

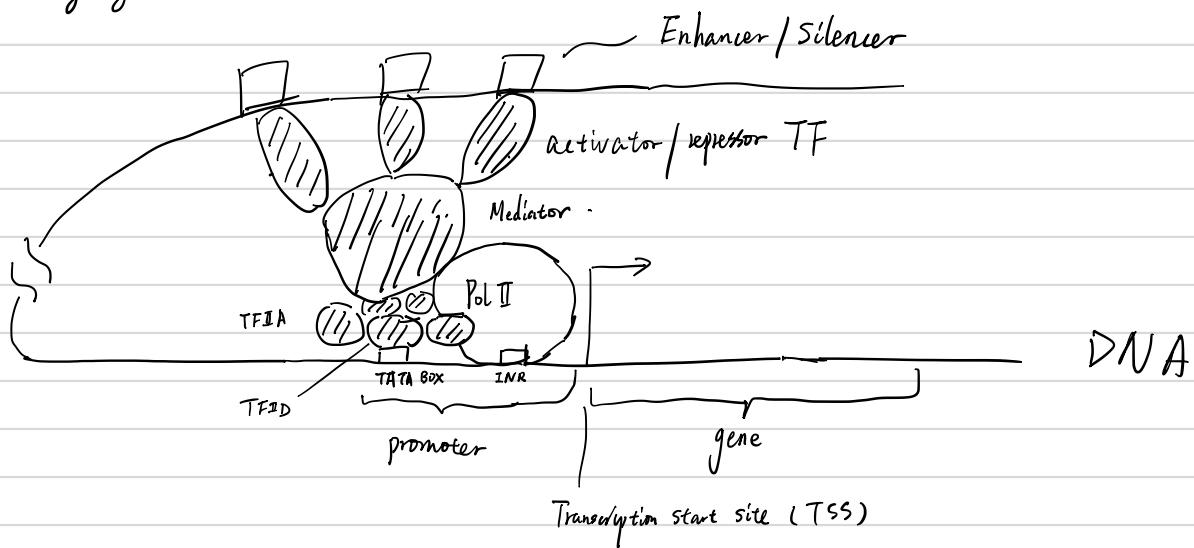
Lecture 6 Mar 7 Characterizing Protein-DNA interactions using ChIP-seq and Motif Disc.

Agenda

- Landscape of gene regulation
- Assay *in vivo* TF binding via ChIP-seq
 - Binding events determination ↴ EM.
 - Motif discovery

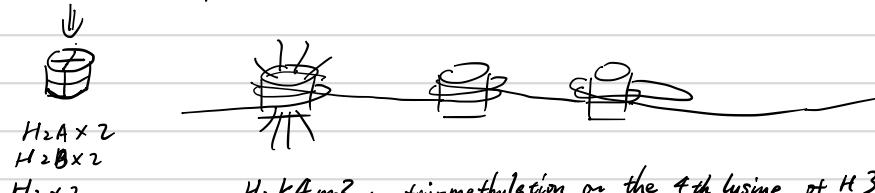
Gene Regulation

- Gene < 1% of the genome
 - Different cells share the same DNA sequence
 - Most of disease-associated mutations reside in non-coding region
- } 99% of the genome regulate the expression of gene in a cell-type specific way.



- Transcription factors
 - regulates gene regulation
 - Binding has sequence specificity

