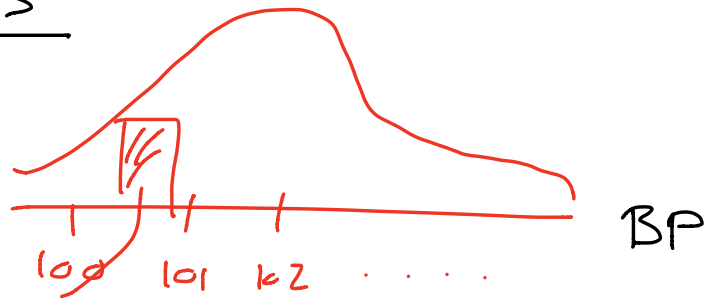


Cutflinks

Notation:



$$F(l_0) \quad \sum_{i=1}^{\infty} F(i) = 1$$

Our data



Task: Given this genome we try
and figure out the attribution
to some transcript assignments

$$P = \{P_{t_1}, \dots, P_{t_{|T|}}\}$$

"abundances of transcripts"

$$\begin{matrix} r_1 \\ \vdots \\ r_n \end{matrix} \left| \begin{matrix} t_1 & \dots & t_{|T|} \end{matrix} \right.$$

$$H_{RT} = \begin{matrix} & i & \\ r_R & \left| \begin{array}{c} \square \end{array} \right. & \\ & \text{---} \end{matrix}$$

1 if fragment is completely
contained in t

0 o.w.

• $I_t(r)$ implied length in
fragment t

• f_i = i th frag. alignment.

$$L(g | R) = \prod_{i=1}^R P(f_i = r_i)$$

$$= \prod_{i=1}^R \sum_{t \in T} P(f_i = r_i | T_i = t) P(T_i = t)$$

$$= \prod_{i=1}^R \sum_{t \in T} P(f_i = r_i | T_i = t) P_t \tilde{I}(t)$$

$$= \prod_{i=1}^R \sum_{t \in T} \frac{F(i) (l(t) - l_t(i) + 1)}{\sum_{u \in T} p_u \tilde{l}(u)}$$

with $\tilde{l}(t_i) = \sum_{i=1}^{\infty} F(i) (l(t_i) - i + 1)$

called the "adjusted length"

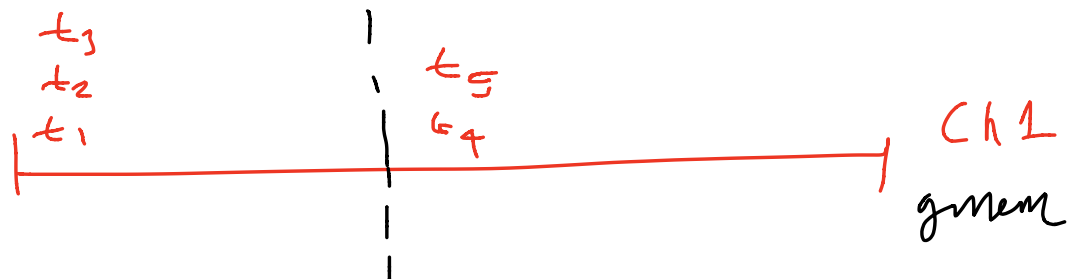
$$= \prod_{i=1}^R \sum_{t \in T} \frac{F(l_t(r_i))}{l(t) - l_t(r_i) + 1} \alpha_t$$

need to estimate

We could optimize but $|T| = 60,000$

Too big to optimize.

Idea partition the genome.



Introduce L_i a region where f_i
 corresponds to and $X_G = \#$
 of fragments falling into region G_j .

$$\begin{aligned}
 \beta_g &\equiv P(L_i = g) = \sum_{t \in G_g} \alpha_t \\
 &= \sum_{t \in G_g} \sigma_t \tau_t \tilde{l}(t) \\
 &= \frac{\sum_{n=1}^{|G|} \sum_{u \in G_n} \sigma_n \tau_u \tilde{l}(u)}{\sum_{n=1}^{|G|} \sigma_n \tilde{l}(G_n)} \\
 &= \frac{\sigma_g \tilde{l}(G_g)}{\sum_{n=1}^{|G|} \sigma_n \tilde{l}(G_n)}
 \end{aligned}$$

$$\sigma_g = \sum_{t \in G_g} p_t \qquad \tau_t = \frac{p_t}{\sum p_u}$$

$$u \in G_g$$

So the full likelihood

$$L(p|R) = \prod_{i=1}^n P(f_i = r_i | L_i = g) \beta_g I(r_i \in G_g)$$

$$= \left\{ \prod_{g=1}^{|G|} \prod_{r_i \in G_g} \sum_{t \in G_g} \gamma_t \frac{F(L_t(r_i))}{L(t) - I_t(r_i) + 1} \right\}$$

$$\propto \prod_{g=1}^{|G|} \beta_g^{X_g}$$

So $\vec{\beta} = \frac{X_g}{R}$ γ_t can be

found via constrained optimization.

• We even have variance estimate

• Unstable

— identifiability

- Use importance sampling from likelihood. for mean & variance estimate. $\psi \sim \text{Gov}(\gamma)$.

- Using this we can define the FPKM statistic with its corresp. Variance.

Testing Diff. Expression

$$\log \left(\frac{x_g^a r_t^a R^b}{x_g^b \hat{\gamma}_t^a R^a} \right) \quad \text{and we can estimate the variance as well.}$$

Then the test stat. is

$$\sqrt{n} \frac{\log(\hat{\text{ratio}}) - 0}{\quad} \sim N(0, 1)$$

test var

- HW: how can we test differences between multiple groups.
- Differential gene transcript expression analysis of ...
- Question HW: how do we do multiple testing in Cufflinks.