

# EDA: Principal Component Analysis (PCA)

---

RNA-seq data analysis

**Paulo Czarnewski**  
**Julie Lorent**



# Why PCA?

Simplify complexity, so it becomes easier to work with.

*Reduce number of features (genes)*

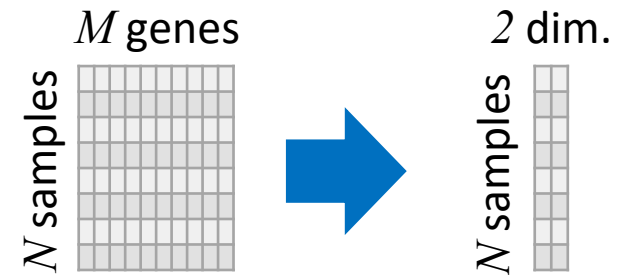
“Remove” redundancies in the data

Identify the most relevant information

*Find and filter noise*

*Detect data quality outliers or batches*

Data visualization

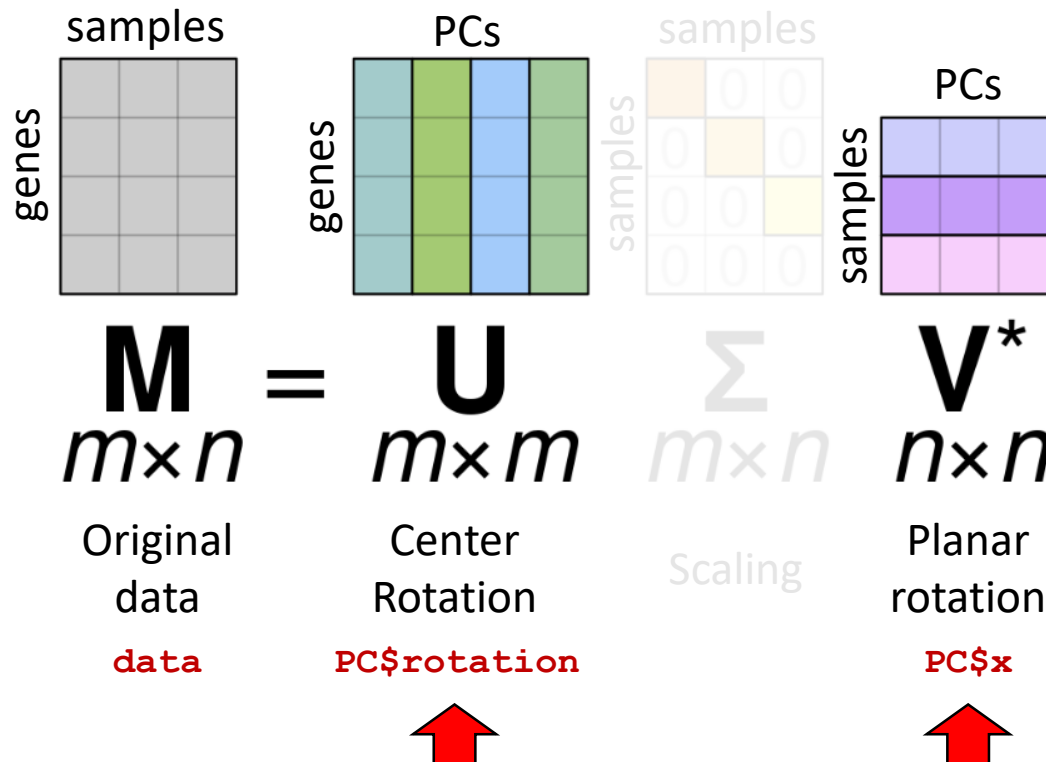


# How PCA works

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)



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

- sdev
- rotation
- center
- scale
- x

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)

