

Lecture 7: Expectation-Maximization (EM) Algorithm

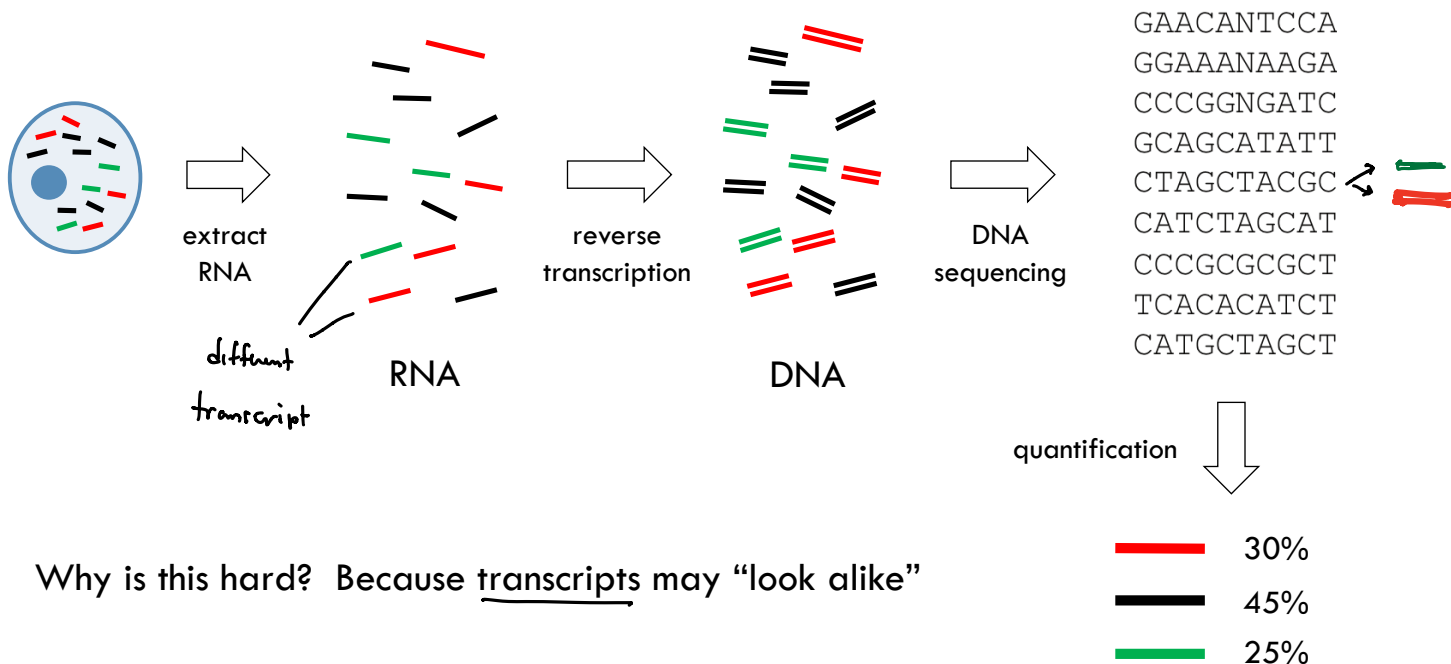


ECE 365

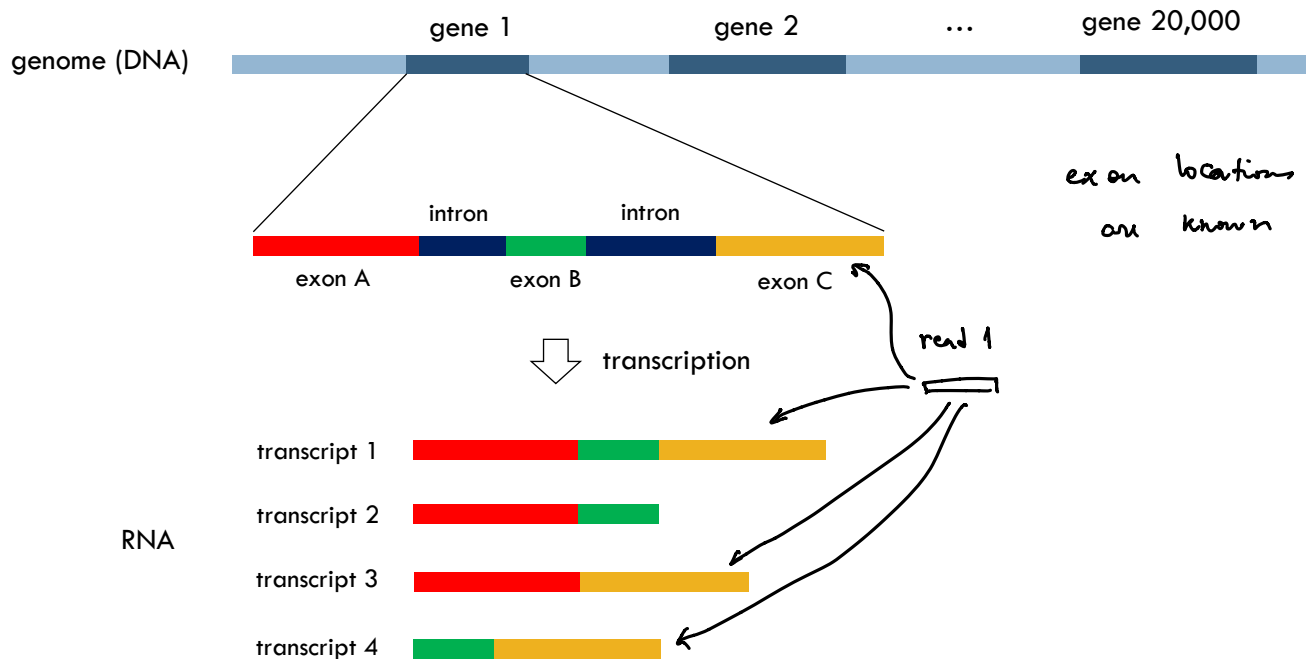
Announcements:

- Lab 3 due on Friday
- Lab 4 released tonight (due April 2)

From previous lecture: RNA quantification



From genes to transcripts



General RNA quantification problem

□ Input: **aligned read data**

	tr. 1	tr. 2	...	tr. K		
read 1	0	1	1	0	0	0
read 2	0	0	0	1	0	0
	0	1	1	0	0	0
⋮	0	0	0	0	1	0
	1	0	0	0	0	0
	1	0	0	0	1	0
read N	0	1	0	0	0	1

means: read N maps to tr. 2

relative abundances

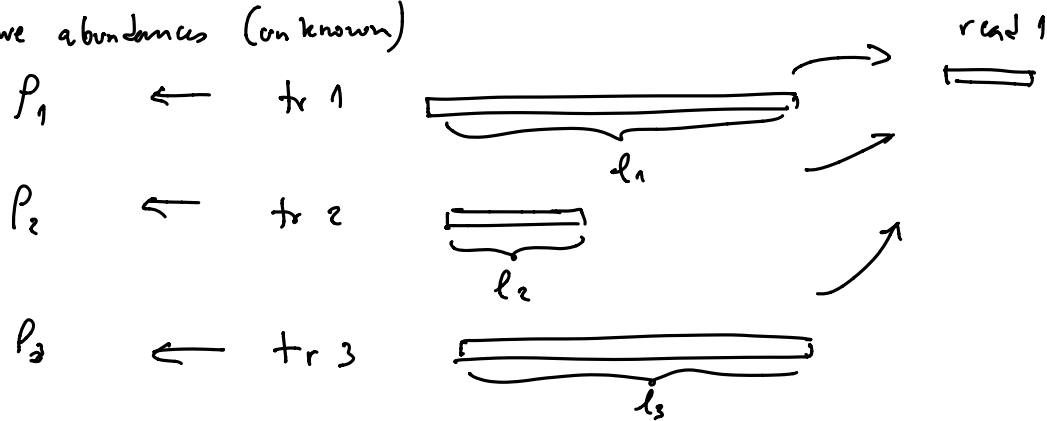
$$\Rightarrow \left. \begin{array}{l} \hat{p}_1 = 0.01 \\ \hat{p}_2 = 0.04 \\ \vdots \\ \hat{p}_K = 0.02 \end{array} \right\} \text{add up to 1}$$

General approach to solve this:

Expectation - Maximization Algorithm (EM)

Model for read data generation:

relative abundances (unknown)



$$\theta_k = P(\text{read comes from tr. } k) = \frac{p_k l_k}{\sum_{j=1}^K p_j l_j}$$

True read origin/assignment (unknown)

$$Z_{ik} = \begin{cases} 1 & \text{if read } i \text{ comes from transcript } k \\ 0 & \text{otherwise} \end{cases}$$

form a matrix

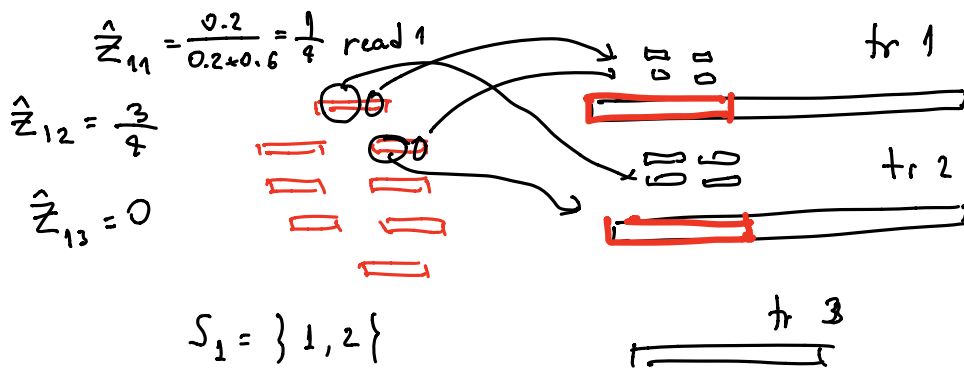
$$N \text{ reads} \left\{ \begin{matrix} \text{K transcripts} \\ \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix} \end{matrix} \right. \begin{matrix} \text{(one 1 per row)} \\ \Rightarrow \text{"cleaned up" version} \\ \text{of data matrix} \end{matrix}$$

