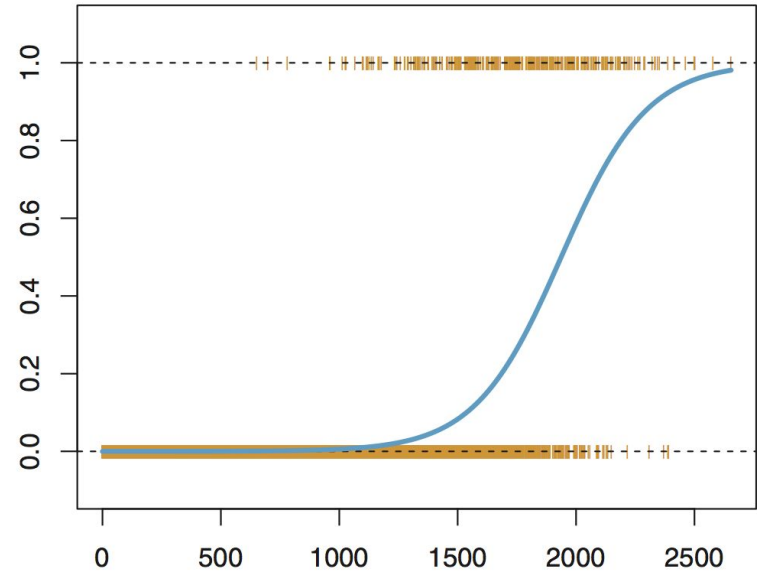


# Logistic Regression 1/2

Classification, metrics  
and ROC curves

DSI, jf.omhover, Dec 6, 2016



# Logistic Regression 1/2

Classification, metrics  
and ROC curves

DSI, jf.omhover, Dec 6, 2016

## OBJECTIVES (morning)

- **Relate** Regression to Classification in the context of supervised learning
- **Compare** Logistic Regression to Linear Regression
- **Define** and **compute** metrics for evaluating classifiers

## OBJECTIVES (afternoon)

- **Describe** the process for computing parameter values in LogReg
- **Use** the parameters of a LogReg model to **compute** the class of an observation





# Supervised Learning

Learning / Estimating FUNCTIONS based on examples

# Reality VS Model : assumptions and learning



## REALITY

	type	income	education	prestige
accountant	prof	62	86	82
pilot	prof	72	76	83
architect	prof	75	92	90
author	prof	55	90	76
chemist	prof	64	86	90
minister	prof	21	84	87
professor	prof	64	93	93
dentist	prof	80	100	90
reporter	wc	67	87	52
engineer	prof	72	86	88
undertaker	prof	42	74	57
lawyer	prof	76	98	89

data

$(x_1, y_1)$

...

$(x_n, y_n)$

$x \ y$

**OBJECTIVE:**  
descriptive  
predictive  
normative

...

$$\sum (y_i - \hat{f}(x_i))^2$$

COST FUNCTION

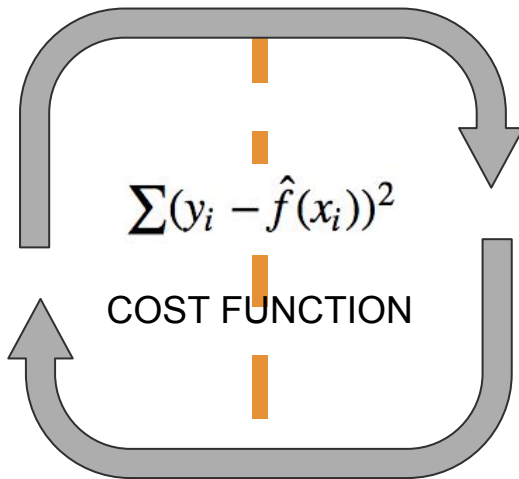
## MODEL

$$y = f(x) + \epsilon$$

take a function as  
an assumption

$$\hat{y} = \hat{f}(x)$$

Estimator  
of the function



# Linear Regression



## REALITY

	type	income	education	prestige
accountant	prof	62	86	82
pilot	prof	72	76	83
architect	prof	75	92	90
author	prof	55	90	76
chemist	prof	64	86	90
minister	prof	21	84	87
professor	prof	64	93	93
dentist	prof	80	100	90
reporter	wc	67	87	52
engineer	prof	72	86	88
undertaker	prof	42	74	57
lawyer	prof	76	98	89

data

$(x_1, y_1)$

...

$(x_n, y_n)$

$x \ y$

## OBJECTIVE:

descriptive  
predictive  
normative

...

