# EDA: Principal Component Analysis (PCA)

RNA-seq data analysis

**Paulo Czarnewski**

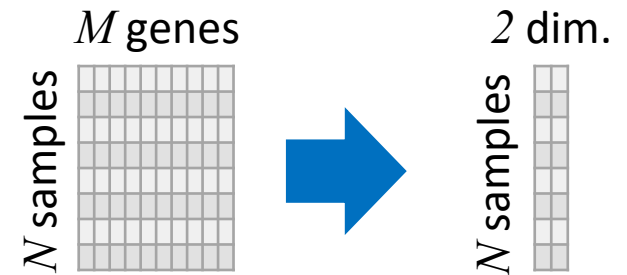https://czarnewski.github.io/czarnewski/index.html

# Why PCA?
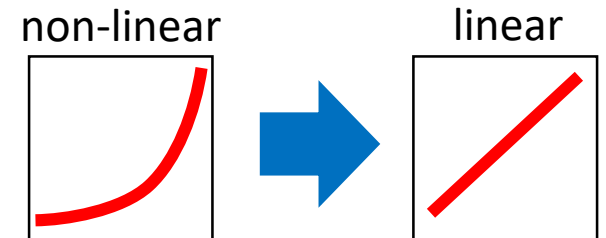
Simplify complexity, so it becomes easier to work with.
- Reduce number of features (genes)
- Need to transform non-linear relationships to linear

"Remove" redundancies in the data

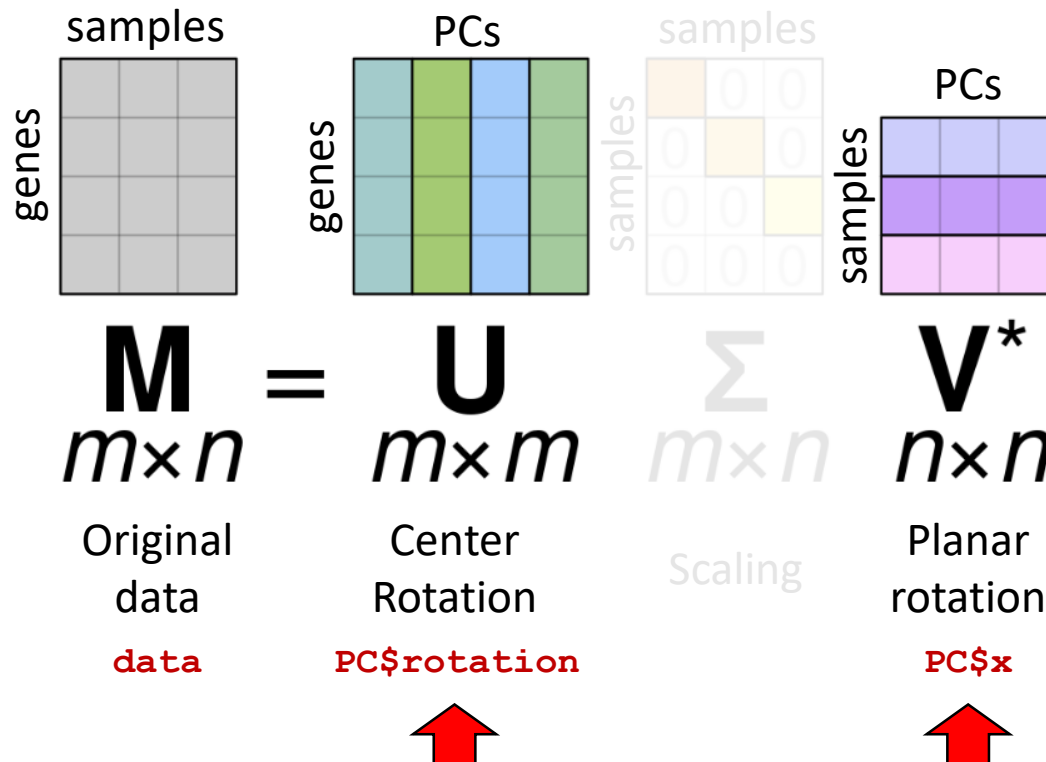Identify the most relevant information
*Find and filter noise*

Data visualization

# How PCA works

It is a <u>LINEAR</u> algebraic method of dimensionality reduction.

