

# Logistic Regression

Luca Martino

# Introduction

## DATA:

$x_k \in \mathbb{R}$  it could be easily generalized for  $\mathbf{x}_k \in \mathbb{R}^D$

$y_k \in \{0, 1\}$

$\{x_k, y_k\}_{k=1}^N$

**We model directly the probability of an input  $x$  belonging to a class.**

# Introduction

**We model directly the probability of an input  $x$  belonging to a class. In this binary classification case, we focus on the probability of belonging to the class labelled with the label “ $y_k=1$ ”.**

## Model [\[edit\]](#)

The **logistic function** is of the form:

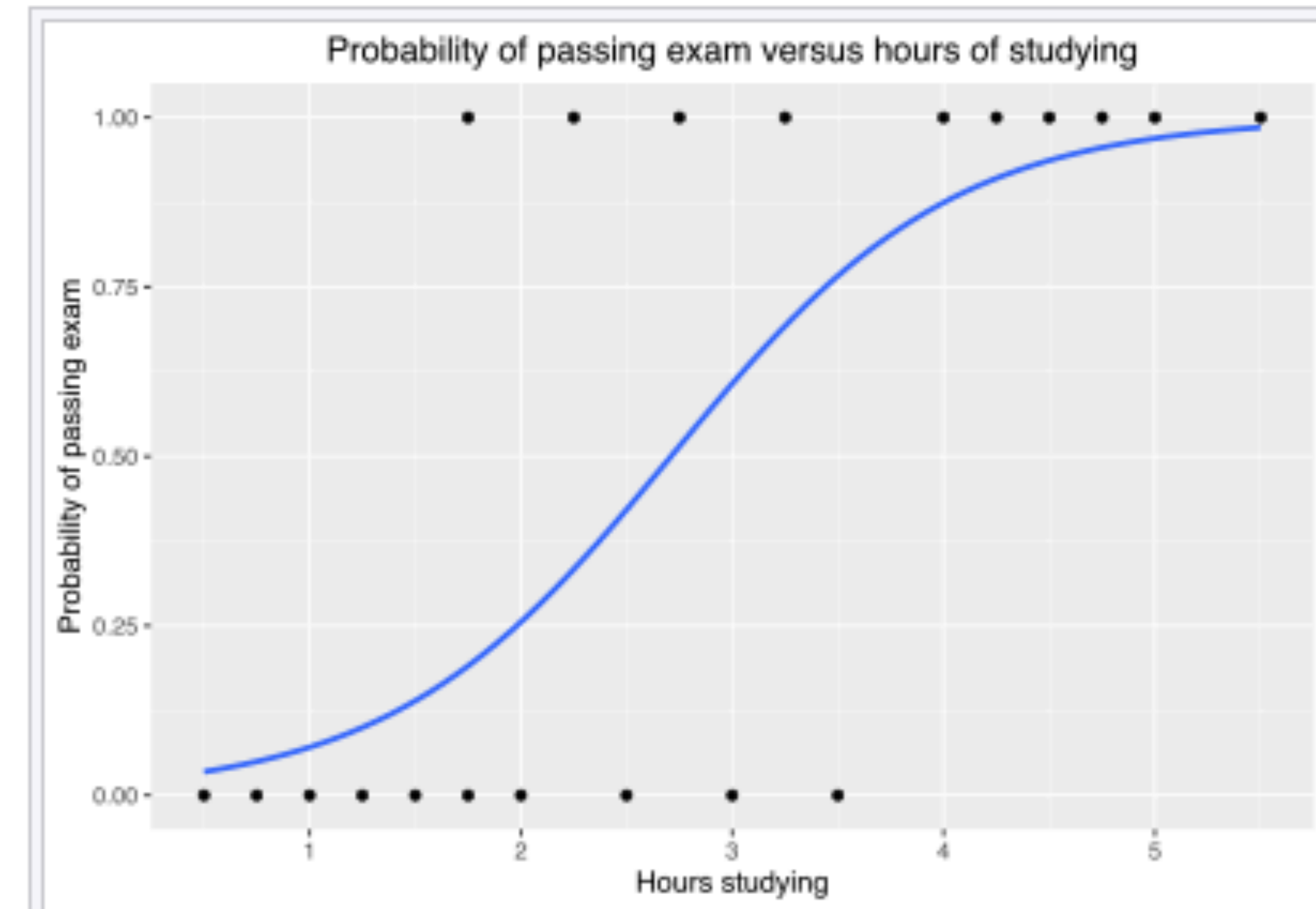
$$p(y = 1|x) = \frac{1}{1 + e^{-(x-\mu)/s}}$$