## MODEL

$$y \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

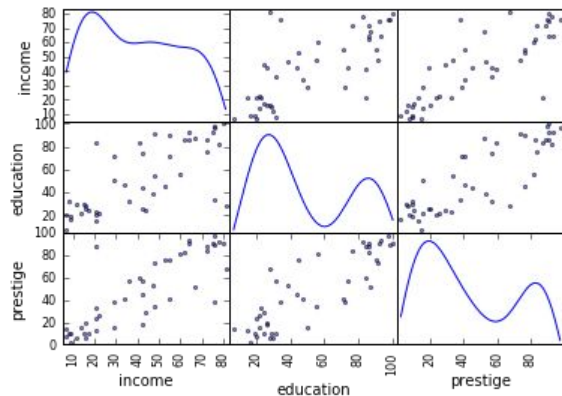
We make the assumption that we  
have a linear relation

# Linear Regression - General Process



## REALITY

- 1) Having a data sample  
Observing an underlying behavior



## MODEL

- 2) Make an assumption  
on the model underlying the data

$$y \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

linear relation  
(+ assumptions)

- 3) Find the instance of the model  
that fits with data sample

# Multi-Linear Regression



**COST FUNCTION (Residual Sum of Squares)**

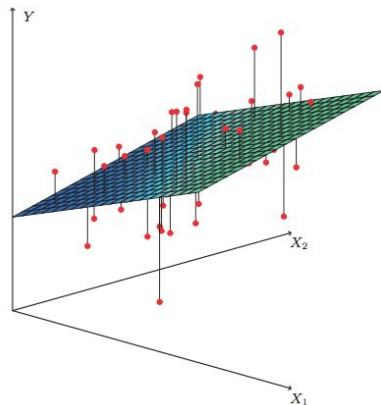
$$RSS(\beta) = (y - X\beta)^T(y - X\beta)$$

O.L.S.

**REALITY**

**DATA**

$$X = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,p-1} \\ 1 & x_{2,1} & \cdots & x_{2,p-1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,p-1} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = y$$



**MODEL**

**model class**

$$y \approx X\beta$$

**PROBLEM**

$$\hat{y} = X\hat{\beta}$$

**model instance  
estimator  
parameters**

**SOLUTION**

$$\hat{\beta} = (X^T X)^{-1} X^T y$$



# Classification

Learning / Estimating “models of classes” based on examples



# Reality vs Model : assumptions and learning



**REALITY**

$$X = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,p-1} \\ 1 & x_{2,1} & \cdots & x_{2,p-1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,p-1} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$= y$