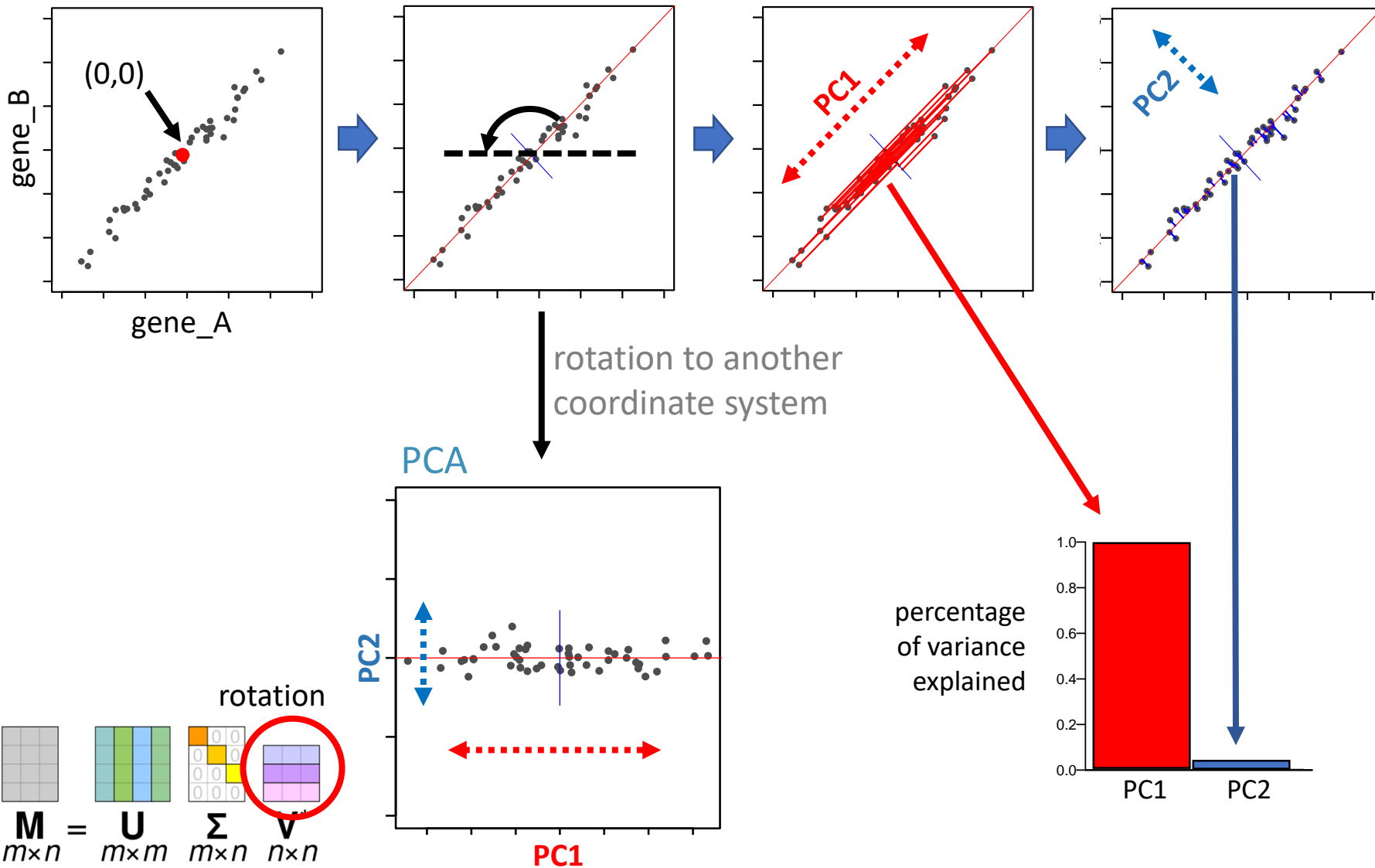
*Among the  $n$ -dimensional space, PCA will find the direction with highest variance*

*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)

