
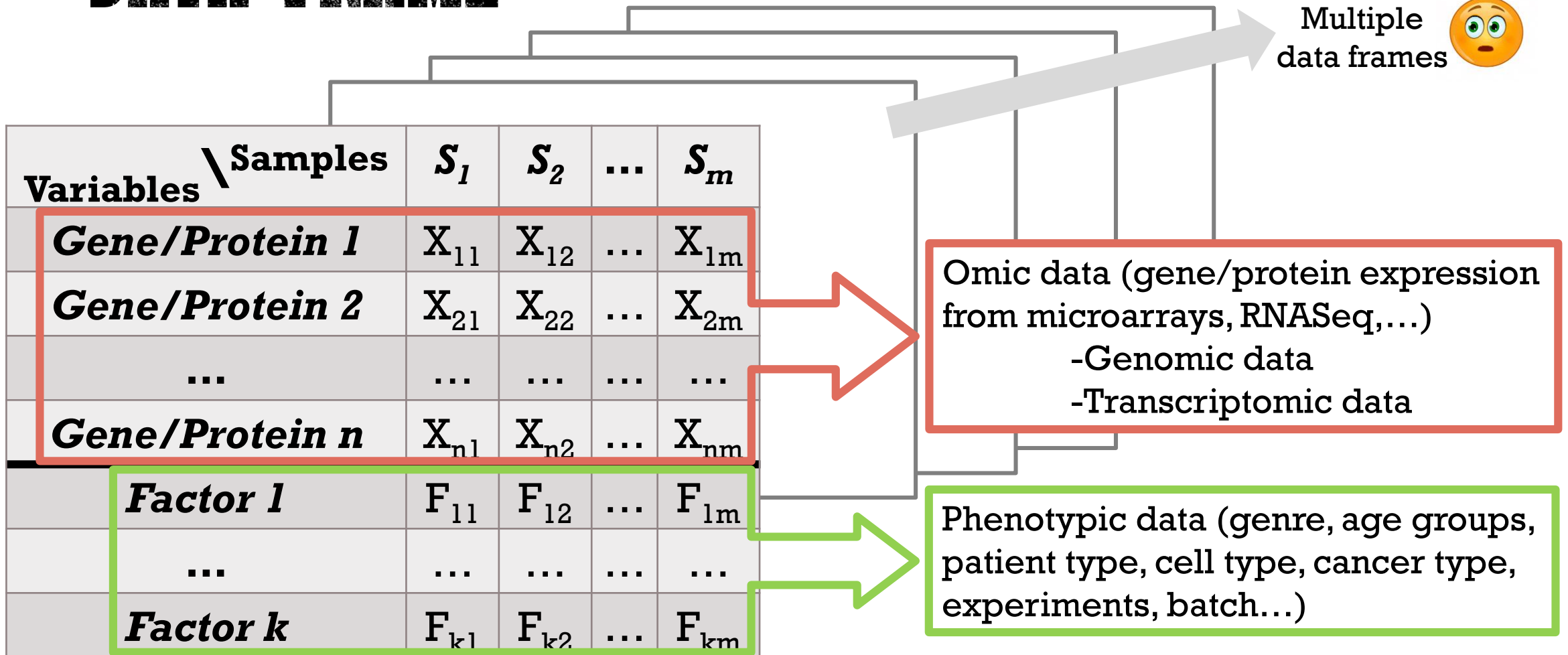


**HUFF POST**



CENTRO DE INVESTIGACIÓN  
DEL CÁNCER

# DATA FRAME



# OMIC DATA

- Gene expression values up to  $2^{20}$ , then log-scale.



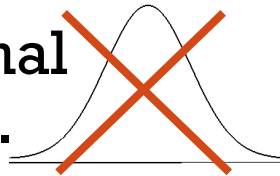
- Lots of genes with low expression.



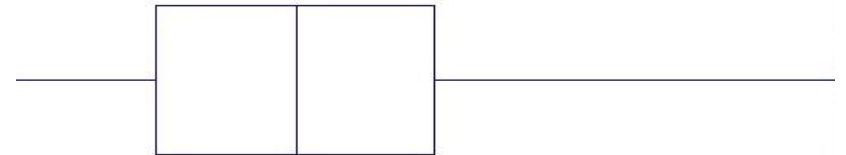
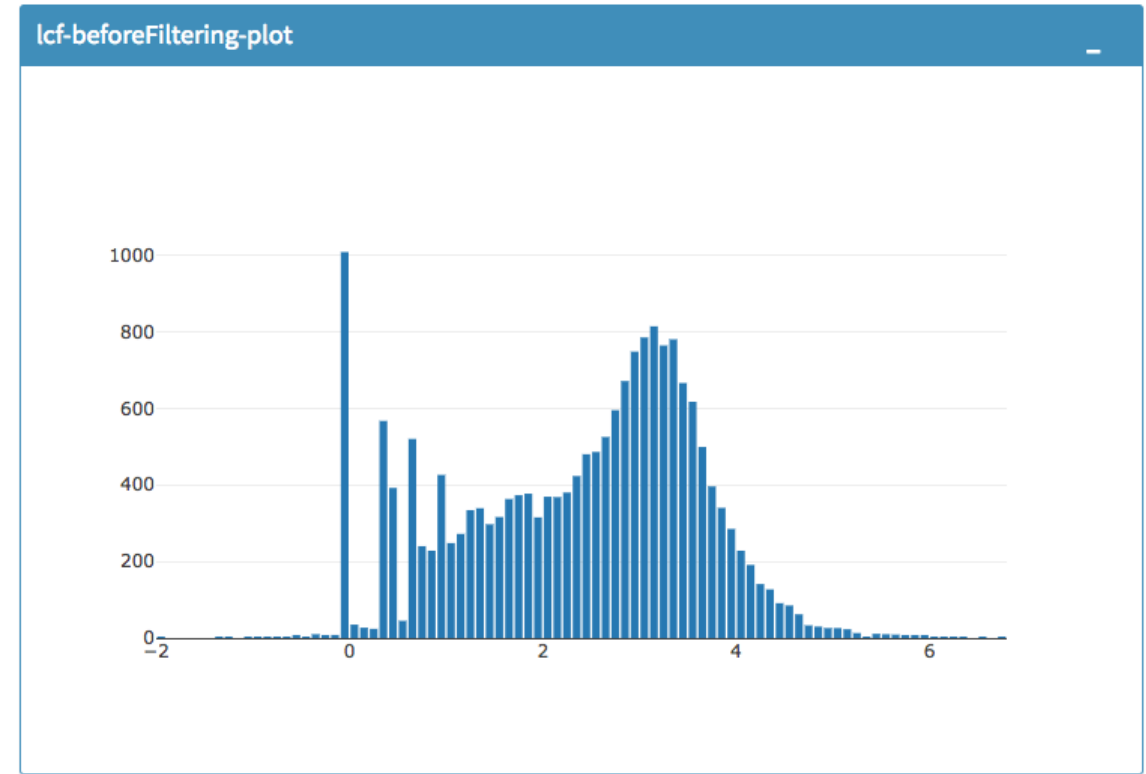
- Very few genes with high expression.