It is a case inside Singular Value Decomposition (SVD) method (data compression)
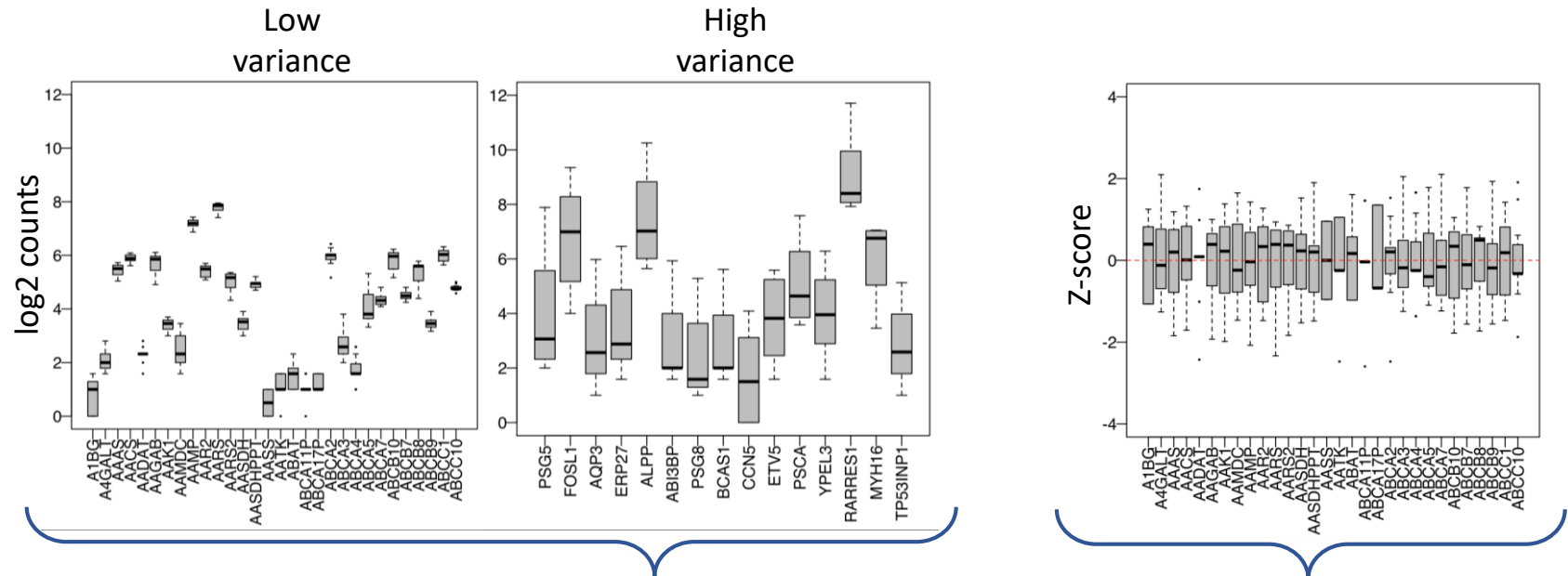*Any matrix can be decomposed as a multiplication of other matrices (Matrix Factorization).*



$$\underset{m\times n}{\mathbf{M}} = \underset{m\times m}{\mathbf{U}} \quad \underset{m\times n}{\Sigma} \quad \underset{n\times n}{\mathbf{V^*}}$$

| Original data | Center Rotation | Scaling | Planar rotation |
|---|---|---|---|
| **data** | **PC$rotation** | | **PC$x** |

```
PC <- prcomp( data )
PC$|
  ◇ sdev
  ▦ rotation
  ◆ center
  ◆ scale
  ▦ x
```
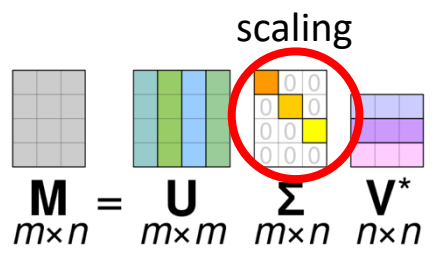
# Data transformation and scaling

Before applying PCA, the data should be first transformed to a <u>linear</u> scale (i.e. log)

