

Lecture 10: Clustering Single-Cell RNA-seq data

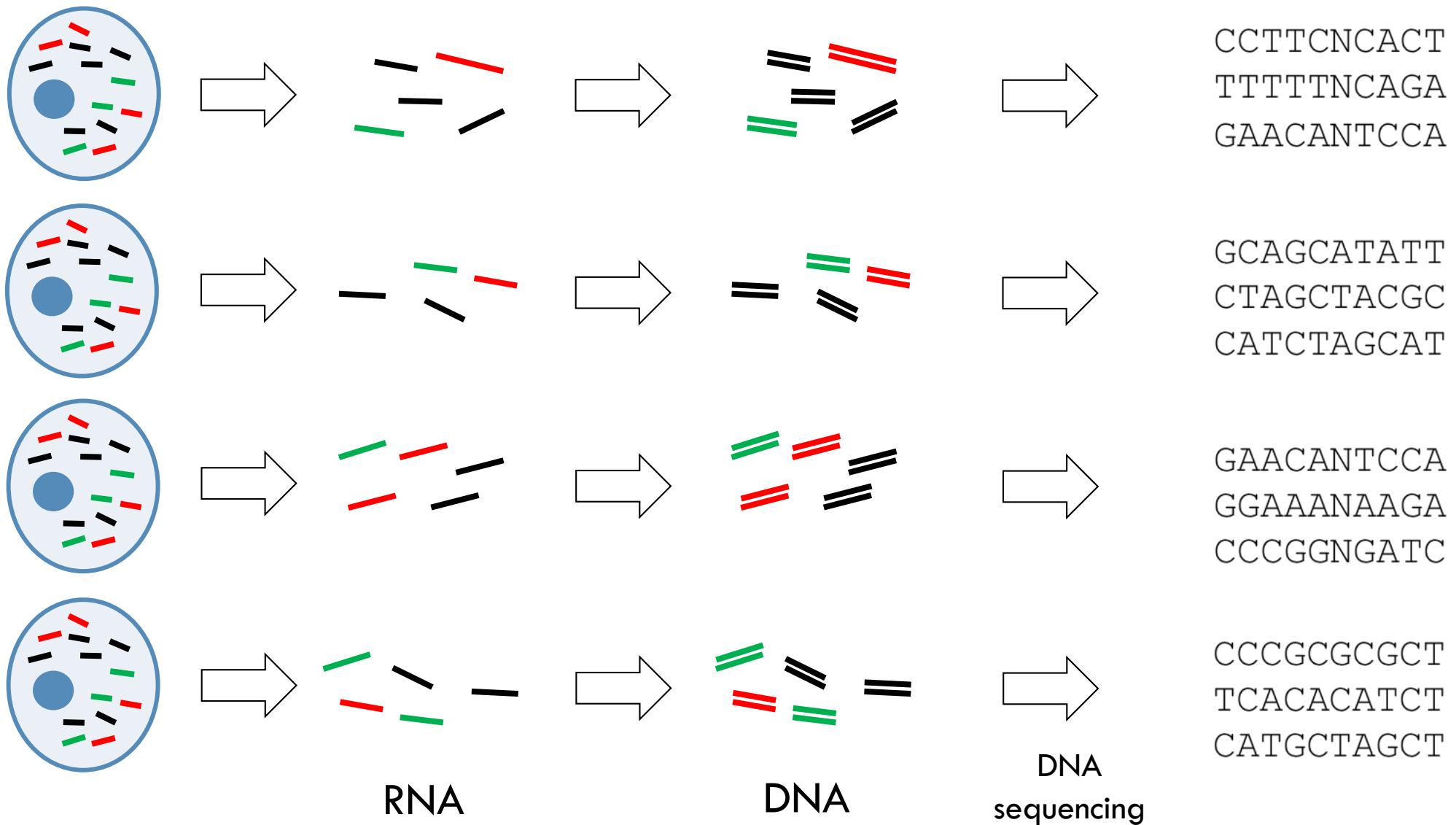


ECE 365

Announcements:

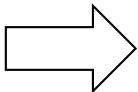
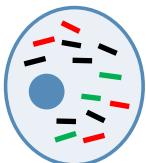
- Lab 4 (RNA-seq) due tomorrow
- Quiz 2 today

From previous lecture: Single-Cell RNA-seq

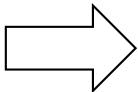
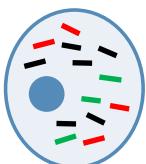


From previous lecture: Single-Cell RNA-seq

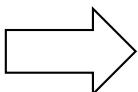
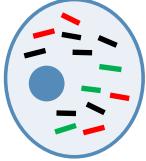
- We obtain reads for each cell



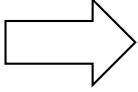
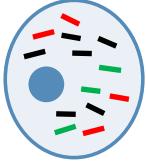
GCAGCATATT
CTAGCTACGC
CATCTAGCAT



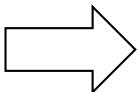
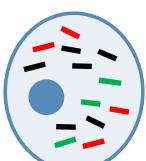
GAACANTCCA
GGAAANAAGA
CCCGGNGATC