Histone modification

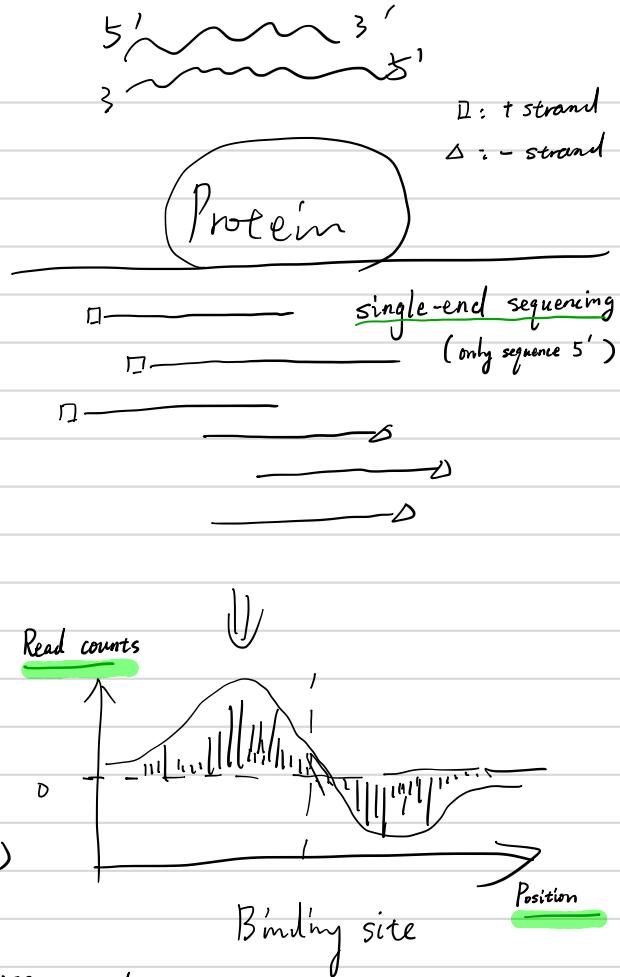
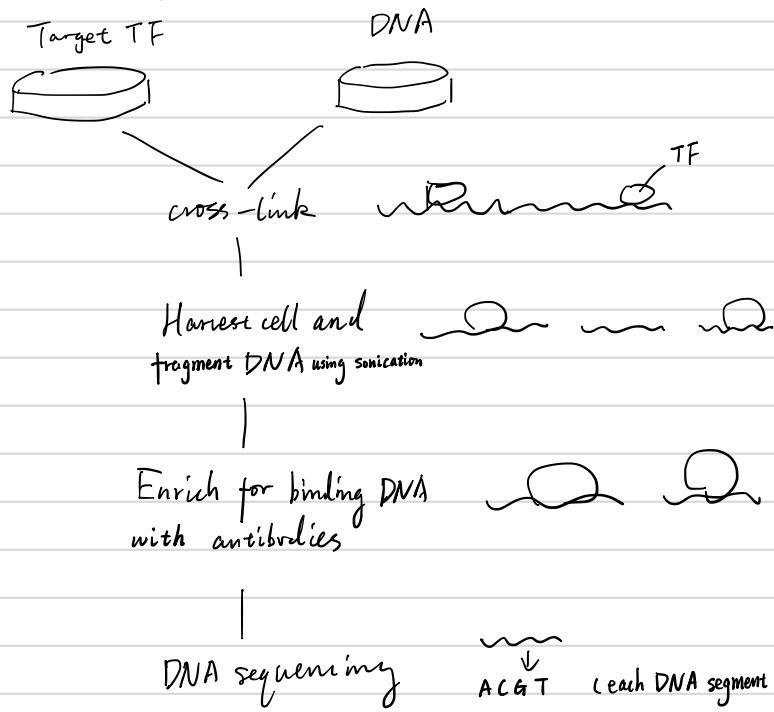


CTCF and cohesins



H₃K4me3 : tri-methylation on the 4th lysine of H3 ⇒ promoters
 H₃K27ac : acetylation on the 27th lysine of H3 ⇒ active regulatory elements

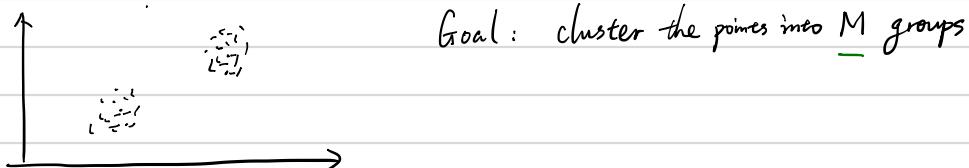
ChIP-seq



Computation problem

- Predict binding events : predict where the TF bind from the noisy ChIP-seq read counts
- Motif discovery : given a set of binding sequences (e.g. 10 base pair around the predicted binding sites), model the sequence specificity (sequence "motif")

Expectation - Maximization (Mixture Model as example)



Assumption : Points in the i th group distribute $\sim N(\mu_i, \Sigma_i) \triangleq N(\theta_i)$
 (Gaussian Mixture)

Notation :

- x_1, x_2, \dots, x_n samples $\in \mathbb{R}^2$
- z_1, z_2, \dots, z_n group assignments $\in \{1, 2, \dots, M\}$, $M = \# \text{ of groups}$
- θ_1, θ_2 gaussian parameters
- π_1, π_2 prior prob. of each group $\pi_k = P(z_i = k), \forall i$

$$\text{likelihood of observe } p(x_i; \theta, \pi) = \sum_{k=1}^M p(z_i=k) p(x_i | z_i=k; \pi, \theta) = \sum_{k=1}^M \pi_k N(x_i; \theta_k)$$

$$\Rightarrow \ln p(x; \theta, \pi) = \ln \prod_{i=1}^n p(x_i; \theta, \pi) = \sum_{i=1}^n \ln \sum_{k=1}^M \pi_k N(x_i; \theta_k)$$

