

EDA: Principal Component Analysis (PCA)

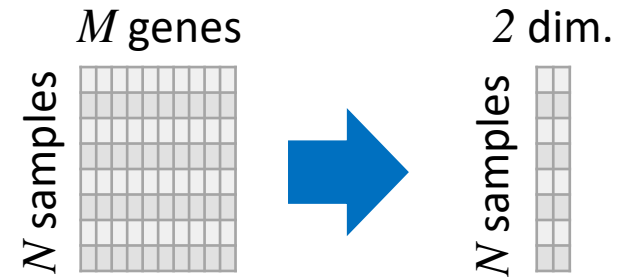
RNA-seq data analysis

Paulo Czarnewski | 30-Nov-2020

Why PCA?

Simplify complexity, so it becomes easier to work with.

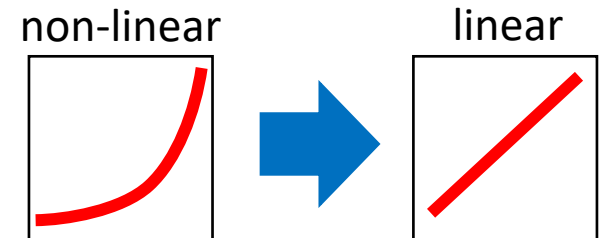
- Reduce number of features (genes)
- Need to transform non-linear relationships to linear



“Remove” redundancies in the data

Identify the most relevant information

Find and filter noise

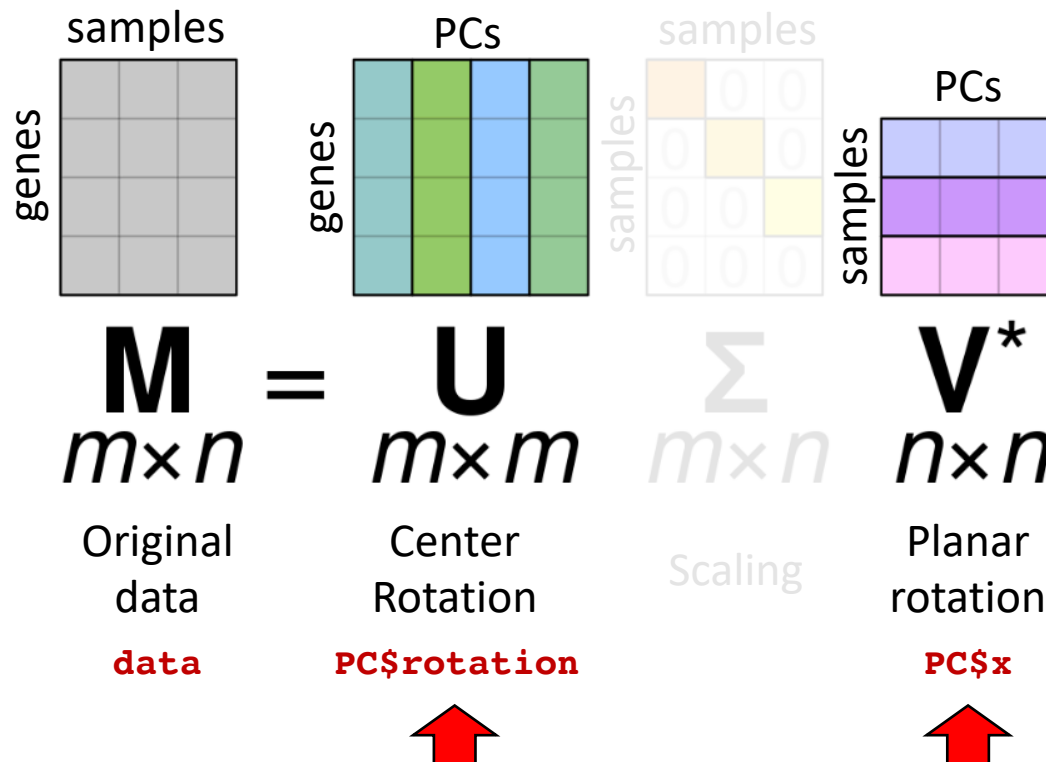


Data visualization

How PCA works

It is a LINEAR 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).