**OBJECTIVE:**  
descriptive  
predictive  
normative  
...

$$\sum (y_i - \hat{f}(x_i))^2$$

COST FUNCTION

**MODEL**

$$y = f(x) + \epsilon$$

take a function as  
an assumption

$$\hat{y} = \hat{f}(x)$$

Estimator  
of the function

# Mapping // Classification algorithms



**Logistic Regression**

**k-NN**

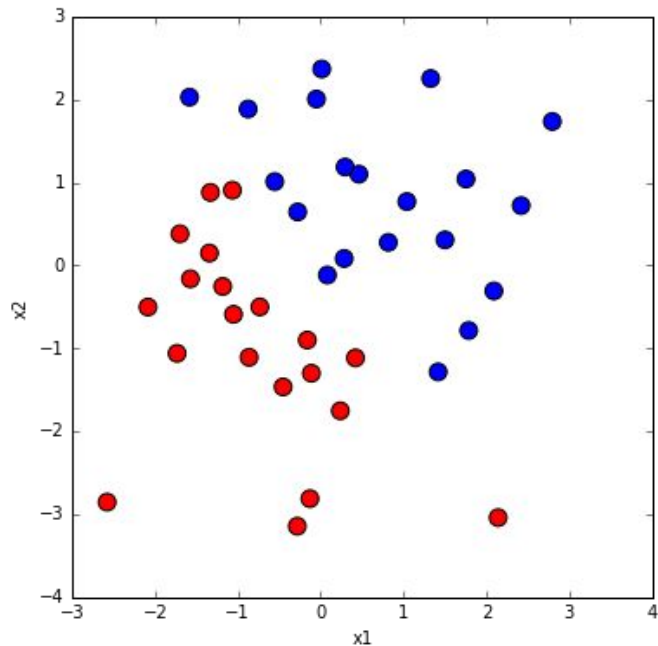
**Decision Trees**

**Random Forest, Boosting**

**Support Vector Machines (SVM)**

**Neural Networks**

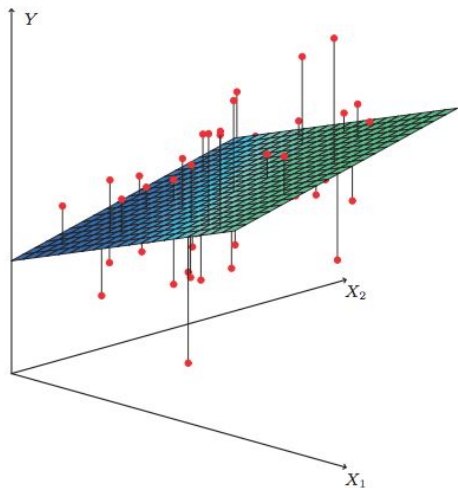
...