- Expression is not Normal (Gaussian) distributed.



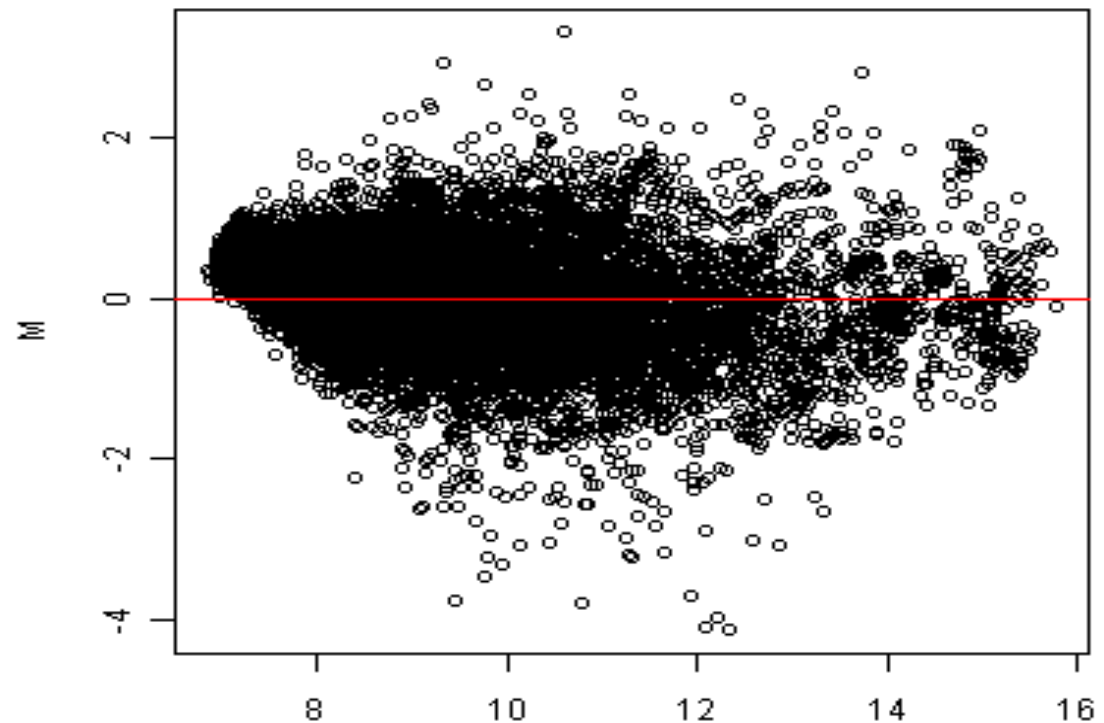
- Plots: histogram, density, boxplot, MA-plot...