CCCGCGCGCT
TCACACATCT
CATGCTAGCT



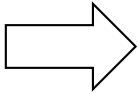
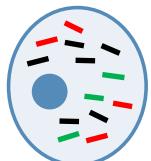
GAACANTCCA
GGAAANAAGA
CCCGGNGATC



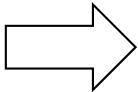
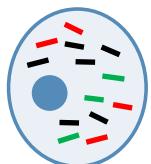
CCTTCNCACT
TTTTNCAGA
GAACANTCCA

From previous lecture: Single-Cell RNA-seq

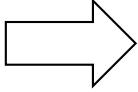
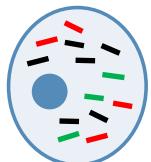
- We obtain reads for each cell



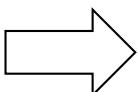
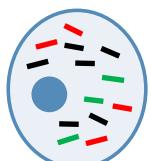
GCAGCATATT
CTAGCTACGC
CATCTAGCAT



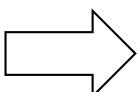
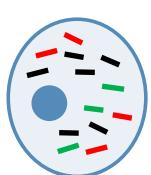
GAACANTCCA
GGAAANAAGA
CCCGGNGATC



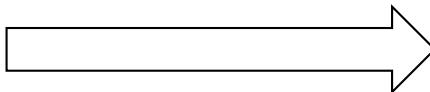
CCCGCGCGCT
TCACACATCT
CATGCTAGCT



GAACANTCCA
GGAAANAAGA
CCCGGNGATC



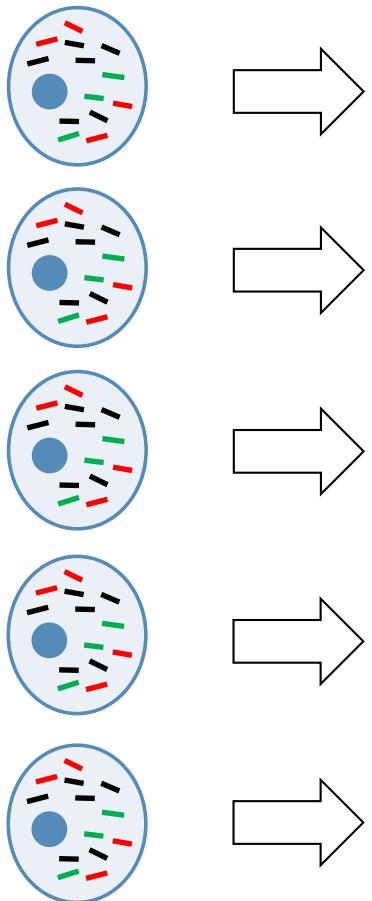
CCTTCNCACT
TTTTTNCAGA
GAACANTCCA



align to genes
(or transcripts)

