# Computational genomics: hands on course
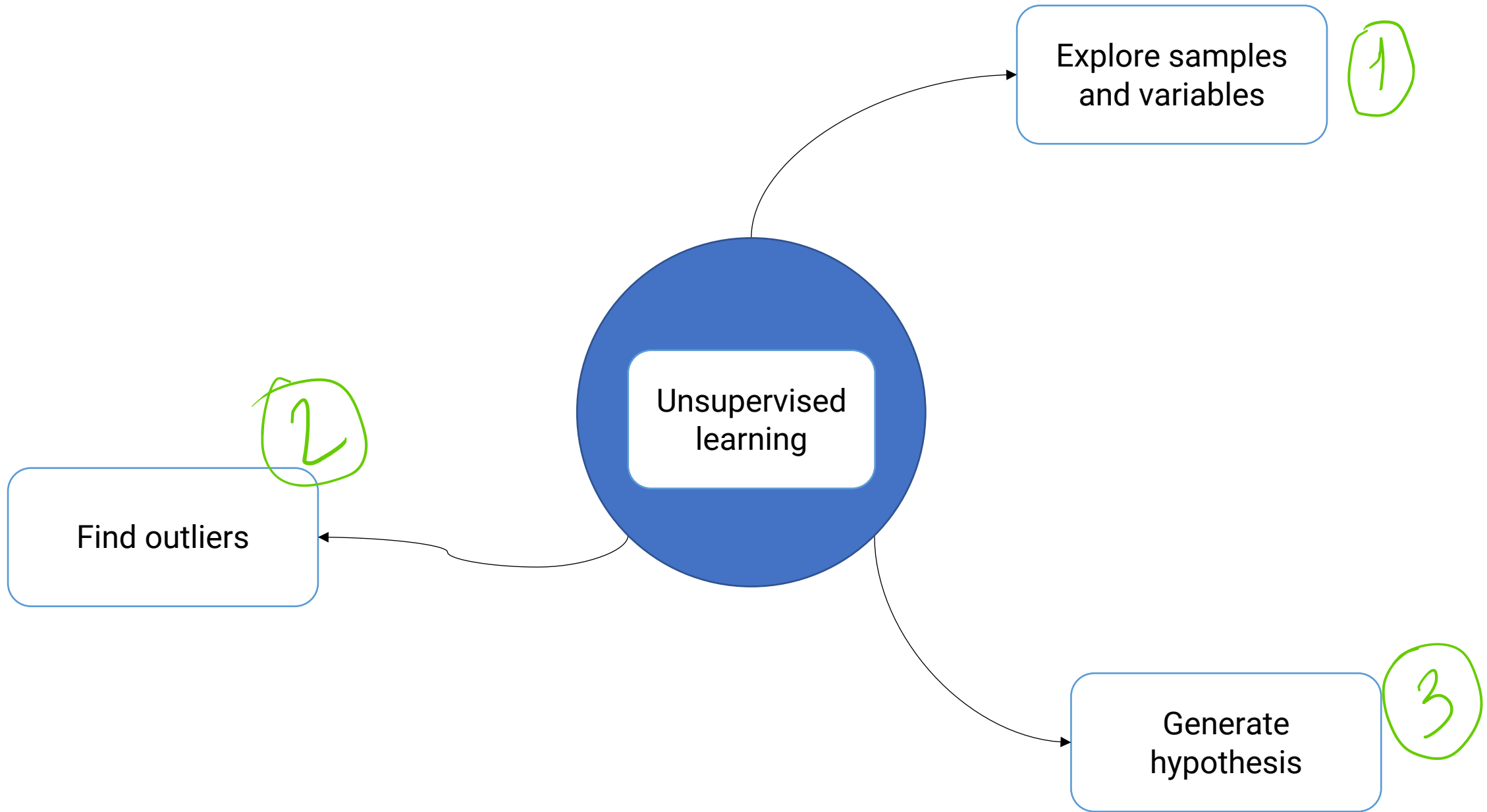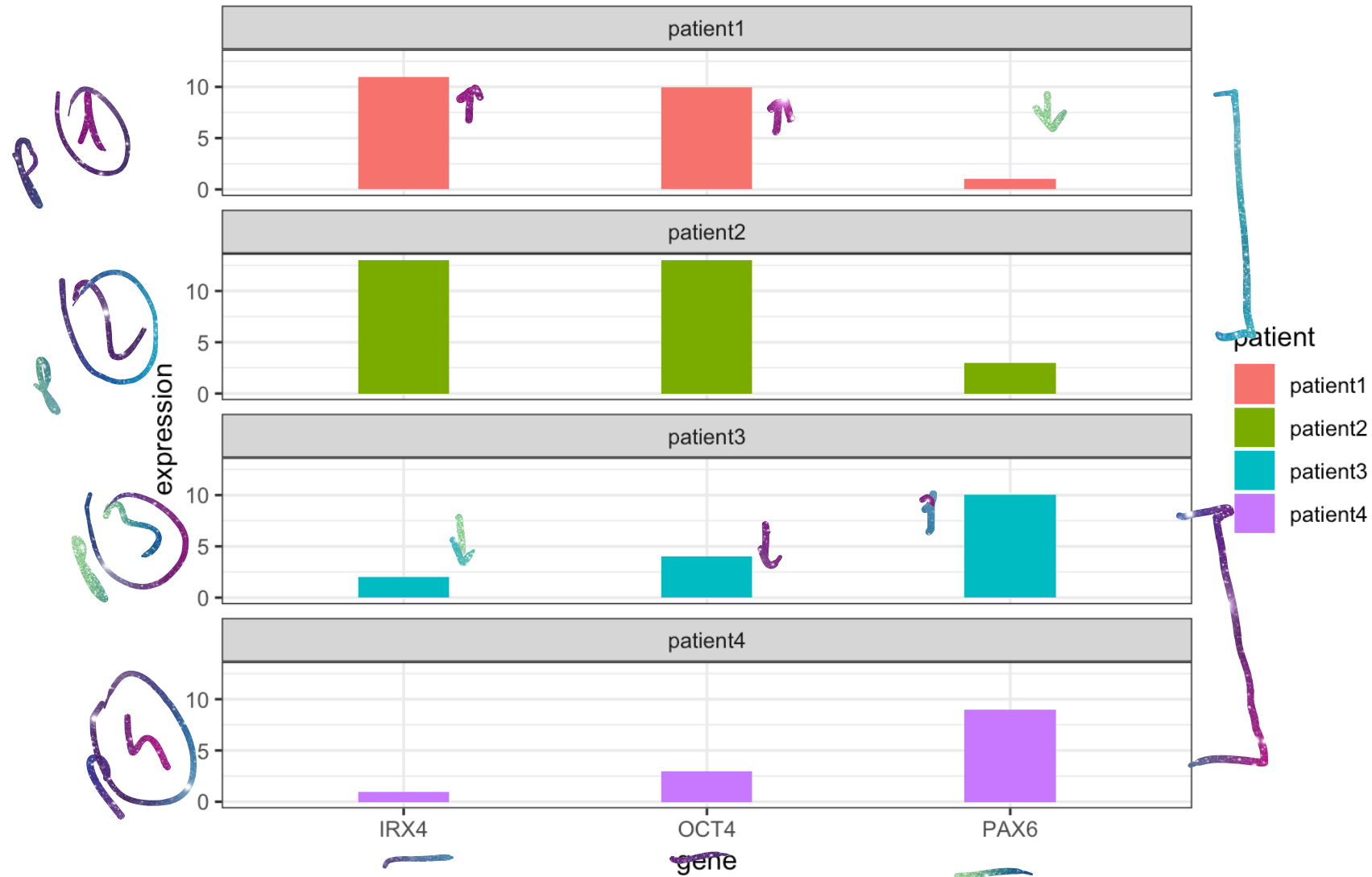
## Unsupervised learning