where  $\mu$  is a **location parameter** (the midpoint of the curve, where  $p(\mu) = 1/2$ ) and  $s$  is a **scale parameter**. This expression may be rewritten as:

$$p(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

where  $\beta_0 = -\mu/s$  and is known as the **intercept** (it is the *vertical* intercept or  $y$ -intercept of the line  $y = \beta_0 + \beta_1 x$ ), and  $\beta_1 = 1/s$  (inverse scale parameter or **rate parameter**): these are the  $y$ -intercept and slope of the log-odds as a function of  $x$ .

Conversely,  $\mu = -\beta_0/\beta_1$  and  $s = 1/\beta_1$ .



# Equivalent expressions

$$p(y_k = 1|x_k) = p_k = \frac{1}{1 + e^{-(x_k - \mu)/s}}$$

where  $\beta_0 = -\mu/s$  and is known as the **intercept** (it is the *vertical* intercept or y-intercept of the line  $y = \beta_0 + \beta_1 x$ ), and  $\beta_1 = 1/s$  (inverse scale parameter or **rate parameter**): these are the y-intercept and slope of the log-odds as a function of x. Conversely,  $\mu = -\beta_0/\beta_1$  and  $s = 1/\beta_1$ .

$$p(y_k = 1|x_k) = p_k = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_k)}}$$

$$p(y_k = 1|x_k) = p_k = \frac{e^{\beta_0 + \beta_1 x_k}}{1 + e^{\beta_0 + \beta_1 x_k}}$$

# Binary case: likelihood function

$$p(y_k = 1|x_k) = p_k = \frac{1}{1 + e^{-(x_k - \mu)/s}}$$

$$p(\mathbf{y}|\mathbf{x}) = p(y_1, \dots, y_N|x_1, \dots, x_N) = \prod_{k:y_k=1} p_k \prod_{k:y_k=0} (1 - p_k)$$

**If  $y_k=0,1$ , we can rewrite:**

$$p(\mathbf{y}|\mathbf{x}) = p(y_1, \dots, y_N|x_1, \dots, x_N) = \prod_{k=1}^N p_k^{y_k} (1 - p_k)^{1-y_k}$$

# Binary case: likelihood function

If  $y_k=0,1$ , we can rewrite:

$$p(\mathbf{y}|\mathbf{x}) = p(y_1, \dots, y_N | x_1, \dots, x_N) = \prod_{k=1}^N p_k^{y_k} (1 - p_k)^{1-y_k}$$

$$\log p(\mathbf{y}|\mathbf{x}) = \sum_{k=1}^N \left[ \log [p_k^{y_k}] + \log [(1 - p_k)^{1-y_k}] \right]$$

In the code that I sent in Studium, I used this formula in order to avoid numerical issues, such as NaN for “0 times -Inf” = “0 x -Inf”.

$$\log p(\mathbf{y}|\mathbf{x}) = \sum_{k=1}^N [y_k \log (p_k) + (1 - y_k) \log (1 - p_k)]$$

# **First generalization: With multidimensional inputs (straightforward)**

**Luca Martino**

# With multimensional inputs

**DATA:**

$$\mathbf{x}_k = [x_{k,1}, \dots, x_{k,D}]^\top \in \mathbb{R}^D$$

$$y_k \in \{0, 1\}$$

$$\{\mathbf{x}_k, y_k\}_{k=1}^N$$



# With multidimensional inputs

**Then we consider:**

$$p(y_k = 1 | \mathbf{x}_k) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{k,1} + \beta_2 x_{k,2} + \dots + \beta_D x_{k,D})}}$$

$$p(y_k = 1 | \mathbf{x}_k) = \frac{1}{1 + e^{-(\beta_0 + \sum_{d=1}^D \beta_d x_{k,d})}}$$

