

Welcome to the practical session **"big data & predictive models"**

Loic Verlingue
ISFBM module 12
January 23th 2019

Running instructions

- Open R studio
- Open New File -> R Notebook
- Go to the web https://github.com/gustaveroussy/IFSBM-bigdata/edit/master/TP_IFSBM_module12/
- Copy paste the content of TP_notebookR.Rmd in your new R Notebook
- Follow the notebook instructions
- *Enjoye!*

For very advanced R programmers

- You can erase the R code in the Notebook from the section:

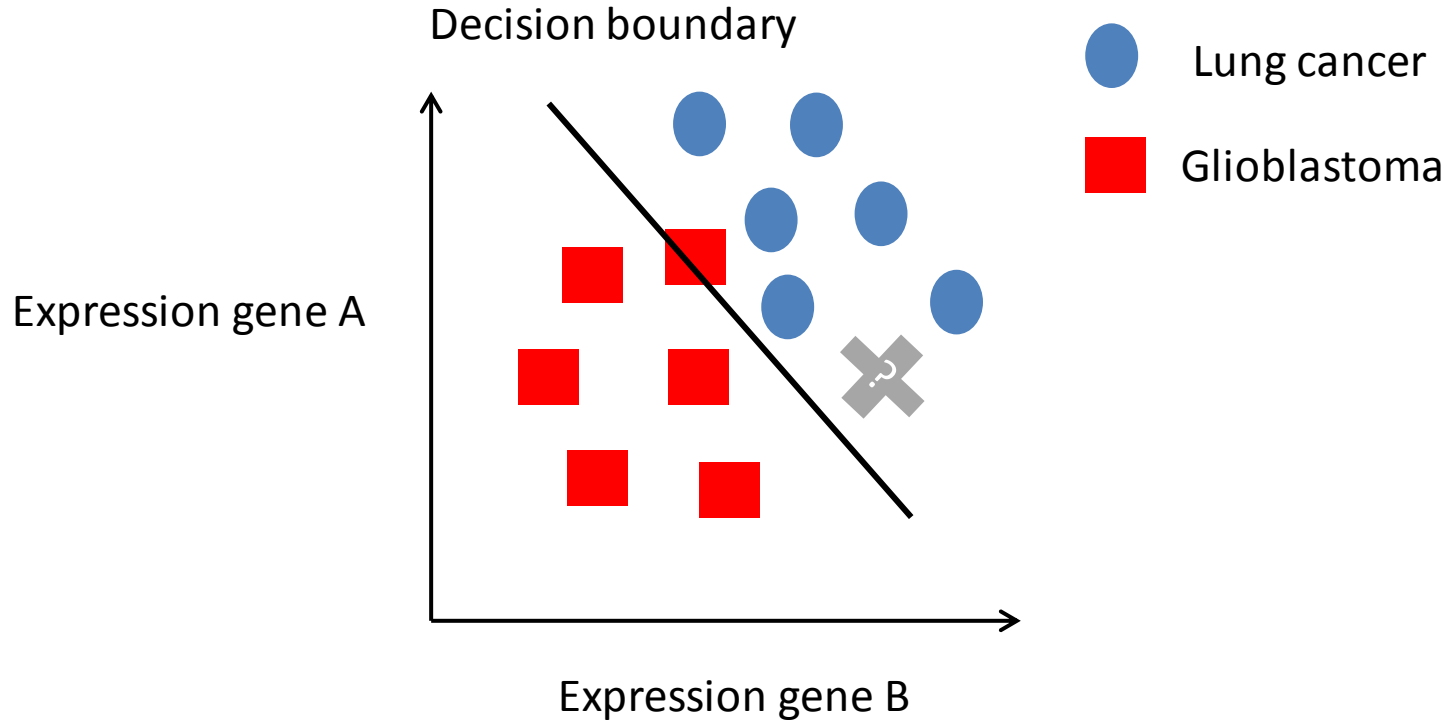
« Description of the data »

To

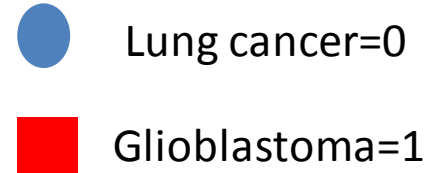
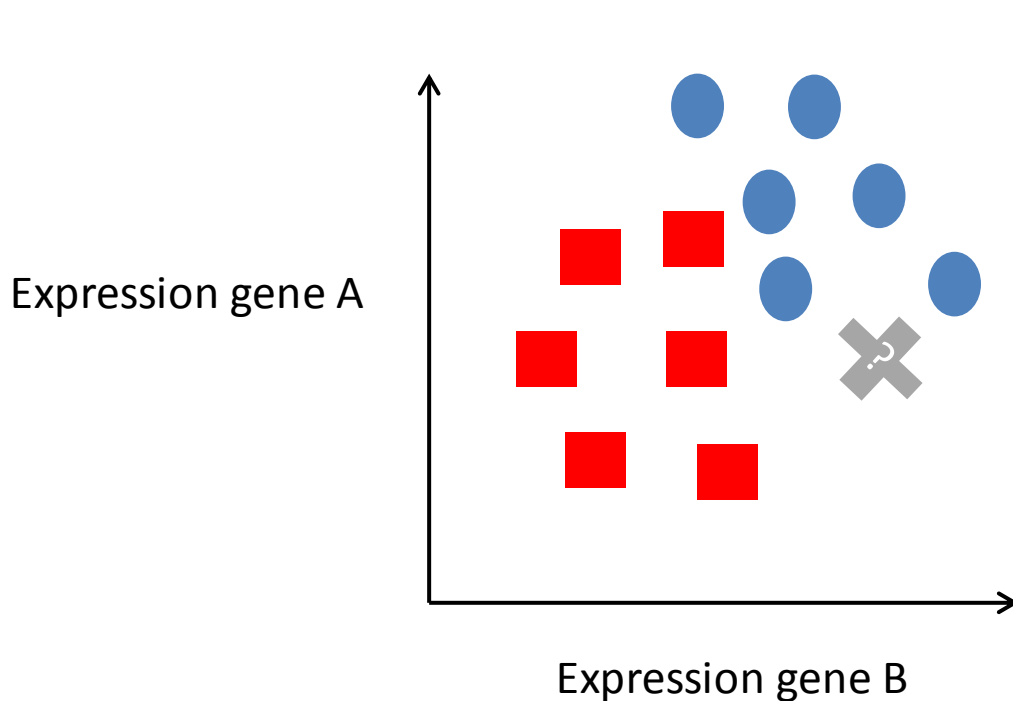
« A glimpse into survival analysis »