# Clustering: grouping samples
## distance metrics

# Clustering: grouping samples
## distance metrics



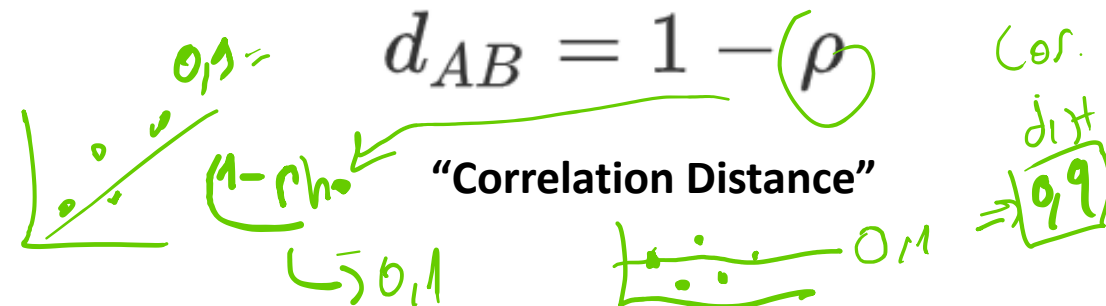$$d_{AB} = \sum_{i=1}^{n} |e_{Ai} - e_{Bi}|$$

**Manhattan distance"** or **"L1 norm"**

$$d_{AB} = \sqrt{\sum_{i=1}^{n} (e_{Ai} - e_{Bi})^2}$$

**"Euclidean Distance"** or **"L2 norm"**

$$d_{AB} = 1 - \rho$$

**"Correlation Distance"**

# Data processing
## Scaling

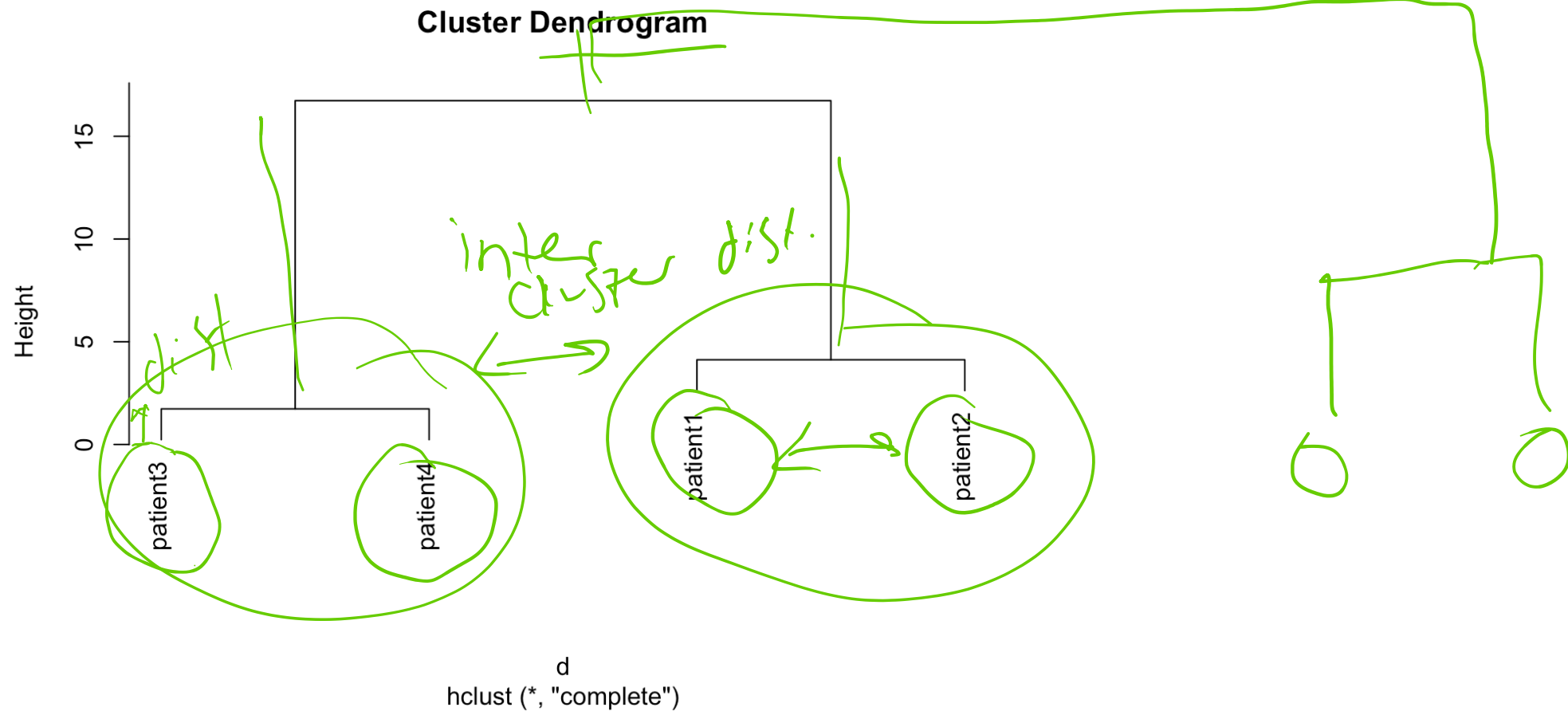**standardization** $(x-\textbf{mean}(x))/sd(x)$ $\Rightarrow \bar{x}=0$, $sd=1$