From previous lecture: Single-Cell RNA-seq

- We obtain reads for each cell



GCAGCATATT
CTAGCTACGC
CATCTAGCAT

GAACANTCCA
GGAAANAAGA
CCCGGNGATC

CCCGCGCGCT
TCACACATCT
CATGCTAGCT

GAACANTCCA
GGAAANAAGA
CCCGGNGATC

CCTTCNCACT
TTTTTNCAGA
GAACANTCCA

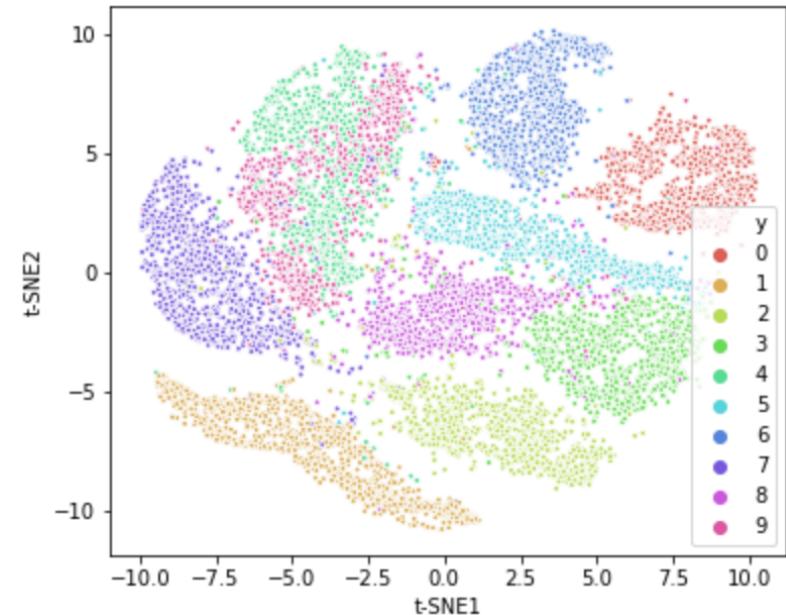
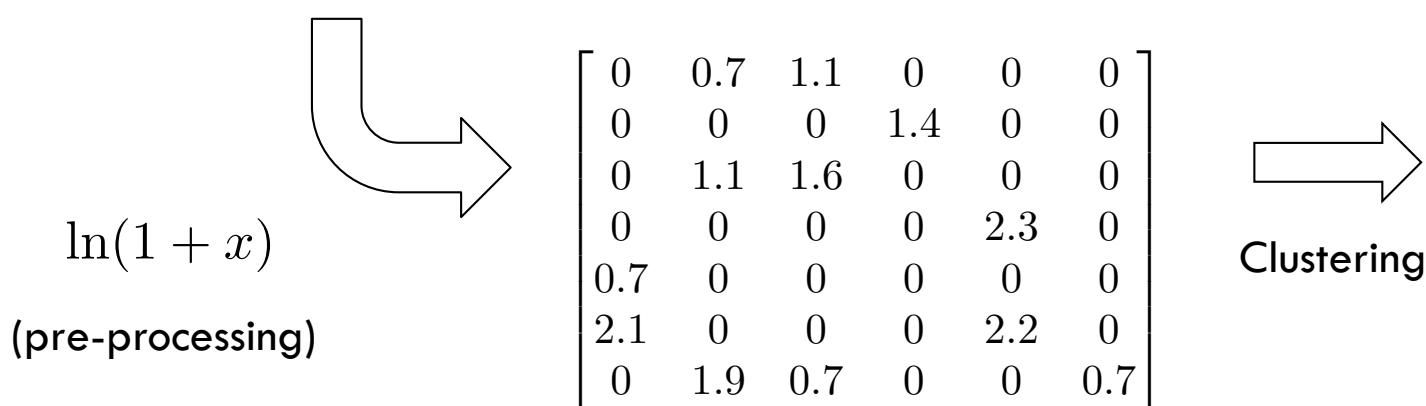
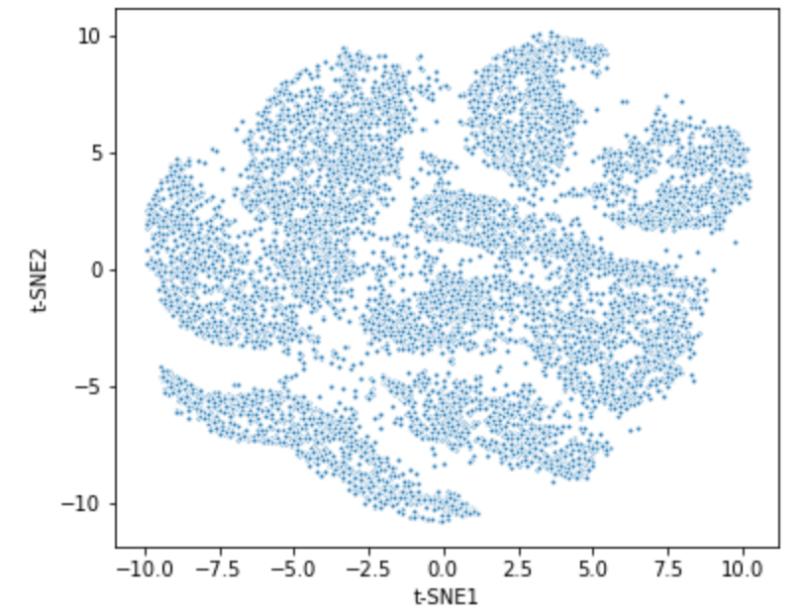
align to genes →

	gene 1	...	gene m			
cell 1	0	1	2	0	0	0
cell 2	0	0	0	3	0	0
⋮	0	2	4	0	0	0
cell n	1	0	0	0	9	0
	7	0	0	0	8	0
	0	6	1	0	0	1

Single-Cell RNA-seq Data Analysis

	gene 1	...		gene <i>m</i>		
cell 1	0	1	2	0	0	0
cell 2	0	0	0	3	0	0
:	0	2	4	0	0	0
	0	0	0	0	9	0
	1	0	0	0	0	0
	7	0	0	0	8	0
cell <i>n</i>	0	6	1	0	0	1

Dimensionality reduction
and visualization



Dimensionality reduction via PCA

Data matrix \mathbf{X} ($n \times m$)

Perform SVD on \mathbf{X} : $\mathbf{X} = \mathbf{U}_{n \times n} \mathbf{S}_{n \times m} \mathbf{V}^T_{m \times m}$

use this to
get transformation \mathbf{W}

$$\mathbf{V}^T_{m \times m}$$

$$\begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots & \\ & & & \sigma_m \end{bmatrix}$$

singular values

$$\sigma_1 > \sigma_2 > \dots > \sigma_m$$

$$\mathbf{V}^T = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_m \end{bmatrix} \quad \text{pick top } k \text{ rows}$$

$$\rightarrow \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_k \end{bmatrix} = \mathbf{W}_{k \times m}$$

