**Machine Learning**
**Regression**

Charlotte Pelletier
*Univ. Bretagne Sud – IRISA Vannes*
Based on C. Friguet's lecture.

28 September 2021

# Content

Recall

- Modelling the relationship between the target variable $y$ and the explanatory variables $X_1, X_2, \cdots, X_d$
    - $y$ is the **qualitative** (binary) variable to be explained
    - $X_i$ are quantitative or qualitative explanatory variables
- → The **linear model** is **not tailored**

Recall

- Modelling the relationship between the target variable $y$ and the explanatory variables $X_1, X_2, \cdots, X_d$
  - $y$ is the **qualitative** (binary) variable to be explained
  - $X_i$ are quantitative or qualitative explanatory variables
- ➜ The **linear model** is **not tailored** because the variable $y$ is not quantitative!
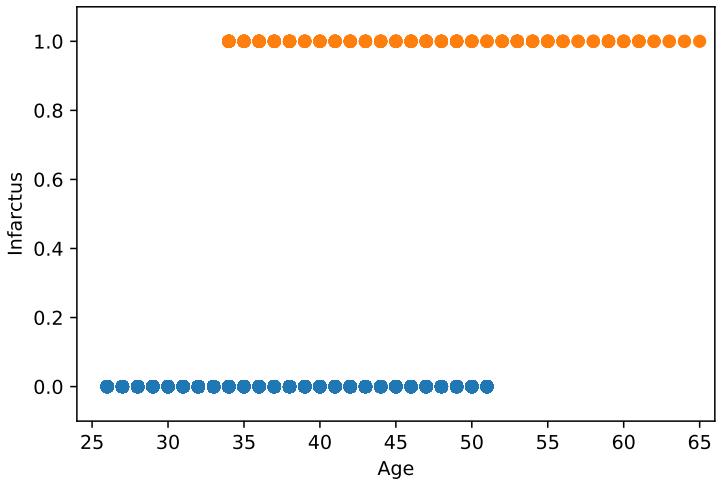
## Example B

We want to evaluate the risk of myocardial infarction (heart attack) on women from a case-study. The collected factors are: taking contraceptive, age, weight, smoking, high blood pressure, a family history of cardiovascular disease, *etc*.
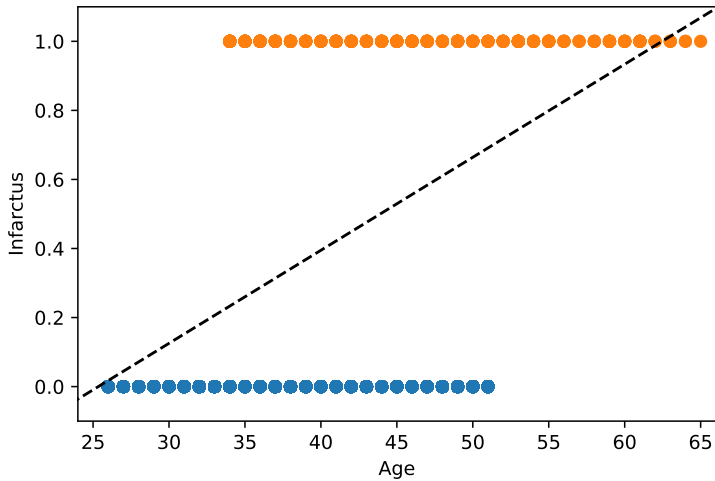
*Data subset:*

| Infarction | Taking contraceptive | Age | Weight | $\cdots$ |
|:---:|:---:|:---:|:---:|:---:|
| no | no | 47 | 48 | $\cdots$ |
| no | no | 35 | 53 | $\cdots$ |
| no | yes | 62 | 41 | $\cdots$ |
| yes | yes | 47 | 45 | $\cdots$ |
| yes | yes | 63 | $\cdots$ | |
| yes | no | 45 | 69 | $\cdots$ |
| yes | no | 90 | 84 | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

*Data from Institut de Santé Publique, d'Epidémiologie et de Développement (ISPED), Bordeaux*
*Source: http://www.biostatisticien.eu/springeR/jeuxDonnees5.html*