Scaling

gene.
exp

scale()

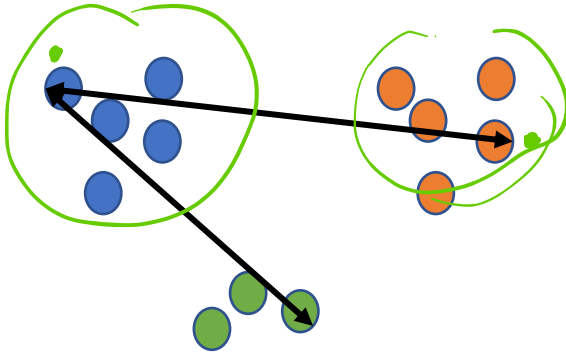More on this later in the supervised learning section.
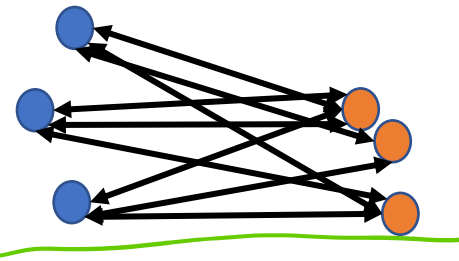
# Clustering: grouping samples
## Hierarchical clustering



**Cluster Dendrogram**

inter cluster dist.

dist

d
hclust (*, "complete")

