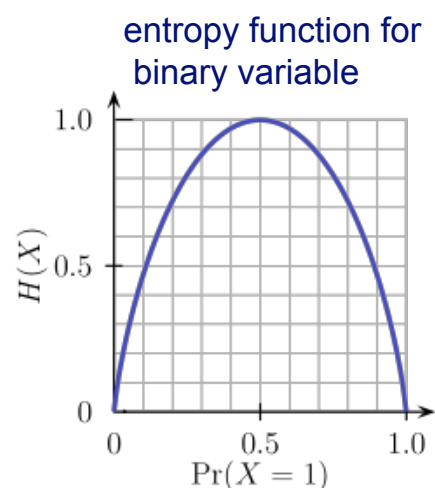
- expected number of bits we have to communicate:
1.75 bits/bike

$$-\sum_{c=1}^{|C|} P(c) \log_2 P(c)$$

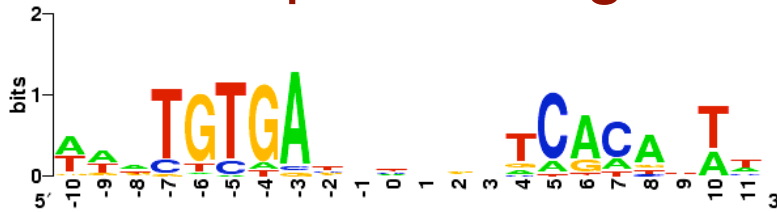
Entropy

- entropy is a measure of uncertainty associated with a random variable
- defined as the expected number of bits required to communicate the value of the variable

$$H(C) = -\sum_{c=1}^{|C|} P(c) \log_2 P(c)$$



Sequence Logos



- based on entropy (H) of a random variable (C) representing distribution of character states at each position
- height of logo at a given position determined by decrease in entropy (from maximum possible)

$$H_{\max} - H(C) = -\log_2\left(\frac{1}{N}\right) - \left(-\sum_c P(c)\log_2 P(c)\right)$$

of characters in alphabet

- height of each character c is proportional to $P(c)$

Mutual Information

- *mutual information* quantifies how much knowing the value of one variable tells about the value of another

$$I(M;C) = H(M) - H(M|C)$$

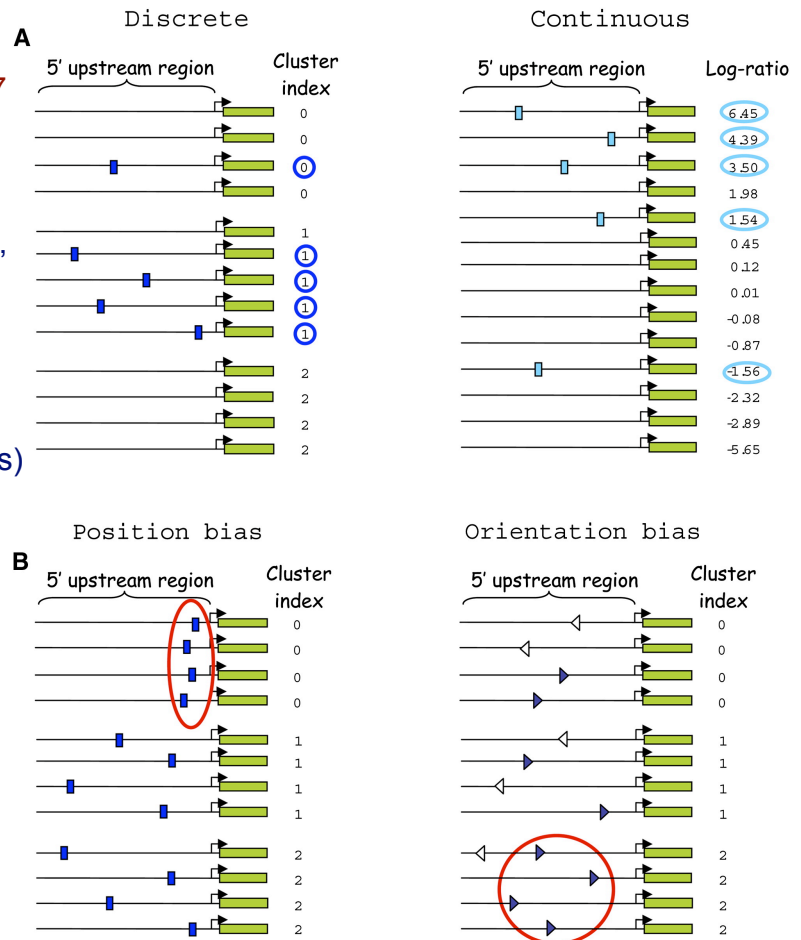
entropy of M
entropy of M conditioned on C