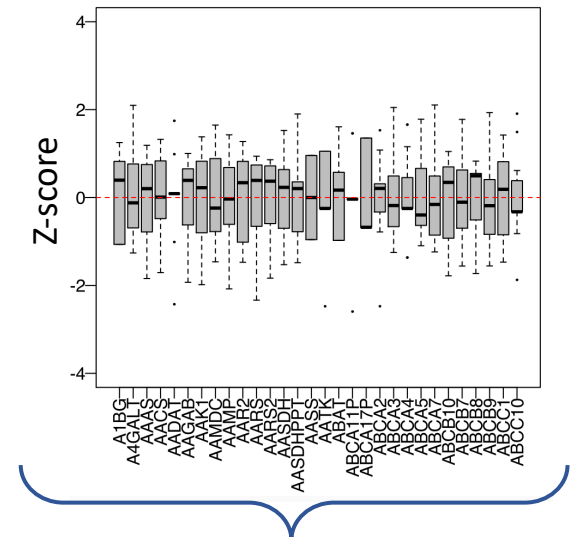
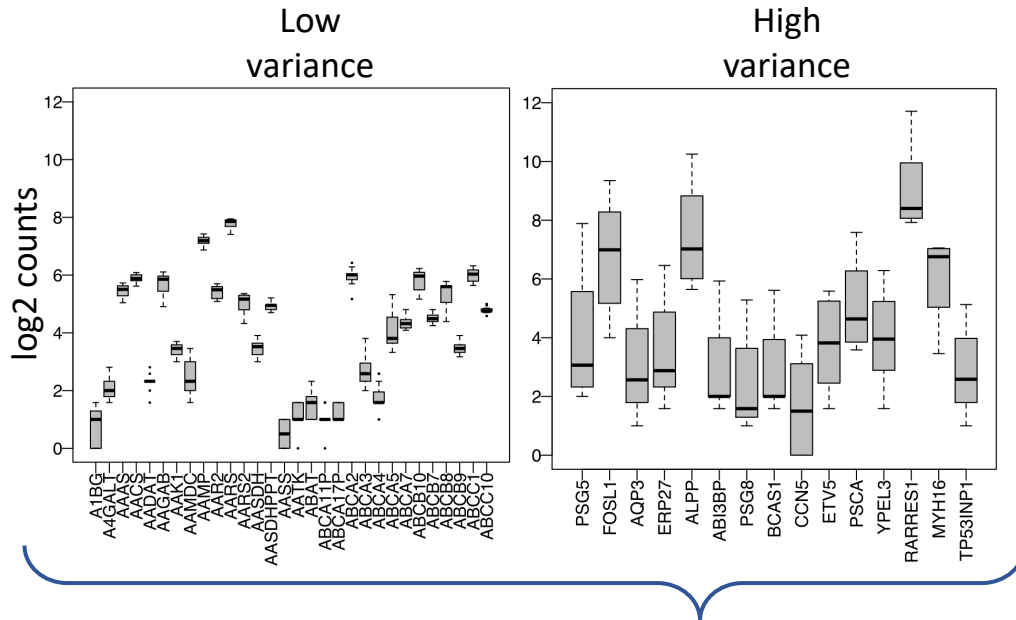


```
PC <- prcomp( data )
PC$|
```

- sdev
- rotation
- center
- scale
- x

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

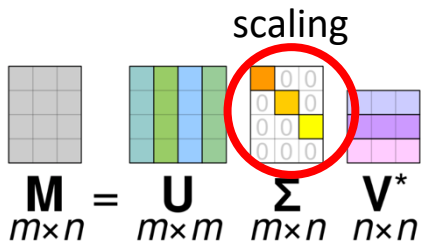
Each feature should be scaled to have a similar center (zero) and similar deviation.