# Clustering: grouping samples
## different approaches for merging clusters

**Complete Linkage**

**Average Linkage**

**Single Linkage**

**Ward's method**

$SS_A$

$SS_B$

$C_{AB} = SS_{AB} - SS_A - SS_B$ → 10

$C_{AC} = 5$

$C_{BC} = 20$
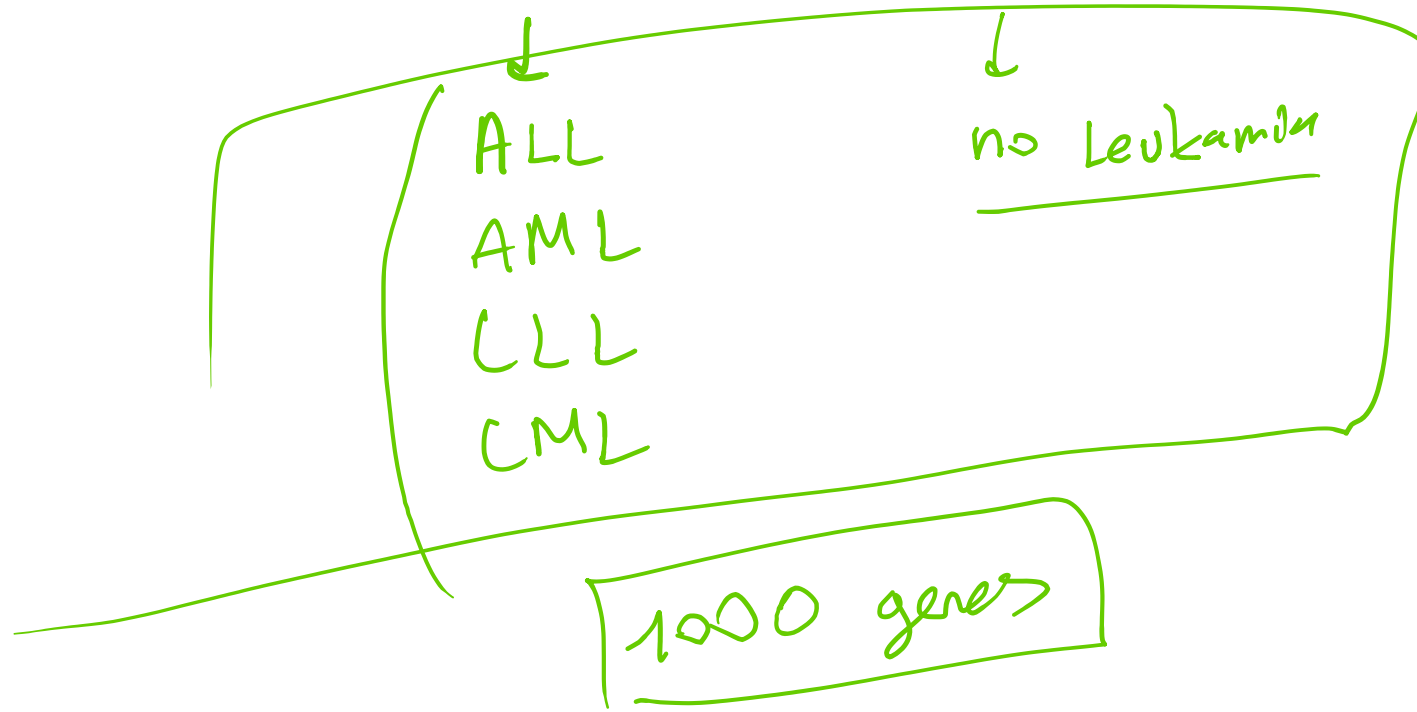
C

# Clustering: grouping samples
## Real world example

**Gene expression profiles from leukaemia patients and healthy controls**

ALL
AML
CLL
CML

no Leukamia

1000 genes

# Clustering: grouping samples
## cutting the tree



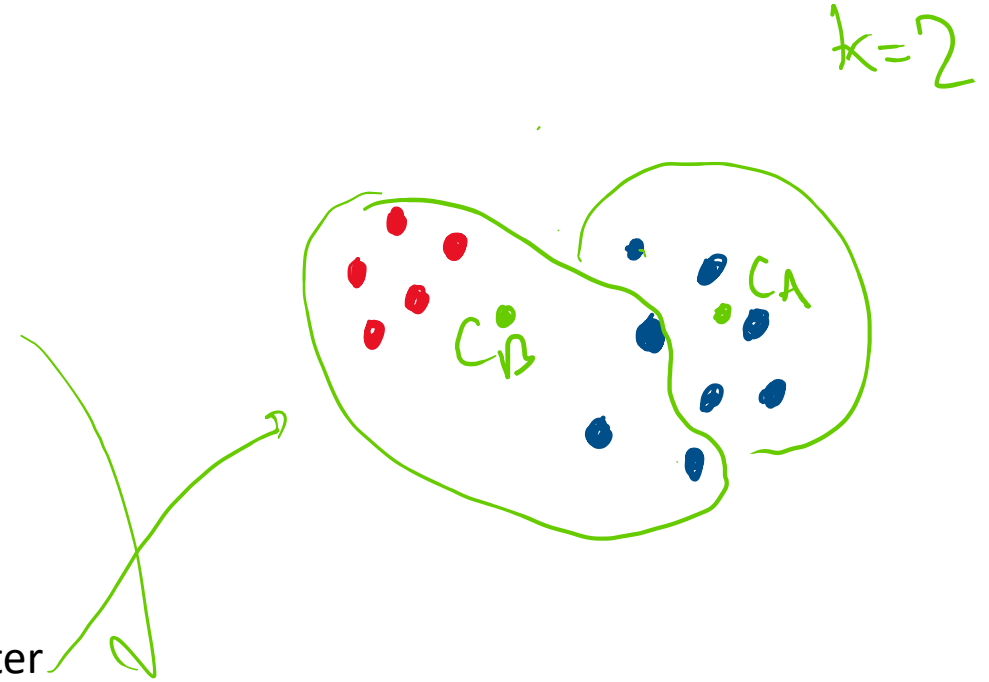**Cluster Dendrogram**

dist(t(mat))
hclust (*, "complete")
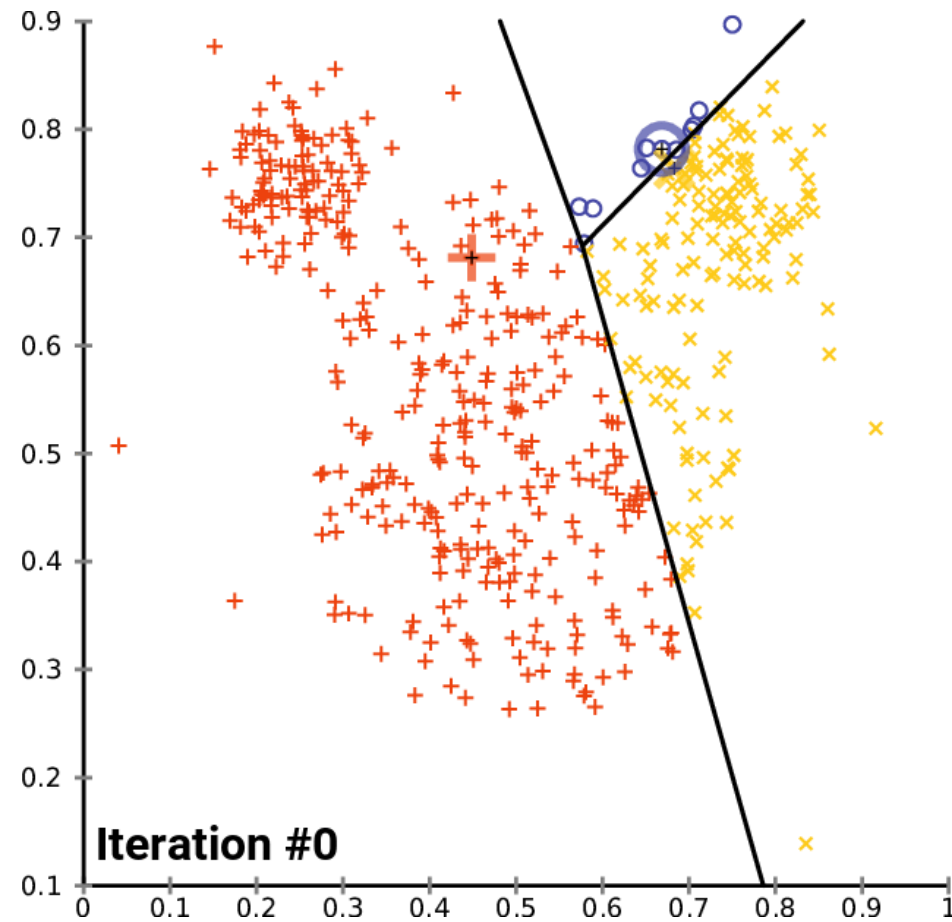
# Clustering: grouping samples
## K-means clustering