# OMIC DATA



One microarray density plot



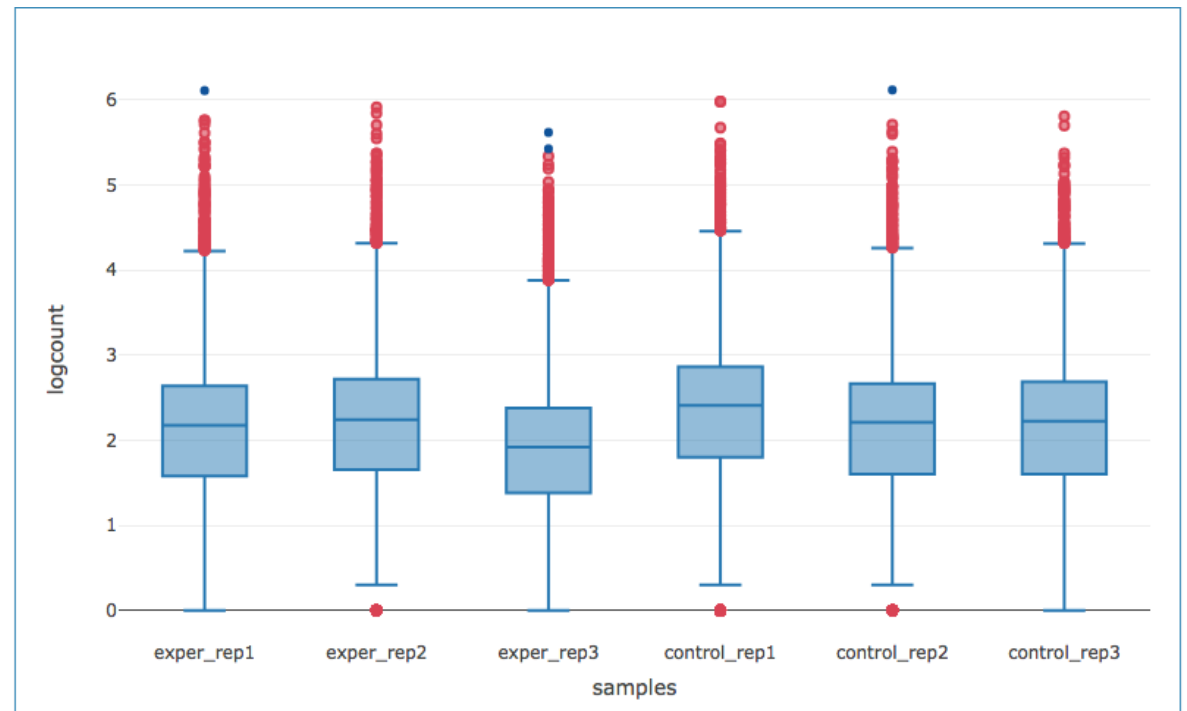
One microarray MA-plot (expresión vs mean)



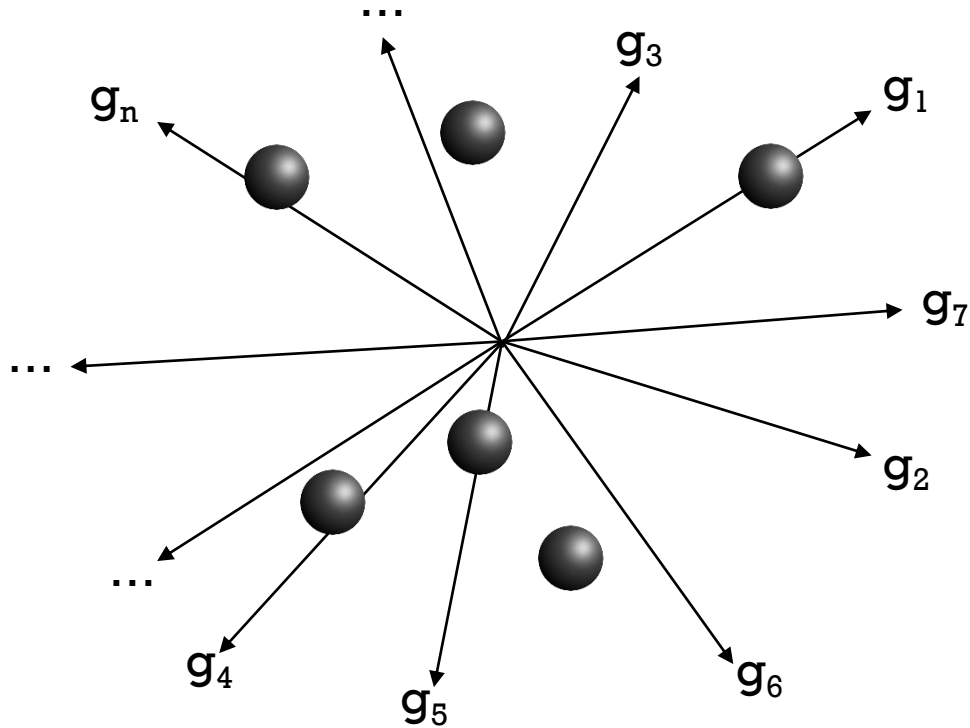
# PHENOTYPIC DATA

WARNING

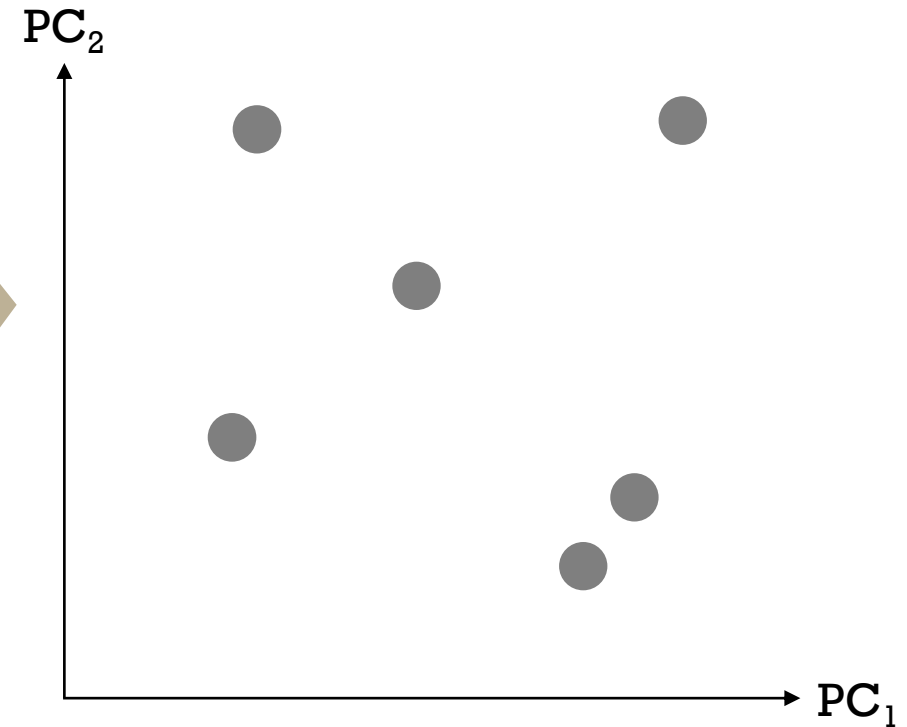
- ❑ Biological variability.
- ❑ Different experimental conditions.
- ❑ Samples from different sources.
- ❑ Join datasets from different experiments or batches.
- ❑ Plots: joint histograms / density / boxplots, PCA plot.