## Content

Adjustment

- we do not modelled $Y$ by a linear relationship
- we look at the probabilities: $\mathbf{P}(Y = 0 | X = x)$ and $\mathbf{P}(Y = 1 | X = x)$
- the knowledge of $\mathbf{P}(Y = 1 | X = x)$ induces the knowledge of $\mathbf{P}(Y = 0 | X = x)$ because

Adjustment

- we do not modelled $Y$ by a linear relationship
- we look at the probabilities: $\mathbf{P}(Y = 0|X = x)$ and $\mathbf{P}(Y = 1|X = x)$
- the knowledge of $\mathbf{P}(Y = 1|X = x)$ induces the knowledge of $\mathbf{P}(Y = 0|X = x)$ because $\mathbf{P}(Y = 0|X = x) = 1 - \mathbf{P}(Y = 1|X = x)$

# Bernoulli's law

Let us write out

$$\mathbf{P}(Y = 1 | X = x) = \pi(x)$$

and

$$\mathbf{P}(Y = 0 | X = x) = 1 - \pi(x),$$

where $\pi(x) \in [0, 1]$.

We can model $Y$ by
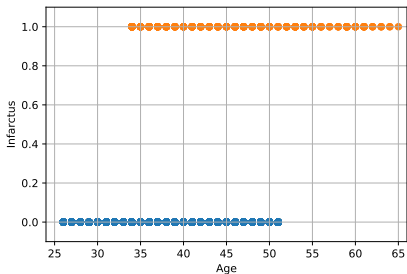
$$Y | X = x \sim \mathcal{B}(\pi(x))$$

(Bernoulli's law)

## Properties (Bernoulli's law)

$$\mathbf{P}(Y = y | X = x) = \pi(x)^y \big(1 - \pi(x)\big)^{1-y}$$

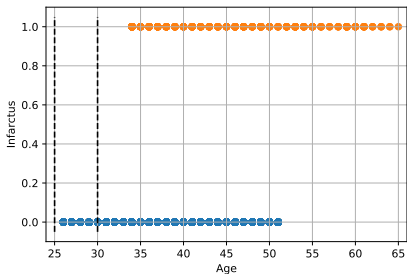$$\mathbb{E}(Y | X = x) = \pi(x) \qquad Var(Y | X = x) = \pi(x)\big(1 - \pi(x)\big)$$

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|----------------|
|        |       |     |     |                |
|        |       |     |     |                |
|        |       |     |     |                |
|        |       |     |     |                |
|        |       |     |     |                |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

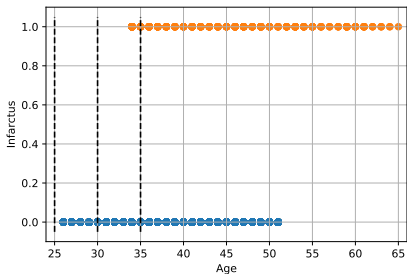| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|----------------|
| [25,30] | 454 | 454 | 0 | 0.00 |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

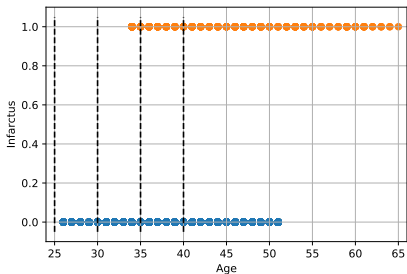| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|----------------|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| | | | | |
| | | | | |
| | | | | |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

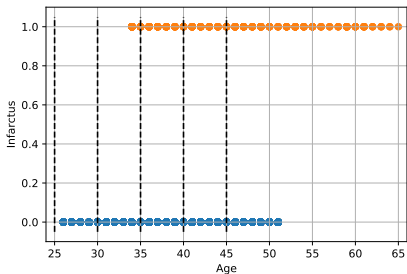| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|----------------|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| (35,40] | 780 | 709 | 71 | 0.10 |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