- And try to do the code by your own with help of the original Notebook

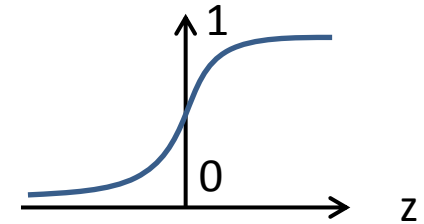
Classification



Logistic regression in 2D

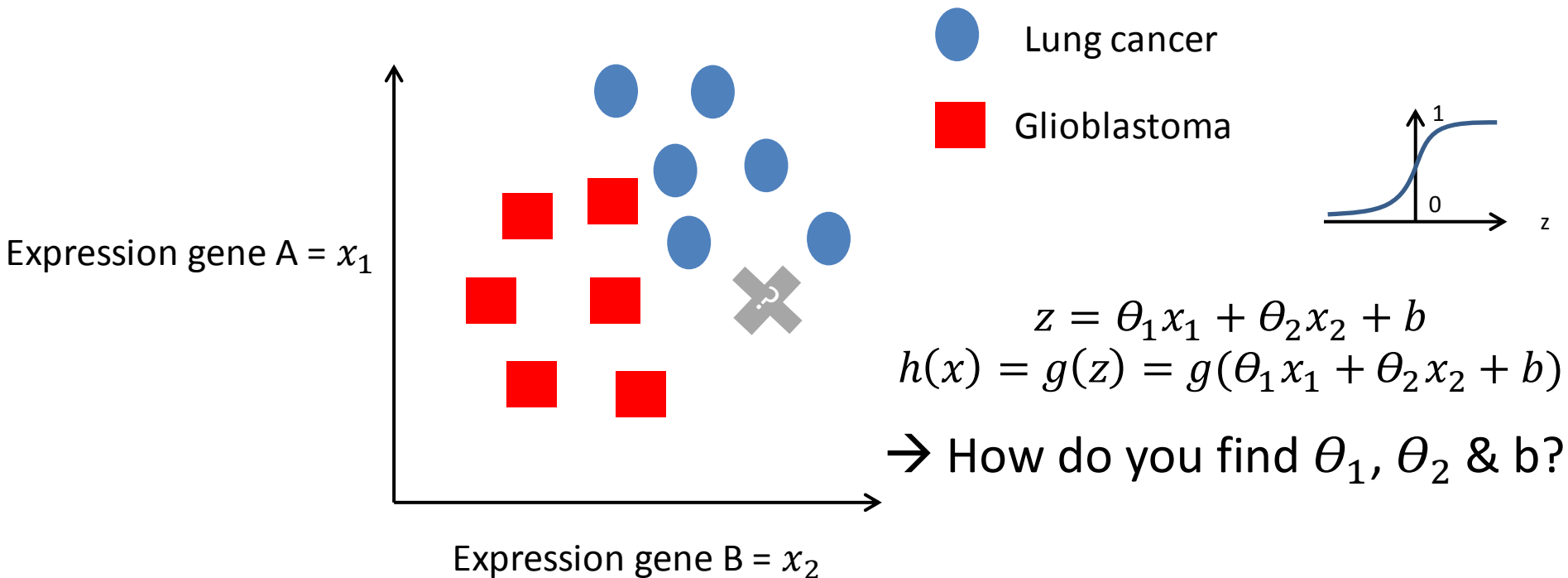


Sigmoid function = logistic function

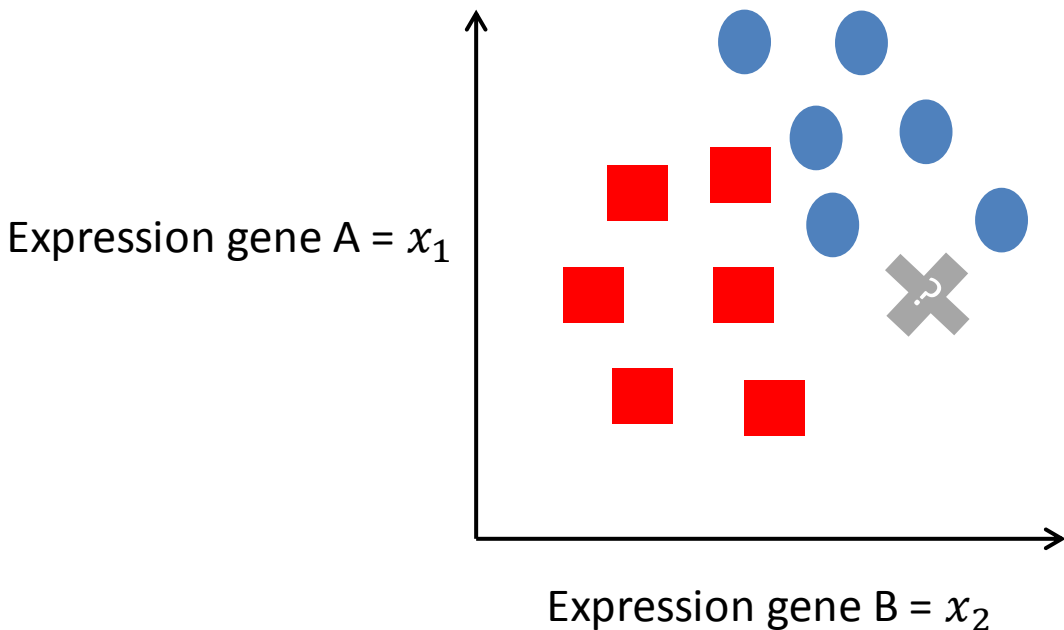


$$g(z) = \frac{1}{1 + e^{-z}}$$

Logistic regression in 2D



Logistic regression in 2D



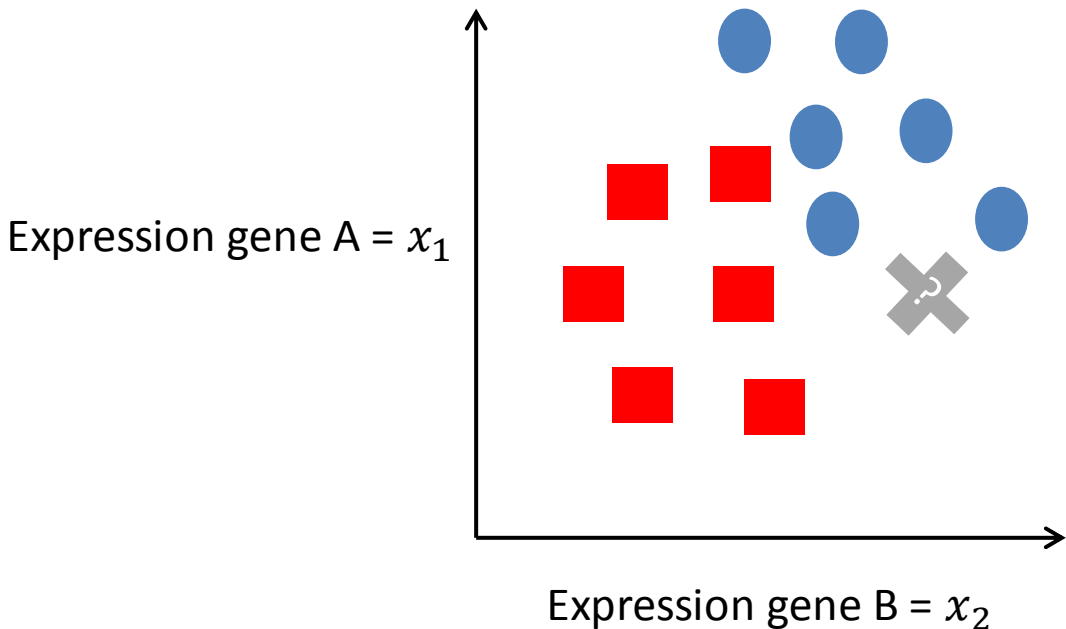
Define a cost function

y^i is your ground truth
 $h(x^i)$ is your prediction

$$\rightarrow \text{cost}(h(x^i), y^i)$$

- If $y^i = 0$ (lung adenocarcinoma) you want $h(x^i)$ to be ~ 0
- If $y^i = 1$ (glioblastoma) you want $h(x^i)$ to be ~ 1