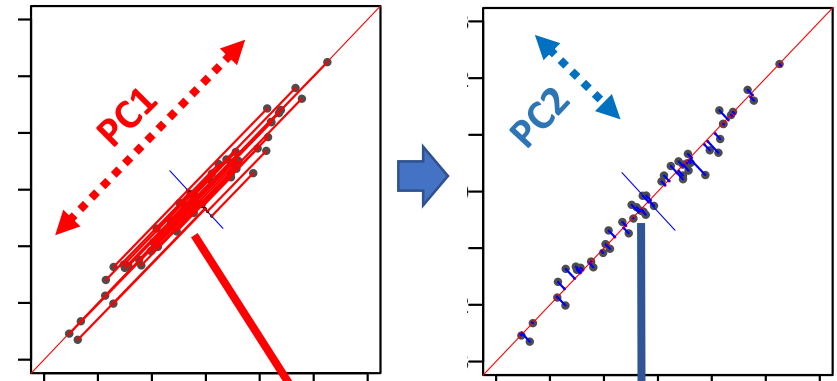


# How PCA works

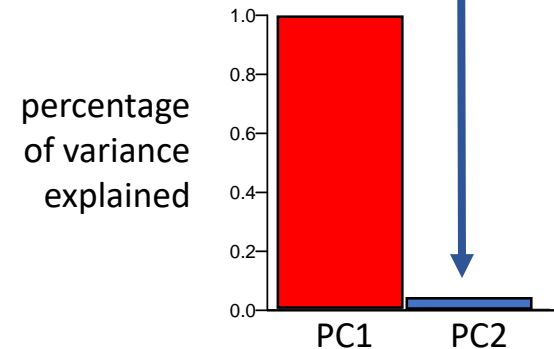
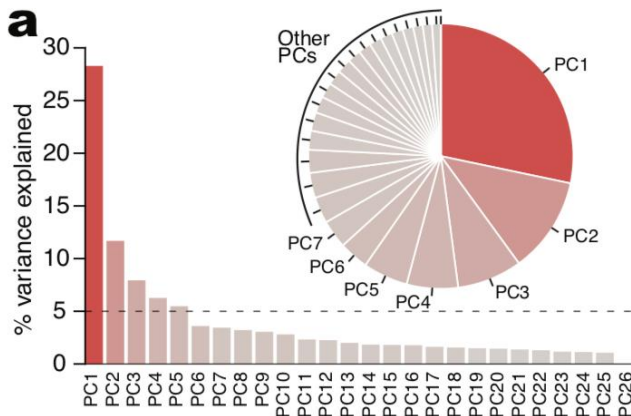
PC1 explains >98% of the variance

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

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

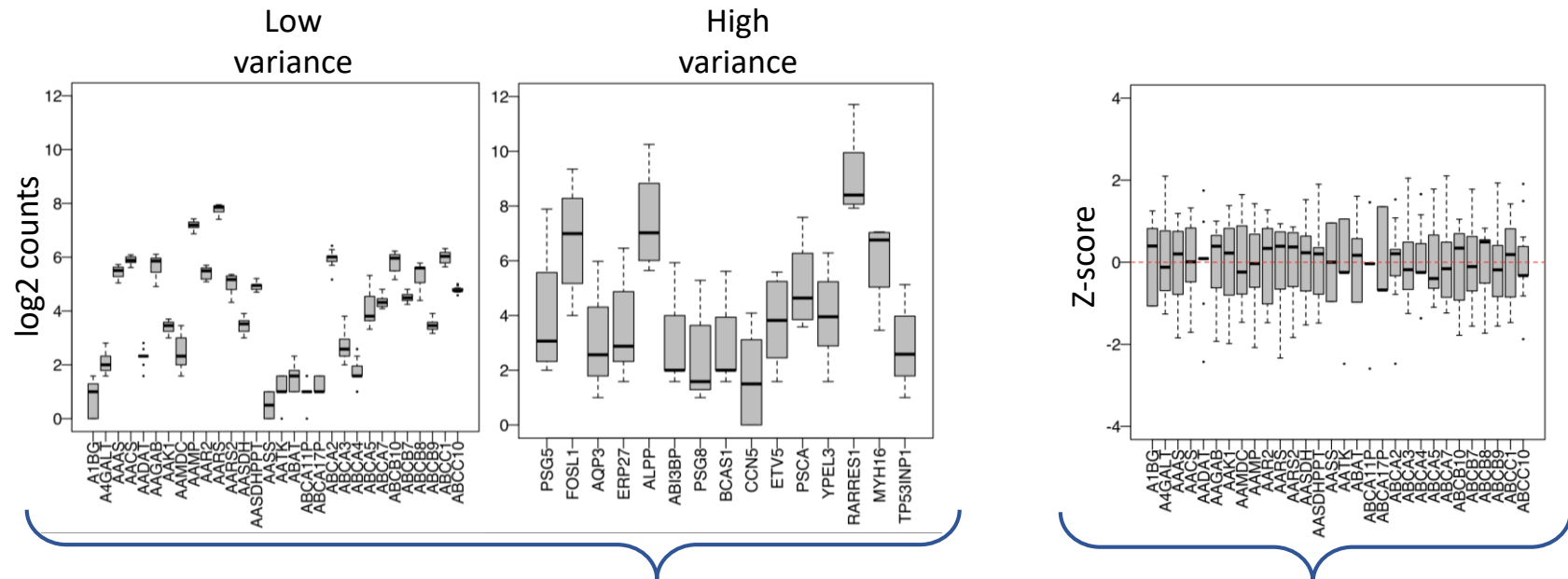


In real life ...