now we can find $\tilde{\mathbf{x}}_i = \mathbf{W} \mathbf{x}_i$

$$\tilde{\mathbf{X}}^T = \mathbf{W} \mathbf{X}^T$$

Non-linear dimensionality reduction via t-SNE

- ***t*-Distributed Stochastic Neighbor Embedding (*t*-SNE)**

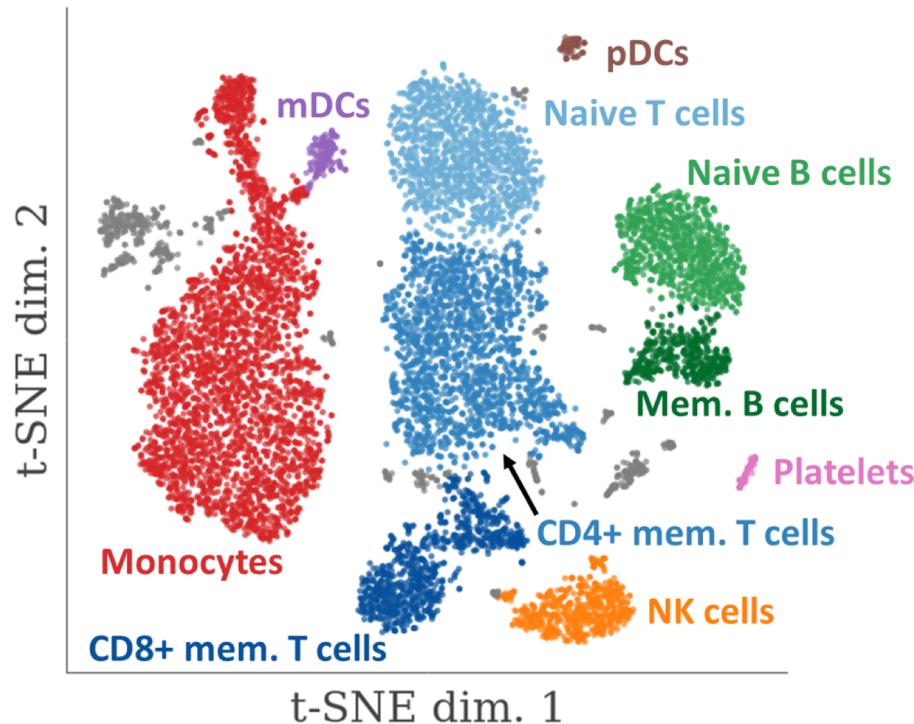
- Popular *non-linear* dimensionality reduction method
 - High-level idea:



- Use gradient descent to move y_i 's to minimize $D(P\|Q) = \sum_{i\neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$

Non-linear dimensionality reduction via t-SNE

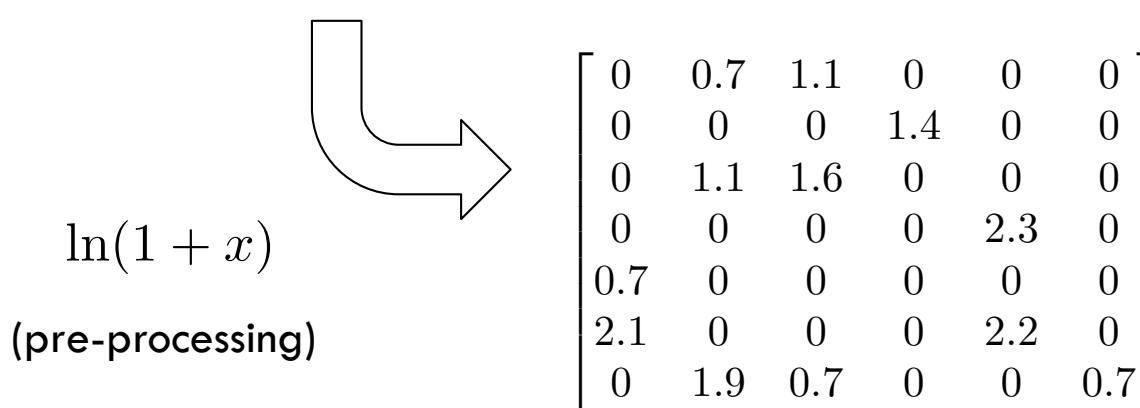
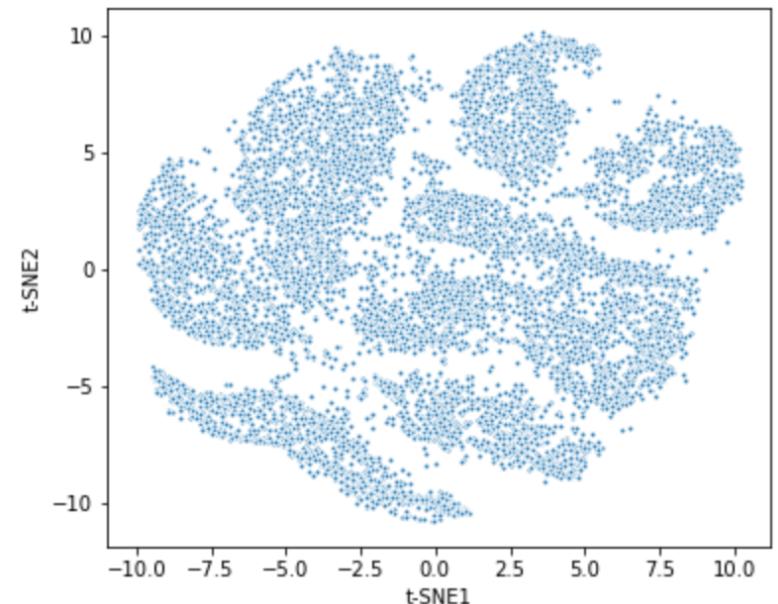
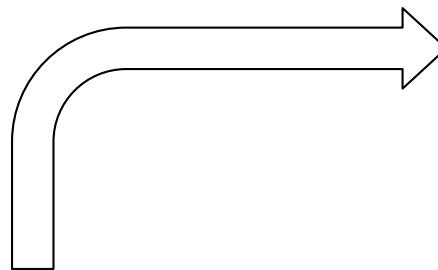
- Very popular for single-cell RNA-seq data