Logistic regression in 2D



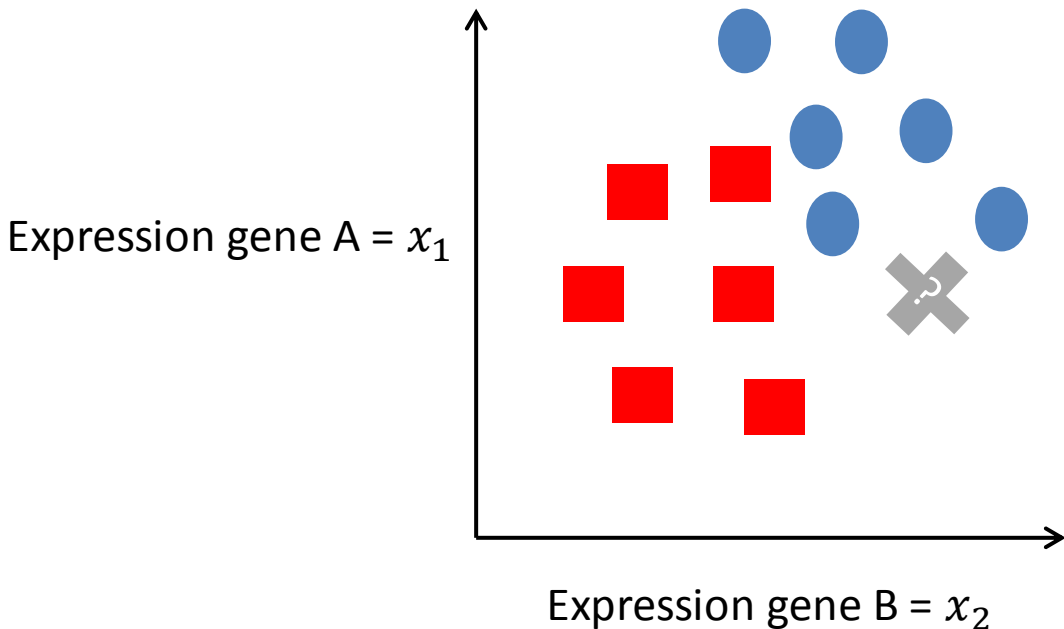
For example

$$\text{cost}(h(x^i), y^i)$$

● $-\log(1 - h(x^i))$ if $y^i = 0$

■ $-\log(h(x^i))$ if $y^i = 1$

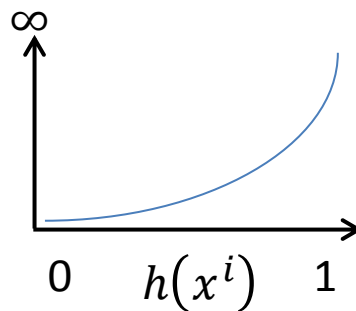
Logistic regression in 2D



For example

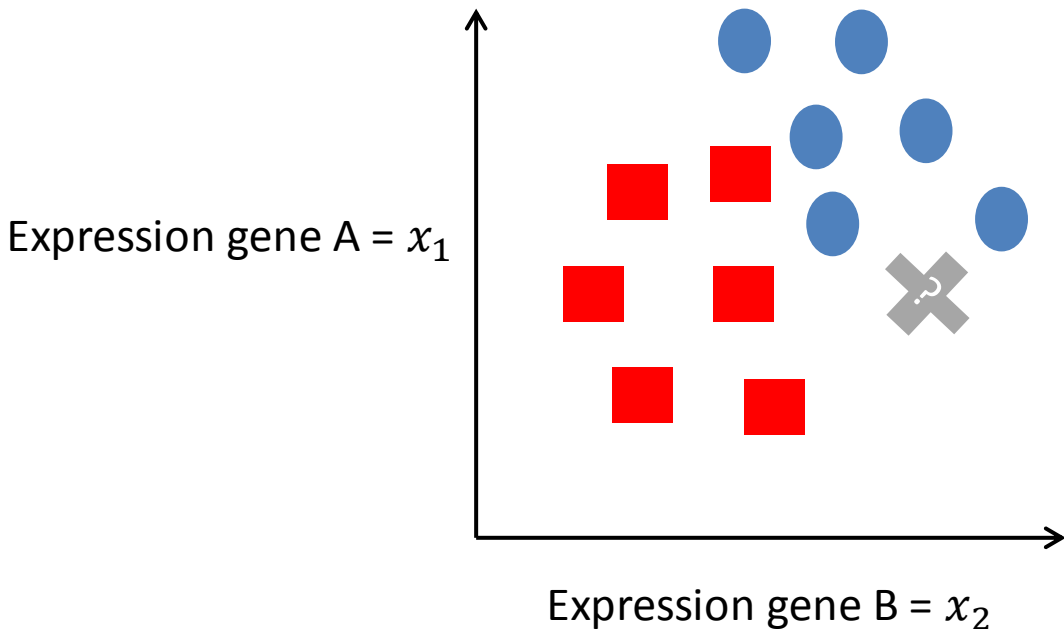
$$\text{cost}(h(x^i), y^i)$$

if $y^i = 0$



$$\text{cost} = -\log(1 - h(x^i))$$

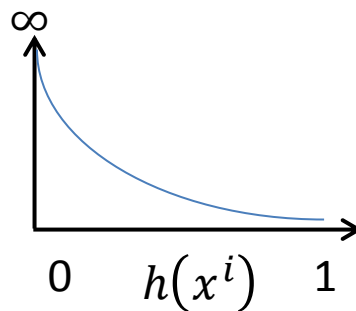
Logistic regression in 2D



For example

$$\text{cost}(h(x^i), y^i)$$

if $y^i = 1$



$$\text{cost} = -\log(h(x^i))$$

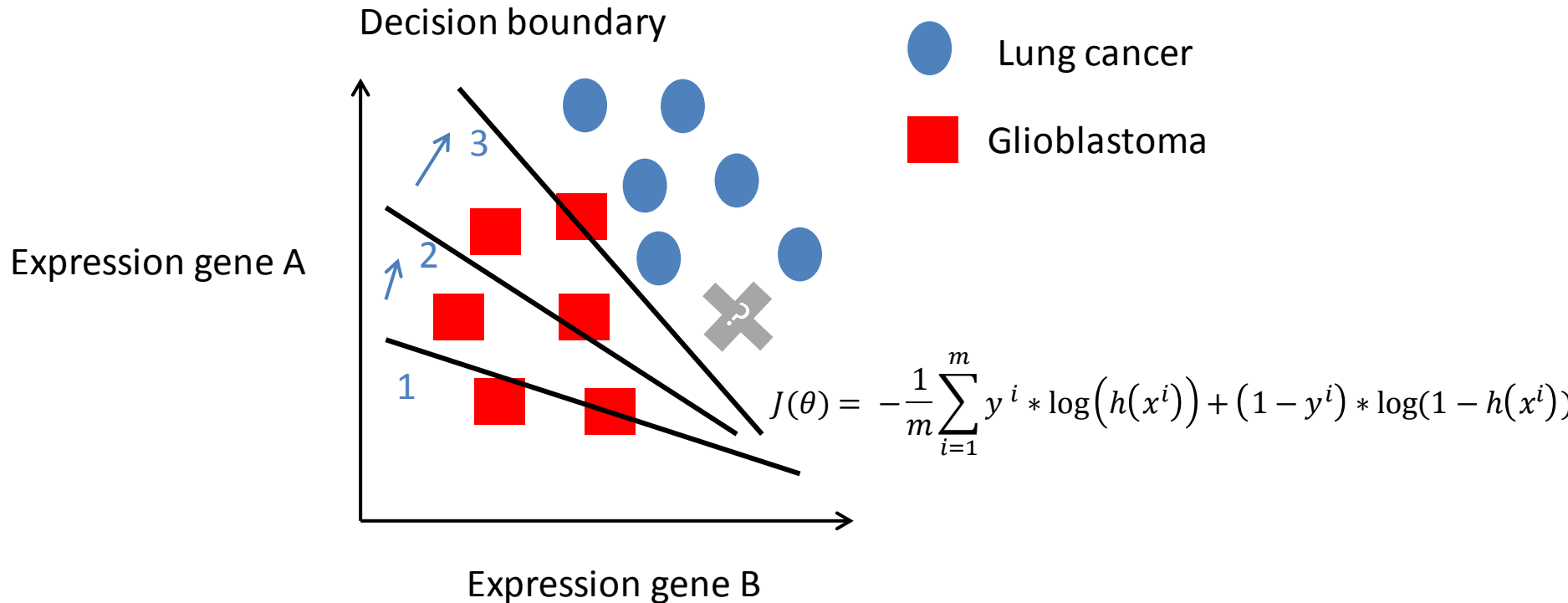
Logistic regression in 2D

Do it on all your data

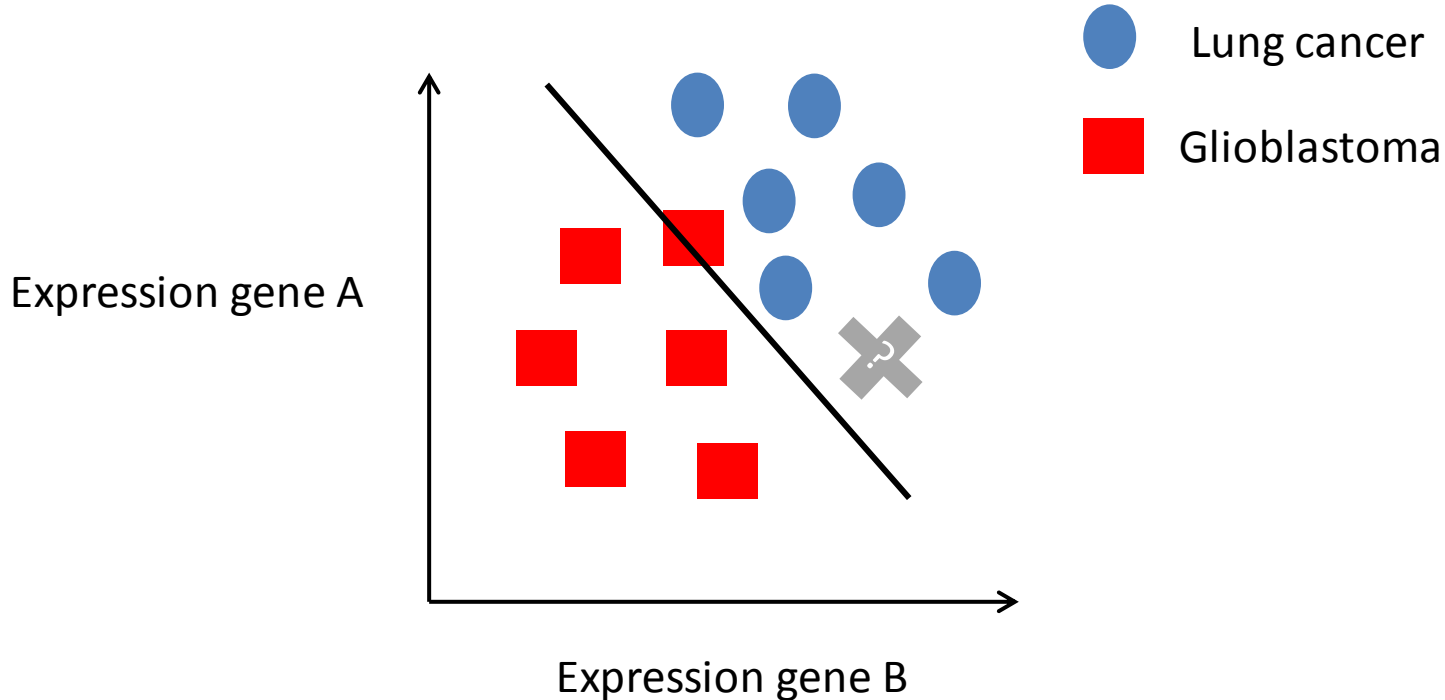
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h(x^i), y^i)$$
$$= -\frac{1}{m} \sum_{i=1}^m y^i * \log(h(x^i)) + (1 - y^i) * \log(1 - h(x^i))$$

And then minimize it!