$$= \sum_m \sum_c P(m,c) \log_2 \left(\frac{P(m,c)}{P(m)P(c)} \right)$$

FIRE

Elemento et al., *Molecular Cell* 2007

- **Given** a set of sequences grouped into clusters
- **Find** motifs, and relationships, that have high *mutual information* with the clusters
- (also can do this when sequences have continuous values instead of cluster labels)



Mutual Information in FIRE

- we can compute the mutual information between a motif and the clusters as follows

$$I(M;C) = \sum_{m=0}^1 \sum_{c=1}^{|C|} P(m,c) \log_2 \frac{P(m,c)}{P(m)P(c)}$$

$m=0, 1$ represent absence/presence of motif

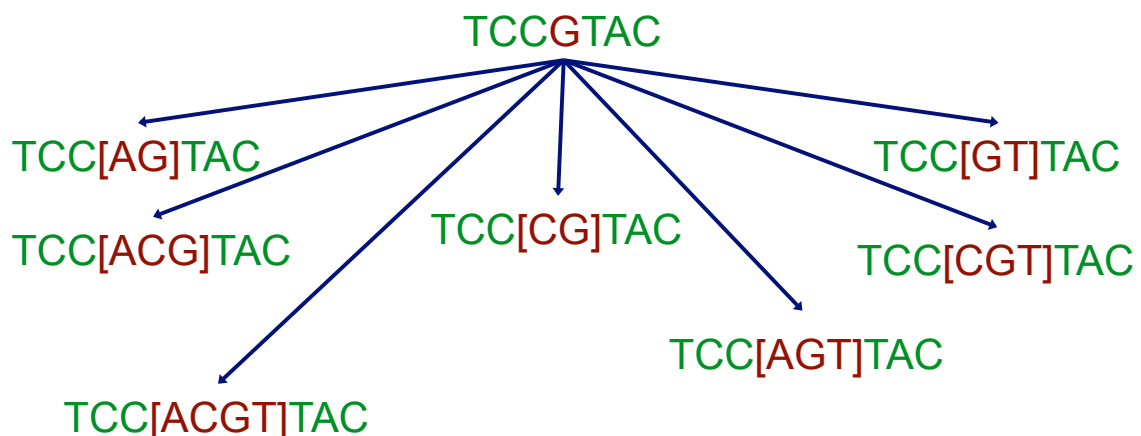
c ranges over the cluster labels

Finding Motifs in FIRE

- motifs are represented by regular expressions; initially each motif is represented by a strict k -mer (e.g. TCCGTAC)
1. test all k -mers ($k=7$ by default) to see which have significant mutual information with the cluster label
 2. filter k -mers using a significance test
 3. generalize each k -mer into a motif
 4. filter motifs using a significance test

Key Step in Generalizing a Motif in FIRE

- randomly pick a position in the motif
- generalize in all ways consistent with current value at position
- score each by computing mutual information
- retain the best generalization