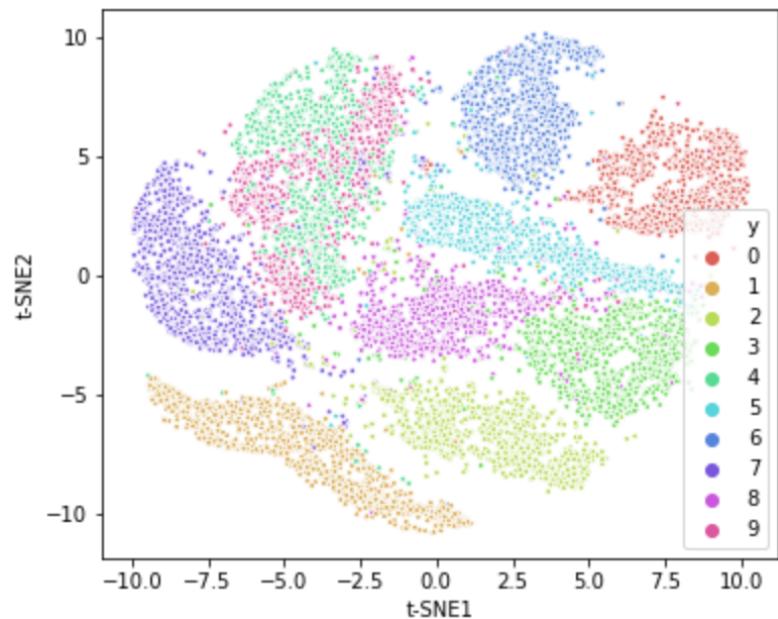
Single-Cell RNA-seq Data Analysis

	gene 1	...	gene m			
cell 1	0	1	2	0	0	0
cell 2	0	0	0	3	0	0
:	0	2	4	0	0	0
	0	0	0	0	9	0
	1	0	0	0	0	0
	7	0	0	0	8	0
cell n	0	6	1	0	0	1

Dimensionality reduction
and visualization



Clustering

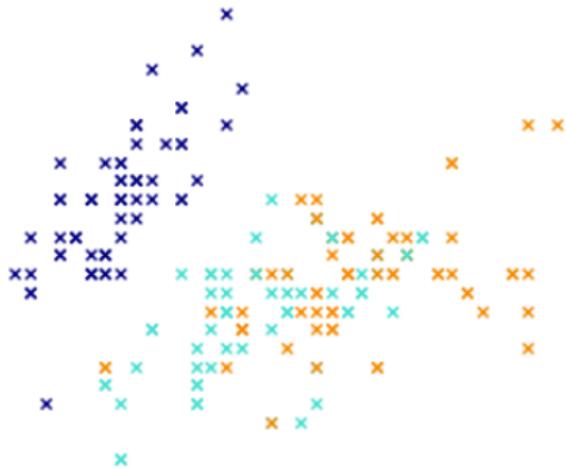


From first module: K-means clustering

- Data points: $x_i \in R^d$ for $i = 1, \dots, N$
- Initialize means μ_ℓ for $\ell = 1, \dots, K$ to K random points
- Repeat until convergence:
 - Assign each x_i to closest mean
 - After all x_i 's are assigned, recompute μ_ℓ for $\ell = 1, \dots, K$ for each cluster

