# Why PCA?

<u>Simplify complexity</u>, so it becomes easier to work with.
*Reduce number of features (genes)*

*M* genes        *2* dim.

*N* samples        *N* samples

"Remove" <u>redundancies</u> in the data

Identify the <u>most relevant</u> information
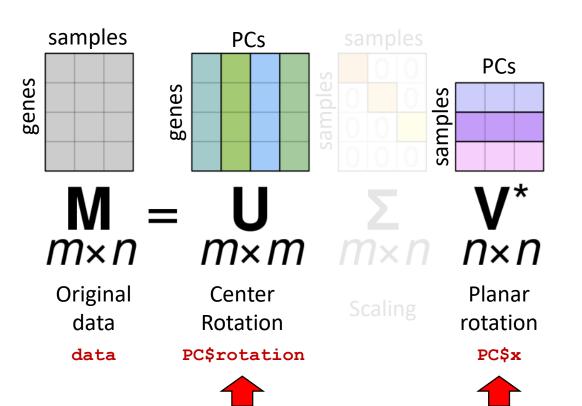*Find and filter noise*
*Detect data quality outliers or batches*

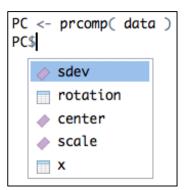Data <u>visualization</u>

# How PCA works

It is an algebraic method of dimensionality reduction (only numerical variables).

Relationships between variables need to be linear to apply PCA. For RNAseq data, it is recommended to have variables on a log-scale with VST or RLOG transformation.

It is a case inside Singular Value Decomposition (SVD) method (data compression)
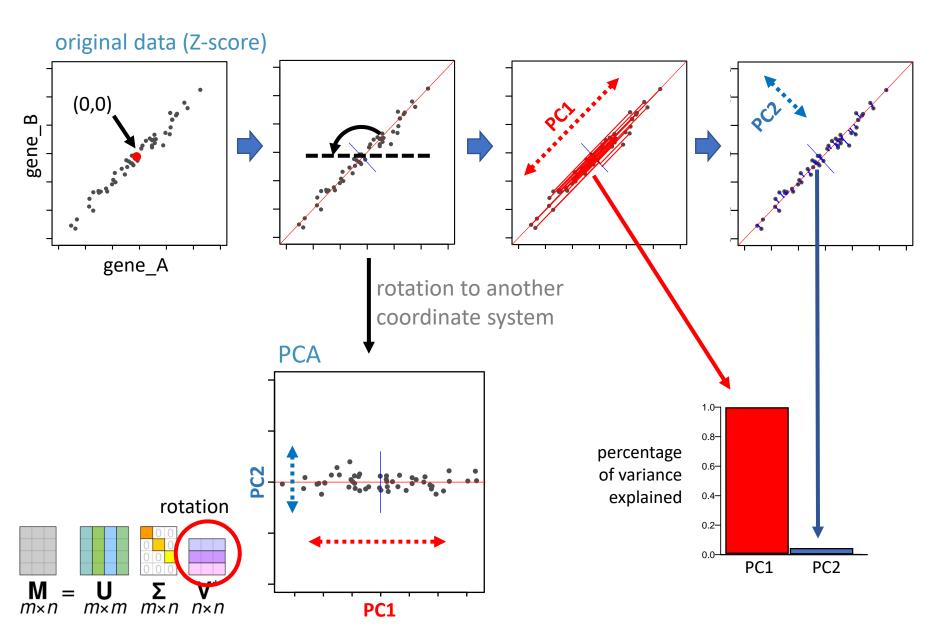
# How PCA works

Transform the data into a space of smaller dimension which would *summarize* the data in the most *relevant* way.

In PCA, *relevance* is measured by the variance (spread)
*Among the n-dimensional space, PCA will find the direction with highest variance*
*Why does PCA consider that higher variance = more relevant?*

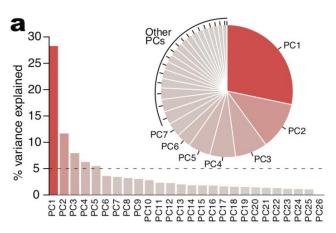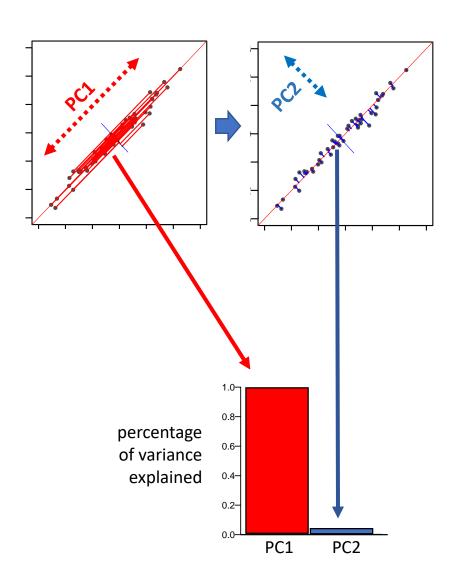| Typical question | PCA perspective |
| --- | --- |
| Is there any significant differences in expression of some genes between my study groups? | Would an "axis" with large spread in my data distinguish samples from one study group to another? |
| Is there any major outlier among my samples? | Would an "axis" with large spread in my data distinguish "failed experiment" samples from successful ones? |
| Is there some batch effect in my data (are technical differences larger than biological differences)? | Would an "axis" with large spread in my data distinguish samples from one batch to another? |

# How PCA works

original data (Z-score)

gene_B

(0,0)

gene_A

rotation to another coordinate system

PCA

PC2

PC1

PC1

PC2

percentage of variance explained

rotation

$$\mathbf{M}_{m \times n} = \mathbf{U}_{m \times m} \; \mathbf{\Sigma}_{m \times n} \; \mathbf{V}_{n \times n}$$

# How PCA works

PC1 explains >98% of the variance

1 PC thus represents 2 genes very well
*"Removing" redundancy*

PC2 is nearly insignificant in this example
*Could be disregarded*

In real life ...