PCA on raw counts will separate genes with higher counts in the first PCs

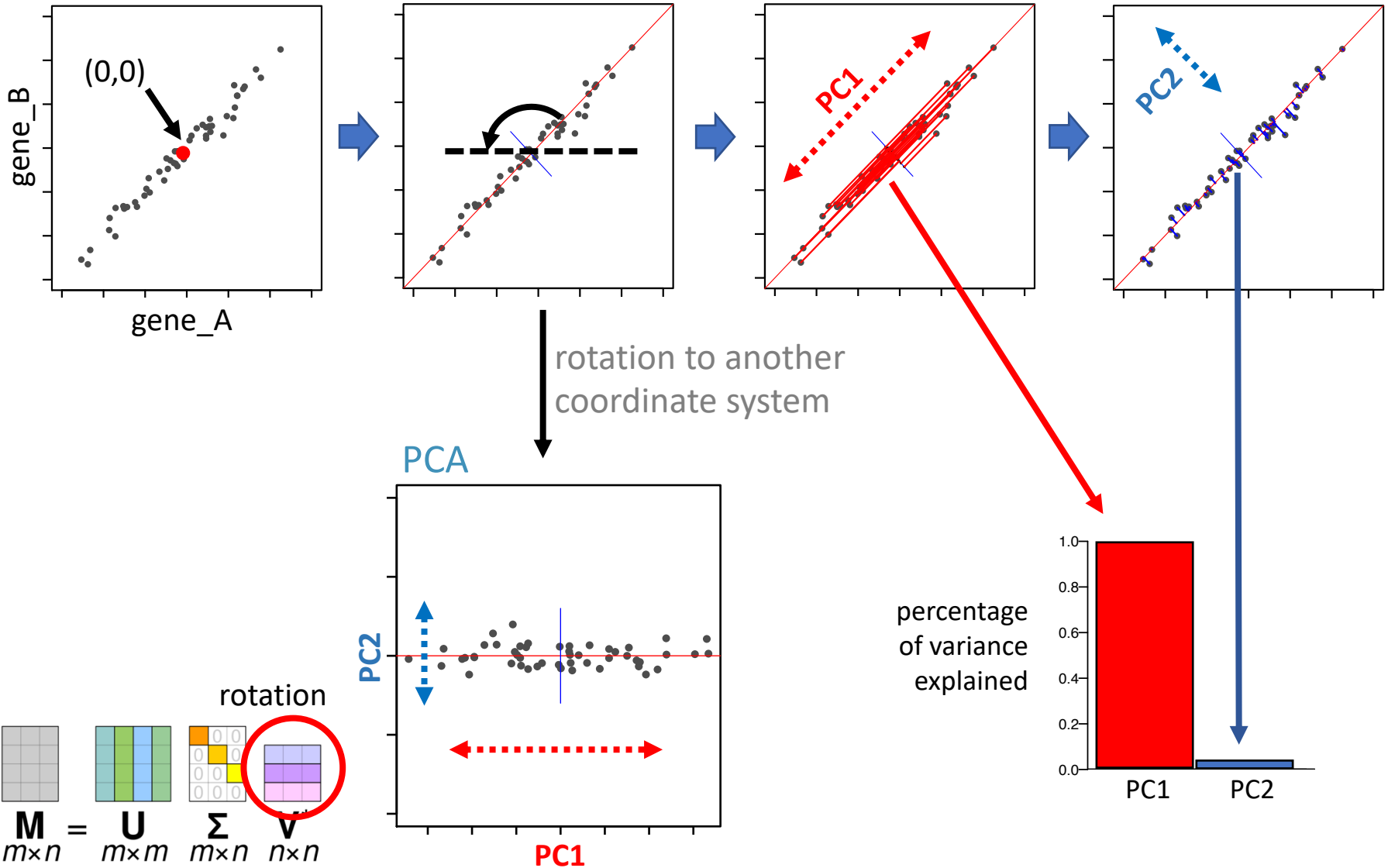
(higher distance to 0)