**Here we have to learn all the betas !!!!**

# Simplifying the previous expressions “using vectors”

**RE-DEFINING:**

$$\mathbf{x}_k = [\mathbf{1}, x_{k,1}, \dots, x_{k,D}]^\top \in \mathbb{R}^{D+1}$$

**AND:**

$$\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_D]$$

**then:**

$$p_k = p(y_k = 1 | \mathbf{x}_k) = \frac{1}{1 + e^{-(\boldsymbol{\beta}\mathbf{x}_k)}} = \frac{e^{\boldsymbol{\beta}\mathbf{x}_k}}{1 + e^{\boldsymbol{\beta}\mathbf{x}_k}}$$

# With multidimensional inputs

The rest remains the same....same likelihood function:

$$p_k = p(y_k = 1 | \mathbf{x}_k)$$

$$p(\mathbf{y} | \mathbf{X}) = p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{k=1}^N p_k^{y_k} (1 - p_k)^{1-y_k}$$

etc....

# Second generalization: More than 2 classes....

“**Multinomial logistic regression:** Many explanatory variables (inputs) and many categories (outputs, more than 2 classes)”

Luca Martino

# Dealing with $M$ classes

**DATA:**

$$\mathbf{x}_k = [x_{k,1}, \dots, x_{k,D}]^\top \in \mathbb{R}^D$$

$$y_k \in \{0, 1, 2, \dots, M - 1\}$$

$$\{\mathbf{x}_k, y_k\}_{k=1}^N$$

# Dealing with $M$ classes

Again, we model directly the probability of an input  $x$  belonging to a class.

**RE-DEFINING:**

$$\mathbf{x}_k = [1, x_{k,1}, \dots, x_{k,D}]^\top \in \mathbb{R}^{D+1}$$

**AND:**

$$\boldsymbol{\beta}_m = [\beta_{m,0}, \beta_{m,1}, \dots, \beta_{m,D}], \quad \text{with } m = 1, \dots, M - 1$$

**then:**

$$p_{k,m} = p(y_k = m | \mathbf{x}_k) = \frac{e^{\boldsymbol{\beta}_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\boldsymbol{\beta}_j \mathbf{x}_k}}$$

$$p_{k,0} = p(y_k = 0 | \mathbf{x}_k) = 1 - \sum_{m=1}^{M-1} \frac{e^{\boldsymbol{\beta}_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\boldsymbol{\beta}_j \mathbf{x}_k}}$$

# Dealing with $M$ classes

then:

$$p_{k,m} = p(y_k = m | \mathbf{x}_k) = \frac{e^{\beta_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}}$$

$$\begin{aligned} p_{k,0} = p(y_k = 0 | \mathbf{x}_k) &= 1 - \sum_{m=1}^{M-1} \frac{e^{\beta_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}} = 1 - \frac{\sum_{m=1}^{M-1} e^{\beta_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}} \\ &= \frac{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k} - \sum_{m=1}^{M-1} e^{\beta_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}} \\ &= \frac{1}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}} \end{aligned}$$

# Dealing with $M$ classes

then, **FINALLY**:

$$p_{k,m} = p(y_k = m | \mathbf{x}_k) = \frac{e^{\beta_m \mathbf{x}_k}}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}}$$

$$p_{k,0} = p(y_k = 0 | \mathbf{x}_k) = \frac{1}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}_k}}$$



# Dealing with $M$ classes

**OR MORE GENERALLY FOR a generic  $\mathbf{x}$  (test input):**

$$p(y_k = m | \mathbf{x}) = \frac{e^{\beta_m \mathbf{x}}}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}}}$$

$$p(y_k = 0 | \mathbf{x}) = 1 - \sum_{m=1}^{M-1} p(y_k = m | \mathbf{x}) = \frac{1}{1 + \sum_{j=1}^{M-1} e^{\beta_j \mathbf{x}}}$$