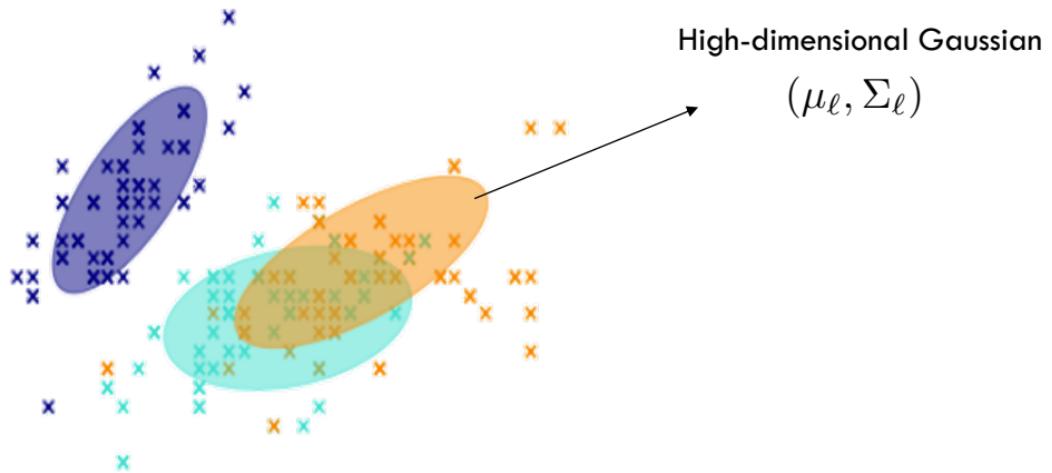
Gaussian Mixture Model

- Generative model: multiple high-dimensional Gaussian distributions



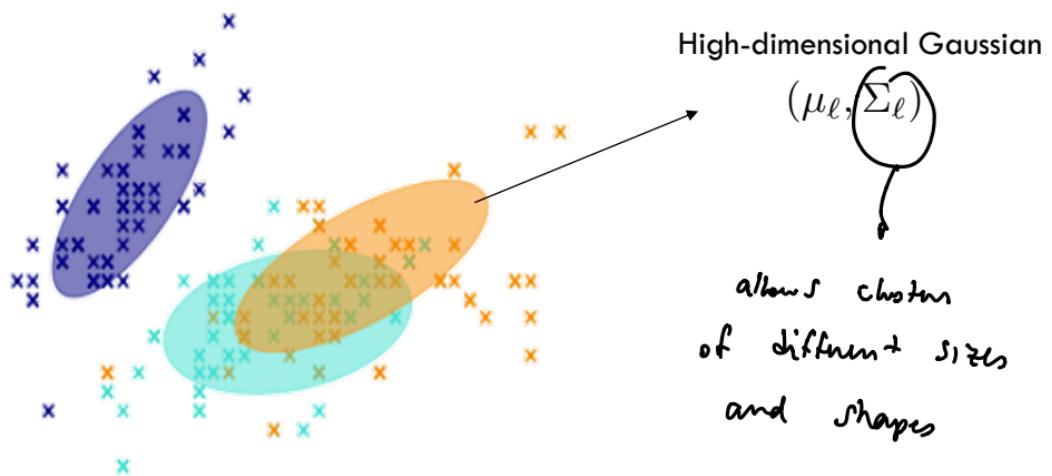
Gaussian Mixture Model

- Generative model: multiple high-dimensional Gaussian distributions



Gaussian Mixture Model

- Generative model: multiple high-dimensional Gaussian distributions



- Model parameters are estimated with the EM algorithm!

Recall: EM algorithm for RNA quantification

- Initialize: $\rho_j = \frac{1}{K}$ for $j = 1, \dots, K$

- Repeat:

E-step

$$Z_{ik} = \begin{cases} \frac{\rho_k}{\sum_{j \in S_i} \rho_j}, & \text{for } k \in S_i \\ 0, & \text{otherwise} \end{cases}$$

↑
membership variable

M-step

$$\rho_k = \frac{\theta_k}{\sum_{j=1}^K \theta_j}, \text{ where } \theta_k = \frac{1}{N} \sum_{i=1}^N Z_{ik}$$

$\theta_k \rightarrow$ length transcript k
 \hookrightarrow prob. of getting a read from tr. k

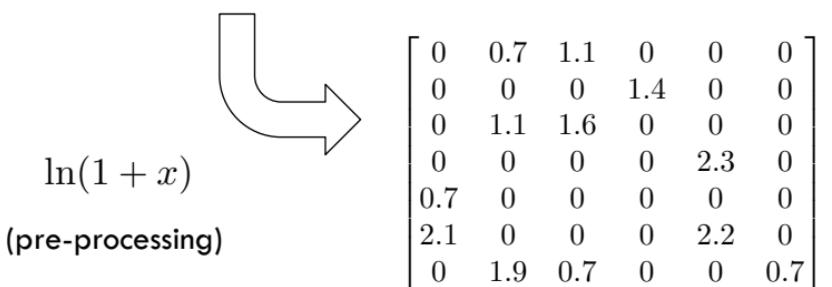
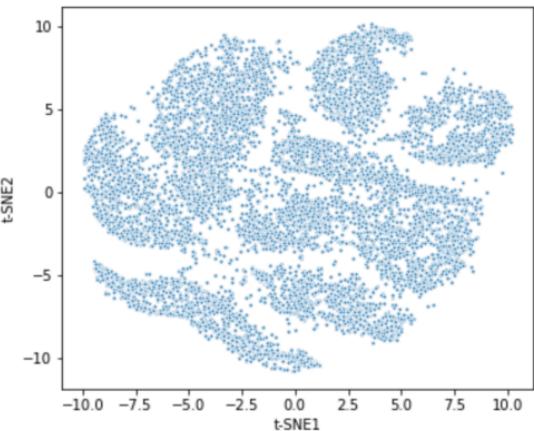
EM algorithm for GMM fitting