# PHENOTYPIC DATA



(imagine) n dimensional space,  
each ● is a sample



PCA: n dimensional space but plot  
of 2 axis explaining most variability



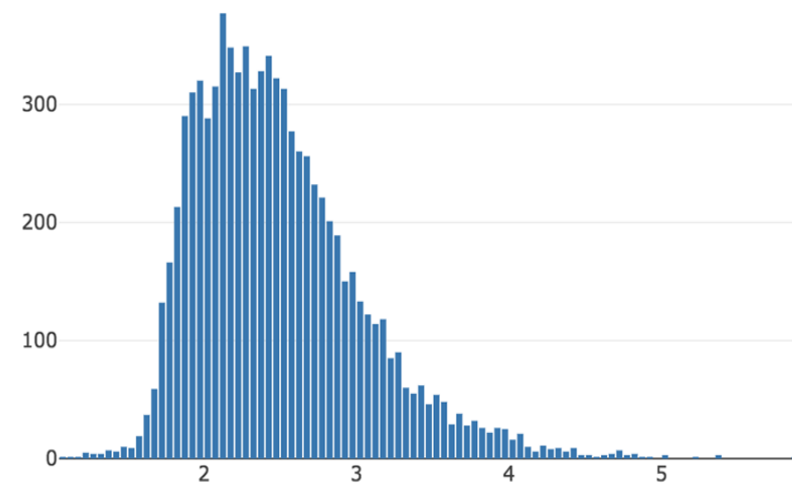
# QUALITY CONTROL AND CORRECTIONS

Usually remove low expressions (low count filtering, background correction...)

**C. Before Filtering**



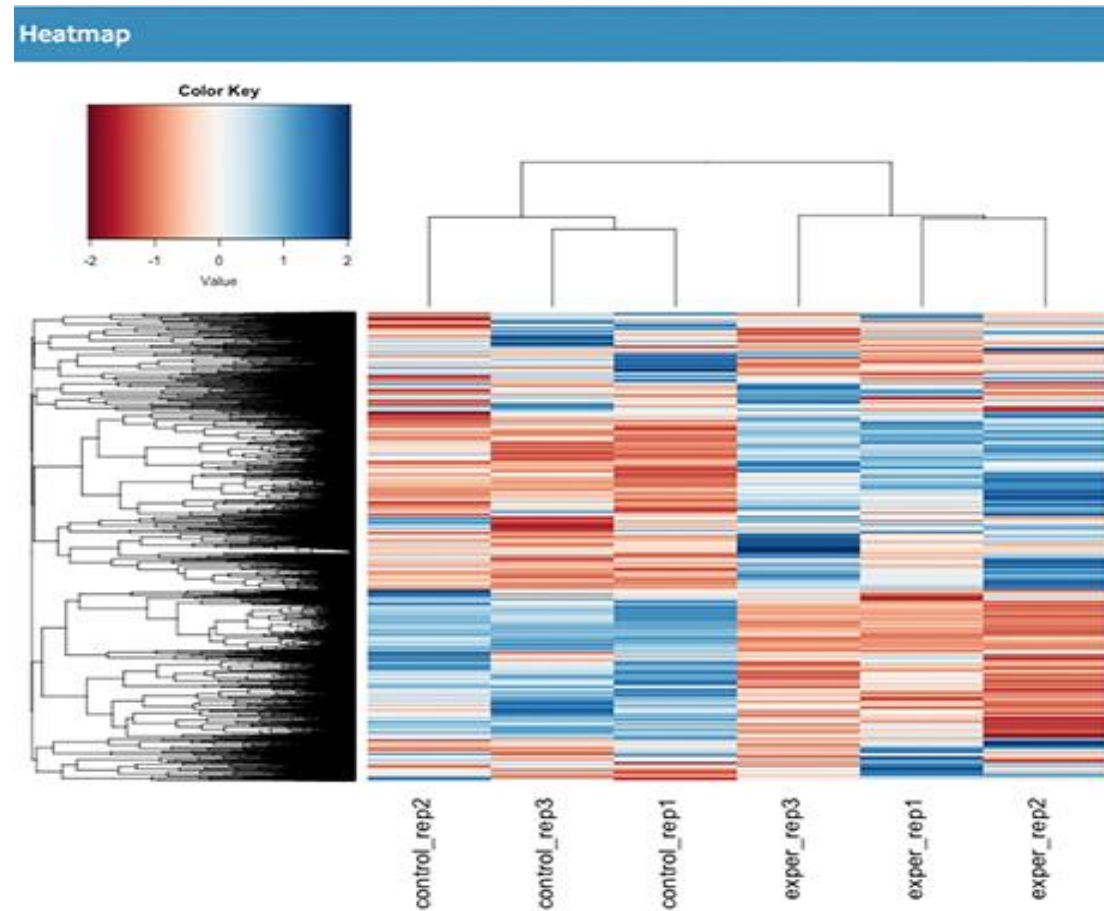
**D. After Filtering**



# QUALITY CONTROL AND CORRECTIONS

Heatmap:

- Search patterns.
- Search batch effects.