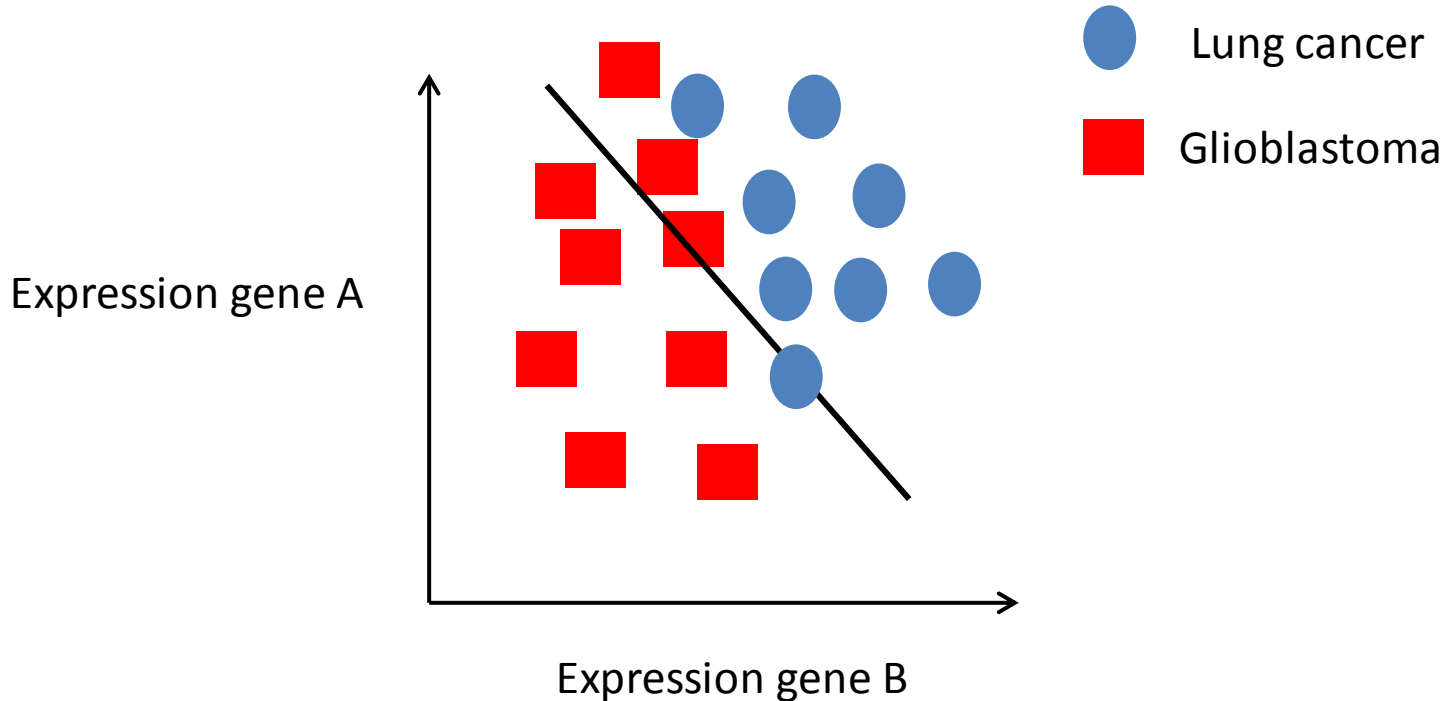
Minimization with gradient descent



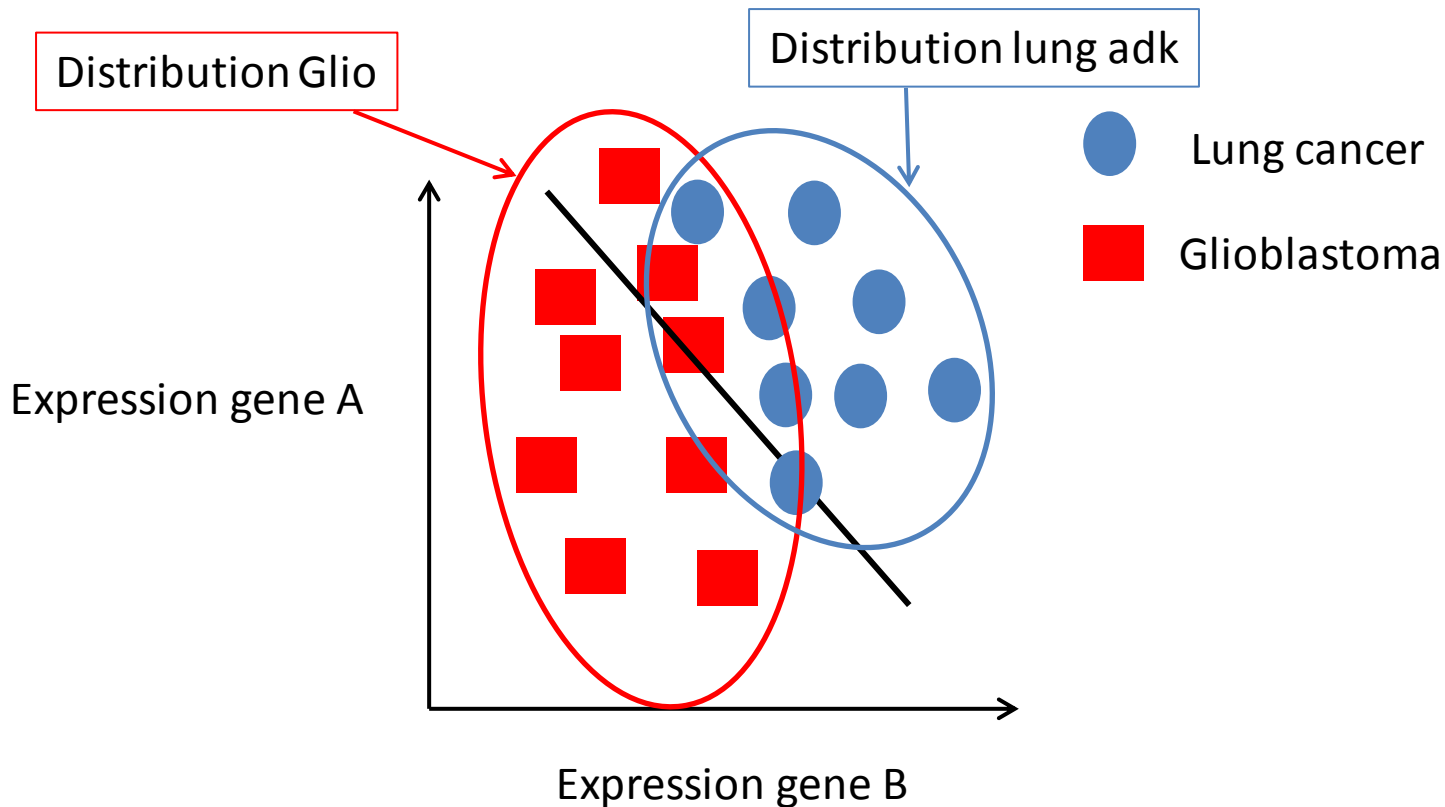
Are you write? Probably yes !



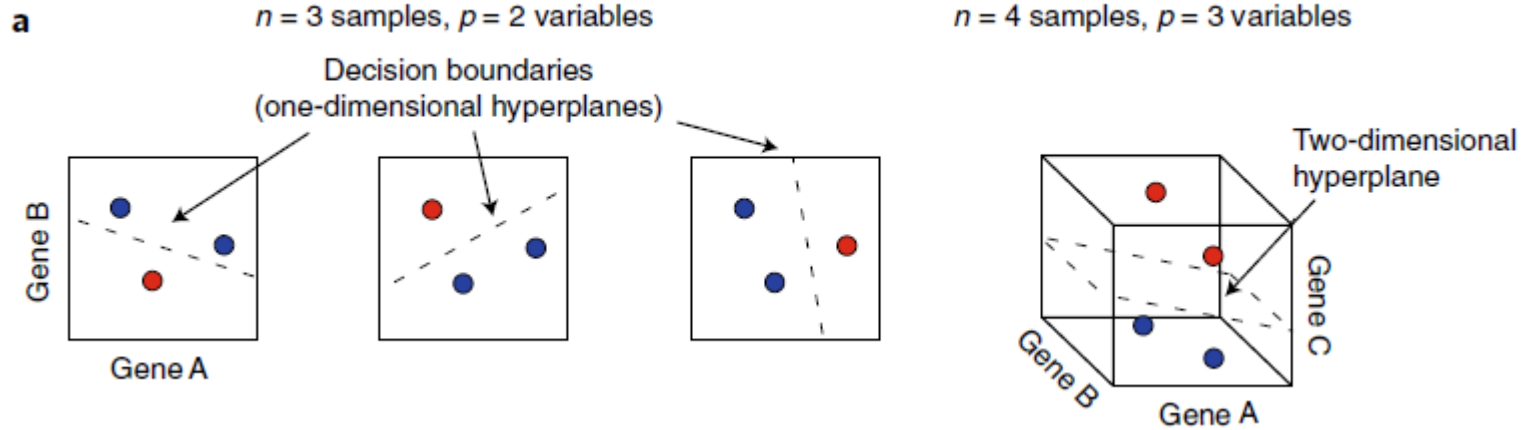
Are you write? No !



Are you write? No !



Generalisation and overfitting



The curse of dimensionality!