- Initialize: (μ_ℓ, Σ_ℓ) and mixture coefficients $\theta_\ell = \frac{1}{K}$ for $\ell = 1, \dots, K$
fraction of each Gaussian
- Repeat:
 - $Z_{ik} = \frac{\mathcal{N}(x_i | \mu_k, \Sigma_k) \cdot \rho_k}{\sum_{j=1}^K \mathcal{N}(x_i | \mu_j, \Sigma_j) \cdot \rho_j}$
membership variable
 - $\theta_k = \frac{1}{N} \sum_{i=1}^N Z_{ik}, \quad \mu_k = \frac{\sum_{i=1}^N Z_{ik} x_i}{\sum_{i=1}^N Z_{ik}}, \quad \Sigma_k = \frac{\sum_{i=1}^N Z_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N Z_{ik}}$

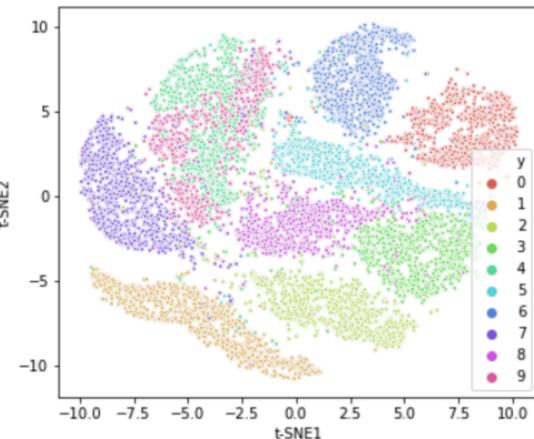
Single-Cell RNA-seq Data Analysis

	gene 1	...	gene <i>m</i>			
cell 1	0	1	2	0	0	0
cell 2	0	0	0	3	0	0
⋮	0	2	4	0	0	0
⋮	0	0	0	0	9	0
⋮	1	0	0	0	0	0
⋮	7	0	0	0	8	0
⋮	0	6	1	0	0	1

Dimensionality reduction
and visualization



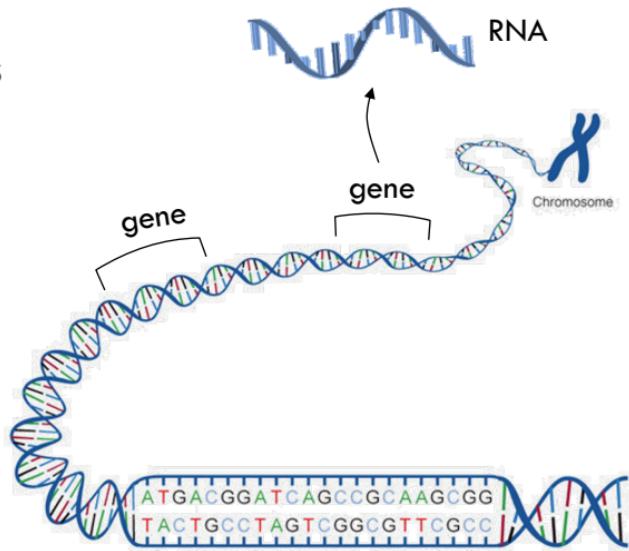
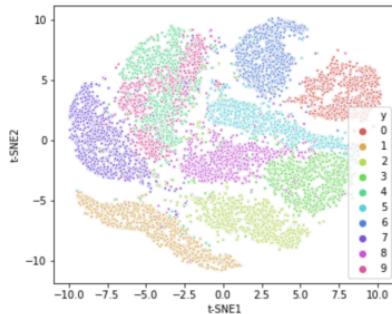
Clustering



Wrapping up

- First part: DNA sequencing data analysis
 - Sequence alignment (dynamic programming, indexing)
 - Genome-wide association studies (GWAS)

- Second part: RNA sequencing data analysis
 - RNA quantification (EM algorithm)
 - Single-cell RNA-seq (dim. reduction, clustering)



Quick review: data types

Genotype data		Bulk RNA-seq (transcript quantification)					Single-cell RNA-seq				
		SNP 1	...	SNP m	tr. 1	...	tr. K	gene 1	...	gene m	
person 1		$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 2 & 1 & 0 & 2 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 1 & 0 \end{bmatrix}$	read 1	$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$	cell 1	$\begin{bmatrix} 0 & 1 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 2 & 4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 9 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 7 & 0 & 0 & 0 & 8 & 0 \end{bmatrix}$					
person n		$\begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 2 \end{bmatrix}$	read 2	\vdots	read N	$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}$	cell 2	\vdots	cell n	$\begin{bmatrix} 0 & 6 & 1 & 0 & 0 & 1 \end{bmatrix}$	

Quick review: genotype data

- VCF file looks like this:

CHR	POS	REF	ALT	P1	P2	P3	P4	P5
7	110	A	G	0 0	1 0	0 1	0 1	0 0
7	112	C	A	0 0	0 1	0 0	0 0	0 1
7	115	G	C	0 0	0 0	0 0	0 0	0 0
7	116	G	T	0 0	1 0	0 0	1 0	1 0
7	118	T	A	0 0	0 0	0 0	0 0	0 0
7	119	(G)	(C)	0 0	0 1	0 1	0 0	1 0
7	121	A	C	0 0	0 0	0 0	0 0	0 0
7	125	C	T	0 0	0 0	0 0	0 1	0 1