$k=2$



1. randomly choose k centers

2. Assign each data point to nearest center

3. Update centroid as the mean value of data points in the cluster

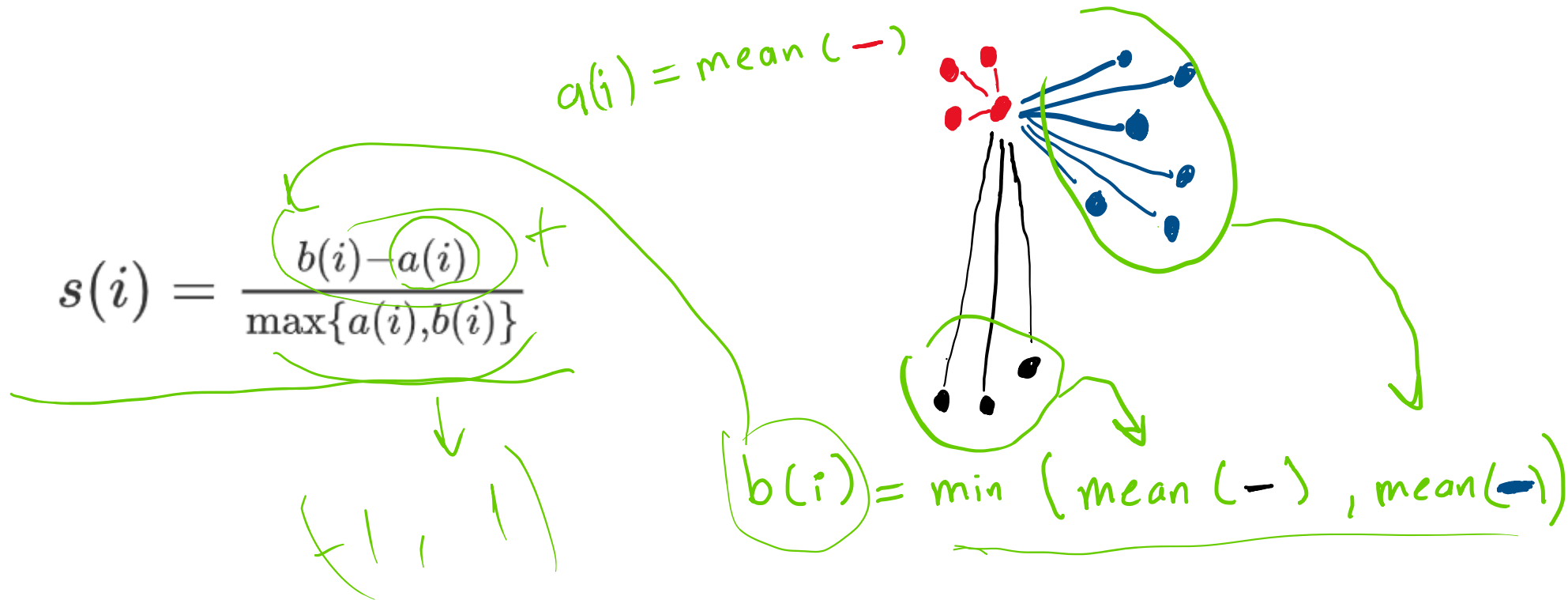4. Repeat steps 2-3 until until sum of squared distances to cluster minimized
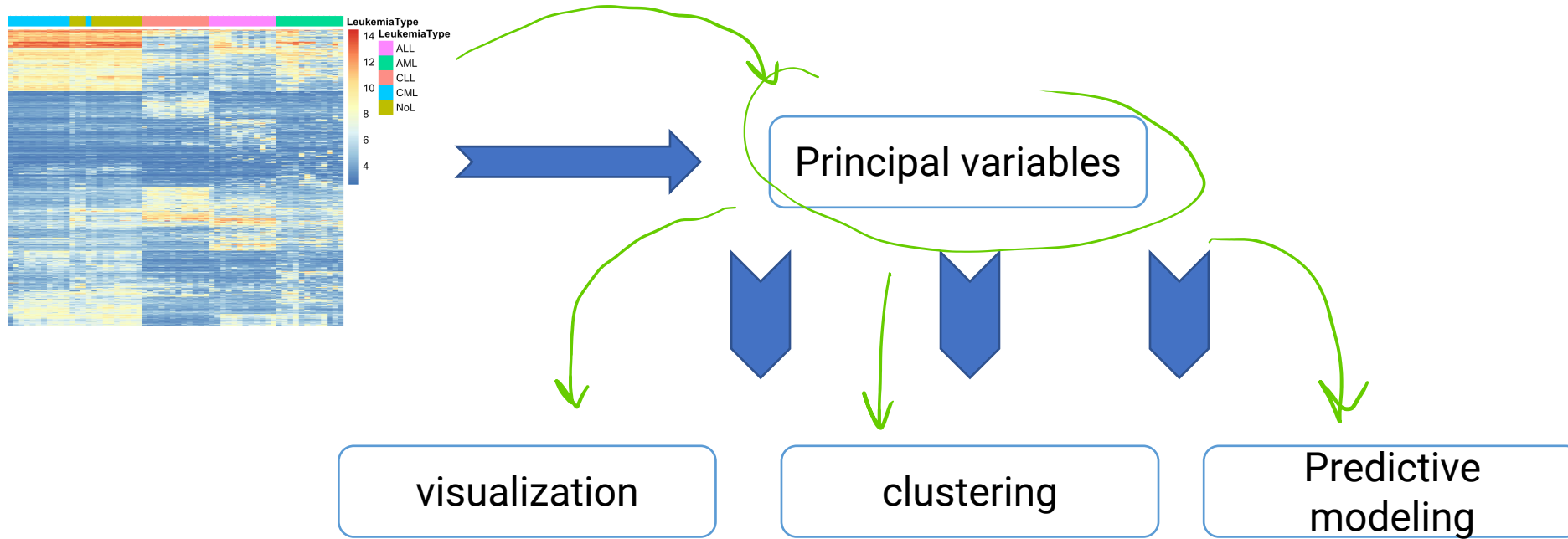
# Clustering: grouping samples
## K-means clustering



Iteration #0

# Clustering: grouping samples
## how to define the best k ? -> Silhouette score



$a(i) = mean(—)$

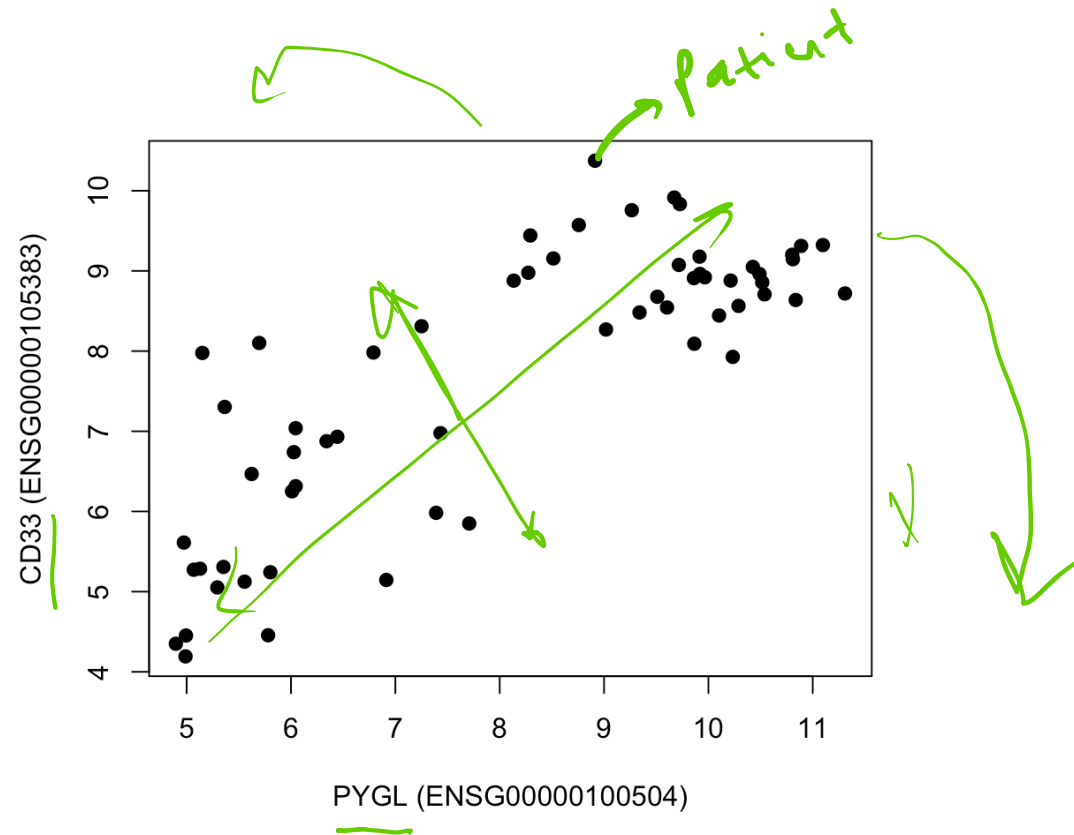$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$b(i) = min (mean (—), mean(—))$