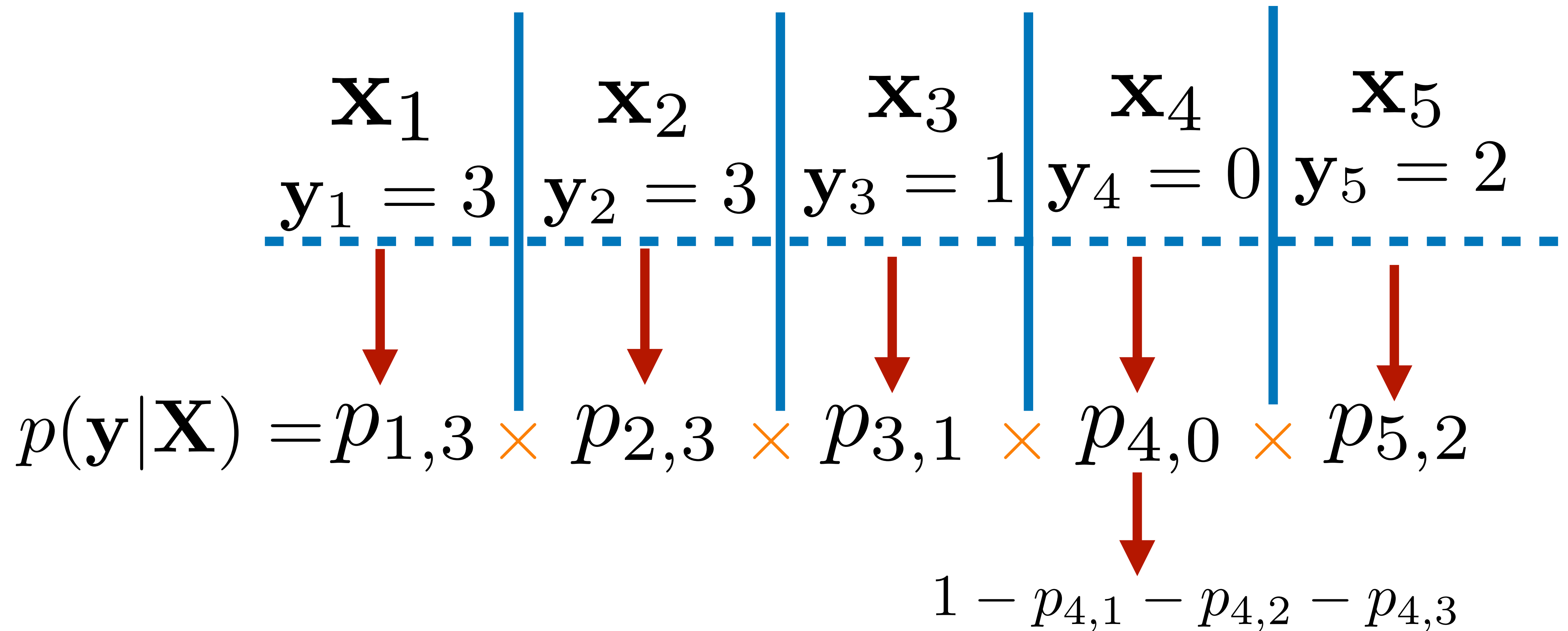
# Dealing with $M$ classes

**WE HAVE TO LEARN  $(M-1)$  VECTORS OF BETAS....  
hence,  $D \times (M-1)$  scalar numbers to learn !!!!**

$$\beta_1, \dots, \beta_{M-1}$$

# Example of construction of the likelihood function

Classes = 0, 1, 2, 3  $\implies M = 4$   $N = 5$



# **Isotopic multi-output logistic regression for parallel classification problems**

**Luca Martino**

# DATA (multioutput - isotopic scenario)

## OUTPUTS - isotopic scenario

$x_1$	$y_{1,1}$	$y_{1,2}$	$y_{1,3}$	...	...	...	$y_{1,N}$
$x_2$	$y_{2,1}$	$y_{2,2}$	$y_{2,3}$	...	...	...	$y_{2,N}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_M$	$y_{M,1}$	$y_{M,2}$	$y_{M,3}$	...	...	...	$y_{M,N}$