EM strategy: Estimate \hat{p}_k and \hat{Z}_{ik} in an alternating fashion and recursively

Initialize : $\hat{p}_k = \frac{1}{K}$ for $k = 1, \dots, K$

(*) Based on \hat{p}_k 's, estimate \hat{z}_{ik} 's :

in the middle
of the
algorithm



$$\hat{p}_1 = 0.2$$

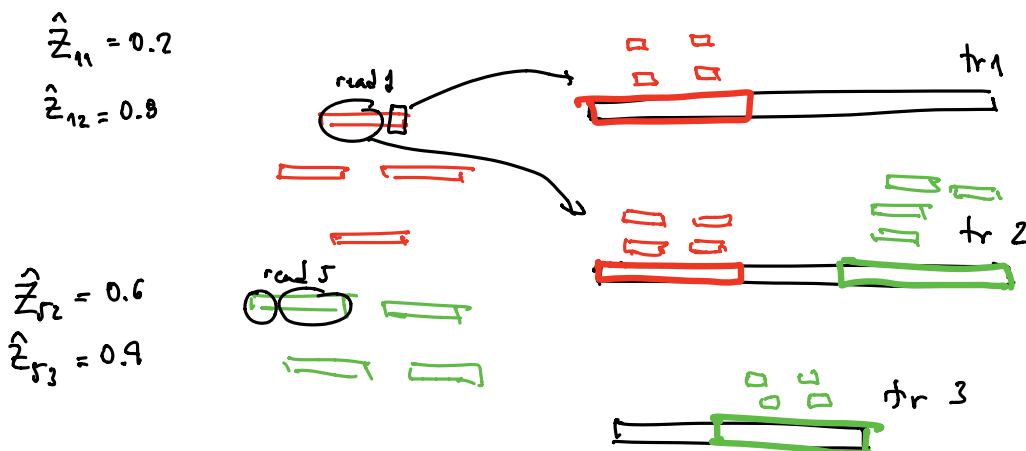
$$\hat{p}_2 = 0.6$$

$$\hat{p}_3 = 0.2$$

Let S_i be the set of transcripts read i maps to

$$\hat{z}_{ik} = \begin{cases} \frac{\hat{p}_k}{\sum_{j \in S_i} \hat{p}_j} & \text{if } k \in S_i \\ 0 & \text{if } k \notin S_i \end{cases} \quad \text{E-step}$$

Based on \hat{z}_{ik} , we want to update \hat{p}_k 's :



$$\hat{\theta}_1 = \frac{0.2 + 0.2 + 0.2 + 0.2}{8} = 0.1$$

$$\hat{\theta}_2 = \frac{4 \times 0.8 + 4 \times 0.6}{8}$$

$$\hat{\theta}_3 = \frac{4 \times 0.4}{8}$$

First, estimate probability θ_k :

$$\hat{\theta}_k = \frac{1}{N} \sum_{i=1}^N z_{ik}, \quad \theta_k \propto p_k l_k$$

$$\hat{p}_k = \frac{\frac{\hat{\theta}_k}{l_k}}{\sum_{j=1}^K \frac{\hat{\theta}_j}{l_j}} \quad \left. \vphantom{\frac{\hat{\theta}_k}{l_k}} \right] \text{M-step}$$

Repeat. Go back to $\textcircled{*}$

EM algorithm for RNA quantification

□ Initialize: $\hat{\rho}_k^{(1)} = \frac{1}{K}$

□ For $t = 1, 2, \dots$,

E - step □ $\hat{Z}_{ik}^{(t)} = \begin{cases} \frac{\hat{\rho}_k^{(t)}}{\sum_{j \in S_i} \hat{\rho}_j^{(t)}} & \text{for } k \in S_i \\ 0 & \text{otherwise} \end{cases}$

M - step □ $\hat{\rho}_k^{(t+1)} = \frac{\frac{\theta_k^{(t+1)}}{\ell_k}}{\sum_{j=1}^K \frac{\theta_j^{(t+1)}}{\ell_j}} \quad \text{where} \quad \theta_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N \hat{Z}_{ik}^{(t)}$

EM algorithm for RNA quantification

- Does it converge? Yes .
- If so, does it converge to the correct values?

For RNA quant. problem, $\hat{p}_1, \dots, \hat{p}_K$

converge to the maximum likelihood parameters .

- Let's look at a concrete example in Excel