$$\Rightarrow \theta^*, \pi^* = \underset{\theta, \pi}{\operatorname{arg\,max}} \sum_{i=1}^n \ln \sum_{k=1}^M \pi_k N(x_i; \theta_k)$$

hard to optimize

Instead we maximize $\ln p(x, z; \theta, \pi)$

\downarrow unknown! \rightarrow First estimate $p(z|x; \theta, \pi)$ and then use $E_{p(z|x; \theta, \pi)}(\ln p(x, z; \theta, \pi))$ to approximate

$$\text{Expectation step: } r(z) \triangleq p(z|x; \theta, \pi) = \prod_{i=1}^n p(z_i|x_i; \theta, \pi)$$

$$r(z_i=k) \triangleq p(z_i=k|x_i; \theta, \pi) \approx p(x_i|z_i=k; \theta, \pi) p(z_i=k) = N(x_i, \theta_k) \pi_k$$

(Bayes' Rule)

$$\Rightarrow r(z_i=k) = \frac{N(x_i, \theta_k) \pi_k}{\sum_k N(x_i, \theta_k) \pi_k} \Rightarrow (\text{soft assignment of point } i \text{ to group } k)$$

Maximization step:

$$\theta^*, \pi^* = \underset{\theta, \pi}{\operatorname{argmax}} E_{p(z|x; \theta, \pi)}(\ln p(x, z; \theta, \pi))$$

$$= \underset{\theta, \pi}{\operatorname{argmax}} \sum_z r(z) \underbrace{\ln p(x, z; \theta, \pi)}_{\sum_i} \Rightarrow \prod_{i=1}^n N(x_i, \theta_{z_i}) \pi_{z_i}$$

$$= \underset{\theta, \pi}{\operatorname{argmax}} \sum_z r(z) \underbrace{\sum_{i=1}^n [\ln N(x_i, \theta_{z_i}) + \ln \pi_{z_i}]}_L$$

$$\text{Set } \frac{\partial L}{\partial \mu} = \frac{\partial L}{\partial \Sigma} = \frac{\partial L}{\partial \pi} = 0 \Rightarrow \left\{ \begin{array}{l} \mu_k^{\text{new}} = \frac{1}{N_k} \sum_{i=1}^n r(z_i=k) x_i \\ \Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{i=1}^n r(z_i=k) (x_i - \mu_k)(x_i - \mu_k)^T \\ \pi_k^{\text{new}} = \frac{N_k}{n} \\ N_k = \sum_{i=1}^n r(z_i=k) \end{array} \right. \begin{array}{l} \text{(total # of reads assigned)} \\ \text{to } k\text{-th group} \end{array}$$

Binding event prediction

Genome Positioning System (GPS) Guo et al 2010 Bioinformatics

$$\begin{cases} \chi \rightarrow \text{the position of each read} \\ \zeta \rightarrow \text{the assignment of each read to bind events} \end{cases}$$

Special design:

- The property of ChIP-seq: ① Binding events generate reads in a range \Rightarrow segment the genome and perform EM separately
 \Rightarrow in each segment, there is one event on every base pair
- ② Shape of event is the same for any event \Rightarrow Instead of $N(\mu_i, \sigma_i)$ for each event, use emp. distn estimated from genome-wide rough peak-calling (using other software)
 \Rightarrow After performing EM on each segment separately, re-estimate using the updated event position
- We want π_k to be sparse
 \rightarrow negative Dirichlet prior on π_k $P(\pi) \sim \prod_{j=1}^M \frac{1}{(\pi_j)^{\alpha}}, \alpha > 0$
 $\rightarrow \pi_k^{\text{new}} = \frac{N_k - \alpha}{\sum_{j=1}^M (N_j - \alpha)}$ component elimination $\Rightarrow \pi_k^{\text{new}} = \frac{\max(0, N_k - \alpha)}{\sum_{j=1}^M \max(0, N_j - \alpha)}, \alpha \Rightarrow \text{the minimum # of reads an event needs to survive EM iterations}$

Motif discovery

- PWM (position weight matrices)

	1	2	3	4
A	0.1	0.4	0.1	0.1
C	0.3	0.1	0.1	0.4
G	0.4	0.2	0.7	0.4
T	0.2	0.3	0.1	0.1

- Aligned seqs \rightarrow PWM is easy

$$\begin{array}{l} \text{A A C T} \\ \text{A G C T} \\ \text{C T C A} \end{array} \Rightarrow \begin{array}{l} \text{A } 0.66 \ 0.33 \ 0 \ 0.33 \\ \text{C } 0.33 \ 0 \ 1 \ 0 \\ \text{G } 0 \ 0.33 \ 0 \ 0 \\ \text{T } 0 \ 0.33 \ 0 \ 0.66 \end{array}$$