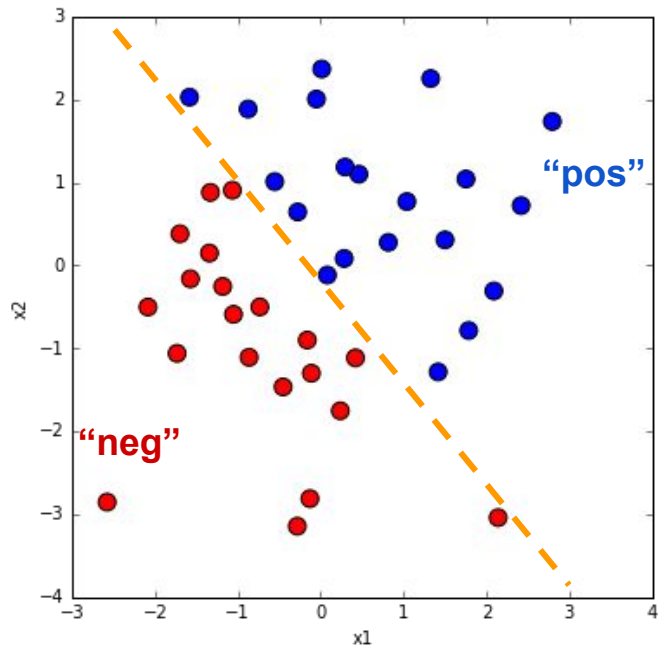
# Regression vs Classification



Quantitative response  $y$  in R



Categorical response  $y$

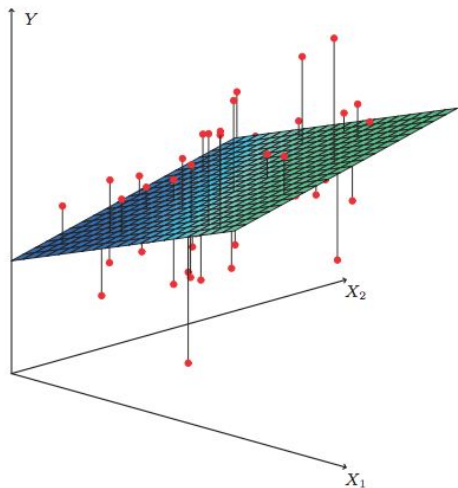


Assigning  $y = 0$  to neg,  $y = 1$  to pos

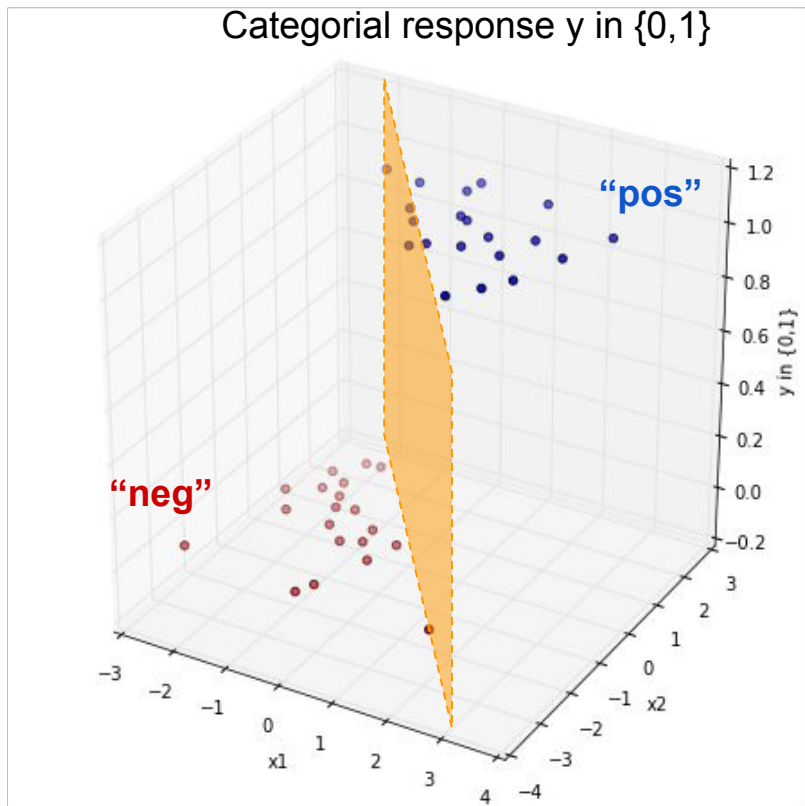
# Regression vs Classification



Quantitative response  $y$  in  $\mathbb{R}$



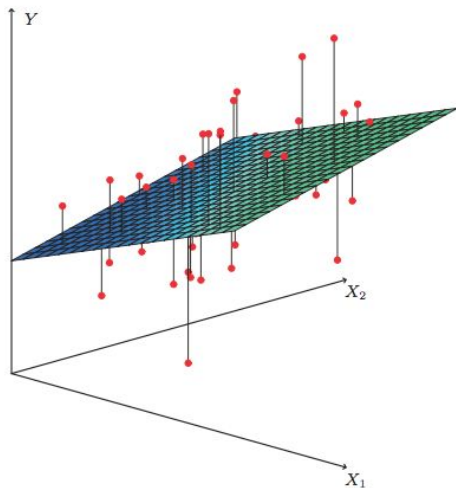
Categorical response  $y$  in  $\{0,1\}$



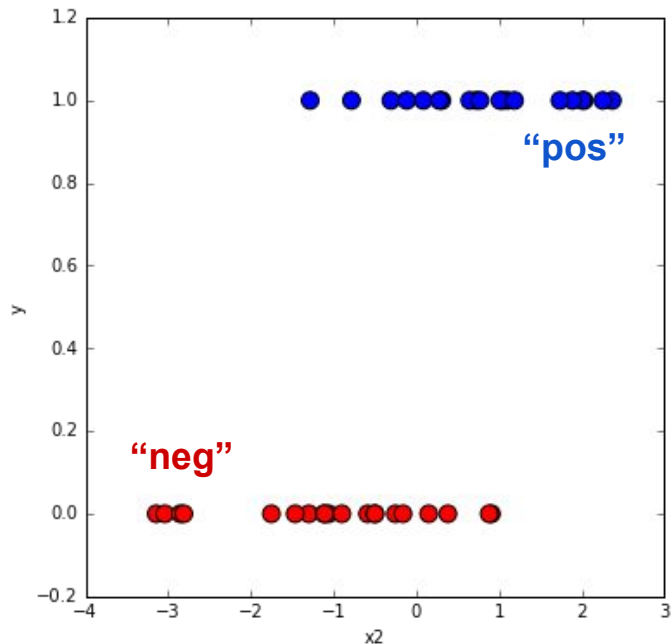
# Regression vs Classification



Quantitative response  $y$  in  $\mathbb{R}$



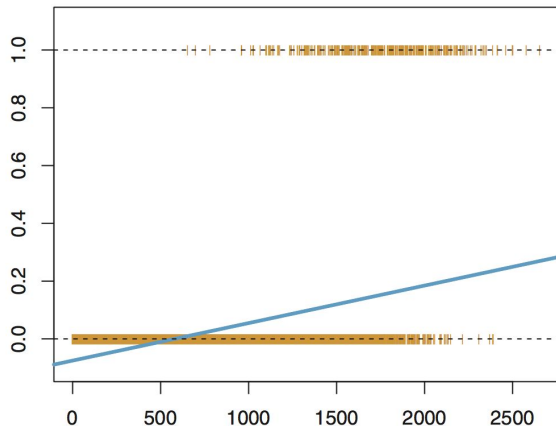
Categorical response  $y$  in  $\{0,1\}$