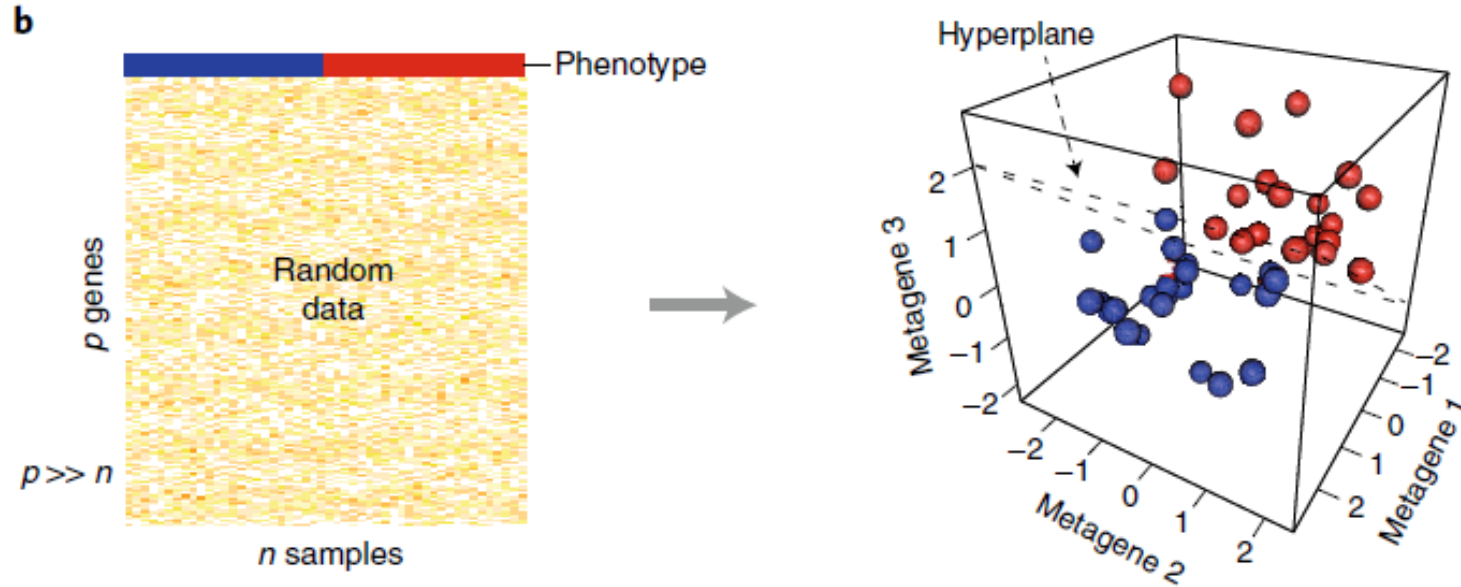
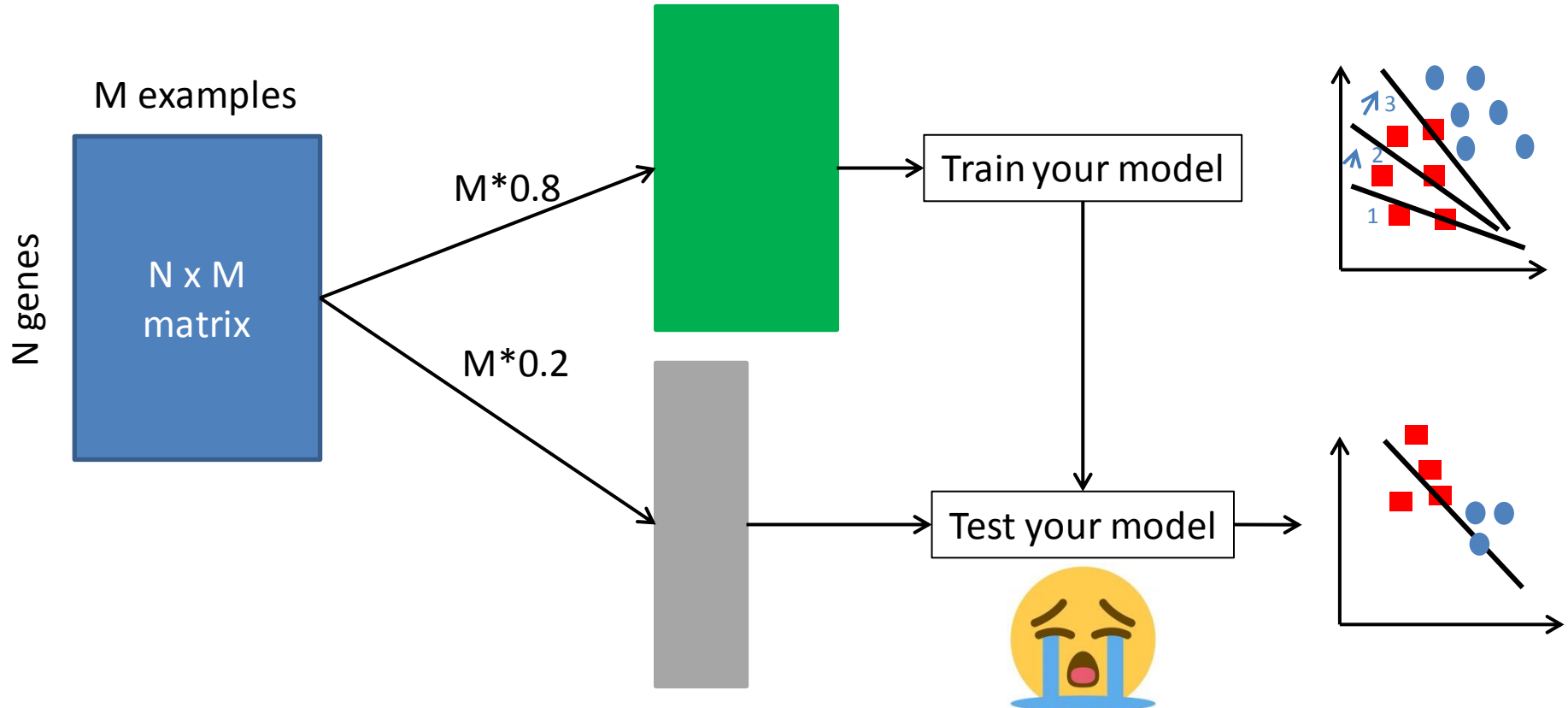


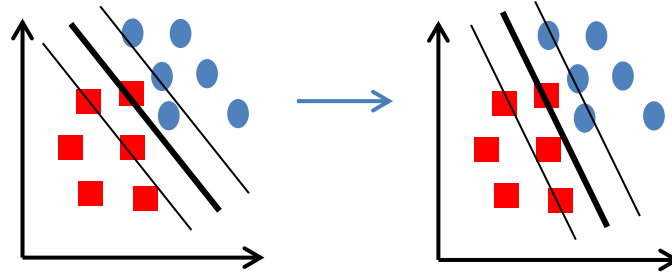
Fig. 1 | The curse of dimensionality and overfitting. a, Low-dimensional examples designed to

How do you evaluate generalization?



Are there strategies not to 🥲 ?

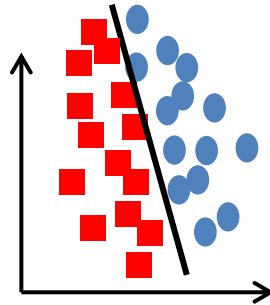
On the modeling / learning:



Eg: SVM,
regularisation /
penalisation..



On the data:

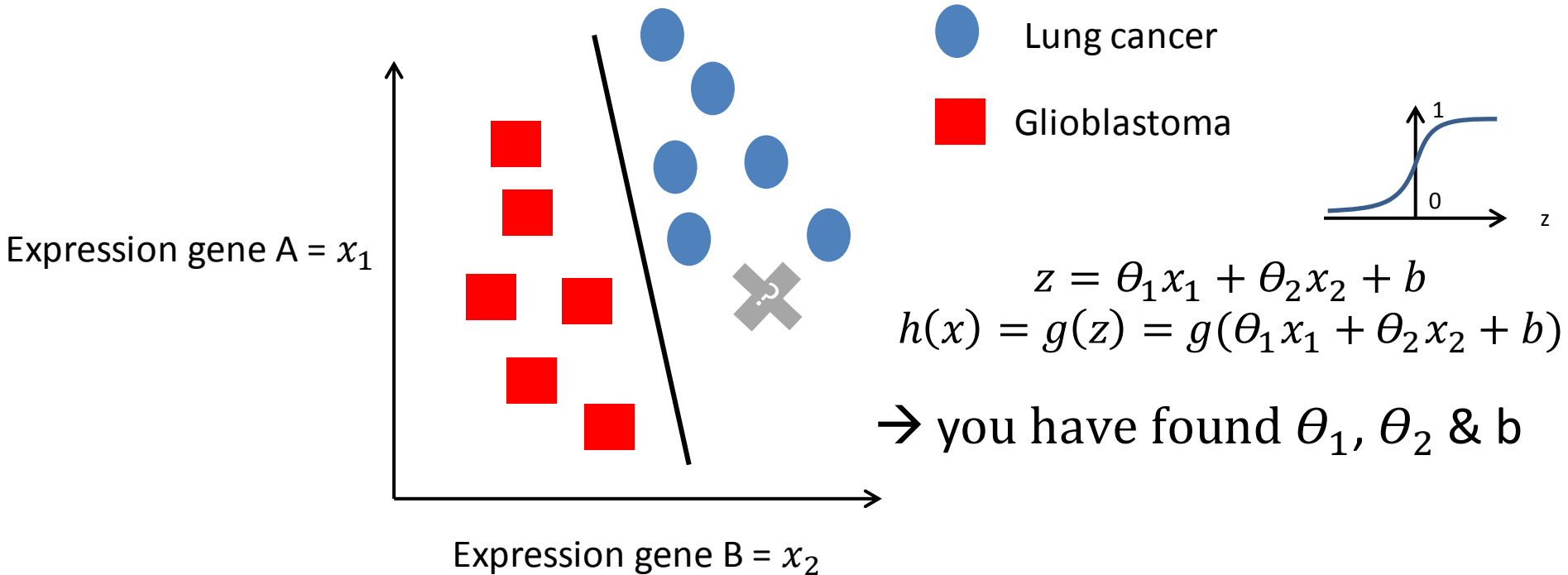


Eg: data augmentation
(real or artificial), work
with distribution...