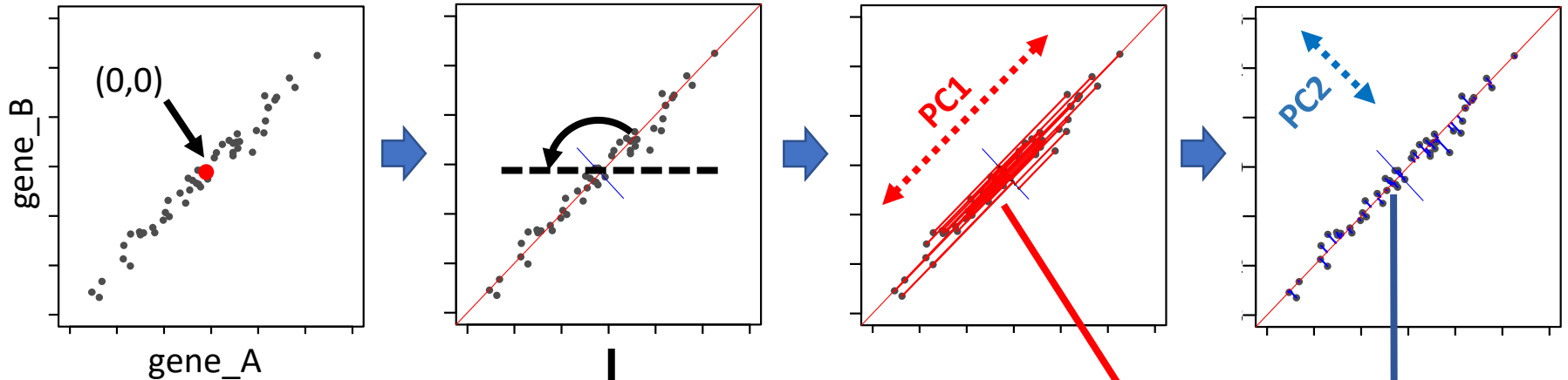Each feature should be scaled to have a similar center (zero) and similar deviation.



Low variance

High variance

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

scaling

PCA on Z-score will separate genes with most <u>common expression trends</u> in the first PCs
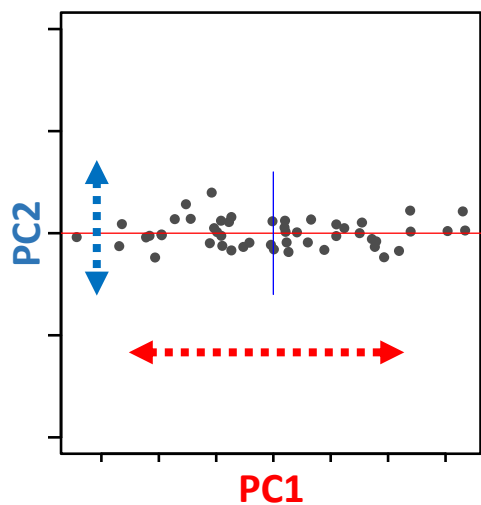
$$M_{m \times n} = U_{m \times m} \, \Sigma_{m \times n} \, V^*_{n \times n}$$

# How PCA works



original data (Z-score)

gene_B

(0,0)

