



The Abdus Salam
International Centre
for Theoretical Physics
50th Anniversary 1964–2014



2584-2

Spring College on the Physics of Complex Systems

26 May – 20 June, 2014

Modeling DNA Motifs

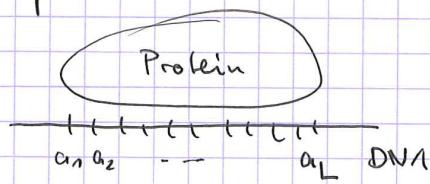
Martin Weigt
Université Pierre et Marie Curie
Paris
France

Modeling DNA Motifs

(1)

transcription factor binding site
 $(TF)BS$

Ingredients



⇒ sequence specific binding

⇒ BS ~ 10 nt

⇒ fuzzy motif \rightarrow statistical modeling

⇒ additive binding energies for closeable neighbors

⇒ statistically independent contribution of sites

⇒ Simplest modeling

Position-specific weight matrix (PSWM)

$$P(a_1, \dots, a_L | w) = \sum_{i=1}^L w_i(a_i)$$

↑
statistical indep.
of sites

↑
position specific
 $= \frac{1}{Z} \exp \left\{ \sum_i \mu_i(a_i) \right\}$

a_i	1	2	3	4	...
A	0.8	0.1	0.01		
C	0.1	0.2	0.01		
G	0.06	0.1	0.3		
T	0.04	0.6	0.68		
	$\Sigma: 1$				

consensus

A T T
 most likely sequence.

Two questions:

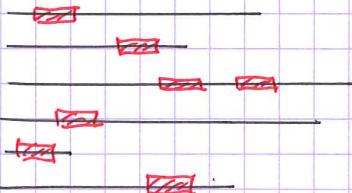
(2)

- 1) How to determine w from sequence data?
 - 2) Given w for some TF, how can we find its BS?
-

For 1: How look data like?

→ collection of sequences which (possibly) contain TFBS

- promoter regions of potential targets of a TF
- ChIP experiments



Here: idealized situation · exactly TFBS seq.

$$\text{①} = \begin{pmatrix} a_1^1 & \dots & a_i^1 & \dots & a_L^1 \\ | & & | & & | \\ a_1^m & \dots & a_i^m & \dots & a_L^m \\ | & & | & & | \\ a_1^M & \dots & a_i^M & \dots & a_L^M \end{pmatrix}$$

a_i^m ... nt in pos. i of seq. m

$n_i(A) + n_i(C) + n_i(G) + n_i(T) = M$

M ... number of seqs.

L ... length of seqs.

a_i^m ... nt in pos. i of seq. m
 $(\in \{A, C, G, T\})$

(3)

For any given ω , we can determine the probability of our data D :

$$P(D | \omega) = \prod_{m=1}^M P(a_1^m, \dots, a_L^m | \omega)$$

$$\begin{aligned} & \stackrel{\text{↑ assume independence of}}{\quad} \\ & \quad \text{sequences} \\ & = \prod_{m=1}^M \prod_{i=1}^L w_i(a_i^m) \end{aligned}$$

$$= \prod_{i=1}^L \underbrace{\prod_{a \in \{ACGT\}}}_{n_i(a)} w_i(a)$$

with

$n_i(a)$ = number of occurrences of nt a in column i of D .

Maximum-likelihood inference



→ choose ω such that ~~P(D)~~ the data are as probable as possible!

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

→ can be solved analytically and gives (no surprise)

$$w_i(a) = \frac{w_i(a) + \mu}{M + 4\mu}$$

pseudo count μ

→ avoids $w_i(a) = 0$ only because a never seen in pos. i → small M .

(4)

But we need

$$P(\omega | D)$$

→ have data, but not the PWM

Bayes theorem:

$$P(\omega | D) =$$

$$\frac{P(D | \omega) \cdot P(\omega)}{P(D)}$$

a priori distribution

$$P(\omega) = \prod_{i=1}^L P(\omega_i)$$

$$= \prod_{i=1}^L c_i \prod_{a \in \{A,C,G,T\}} w_i(a)^{\gamma-1}$$

• $\gamma = 1$ ⇒ flat prior

• $\gamma > 1$ ⇒ favors $\omega_i = (\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$
→ not realistic

• $\gamma < 1$ ⇒ favors $\omega_i \in \left\{ \begin{pmatrix} 1, 0, 0, 0 \\ 0, 1, 0, 0 \\ 0, 0, 1, 0 \\ 0, 0, 0, 1 \end{pmatrix} \right\}$



"consensus" - type ω_i

~~$P(\omega | D) \propto \prod_{i=1}^L \omega_i^{x_i}$~~

Inference of ω from $P(\omega | D)$

Maximum likelihood (max posterior prob)

$$\hat{\omega}_{\text{ML}} = \underset{\omega}{\operatorname{argmax}} P(\omega | D)$$

$$\Rightarrow \hat{\omega}_i(a) = \frac{n_i(a) + \gamma - 1}{M + 2\gamma - 1}$$

But: $P(\omega | D)$ is prob. distr. fct.

e.g. $\langle \omega_i(a) \rangle_p = \frac{n_i(a) + \gamma}{M + 2\gamma}$



difference for small data sets.

Question 2 : How to decide if a Q
NEW sequence is
functional site?