Generalizing a Motif in FIRE

given: k -mer, n

$best \leftarrow \text{null}$

repeat n times

$motif \leftarrow k\text{-mer}$

 repeat

$motif \leftarrow \text{GeneralizePosition}(motif)$ // shown on previous slide

 until convergence (no improvement at any position)

 if $\text{score}(motif) > \text{score}(best)$

$best \leftarrow motif$

return: $best$

Generalizing a Motif in FIRE: Example

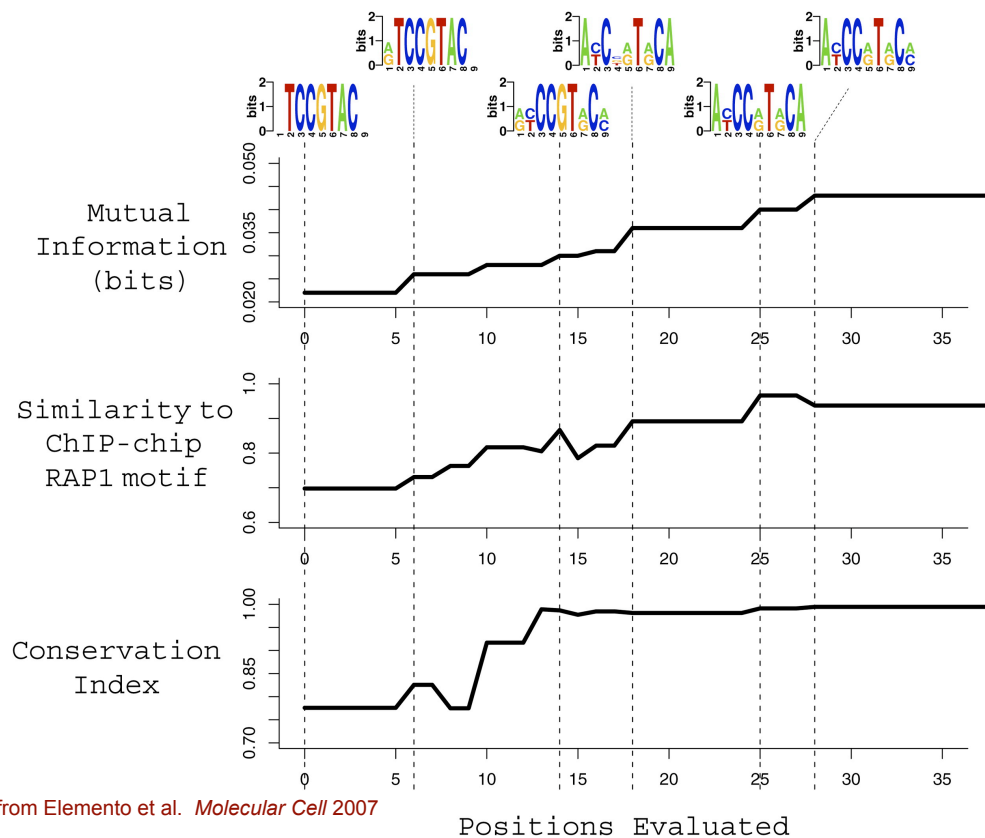


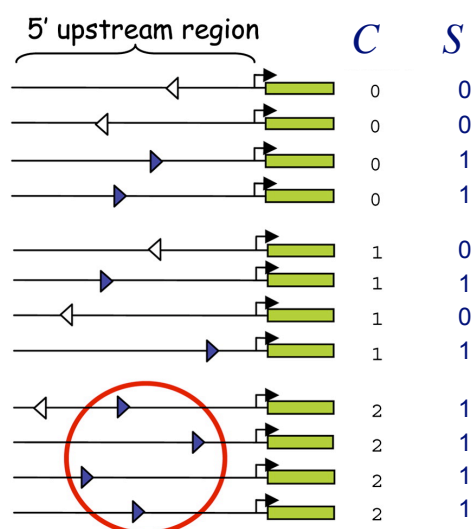
Figure from Elemento et al. *Molecular Cell* 2007

Characterizing Predicted Motifs in FIRE

- mutual information is also used to assess various properties of found motifs
 - orientation bias
 - position bias
 - interaction with another motif