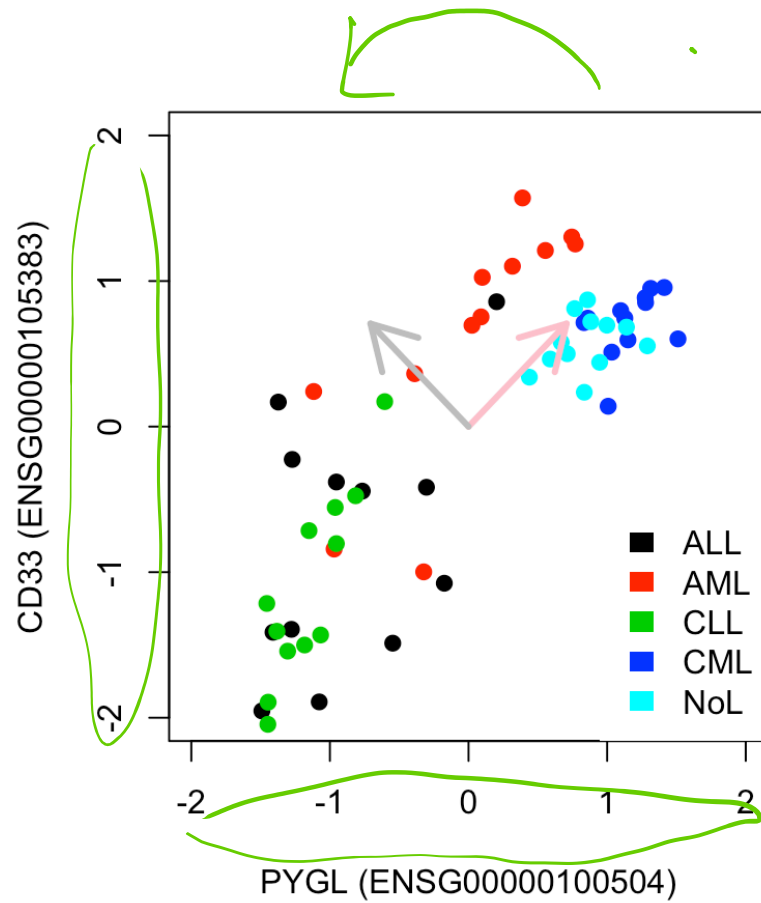# Dimensionality reduction techniques

# Dimensionality reduction techniques
## Principal component analysis (PCA)

# Dimensionality reduction techniques
## Principal component analysis (PCA)
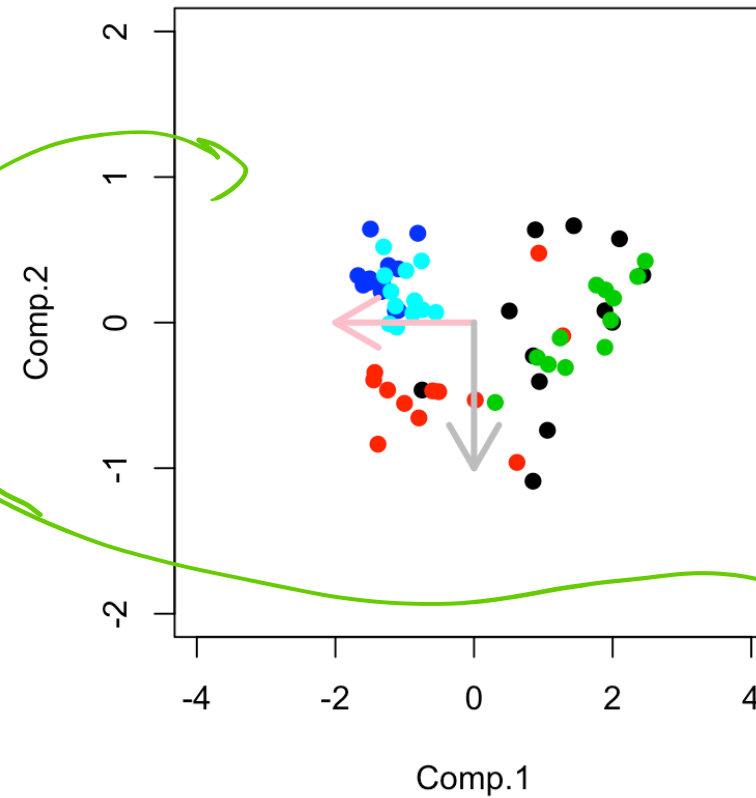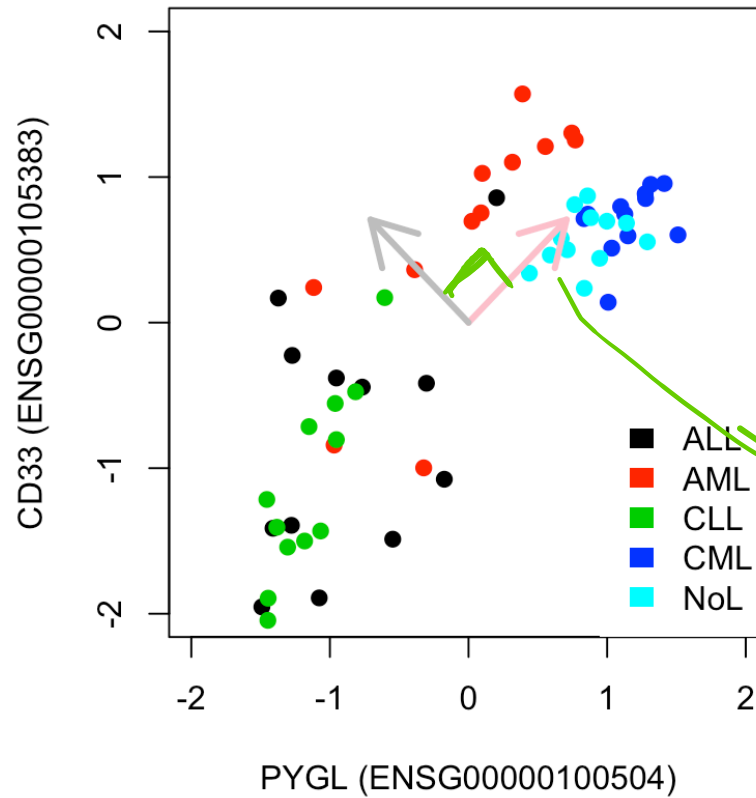
# Dimensionality reduction techniques
## Principal component analysis (PCA)