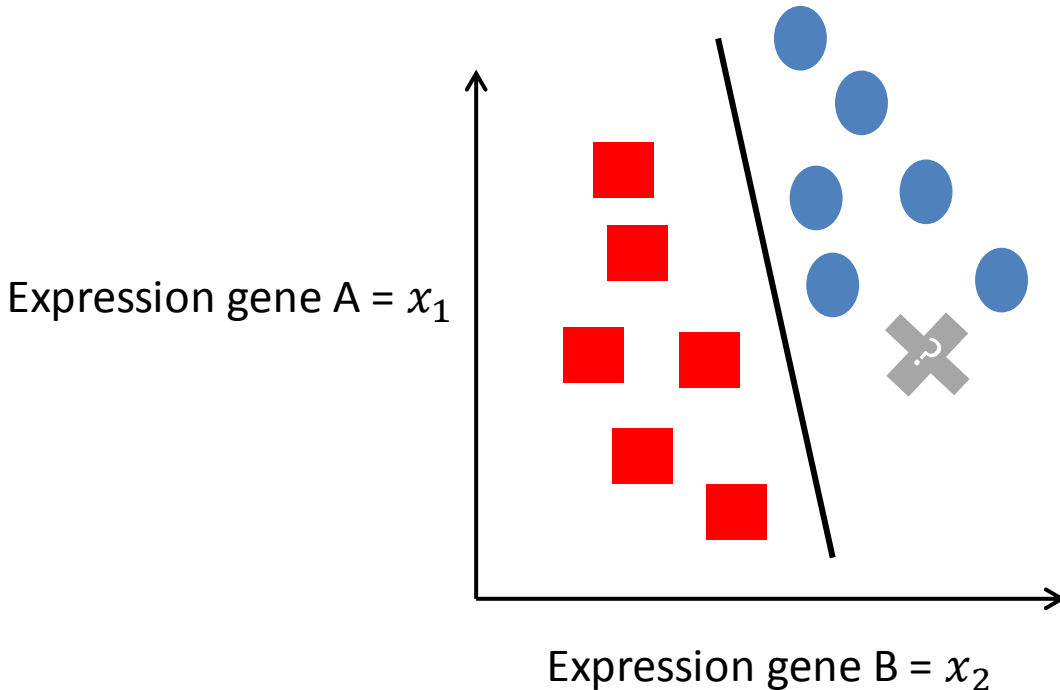


Back to TPBigData.Rmd

Interpret things biologically



Interpret things biologically



Gene A is not very discriminant

- so $abs(\theta_1)$ should be small

Gene B seems discriminant

- so $abs(\theta_2)$ should be larger

Back to TPBigData.Rmd

Lasso regularization or L1N

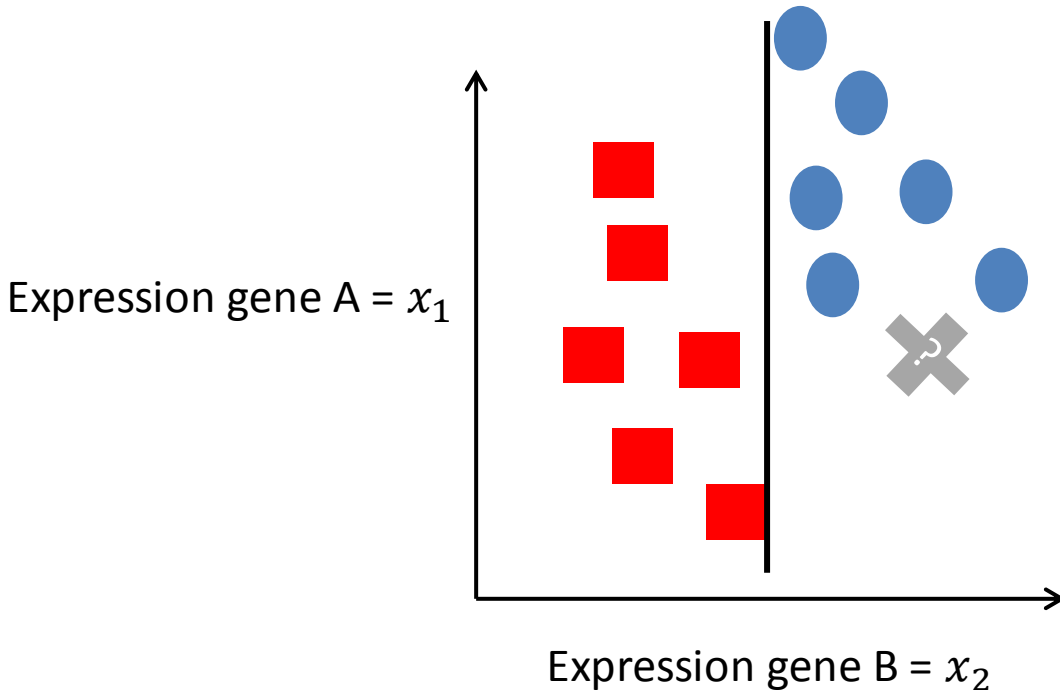
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h(x^i), y^i) + \lambda \sum_{j=1}^n \theta_j$$

→ minimizes $\theta_1, \theta_2 \dots \theta_n$

→ And even set some θ very close to zeros

$$z = \theta_1 x_1 + \theta_2 x_2 + b$$
$$h(x) = g(z) = g(\theta_1 x_1 + \theta_2 x_2 + b)$$

With Lasso regularization

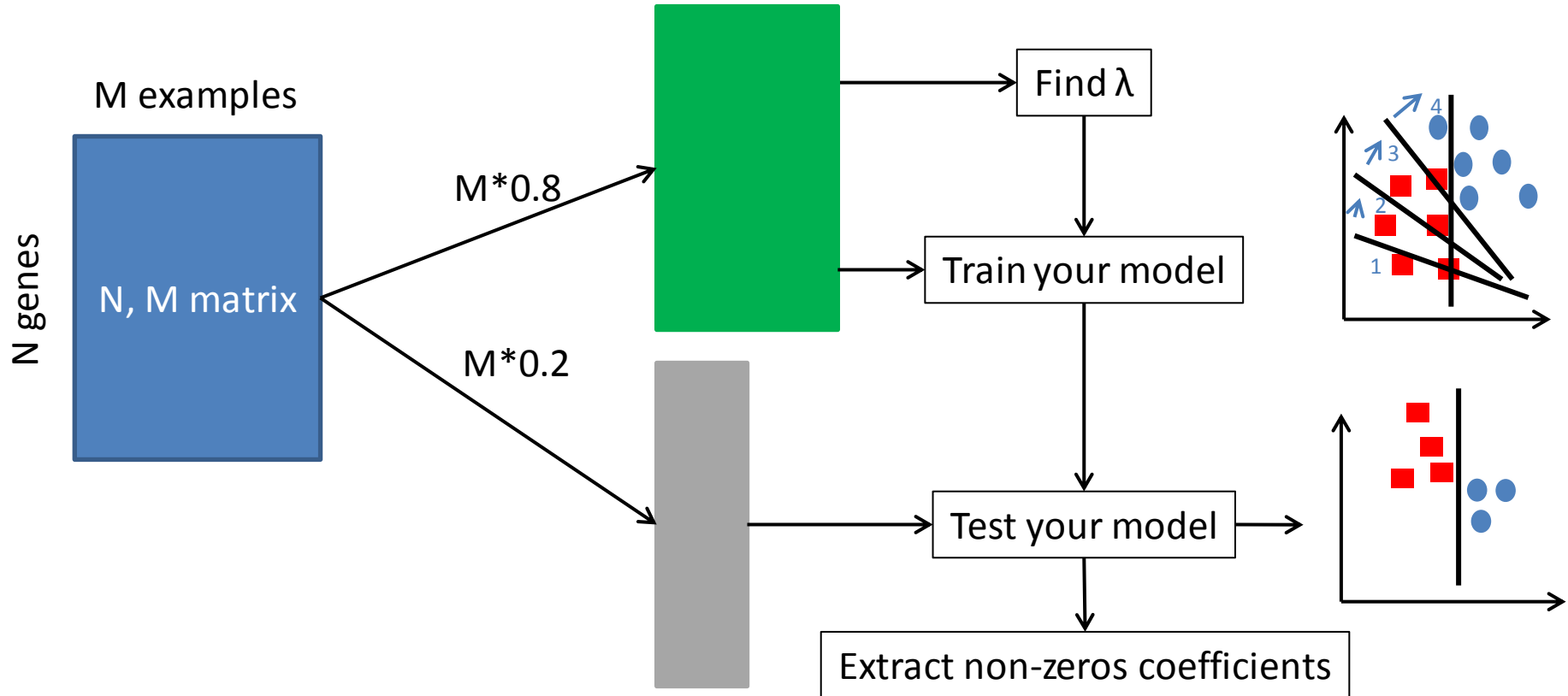


- θ_1 is set to ~ 0
- $abs(\theta_2)$ is kept large

But you have one more
parameter to search = lambda!

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h(x^i), y^i) + \lambda \sum_{j=1}^n \theta_j$$

Train and evaluate with Lasso



Find λ ? Do a K fold cross-validation!