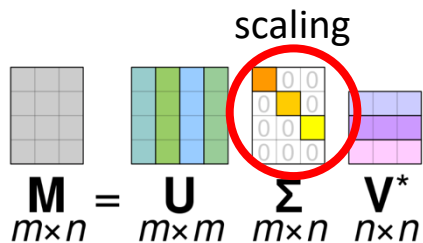
Before applying PCA, the data should be first transformed VST or RLOG

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



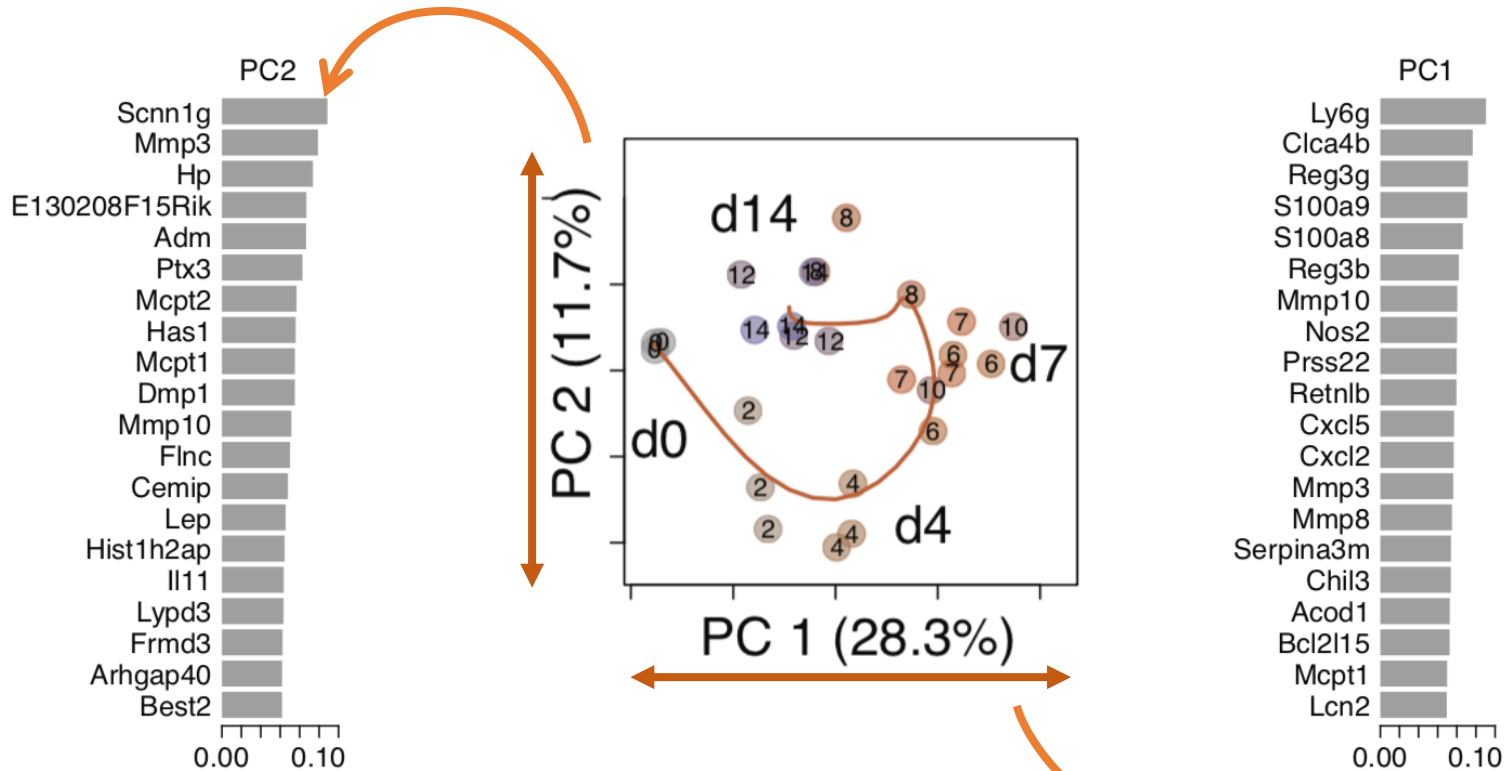
PCA on raw counts will separate genes with higher counts in the first PCs  
(higher distance to 0)