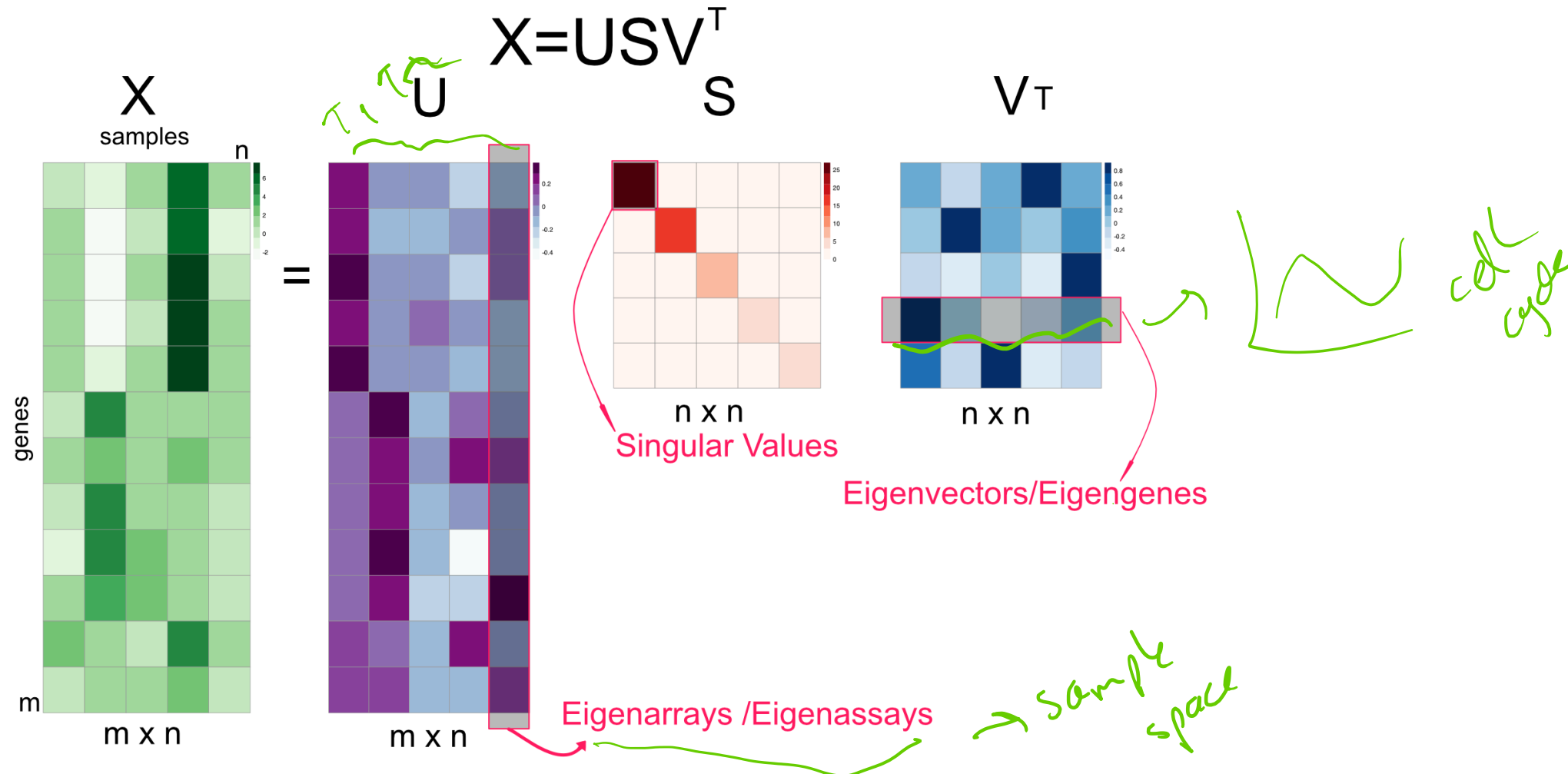$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu_X)(y_i - \mu_Y),$$