$N,$
 $M * 0.8$

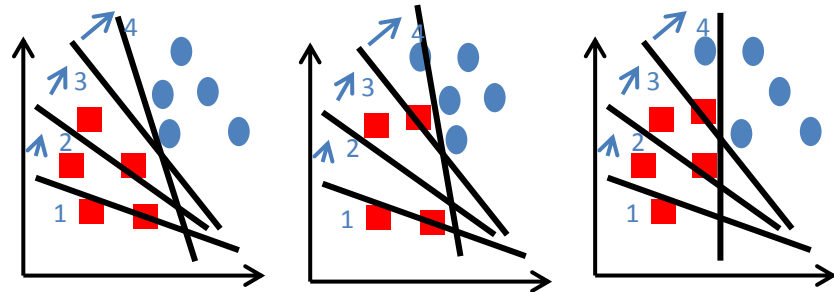
Find λ using
K fold cross-validation

New 80% / 20% split



$\lambda = 0.12$ $\lambda = 1.4$ $\lambda = 2.1$ $\lambda = 4.5 \dots$

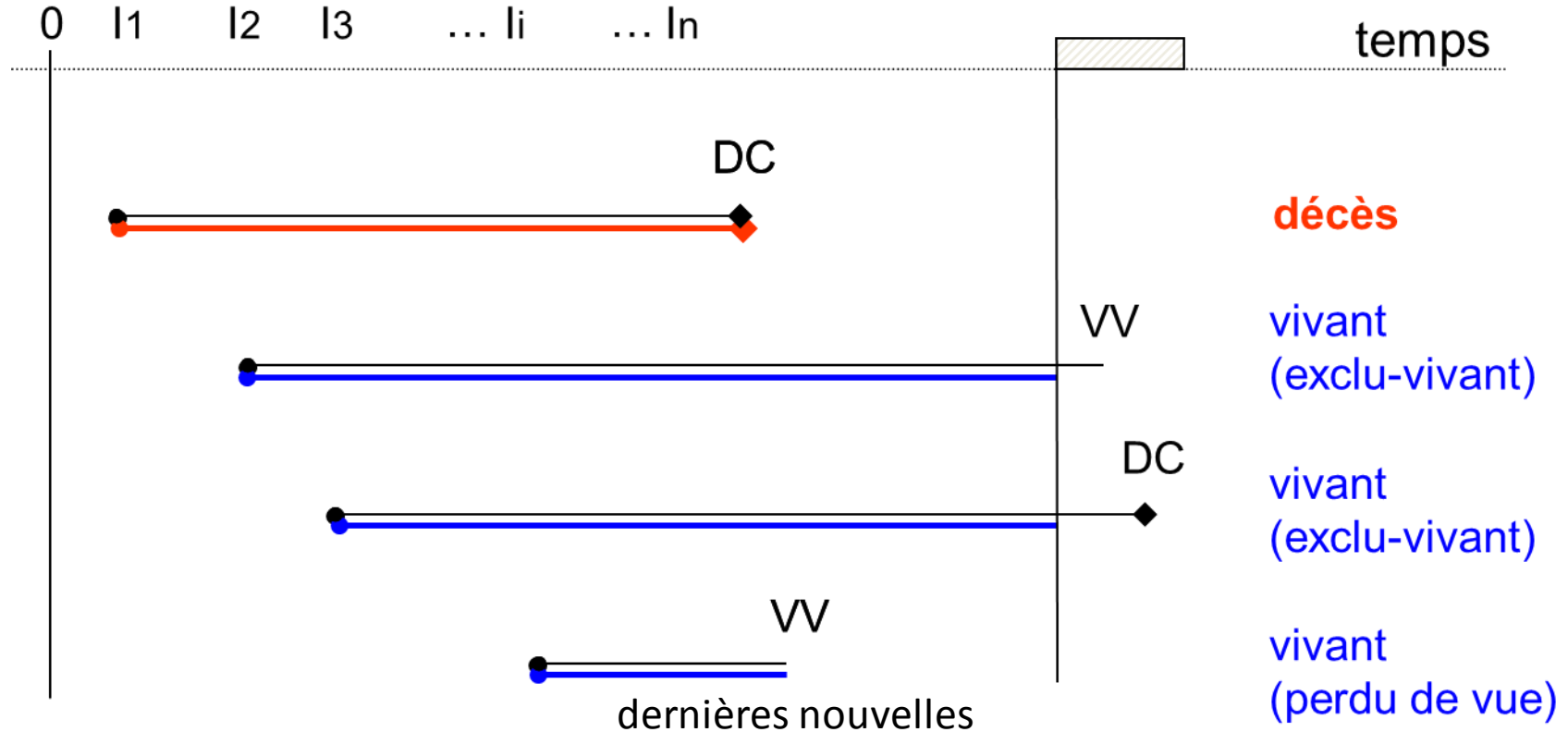
Keep best λ from
cross-validation set



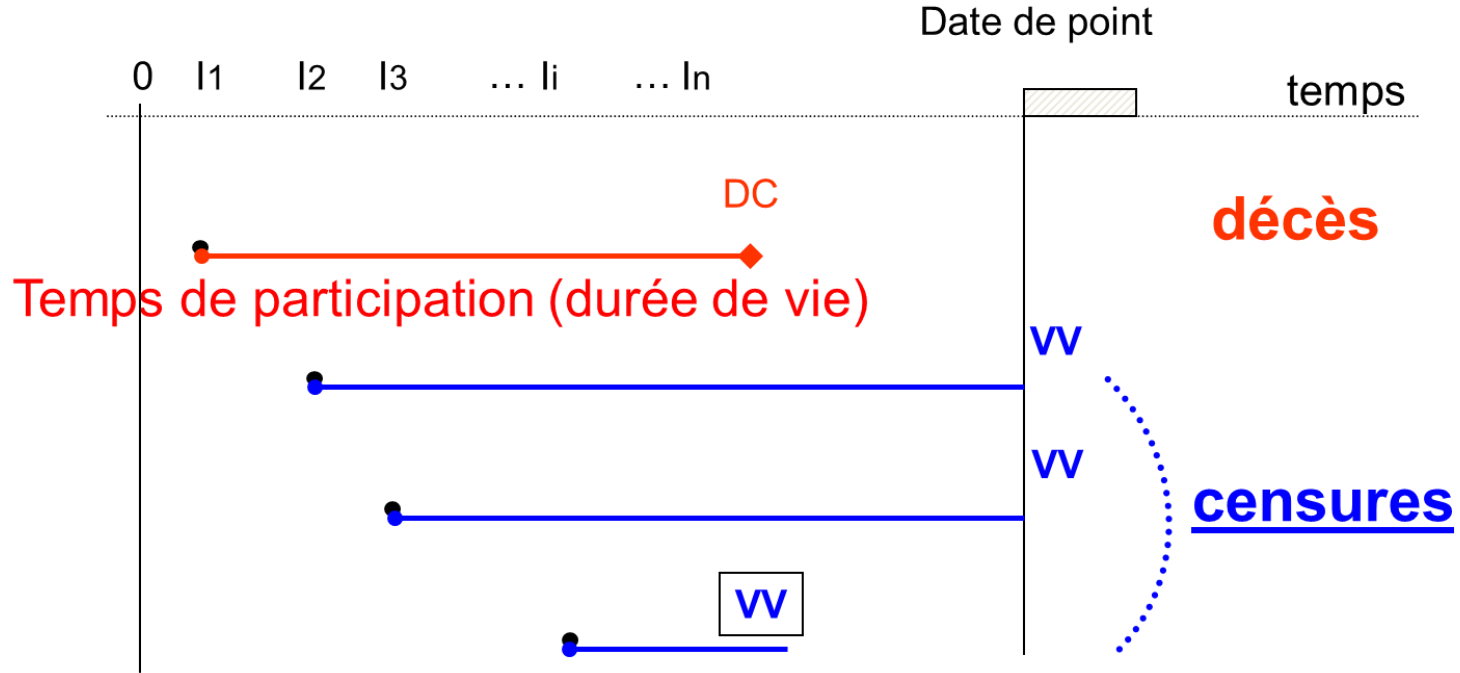
Back to TPBigData.Rmd

La survie : une donnée censurée

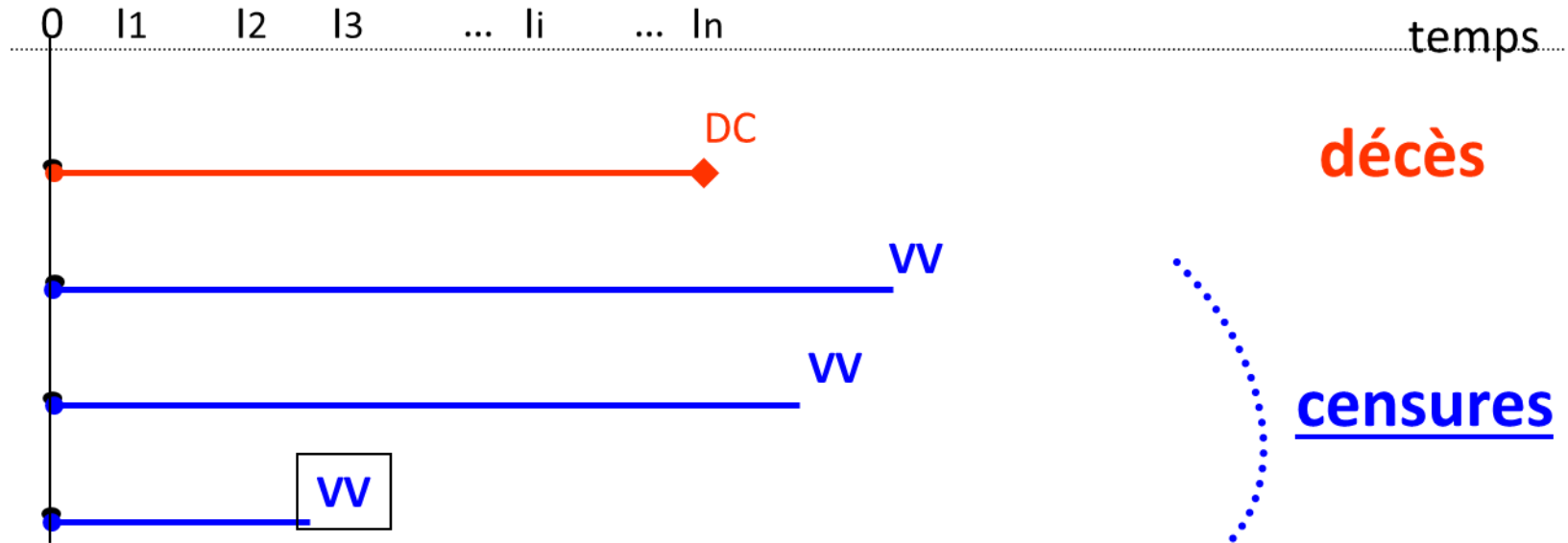
Date de point



La survie : une donnée censurée



La survie : une donnée censurée



Modeling Survival

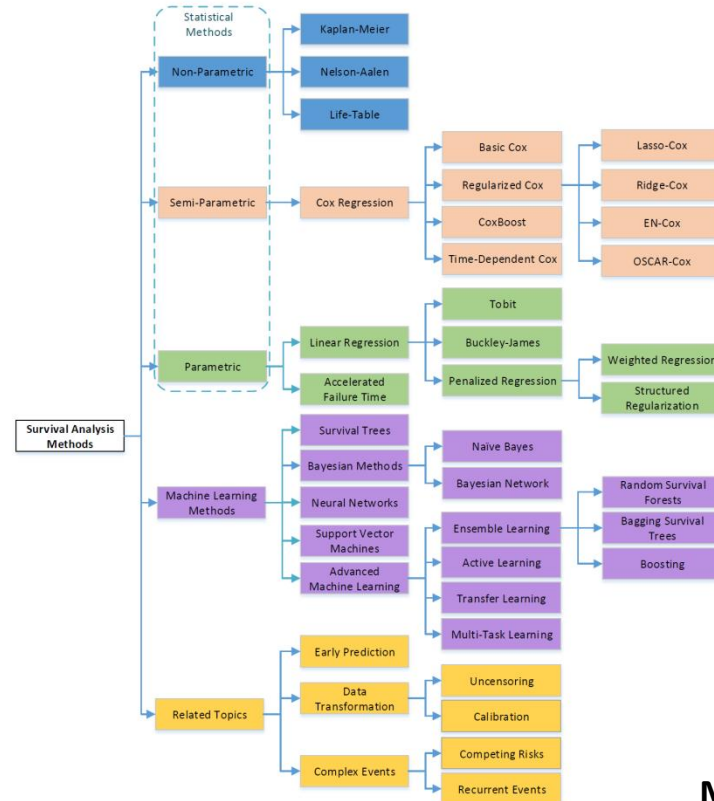
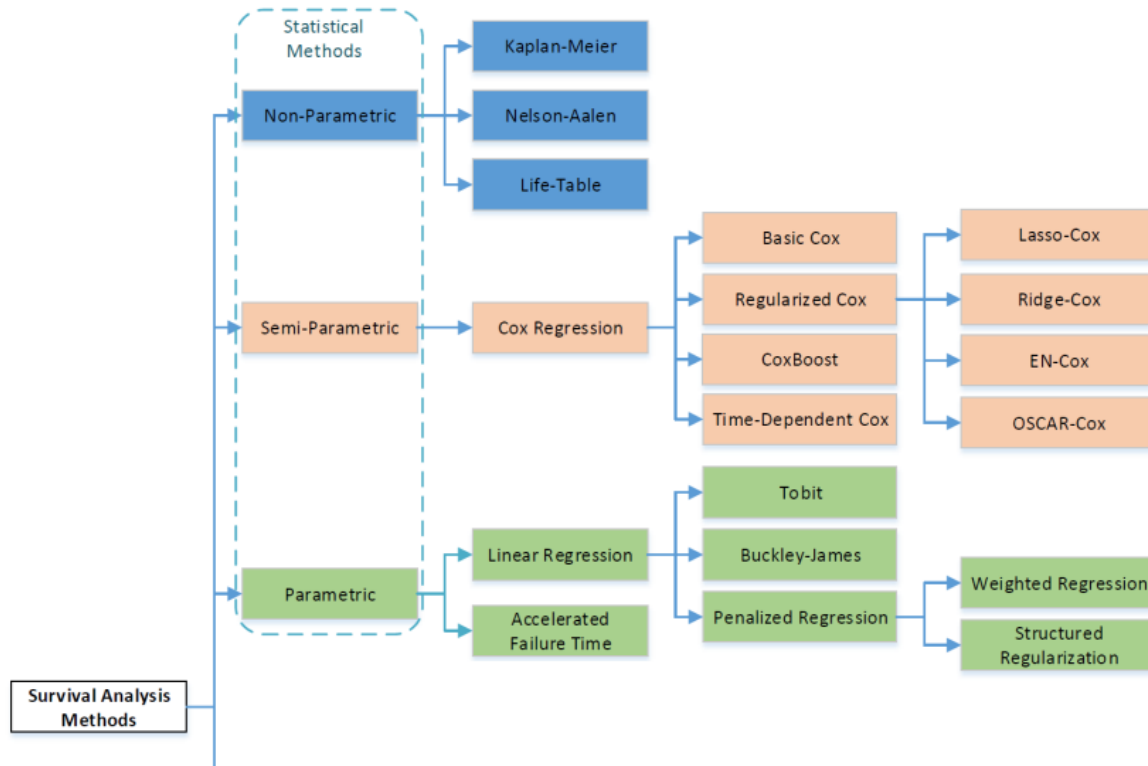


Fig. 3: Taxonomy of the methods developed for survival analysis.

Machine Learning for Survival Analysis: A Survey, PING WANG et al. ArXiv, 2017

Modeling Survival



Modeling Survival

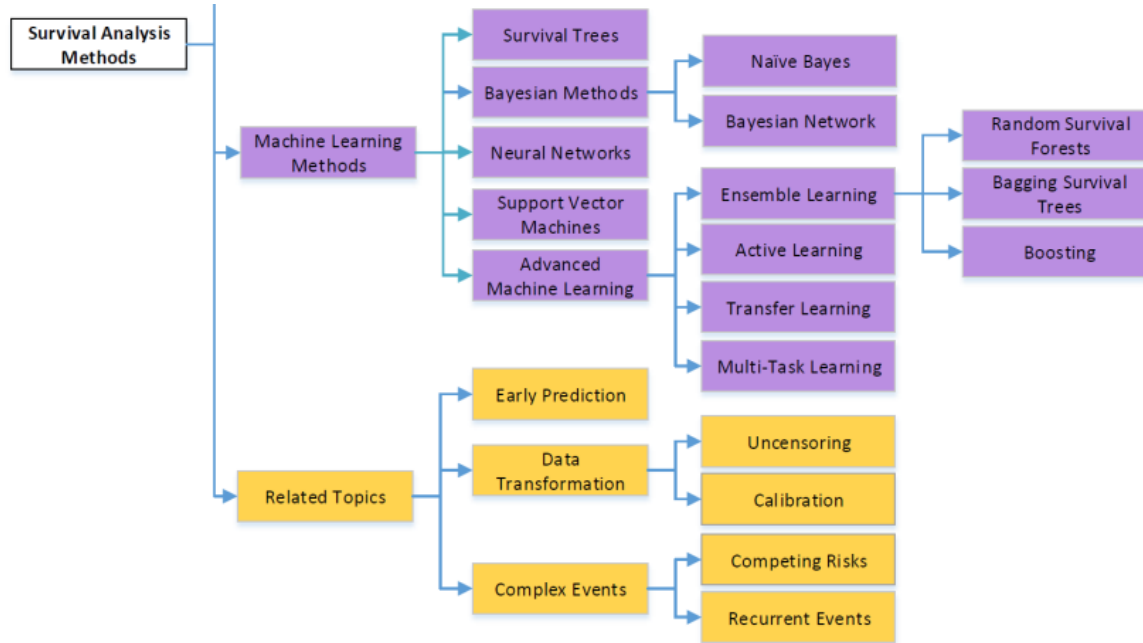


Fig. 3: Taxonomy of the methods developed for survival analysis.

Cox = a semi-parametric model