PCA on Z-score will separate genes with most common expression trends in the first PCs



How PCA works

original data (Z-score)

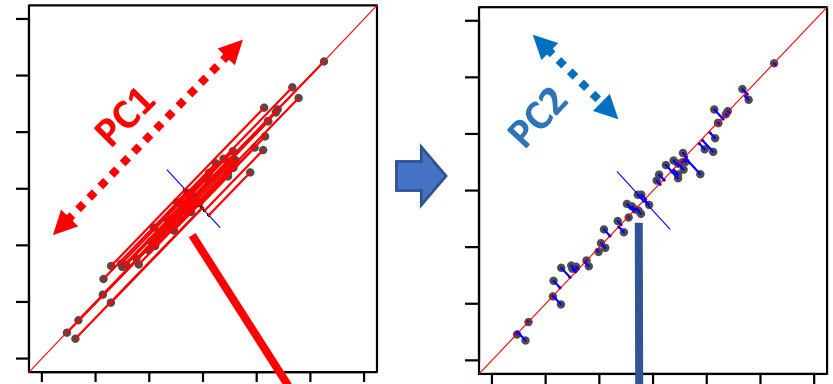


How PCA works

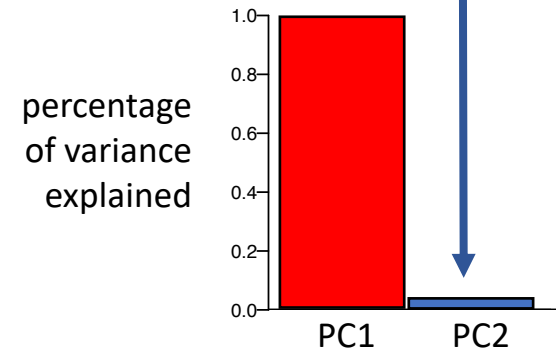
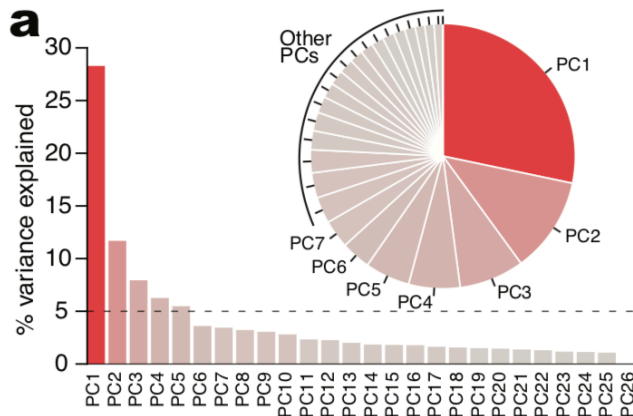
PC1 explains >98% of the variance