# Trying to apply LinReg to y



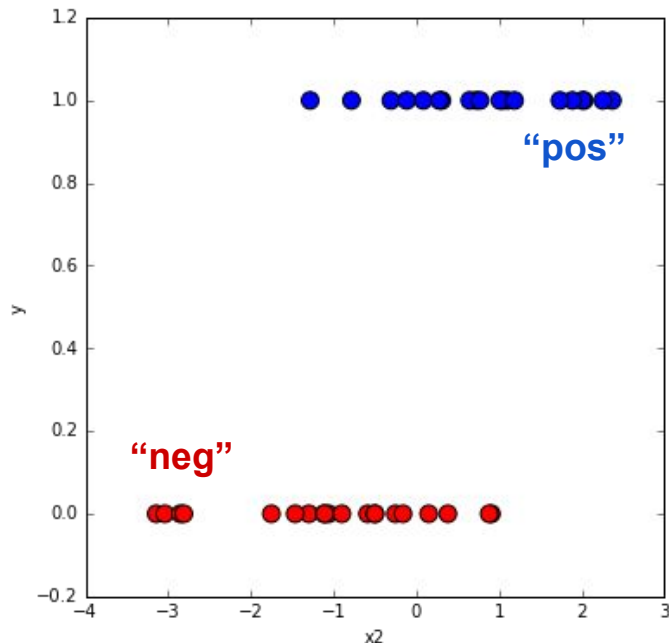
Quantitative response y in R



$$y \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

Negative probabilities ?  
How to cut-off ?

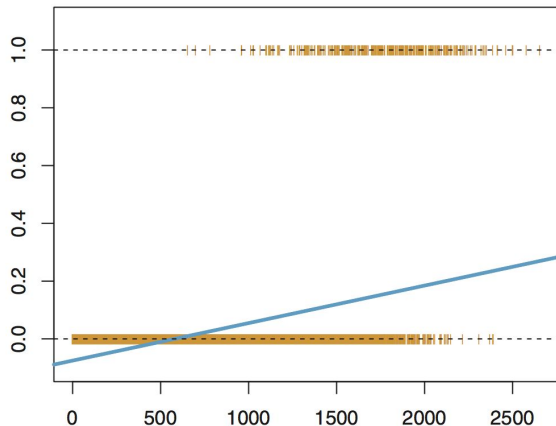
Categorical response y in {0,1}



# LogReg as model of probability



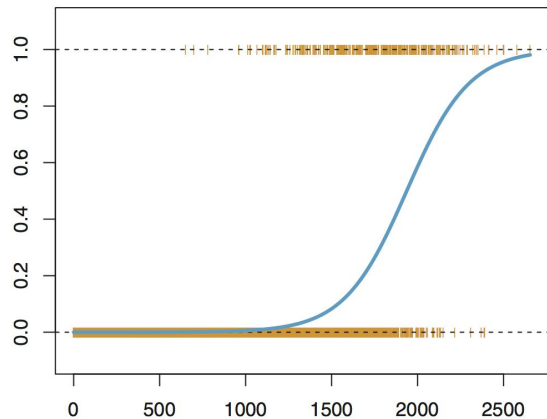
Quantitative response  $y$  in  $\mathbb{R}$



$$y \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

Negative probabilities ?  
How to cut-off ?

Categorical response  $y$  in  $\{0,1\}$



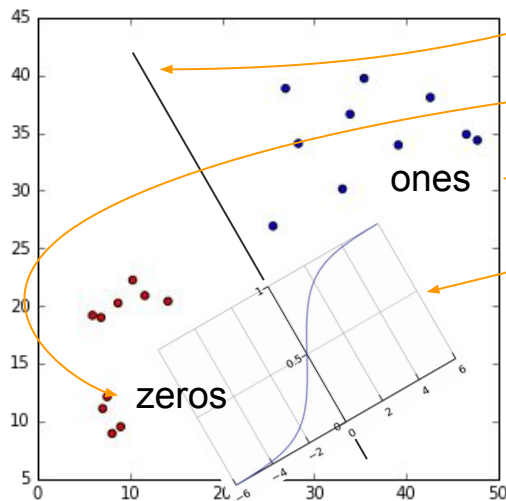
$$p(X) = h(\beta_0 + \beta_1 \cdot x_1 + \dots + \beta_p \cdot x_p)$$

Idea : model probability of being positive  
as a function of a linear model

# LogReg in a nutshell



## REALITY



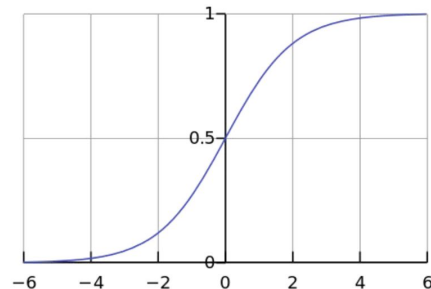
*It (badly) translates as :  
computes the probability  
of being in one of the two  
classes  
depending on of the side  
and distance of the plan*