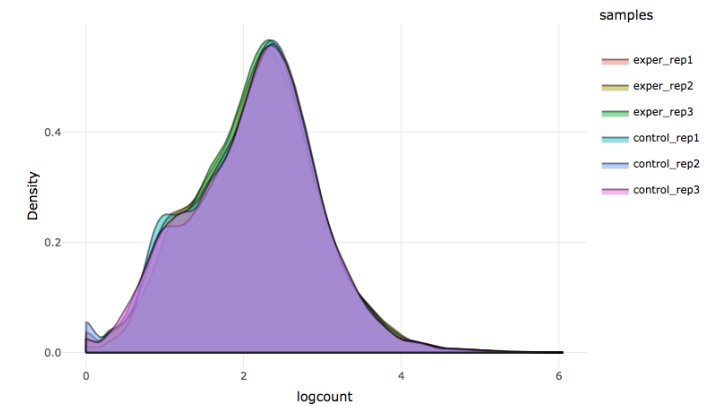
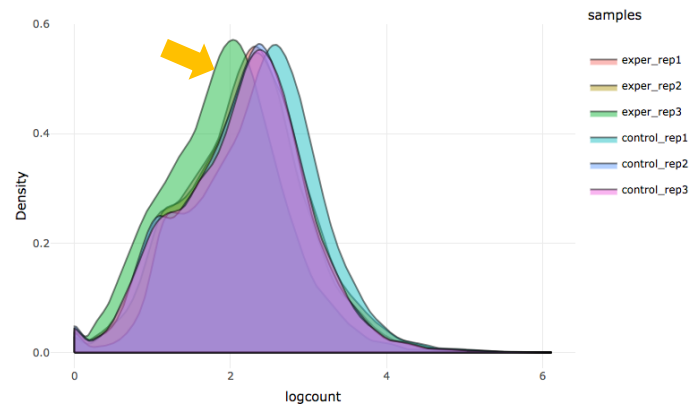
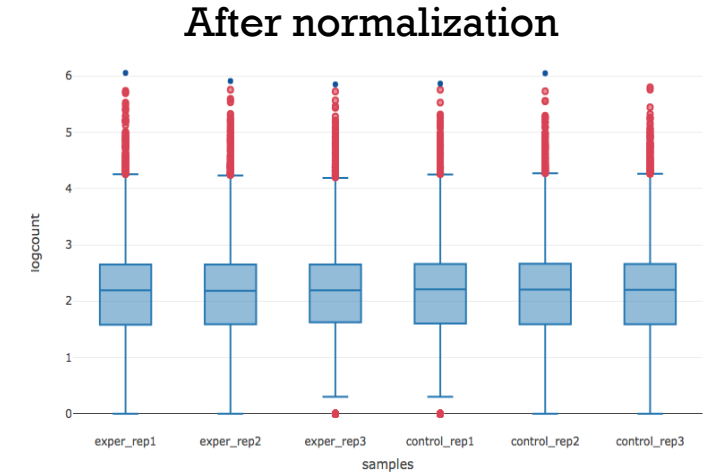
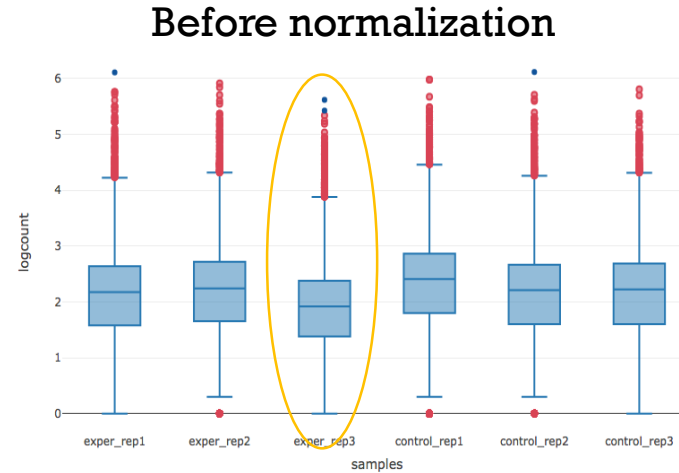




# QUALITY CONTROL AND CORRECTIONS

Normalize:

- Make the expression signal comparable.
- RMA, MAS5, fRMA...