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

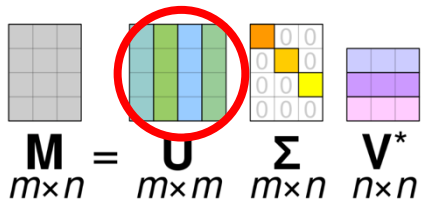


# How PCA works

## Interpretability of principal components



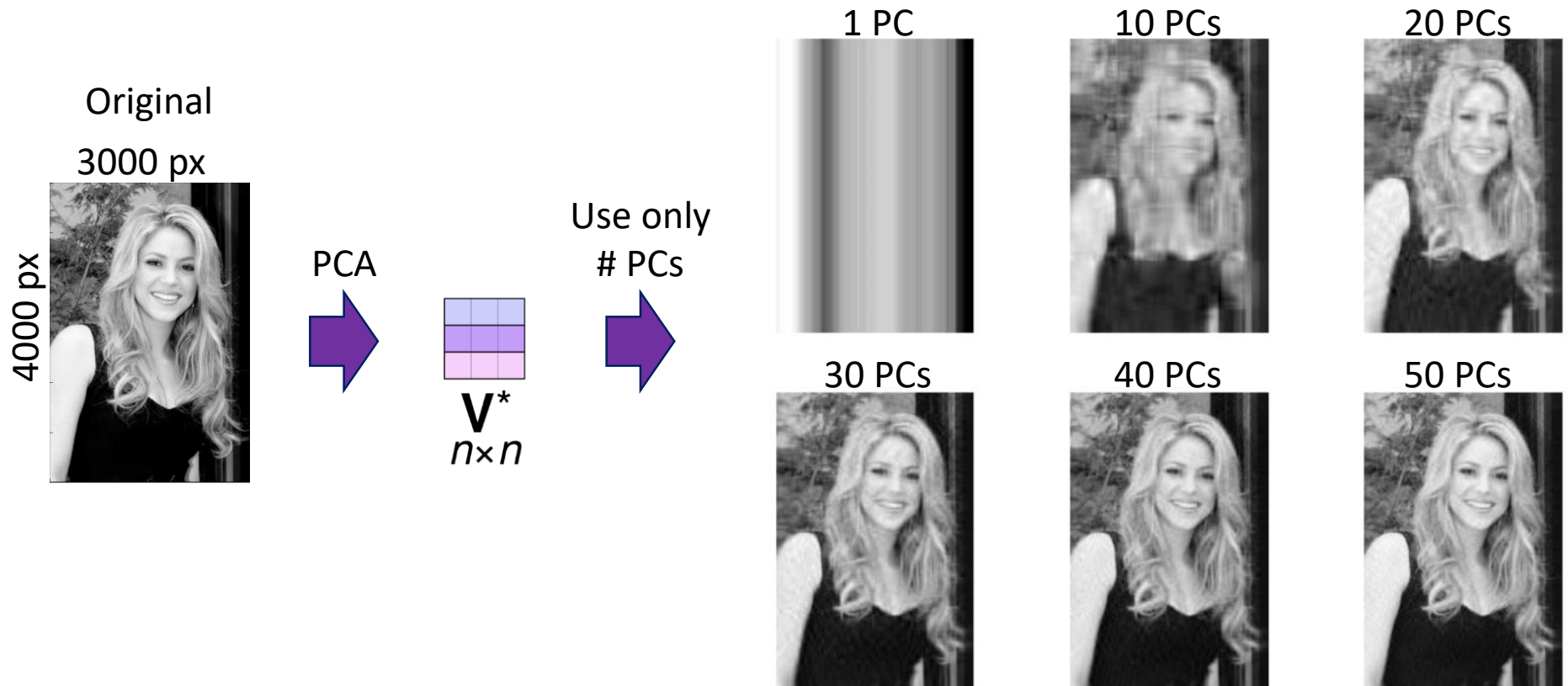
rotation





# A visual intuition of PCA

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



PCA is a dimensionality reduction method for numerical variables

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

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

It is an interpretable dimensionality reduction

The top principal components contain higher variance from the data and PCA preserves the whole variance

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

# To know more about PCA

If you are interested in a more comprehensive description of PCA, I would recommend reading Payam Emami's chapter on the topic:  
[https://payamemami.github.io/pca\\_basics/](https://payamemami.github.io/pca_basics/)



**Thank you. Questions?**

---