CHR	POS	REF	ALT	P1	P2	P3	P4	P5
7	110	A	G	0	1	2	1	0
7	112	C	A	0	1	0	0	1
7	115	G	C	0	0	0	0	0
7	116	G	T	0	1	0	1	1
7	118	T	A	0	0	0	0	0
7	119	G	C	0	1	1	0	1
7	121	A	C	0	0	0	0	0
7	125	C	T	0	0	0	1	1

phenotype (binary):

0	1	0	1	0
---	---	---	---	---

Used to obtain a GWAS risk model: $\ln\left(\frac{p(1|x)}{1-p(1|x)}\right) = \beta_0 + \beta_{112}x_{112} + \beta_{121}x_{121}$

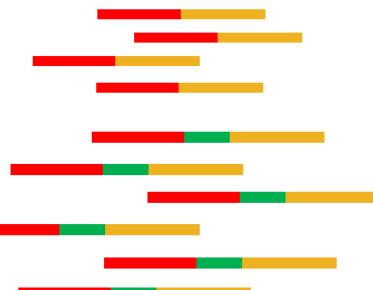
log-odds ratio

log-odds-ratio for general population

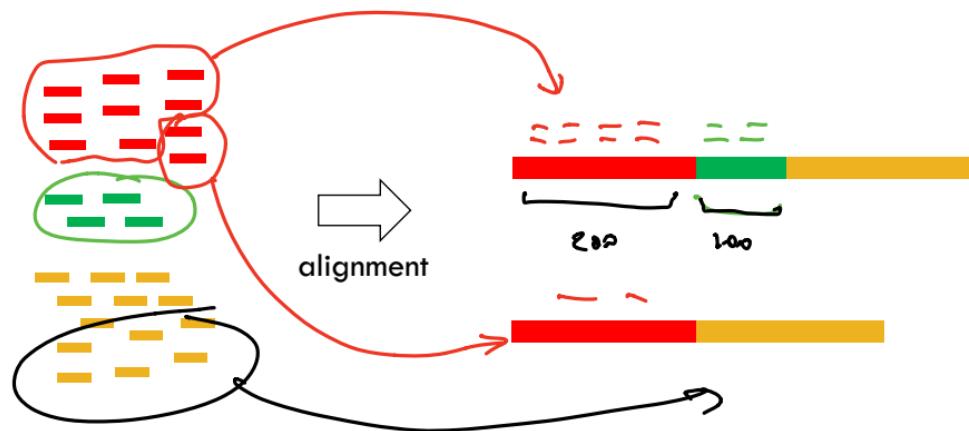
Quick review: transcript quantification

transcript 1 

transcript 2 

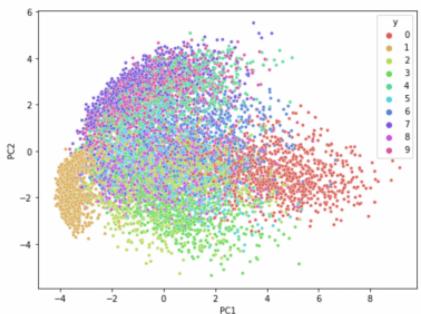
RNA 

RNA-seq

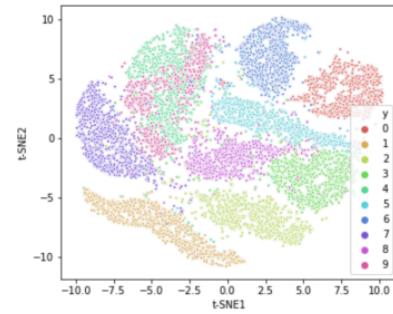


Quick review: single-cell RNA-seq

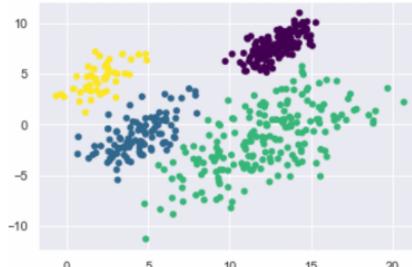
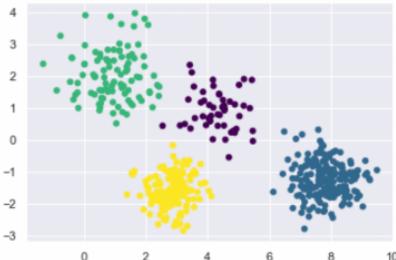
Linear dimensionality reduction (PCA)



Non-linear dimensionality reduction (t-SNE)



Clustering: K-means vs GMM



Better algorithms exist
for single-cell RNA-seq data

Wrapping up

- First part: DNA sequencing data analysis
 - Sequence alignment (dynamic programming, indexing)
 - Genome-wide association studies (GWAS)

- Second part: RNA sequencing data analysis
 - RNA quantification (EM algorithm)
 - Single-cell RNA-seq (dim. reduction, clustering)