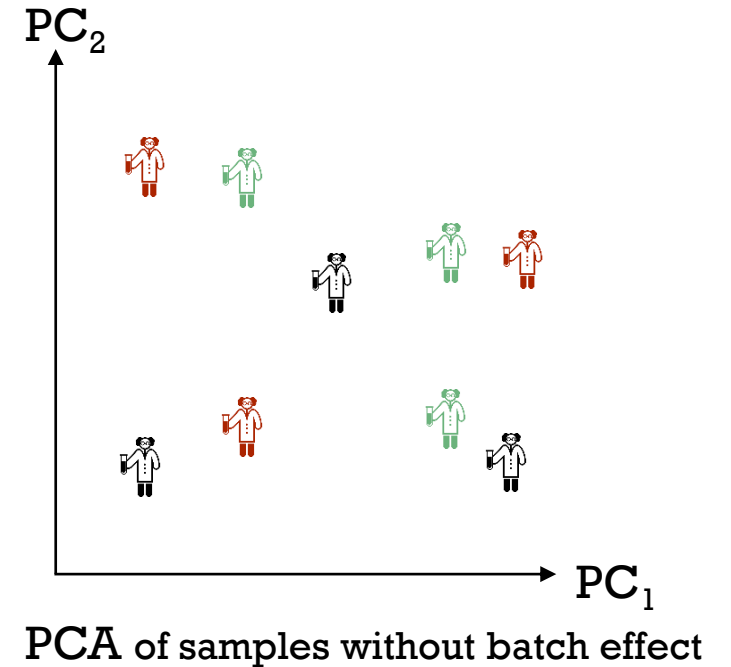
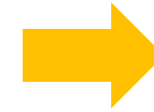
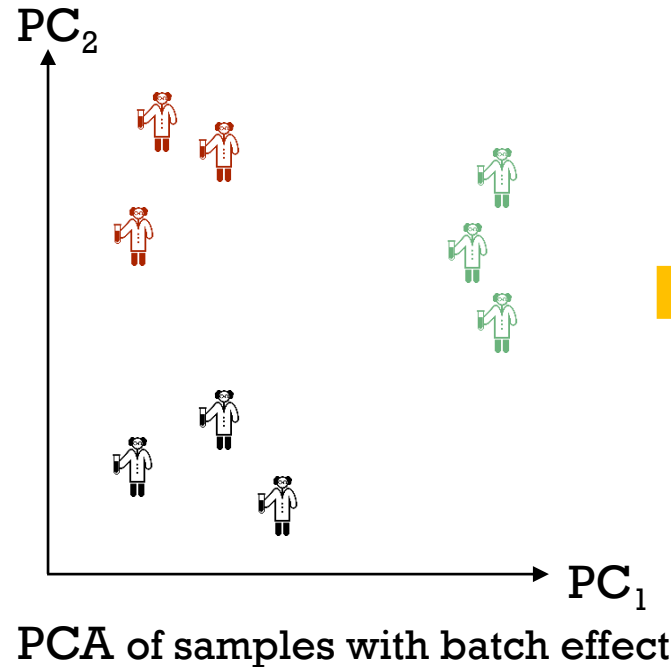
# QUALITY CONTROL AND CORRECTIONS

Batch effects correction:

- Correct the technical sources of variation, such as different processing times or different handlers.



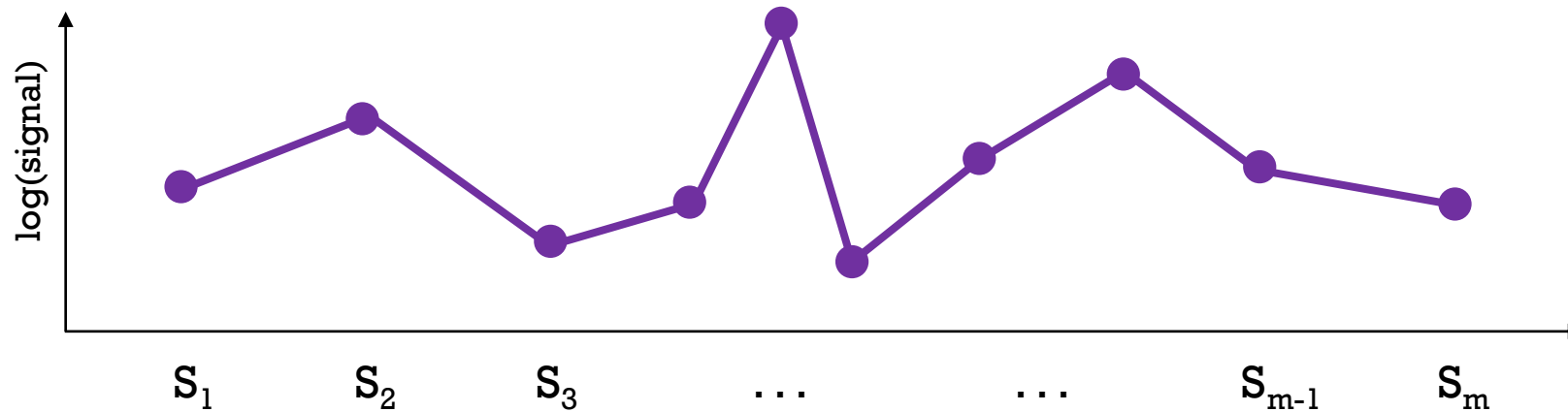
- Combat, f-RMA, TMM, RLE, upperquartile...