Goal:

Given a set of binding sequences (same length, not aligned), find a PWM (of length W) that describe the binding motif

$\{ \chi \rightarrow \text{sequences} \}$

$\{ \zeta \rightarrow \text{the starting position of each seq} \}$

$p(X_i | \zeta_i = k)$ is parameterized by motif PBM $M \in \mathbb{R}^{4 \times W}$ and background PBM $B \in \mathbb{R}^{4 \times 1}$

$$\Rightarrow p(X_i | \zeta_i = k; M, B) = \prod_{j=1}^{k-1} B(X_i^j) \prod_{j=k}^{k+W} M(j-k+1, X_i^j) \prod_{j=k+W}^{S(X_i)} B(X_i^j) \quad (1)$$

↓
The prob. that the motif sequence in X_i starts at the k th position

$$\text{E step: } p(\zeta | X; M, B) = \prod_{i=1}^n p(\zeta_i | X_i; M, B) \stackrel{?}{=} \prod_{i=1}^n Y(\zeta_i)$$

$$\Rightarrow Y(\zeta_i = k) \sim p(X_i | \zeta_i = k; M, B) p(\zeta_i = k)$$

\Rightarrow If we assume $p(\zeta_i)$ is uniform (all the starting position are equally likely)

$$Y(\zeta_i = k) = \frac{\prod_{j=1}^{k-1} B + \prod_{j=k}^{k+W-1} M + \prod_{j=k+W}^{S(X_i)} B}{\sum_{k=1}^{S(X_i)-L+1} \prod_{j=1}^{k-1} B \prod_{j=k}^{k+W-1} M \prod_{j=k+W}^{S(X_i)} B}$$

$$M \text{ step: } M^*, B^* = \underset{M, B}{\operatorname{argmax}} E_{P(Z|X; M, B)} \left(\ln P(X, Z; M, B) \right) + \sum_{l=1}^W \lambda_l (\sum_c M(l, c) - 1)$$

\downarrow

$$p(X|Z) p(Z) = \pi_B \pi_M \pi_B \times p(Z)$$

$$\frac{\partial L}{\partial M} = \frac{\partial L}{\partial B} = 0$$

For all sequences and all the possible starting position, sum the prob. that the k-th nucleotide in the motif is C

$$\Rightarrow M^{\text{new}}(k, c) = \frac{n_{k,c} + d}{\sum_c (n_{k,c} + d)}, \quad n_{k,c} = \sum_i \sum_{\{Z_i | X_i^{z_i+k-1} = c\}} P(Z_i | X_i), \quad d \text{ is a pseudo-count}$$

reflects how much we believe a priori that PWM should be uniform

$$B^{\text{new}}(c) = \frac{n_{0,c} + d}{\sum_c (n_{0,c} + d)}, \quad n_{0,c} = n_c - \sum_k n_{k,c}$$

Example of M step:

Assume $W=3$, after E step we have the following $p(Z|X)$ (prob. of motif starting position in each sequence)

$$X_1 = A \ C \ A \ G \ C \ A \quad X_2 = A \ G \ G \ C \ A \ G \quad X_3 = T \ C \ A \ G \ T \ C$$

$$p(Z_1|X_1) \ 0.1 \ 0.1 \ 0.1 \ 0.1 \quad p(Z_2|X_2) \ 0.4 \ 0.1 \ 0.1 \ 0.4 \quad p(Z_3|X_3) \ 0.2 \ 0.6 \ 0.1 \ 0.1$$

\Rightarrow M step

$$\begin{aligned} n_{1,A} &= 0.1 + 0.1 + 0.4 + 0.1 = 0.7 \\ n_{1,C} &= 0.7 + 0.4 + 0.6 = 1.7 \\ n_{1,G} &= 0.4 \\ n_{1,T} &= 0.2 \end{aligned} \quad \left. \begin{array}{l} \sum_c n_{k,c} = 3 \end{array} \right.$$

$$\Rightarrow \text{Let } d=1, \quad M(1, A) = \frac{0.7+1}{3+4} = 0.24$$

$$M(1, C) = (1.7+1) / (3+4) = 0.39$$

$$M(1, G) = (0.4+1) / (3+4) = 0.2$$

$$M(1, T) = (0.2+1) / (3+4) = 0.17$$

Similarly we can calculate the rest of M^{new}

	1	2	3	
$M(k, c)$	A	0.24	0.39	0.21
	C	0.39	0.21	0.18
	G	0.2	0.24	0.44
	T	0.17	0.16	0.16