## MODEL

$$p(X) = h(\beta_0 + \beta_1 \cdot x_1 + \dots + \beta_p \cdot x_p)$$

$$h : \mathbb{R} \rightarrow [0, 1]$$

$$h(t) = \frac{1}{1+e^{-t}}$$







# How to evaluate a classifier ?

# Conforming a classifier to the actual response



$x_1, x_2, x_3, x_4$	$y$	$\hat{y} = \hat{f}(x)$	
...	1	1	0.95
...	0	0	0.21
...	0	1	0.55
...	1	0	0.43
...	1	1	0.77
...	1	0	0.44
...	0	0	0.15
...	1	1	0.81

		$\hat{y} = \hat{f}(x)$		
		Pred P	Pred N	
$y$	Actual P	3	2	P = 5
	Actual N	1	2	N = 3
		P*	N*	

# Confusion Matrix



$\hat{y} = \hat{f}(x)$

		$\hat{y} = \hat{f}(x)$	
		Pred P	Pred N
$y$	Actual P	True Positive	False Negative
	Actual N	False Positive	True Negative
		P*	N*

P = 5  
N = 3

# Confusion Matrix - Metrics



The proportion of observations that are correctly classified ?

**Accuracy :**

The proportion of positives that are correctly identified as such ?

**True Pos Rate :**

(aka recall, sensitivity)

The proportion of negatives that are correctly identified as such

**True Neg Rate :**

(aka specificity)

$\hat{y} = \hat{f}(x)$

		$\hat{y} = \hat{f}(x)$	
		Pred P	Pred N
$y$	Actual P	True Positive	False Negative
	Actual N	False Positive	True Negative
		P*	N*

P = 5

N = 3

# Confusion Matrix - Metrics



The proportion of observations that are correctly classified ?

$$\text{Accuracy} : (TN + TP) / (N + P)$$

The proportion of positives that are correctly identified as such ?

$$\text{True Pos Rate} : TP / P$$

(aka recall, sensitivity)

The proportion of negatives that are correctly identified as such

$$\text{True Neg Rate} : TN / N$$

(aka specificity)

$\hat{y} = \hat{f}(x)$

		$\hat{y} = \hat{f}(x)$	
		Pred P	Pred N
$y$	Actual P	True Positive	False Negative
	Actual N	False Positive	True Negative
		P*	N*

P = 5

N = 3

# Confusion Matrix - Metrics



The proportion of observations that are  
NOT correctly classified ?

**Error rate :**

The proportion of positives that are  
NOT correctly identified as such ?

**False Neg Rate :**

(aka fall-out)

The proportion of negatives that are  
NOT correctly identified as such

**False Pos Rate :**

(aka 1-specificity)

$\hat{y} = \hat{f}(x)$

		$\hat{y} = \hat{f}(x)$		
		Pred P	Pred N	
$y$	Actual P	True Positive	False Negative	P = 5
	Actual N	False Positive	True Negative	N = 3
		P*	N*	

# Confusion Matrix - Metrics



The proportion of observations that are  
NOT correctly classified ?

$$\text{Error rate : } (FN + FP) / (N + P)$$

The proportion of positives that are  
NOT correctly identified as such ?

$$\text{False Neg Rate : } FN / P$$

(aka fall-out)

The proportion of negatives that are  
NOT correctly identified as such

$$\text{False Pos Rate : } FP / N$$

(aka 1-specificity)

$\hat{y} = \hat{f}(x)$

		$\hat{y} = \hat{f}(x)$		
		Pred P	Pred N	
$y$	Actual P	True Positive	False Negative	P = 5
	Actual N	False Positive	True Negative	N = 3
		P*	N*	

# Confusion Matrix - Metrics



The proportion of actual positives  
in those identified as such ?

**Precision :  $TP / (FP + TP)$**

The proportion of positives that are  
correctly identified as such ?