# DIFFERENTIAL EXPRESSION

*Gene/Protein expression profile*: expression signal in each sample

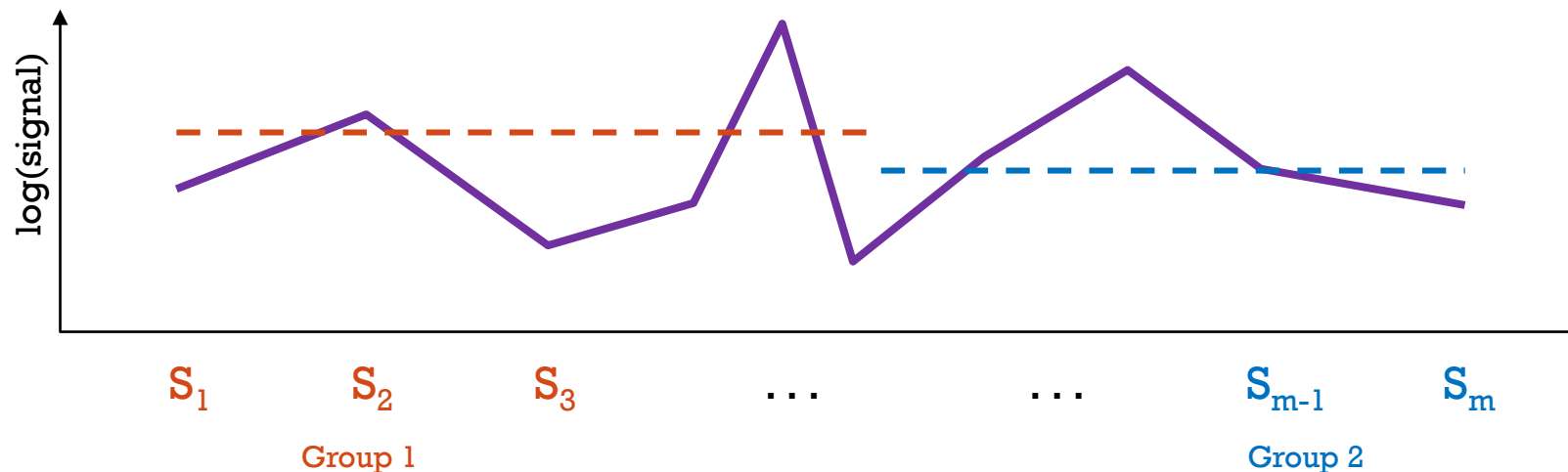
	$s_1$	$s_2$	$s_3$	...	$s_{m-1}$	$s_m$
<b>Gene/Protein</b>	$x_{11}$	$x_{12}$	$x_{13}$	...	$x_{1(m-1)}$	$x_{1m}$



# DIFFERENTIAL EXPRESSION

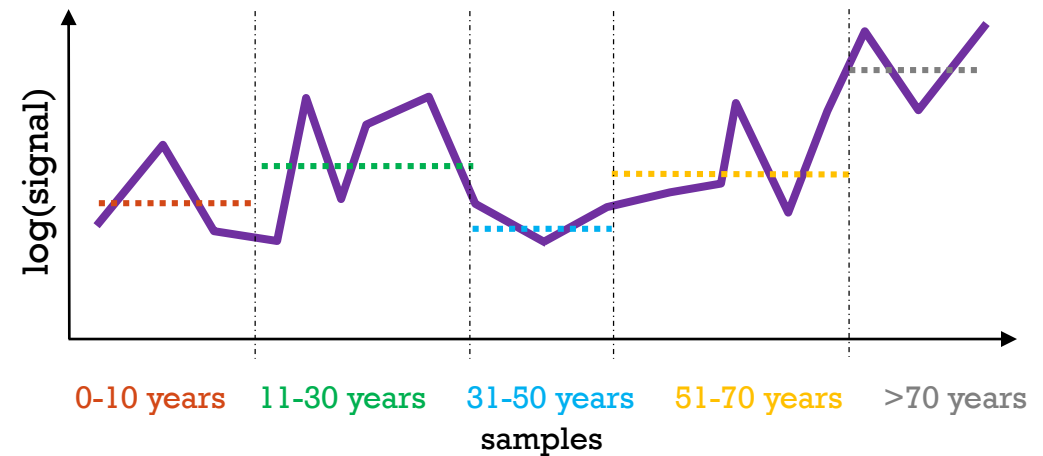
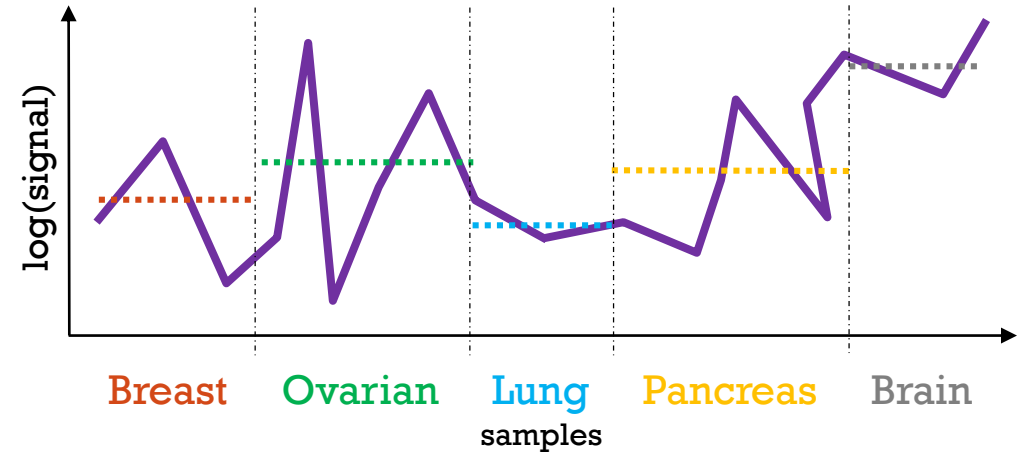
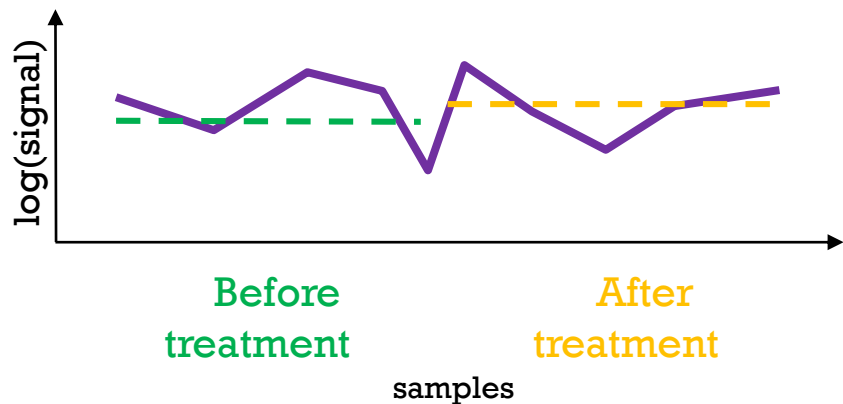
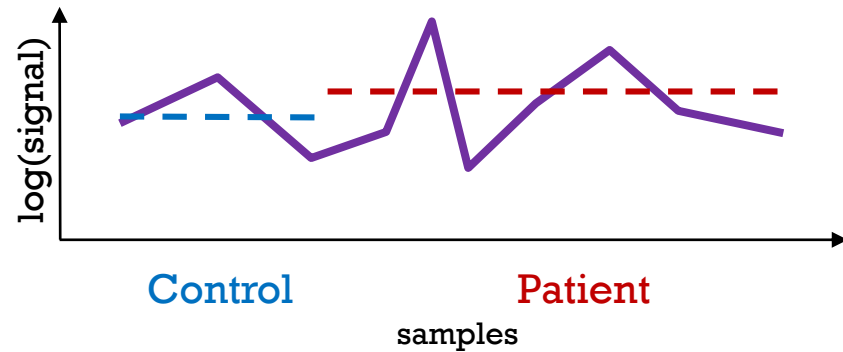
*Differential expression analysis*: expression signal comparison between groups