Czarnewski *et al* 2019



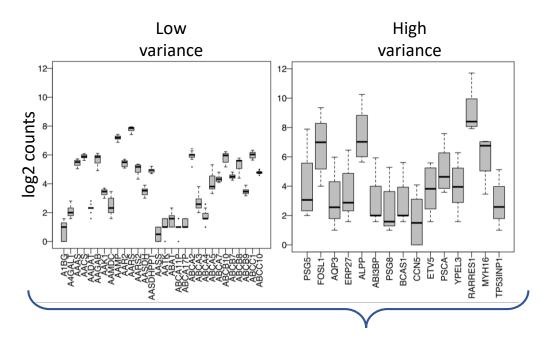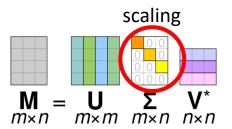percentage
of variance
explained

# Data transformation and scaling

Before applying PCA, the data should be first transformed VST or RLOG
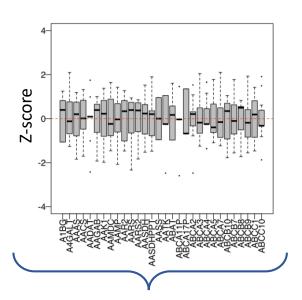
Each feature can be centered and scaled to have a similar center (zero) and similar deviation.



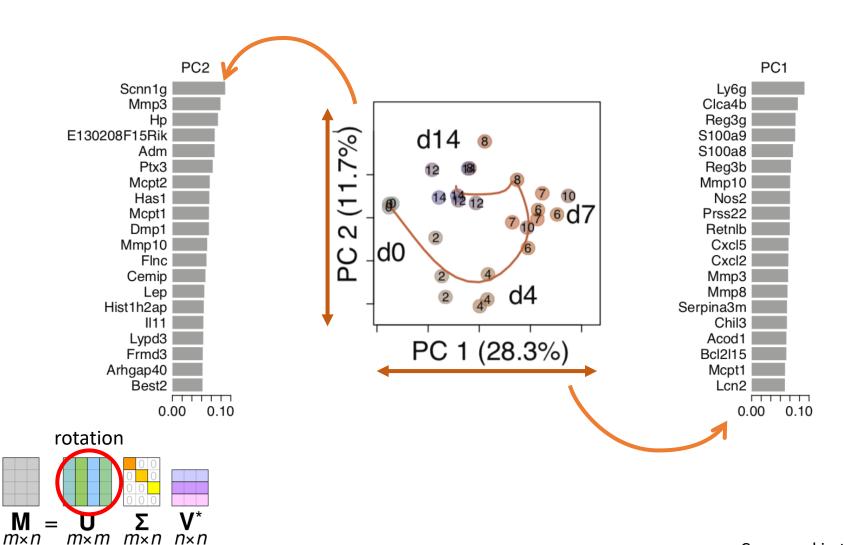PCA on raw counts will separate genes with <u>higher counts</u> in the first PCs
*(higher distance to 0)*

PCA on Z-score will separate genes with most <u>common expression trends</u> in the first PCs
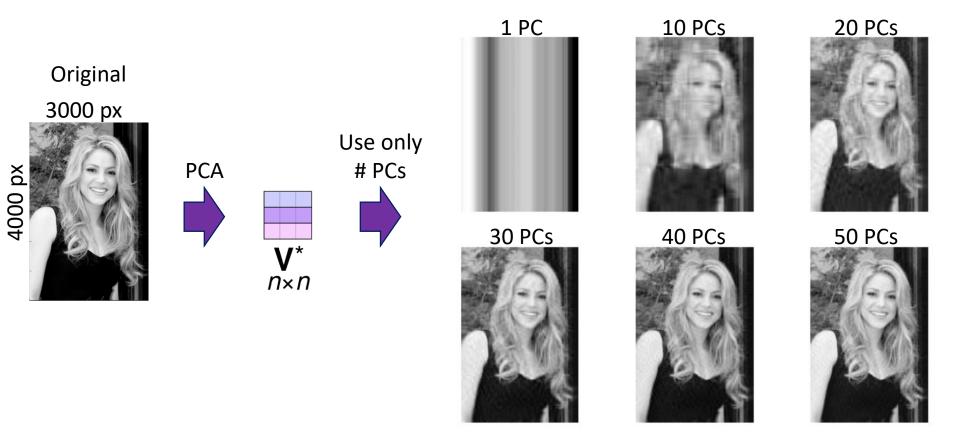
scaling

$$M_{m\times n} = U_{m\times m} \; \Sigma_{m\times n} \; V^*_{n\times n}$$

# How PCA works

Interpretability of principal components

# A visual intuition of PCA

The top principal components store more important ~~Shakira~~ information



Original

3000 px

4000 px

PCA

$V^*$
$n{\times}n$

Use only
# PCs

1 PC          10 PCs          20 PCs

30 PCs          40 PCs          50 PCs

# PCA summary

PCA is a dimensionality reduction method for <u>numerical</u> variables

<u>Linear</u> combinations -> when relationships between variables are non-linear, PCA is not recommended.

The data is usually <u>SCALED</u> and TRANSFORMED (i.e. VST/RLOG) prior to PCA

It is an <u>interpretable</u> dimensionality reduction

The top principal components contain <u>higher variance </u> from the data and PCA preserves the whole variance

Can be used as <u>filtering</u>, by selecting only the top significant PCs
- PCs that explain at least 1% of variance
- The first 5-10 PCs

# Thank you. Questions?