**Recall :  $TP / P$**   
(aka TPR, sensitivity)

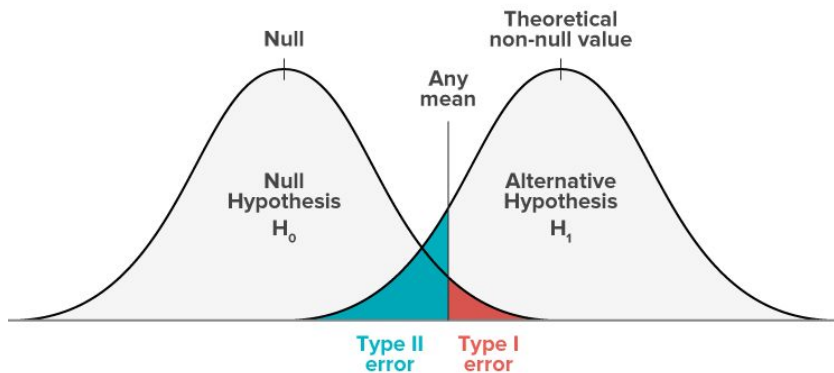
$\hat{y} = \hat{f}(x)$

	Pred P	Pred N	
Actual P	True Positive	False Negative	P = 5
Actual N	False Positive	True Negative	N = 3
	P*	N*	

$y$



# Confusion Matrix - type I and type II error



$y$

$$\hat{y} = \hat{f}(x)$$

	Pred P	Pred N	
Actual P	good	Type II error	P = 5
Actual N	Type I error	good	N = 3
	P*	N*	

# Using response probabilities



$x_1, x_2, x_3, x_4$				$y$	$P > 0.5$	
...	...	...	...	1	1	0.95
...	...	...	...	0	0	0.21
...	...	...	...	0	1	0.55
...	...	...	...	1	0	0.43
...	...	...	...	1	1	0.77
...	...	...	...	1	0	0.44
...	...	...	...	0	0	0.15
...	...	...	...	1	1	0.81

$\hat{y} = \hat{f}(x)$

		Pred P	Pred N	
$y$	Actual P	True Positive	False Negative	P = 5
	Actual N	False Positive	True Negative	N = 3
		P*	N*	

# Cut-offs on probabilities



$x_1, x_2, x_3, x_4$ $y$					$P > 0.5$		$P > 0.6$		$P > 0.7$		$P > 0.8$		$P > 0.9$	
...	...	...	...	1	1	0.95	1	0.95	1	0.95	1	0.95	1	0.95
...	...	...	...	0	0	0.21	0	0.21	0	0.21	0	0.21	0	0.21
...	...	...	...	0	1	0.55	0	<b>0.55</b>	0	0.55	0	0.55	0	0.55
...	...	...	...	1	0	0.43	0	0.43	0	0.43	0	0.43	0	0.43
...	...	...	...	1	1	0.77	1	0.77	1	0.77	0	<b>0.77</b>	0	0.77
...	...	...	...	1	0	0.44	0	0.44	0	0.44	0	0.44	0	0.44
...	...	...	...	0	0	0.15	0	0.15	0	0.15	0	0.15	0	0.15
...	...	...	...	1	1	0.81	1	0.81	1	0.81	1	0.81	0	<b>0.81</b>

Those are sure ones !  $\Rightarrow$  low FPR ! (high precision)  
But we miss so many ones !  $\Rightarrow$  low TPR ! (low recall)