	$S_1$	$S_2$	$S_3$	...	$S_{m-1}$	$S_m$
<b>Gene/Protein</b>	$X_{11}$	$X_{12}$	$X_{13}$	...	$X_{1(m-1)}$	$X_{1m}$
<b>Factor/Group</b>	$F_1$	$F_2$	$F_3$	...	$F_{m-1}$	$F_m$



# DIFFERENTIAL EXPRESSION

Some scenarios:

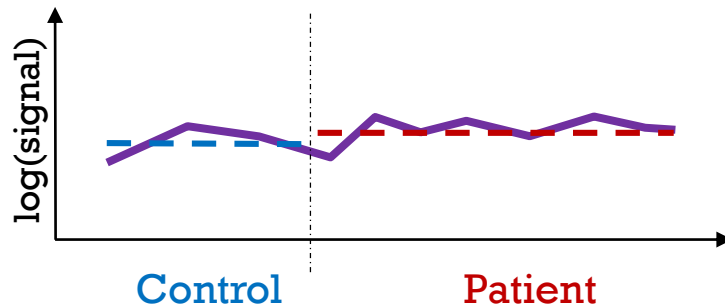


# DIFFERENTIAL EXPRESSION

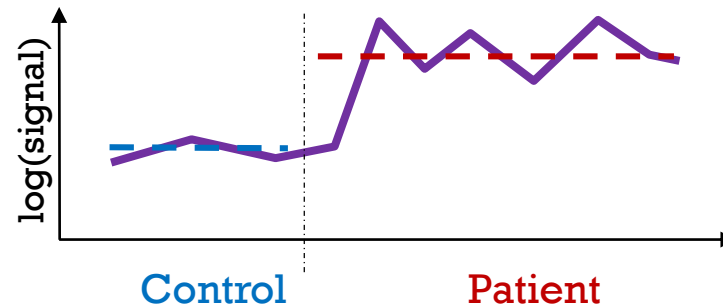
In statistics it is a hypothesis contrast:

$H_0$ : The gene does not change (not DE)

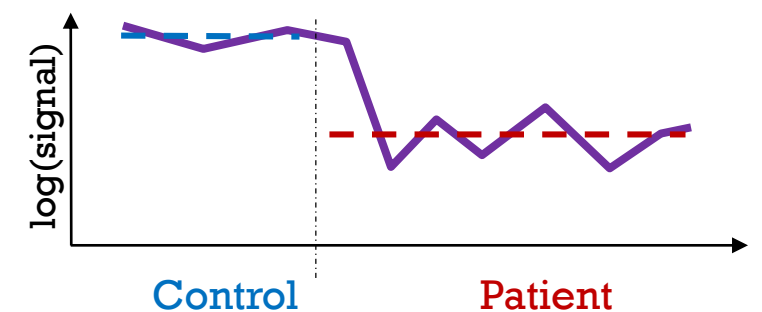
$H_1$ : The gene changes (there are significant differences between groups)



$\text{signal(C)} \approx \text{signal(P)}$   
no significant gen



$\text{signal(C)} < \text{signal(P)}$   
over-expressed gen



$\text{signal(C)} > \text{signal(P)}$   
under-expressed gen



# DIFFERENTIAL EXPRESSION

Objective:

- Find genes that change between groups.



Over-expressed genes → the disease activates it



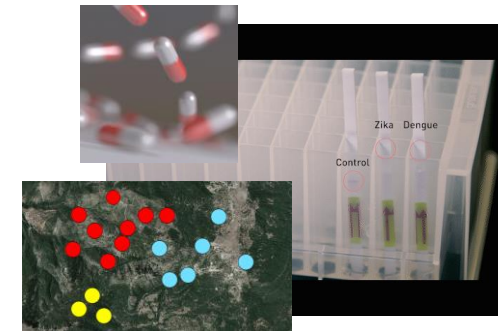
Under-expressed genes → the disease represses it



OFF

\* Examples:

- Use this information to activate or repress genes with drugs to fight the disease.
- Find genetic signatures, groups of genes that represent a disease, a type of cancer, a group of individuals...



# ACKNOWLEDGMENTS

## *Group CIC-IBMCC-USAL*

Javier De Las Rivas

Diego Alonso

Alberto Berral

Fernando Bueno

Santiago Bueno

Óscar González

Elena Sánchez

José Manuel Sánchez