Using MI to Determine Orientation Bias

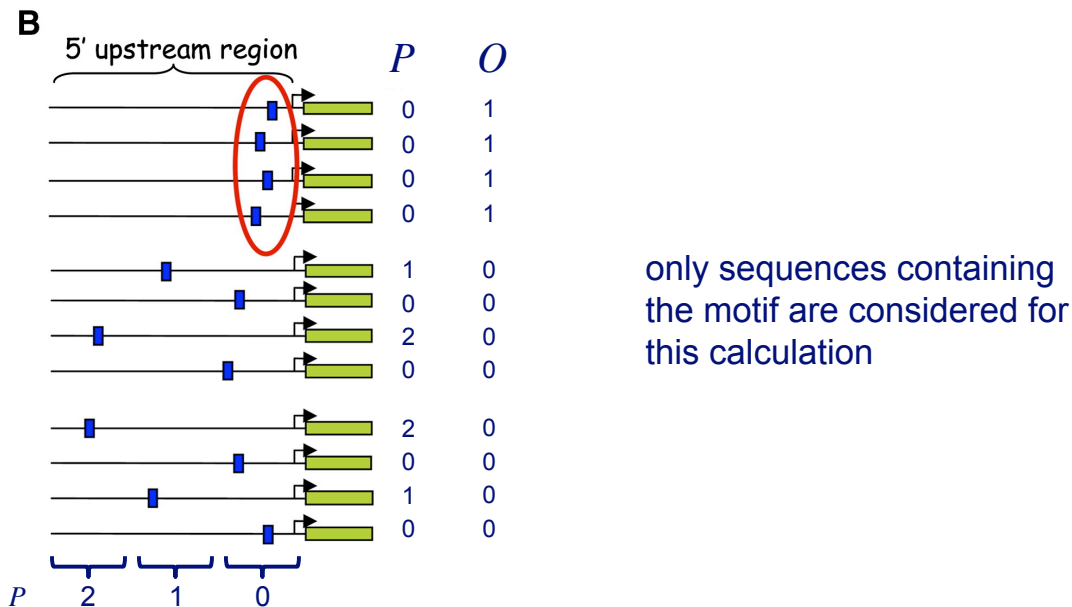
$I(S;C)$ C indicates cluster
 $S=1$ indicates motif present on transcribed strand
 $S=0$ otherwise (not present or not on transcribed strand)



also compute MI where $S=1$
 indicates motif present on
 complementary strand

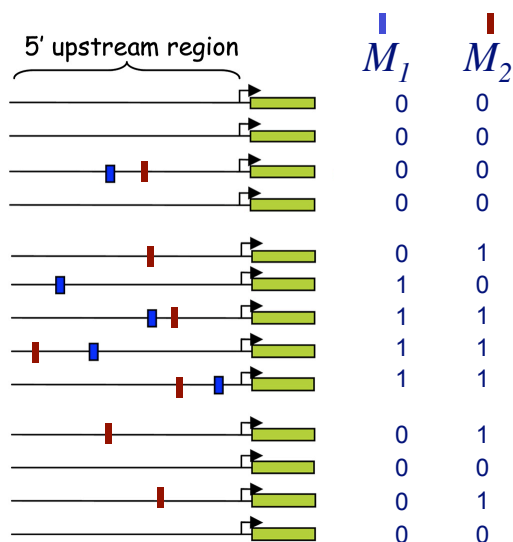
Using MI to Determine Position Bias

$I(P;O)$ P ranges over position bins
 $O=0, 1$ indicates clusters in which the motif is overrepresented or not



Using MI to Determine Motif Interactions

$I(M_1;M_2)$ $M_1=0, 1$ indicates clusters in which motif 1 is overrepresented or not; similarly for M_2



Discussion of CRM Finding Methods

- Noto & Craven
 - HMM structure search to find CRM model
 - search operators apply to compact, logical representation instead of directly to HMM
 - employs generalized (a.k.a. semi-Markov) HMM approach to model *background* sequence lengths
- FIRE
 - mutual information used to identify motifs and relationships among them
 - motif search is based on generalizing informative *k*-mers
- in contrast to many motif-finding approaches, both CRM methods take advantage of *negative* sequences
- FIRE returns all informative motifs/relationships found, whereas the Noto & Craven approach returns single discriminative model