⇒ compare to null model

- indept. nucleotides
- no site specificity
- background nucleotide freqs.

$$P_o(\vec{s}) = \prod_{i=1}^L p_o(s_i)$$

⇒ log odds ratio

$$S(\vec{s}) = \log$$

$$\frac{P_f(\vec{s} | w)}{P_o(\vec{s})}$$

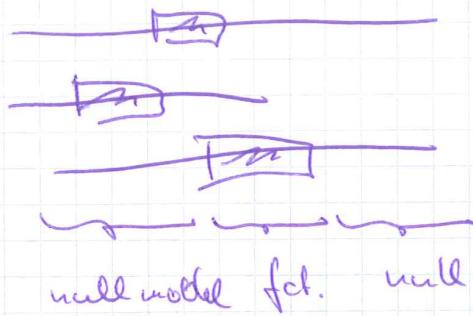
$$= \sum_{i=1}^L \ln \frac{w_i(s_i)}{p_o(s_i)}$$

more likely under functional or null model?

remark: threshold zero unreliable because of much better sampling of background sites, many random sites look functional...

more general data situation.

(6)



$$\vec{s} = (s_1, \dots, s_N) \quad N > L$$

null fact. null

$$\Rightarrow P(\vec{s} | \omega, \ell) \rightarrow \prod_{i=1}^l p_0(s_i) \prod_{j=1}^L w_j(s_{j+l}) \\ \cdot \prod_{i=l+1}^{N-l} p_0(s_i) \\ = \prod_{i=1}^N p_0(s_i) \prod_{j=1}^L \frac{w_j(s_{j+l})}{p_0(s_{j+l})}$$

Data set: $D = \{\vec{s}^n \mid n = 1..M\}$

$$\hookrightarrow \vec{s}^n = (s_1, \dots, s_{N_n})$$

$$\Rightarrow P(D \mid \omega, (\ell^n)_{n=1..M})$$

$$= \underbrace{\prod_{n=1}^M P(\vec{s}^n \mid \omega, \ell^n)}$$

independence of different segs.

\Rightarrow estimate w and (e^u) !

(2)

two strategies

$$1) \quad P(D|w) = \sum_{(e^u)} P(D|w, (e^u)) P(e^u)$$

+ max w.r.t. w (MEME)

$$2) \quad P(D|(e^u)) = \int dw \quad P(D|w, (e^u)) P(w)$$

↑
Dirichlet

+ MC sampling of (e^u)

Strategy 1: Max. of $P(D|w)$

by expectation-maximization (EM)

For given (e^u)

$$\Rightarrow w_i(a) = \frac{1}{M} \sum_{m=1}^M \langle \delta_a, s_{i+m} \rangle_{e^u}$$

→ averaging over (e^u)

$$\begin{aligned} w_i(a) &= \frac{1}{M} \sum_{m=1}^M \langle \delta_a, s_{i+m} \rangle_{e^u} \\ &= \frac{1}{M} \sum_{m=1}^M \frac{\sum_{e^u} \delta_a, s_{i+m} P(\vec{s}^m | w, e^u)}{\sum_{e^u} P(\vec{s}^m | w, e^u)} (*) \end{aligned}$$

Problem: w on both sides !

(P)

- Algo:
- initialize ω randomly
 - Plug in rhs of (*) , update ω
 - iterate until ω remains stationary
 - local max of $P(D|\omega)$
 - reinitialize for closing global local max (if not global)

Strategy 2

- calculate $P(D|(\ell^n)) = \int d\omega P(D|\omega, (\ell^n)) P(\omega)$ explicitly
- sample over (ℓ^n)

$$P(D|\omega, (\ell^n)) = \prod_{a \in \{A, C, G, T\}} p_0(a)^{w(a)} \times \prod_{i=1}^L \prod_{b \in \{A, C, G, T\}} \frac{(w_i(b))}{\cancel{p_0(b)}}$$

Attention: $w(a), w_i(b)$ depends on (ℓ^n) !

With Dirichlet prior $P(\omega) \sim \prod_{i,a} w_i(a)^{\gamma-1}$

$$P(D|(\ell^n)) = \prod_a p_0(a)^{w(a)} \times \prod_{i=1}^L \left[\frac{\Gamma(u_i + \gamma)}{\Gamma(u_i + \gamma_i)} \frac{\prod_b \frac{\Gamma(u_i(a) + b)}{\Gamma(u_i(a))}}{\prod_b \frac{\Gamma(u_i(a) + b)}{\Gamma(u_i(a))}} \right]$$

⇒ complicated (ℓ^n) -dependence

(5)

Gibbs Sampling

- initialize (l^0)



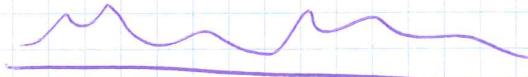
rand.

- select seq. $\star \in \{1, \dots, m\}$ ~~l^0~~

- calculate

$$P(l^{\star} = l \mid D, (l^{-\star}))$$

for each $l = 0, \dots, N_v - L$



- select new value $l^{\star} = l$ with this proba.



+ iteration.

Most probable positions
→ simulated annealing

$$P(D \mid (l^{\star}))^{\beta}$$

$$\beta = 1 \rightarrow +\infty$$