1 PC thus represents 2 genes very well
“Removing” redundancy

PC2 is nearly insignificant in this example
Could be disregarded

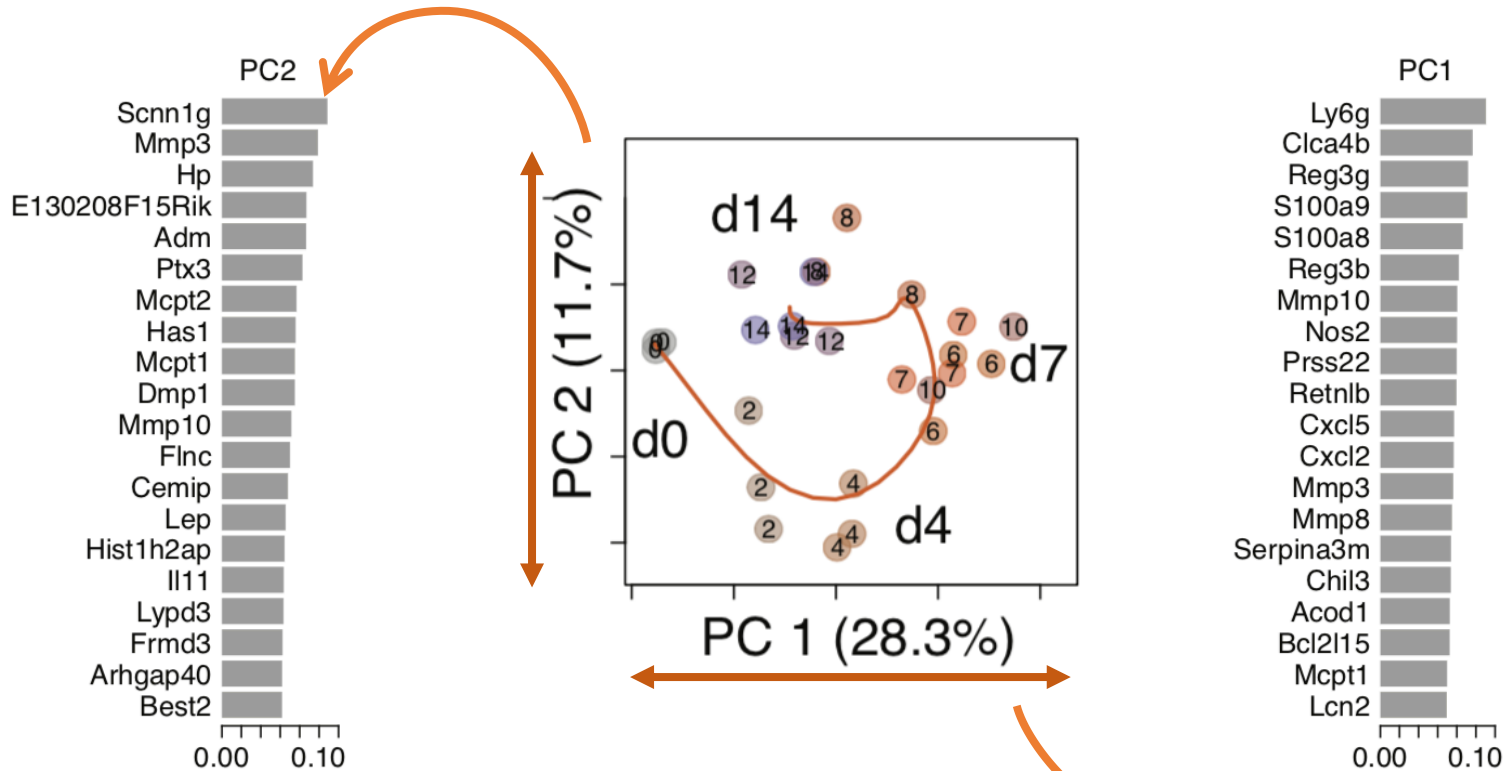


In real life ...