The knowledge of the underlying distribution of time to event of interest is not required, but the attributes are assumed to have an exponential influence on the outcome.

Individual hazard function:

$$h(t, X_i) = h_0(t) \exp(X_i \beta)$$

Negative partial log likelihood to minimize:

$$LL(\beta) = - \sum_{j=1}^N \delta_j \{X_j \beta - \log[\sum_{i \in R_j} \exp(X_i \beta)]\}.$$

N = number of subjects

X = variables

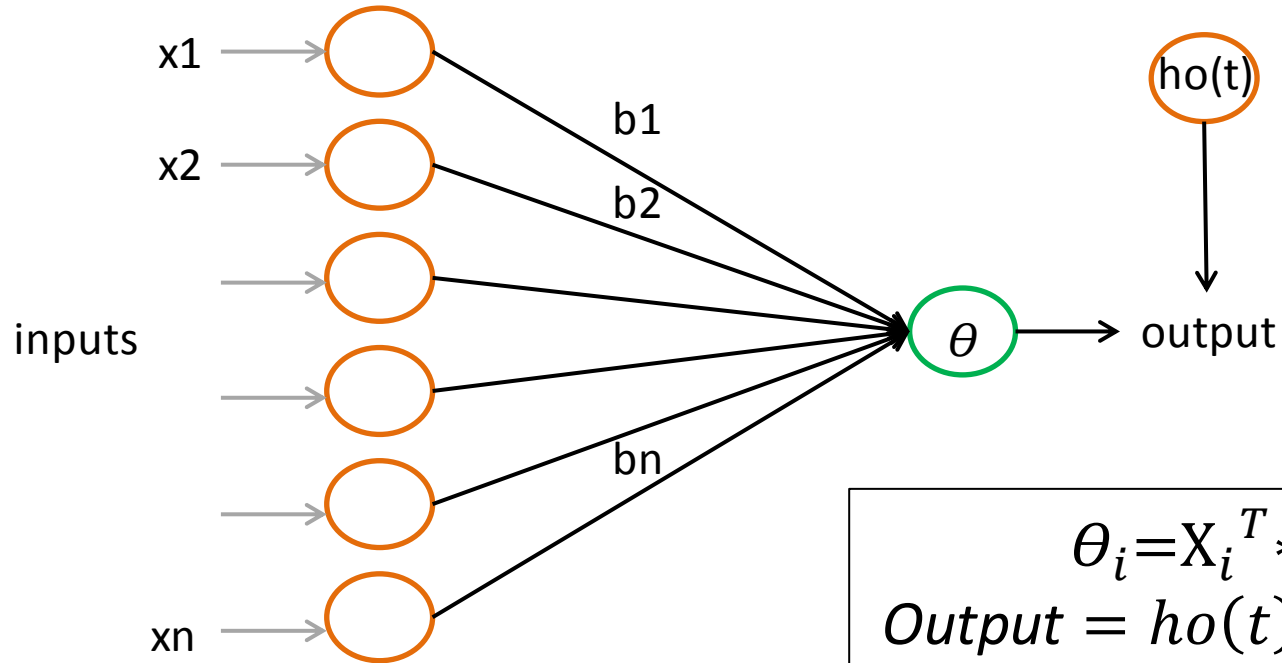
β = parameter to estimate

δ = censoring status

R = the set of risk subjects at

Time to event of interest

Cox model



$$\theta_i = X_i^T * \beta$$
$$Output = ho(t)_i * exp^{\theta_i}$$

C index

$$c = \Pr(\hat{y}_1 > \hat{y}_2 \mid y_1 > y_2)$$

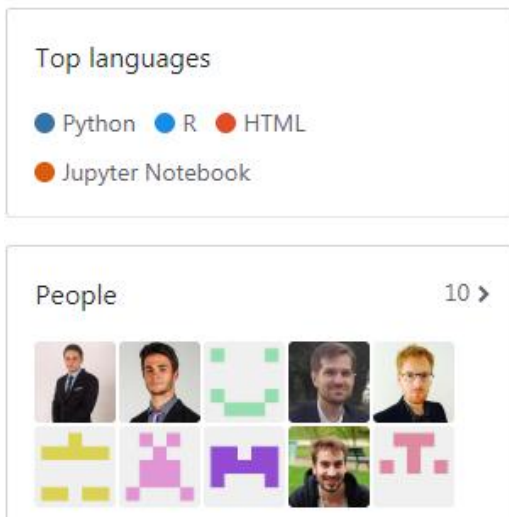
Data Science at DITEP

Info <https://github.com/DITEP>

Medical team



Data Science team



- 1 Senior Bioinformatician
Leo Colmet Daage
- Students from CentraleSupélec
- 1 Student from Telecom Paristech
- For 2019: MSc, PhDc

loicverlingue@yahoo.fr