gene_A

PC1

PC2

rotation to another coordinate system

PCA

PC2

PC1

rotation

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

percentage of variance explained

1.0
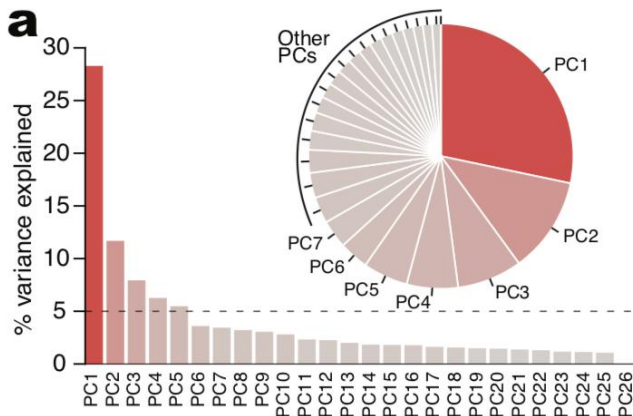
0.8

0.6

0.4

0.2

0.0

PC1   PC2

# 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
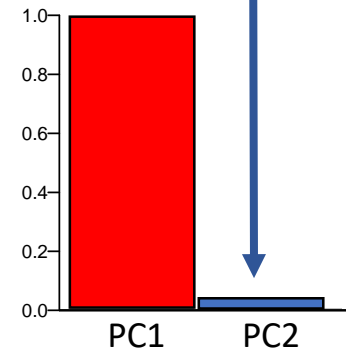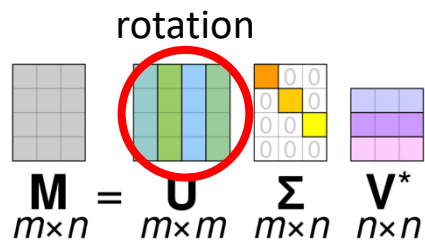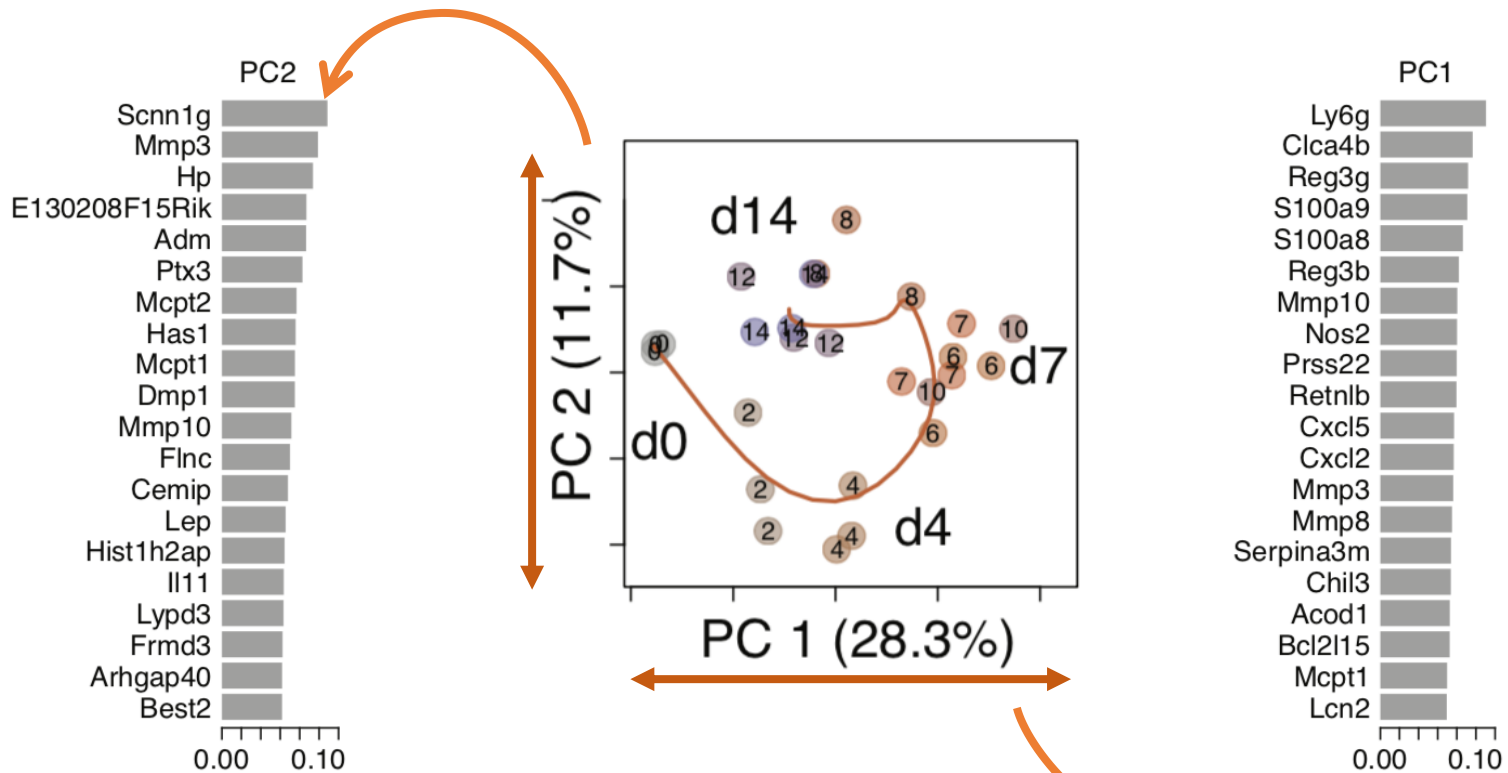*Could be disregarded*

In real life …



Czarnewski *et al* 2019

percentage
of variance
explained

# How PCA works

Each PC has a meaning



Czarnewski *et al* 2019

# A visual intuition of PCA

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



Original
3000 px

4000 px

PCA

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

Use only
# PCs

1 PC

10 PCs

20 PCs

30 PCs

40 PCs

50 PCs

# PCA summary

It is a <u>LINEAR</u> method of dimensionality reduction

The data is usually <u>SCALED</u> (i.e. Z-score ) and TRANSFORMED (i.e. log) prior to PCA

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

The top principal components contain <u>higher variance </u>from the data

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?

**Paulo Czarnewski**