**DATA points ( $M$ ):**

$$\{\mathbf{x}_m, \mathbf{y}_m\}_{m=1}^M$$

$$\mathbf{y}_m = [y_{m,1}, y_{m,2}, \dots, y_{m,N}]$$

$$y_{i,j} \in \{0, 1\}$$

**isotopic scenario:**

**N outputs share the same input x**

$$i = 1, 2, \dots, M$$

$$j = 1, 2, \dots, N$$

# About the notation

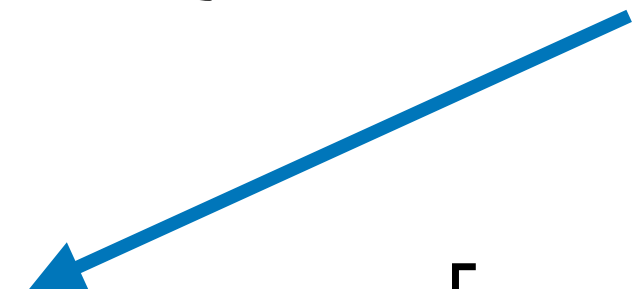
A MORE PROPER NOTATION FOR THESE SLIDES SHOULD BE OBTAINED SWITCHING  $M$  AND  $N$ :

- $N$  should be the number of data points
- and  $M$  should be the number of outputs per each input  $x$

However, we have used this notation for linking this part with the slides on ITEM RESPONSE THEORY (where we use exactly the notation employed here).

**DATA points ( $M$ ):**

$$\{\mathbf{x}_m, \mathbf{y}_m\}_{m=1}^M$$


$$\mathbf{y}_m = [y_{m,1}, y_{m,2}, \dots, y_{m,N}]$$

$$y_{i,j} \in \{0, 1\}$$

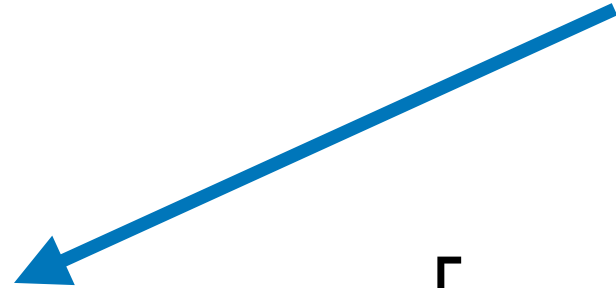
$$i = 1, 2, \dots, M$$

$$j = 1, 2, \dots, N$$

# About the notation

**DATA points ( $M$ ):**

$$\{\mathbf{x}_m, \mathbf{y}_m\}_{m=1}^M$$


$$\mathbf{y}_m = [y_{m,1}, y_{m,2}, \dots, y_{m,N}]$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_M]$$