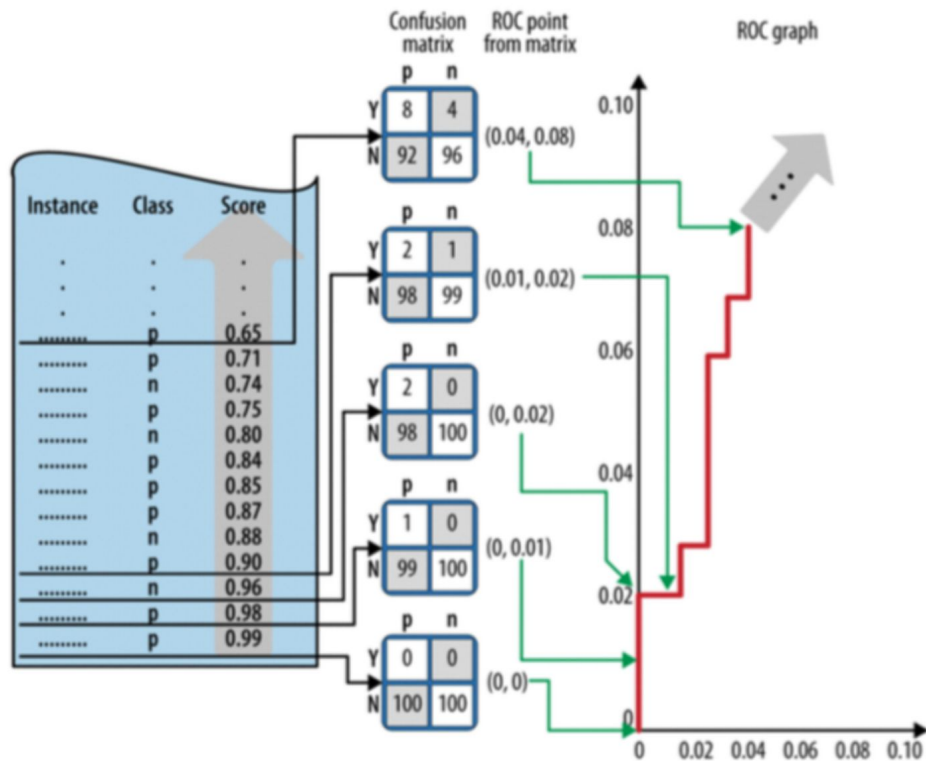
# Cut-offs on probabilities



$x_1, x_2, x_3, x_4$ $y$					$P > 0.5$		$P > 0.4$		$P > 0.3$		$P > 0.2$		$P > 0.1$	
...	...	...	...	1	1	0.95	1	0.95	1	0.95	1	0.95	1	0.95
...	...	...	...	0	0	0.21	0	0.21	0	0.21	1	<b>0.21</b>	1	0.21
...	...	...	...	0	1	0.55	1	0.55	1	0.55	1	0.55	1	0.55
...	...	...	...	1	0	0.43	1	<b>0.43</b>	1	0.43	1	0.43	1	0.43
...	...	...	...	1	1	0.77	1	0.77	1	0.77	1	0.77	1	0.77
...	...	...	...	1	0	0.44	1	<b>0.44</b>	1	0.44	1	0.44	1	0.44
...	...	...	...	0	0	0.15	0	0.15	0	0.15	0	0.15	1	0.15
...	...	...	...	1	1	0.81	1	0.81	1	0.81	1	0.81	1	0.81

We have so many FP !  $\Rightarrow$  high FPR ! (low precision)  
But we capture all the ones !  $\Rightarrow$  high TPR ! (high recall)

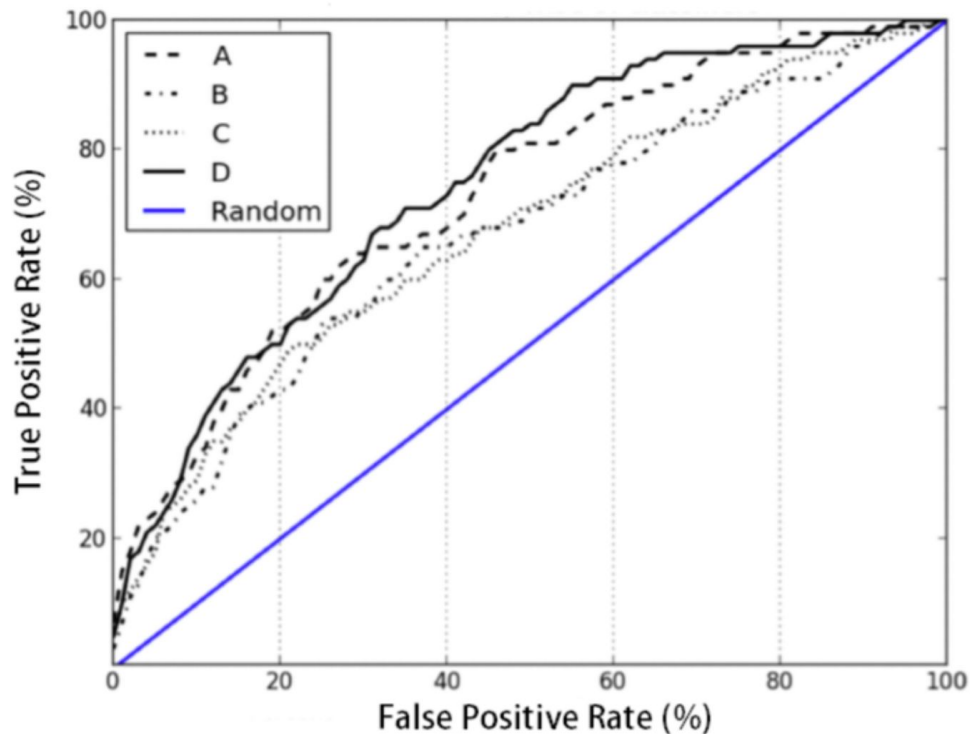
# ROC curve (receiver operating characteristic)



For LogReg, think of it as sliding the purple/red line along the sigmoid function



# Comparing classifiers based on their ROC curve



**Possible metric : AUC**  
**Area-under-curve**

# What is the “ideal” / “worst” classifier ?