How PCA works

Each PC has a meaning

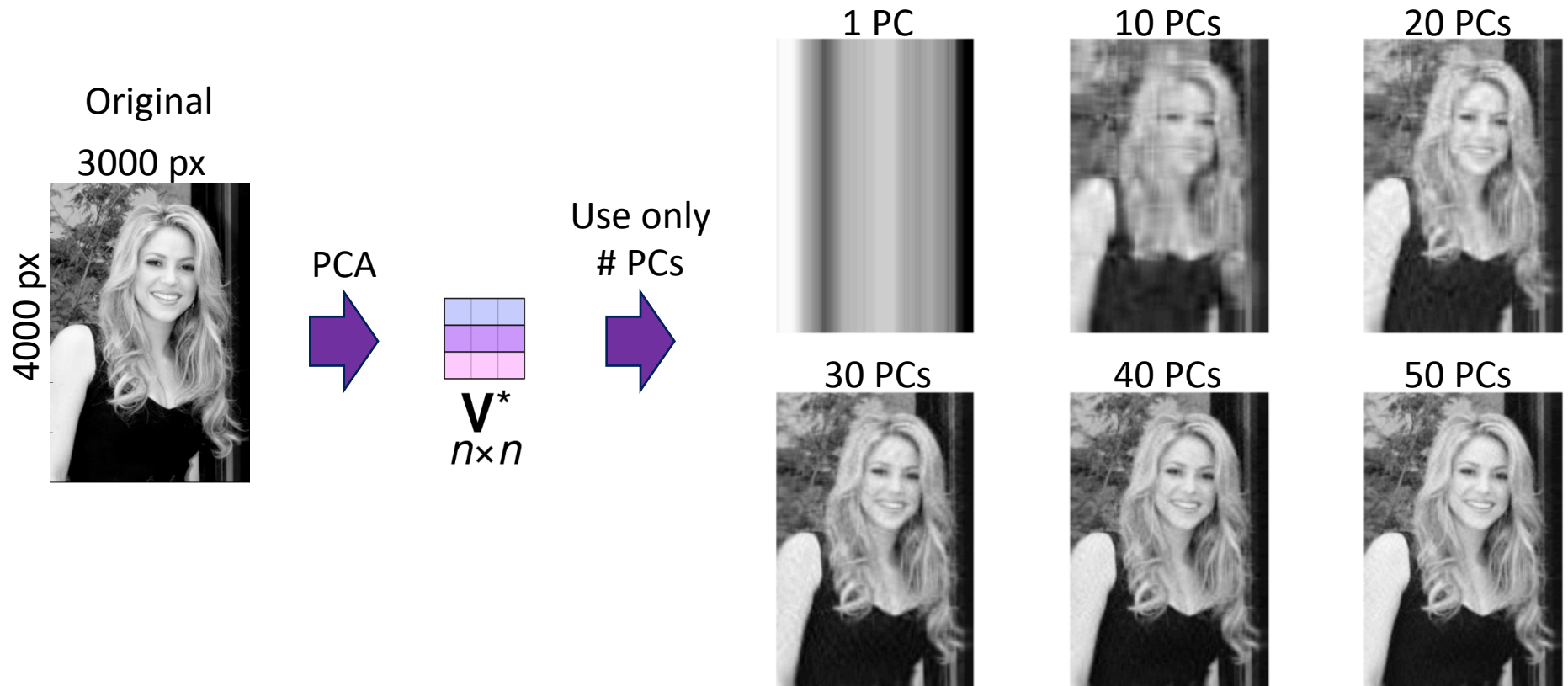


rotation

$$\begin{matrix}
 \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} & = & \begin{matrix} \color{red}\square & \color{green}\square & \color{blue}\square \\ \color{green}\square & \color{blue}\square & \color{green}\square \\ \color{blue}\square & \color{green}\square & \color{blue}\square \end{matrix} & \begin{matrix} \color{orange}\square & \color{yellow}\square & \color{yellow}\square \\ \color{yellow}\square & \color{yellow}\square & \color{yellow}\square \\ \color{yellow}\square & \color{yellow}\square & \color{yellow}\square \\ \color{yellow}\square & \color{yellow}\square & \color{yellow}\square \end{matrix} & \begin{matrix} \color{purple}\square & \color{purple}\square \\ \color{purple}\square & \color{purple}\square \\ \color{purple}\square & \color{purple}\square \\ \color{purple}\square & \color{purple}\square \end{matrix} \\
 \mathbf{M} & = & \mathbf{U} & \mathbf{\Sigma} & \mathbf{V}^* \\
 m \times n & & m \times m & m \times n & n \times n
 \end{matrix}$$

A visual intuition of PCA

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



It is a LINEAR method of dimensionality reduction

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

It is an interpretable dimensionality reduction

The top principal components contain higher variance from the data

Can be used as FILTERING, 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 | 13-May-2019