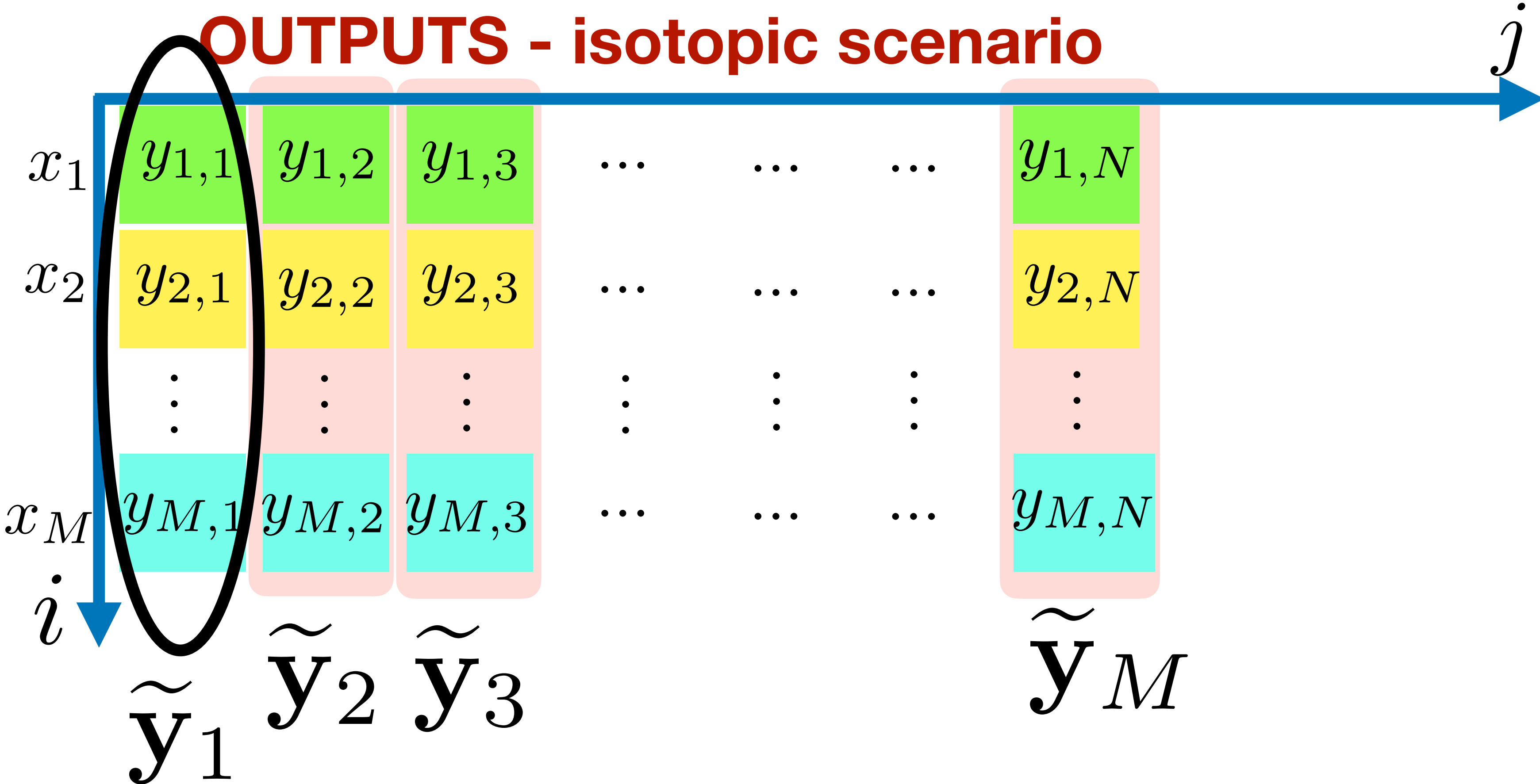
**Other vision of the data (“vertical vision”):**

$$\{x_m, y_{m,j}\}_{m=1}^M$$

$$\tilde{\mathbf{y}}_j = [y_{1,j}, y_{2,j}, \dots, y_{M,j}]^\top$$

# Other vision of the data

Other vision of the data (“vertical vision”):





# But actually what can we do with this data?

For simplicity, considering scalar inputs  $x$ ...just for simplicity and facilitate the comparison with the IRT:

$$x_m \in \mathbb{R} \quad \{x_m, y_{m,j}\}_{m=1}^M \quad \mathbf{X} = [x_1, \dots, x_M]$$

**We have N-parallel classification problems each one with likelihood:**

$$p(y_{m,j} = 1 | x_m) = p_{m,j} = \frac{1}{1 + e^{-(x_m - \mu_j)/s_j}}$$

$$p(\tilde{\mathbf{y}}_j | \mathbf{X}) = p(y_{1,j}, \dots, y_{M,j} | x_1, \dots, x_M) = \prod_{m:y_{m,j}=1} p_{m,j} \prod_{m:y_{m,j}=0} (1 - p_{m,j})$$

$$p(y_{1,j}, \dots, y_{M,j} | x_1, \dots, x_M) = \prod_{m=1}^M p_{m,j}^{y_{m,j}} (1 - p_{m,j})^{(1-y_{m,j})}$$

# But actually what can we do with this data?

We have  $N$ -parallel classification problems each one with likelihood function. **They share the same  $M$  inputs.**

We have to find, in this case with scalar inputs,  $N$  different pairs of  $\mu$  and  $s$ , one for each parallel classification problems.

$$p(y_{m,j} = 1 | x_m) = p_{m,j} = \frac{1}{1 + e^{-(x_m - \mu_j)/s_j}}$$

$$\mu_j, s_j \quad \text{for } j = 1, \dots, N$$