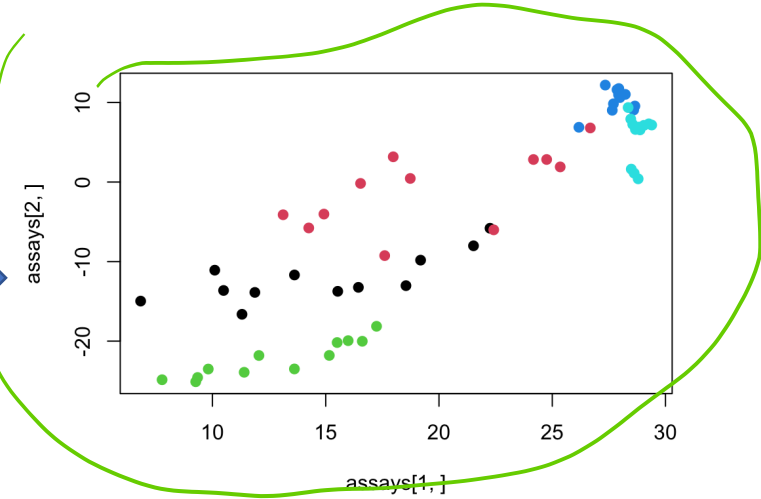# Dimensionality reduction techniques
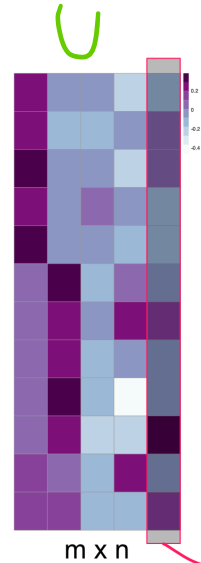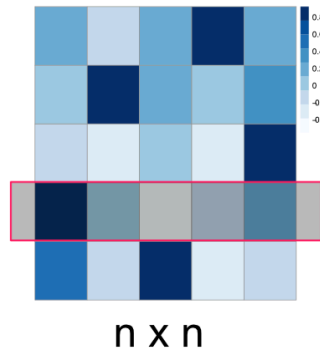## PCA with singular value decomposition (SVD)

# Dimensionality reduction techniques
## PCA with singular value decomposition (SVD)

# Dimensionality reduction techniques
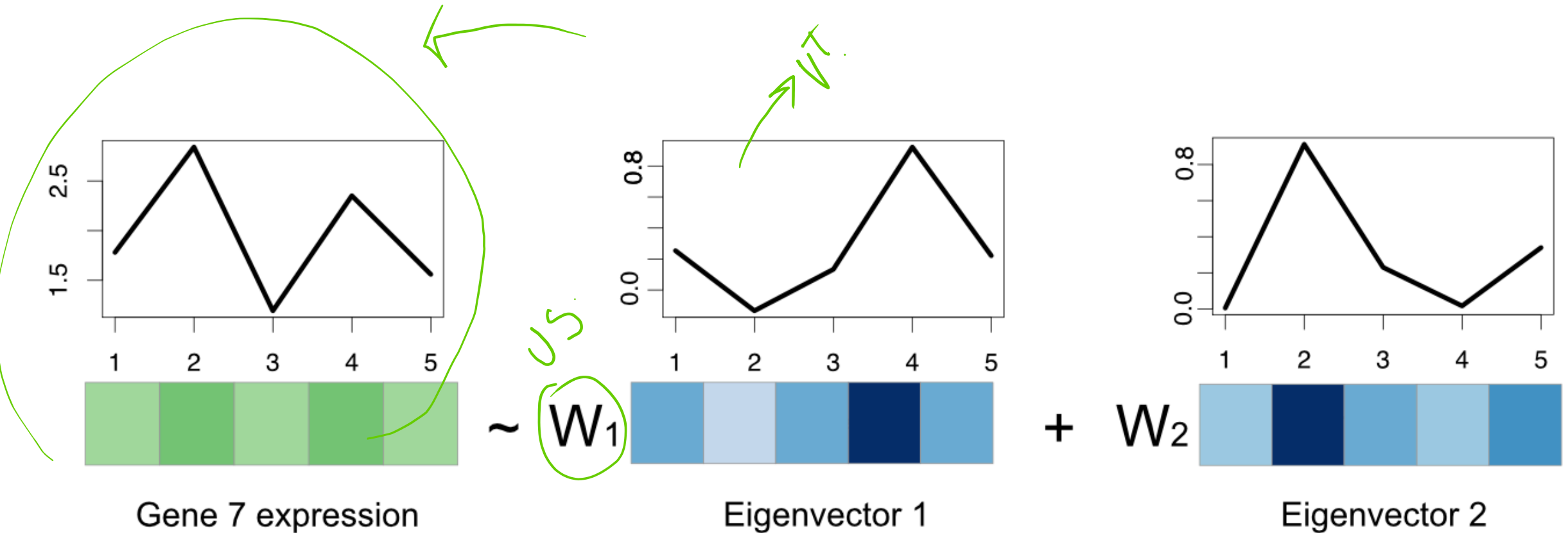## eigenvectors (latent factors) as expression programs



Gene 7 expression ~ W₁ Eigenvector 1 + W₂ Eigenvector 2

# Dimensionality reduction techniques

Matrix factorization

deep learning

| PCA/SVD | ICA | NMF | autoencoders |

Distance compression

| t-SNE | MDS | UMAP |

X=WH

X    W    H

samples

n

=

genes

m

m x n    m x k    k x n

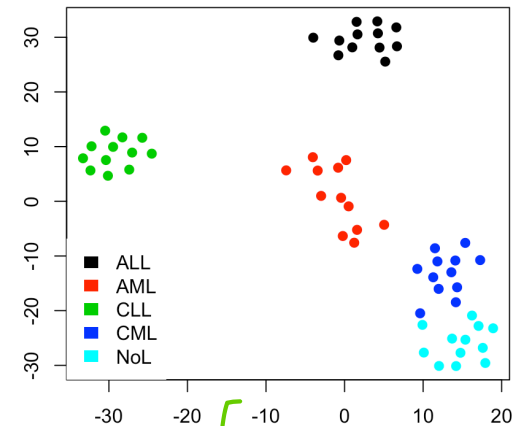metagenes

metaassays/metasamples

Bonus: latent factors

Distance in high dimensions

Distance in 2D/3D

$$\begin{bmatrix} 0 & d_{12}^2 & d_{13}^2 & \dots & d_{1n}^2 \\ d_{21}^2 & 0 & d_{23}^2 & \dots & d_{2n}^2 \\ d_{31}^2 & d_{32}^2 & 0 & \dots & d_{3n}^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ d_{n1}^2 & d_{n2}^2 & d_{n3}^2 & \dots & 0 \end{bmatrix}$$

ALL
AML
CLL
CML
NoL

t-SNE

# Unsupervised learning for Genomics:
## Recap

**Key concepts:**

1) Clustering

Limitations of clustering

Optimal number of clusters

2) Dimension reduction

Visualization

Matrix factorization