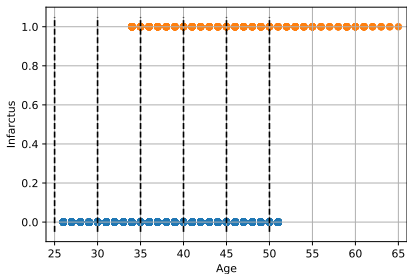| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|----------------|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| (35,40] | 780 | 709 | 71 | 0.10 |
| (40,45] | 608 | 505 | 103 | 0.17 |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

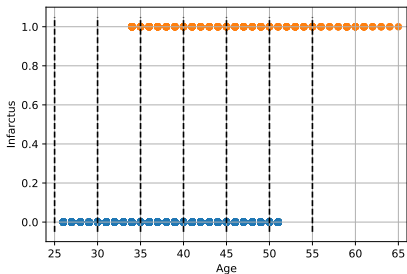| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|---|---|---|---|---|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| (35,40] | 780 | 709 | 71 | 0.10 |
| (40,45] | 608 | 505 | 103 | 0.17 |
| (45,50] | 554 | 355 | 199 | 0.36 |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

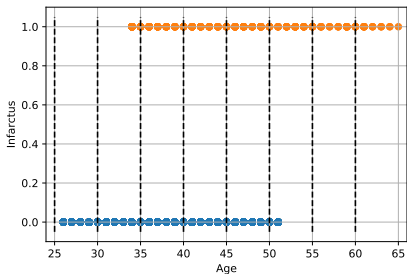| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|--------|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| (35,40] | 780 | 709 | 71 | 0.10 |
| (40,45] | 608 | 505 | 103 | 0.17 |
| (45,50] | 554 | 355 | 199 | 0.36 |
| (50,55] | 134 | 30 | 104 | 0.78 |
| | | | | |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

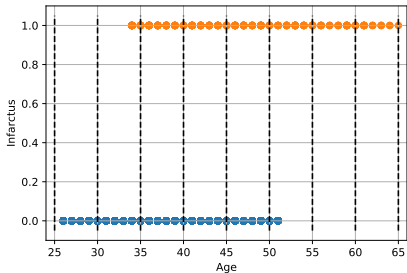| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|-------|-----|-----|----------------|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| (35,40] | 780 | 709 | 71 | 0.10 |
| (40,45] | 608 | 505 | 103 | 0.17 |
| (45,50] | 554 | 355 | 199 | 0.36 |
| (50,55] | 134 | 30 | 104 | 0.78 |
| (55,60] | 33 | 0 | 33 | 1.00 |

- Estimation of $\pi(x)$: $\hat{\pi}(x)$ = proportion of $Y = 1$ for $X$ divided into categories

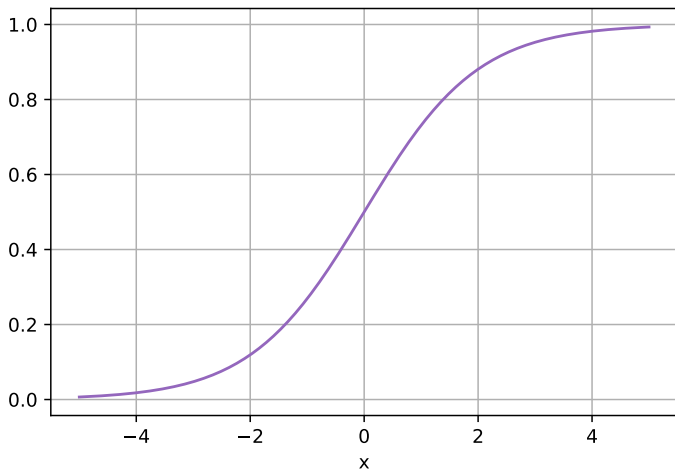| categ. | $n_c$ | #0 | #1 | $\hat{\pi}(x)$ |
|--------|------|-----|-----|------|
| [25,30] | 454 | 454 | 0 | 0.00 |
| (30,35] | 958 | 913 | 45 | 0.04 |
| (35,40] | 780 | 709 | 71 | 0.10 |
| (40,45] | 608 | 505 | 103 | 0.17 |
| (45,50] | 554 | 355 | 199 | 0.36 |
| (50,55] | 134 | 30 | 104 | 0.78 |
| (55,60] | 33 | 0 | 33 | 1.00 |
| (60,65] | 17 | 0 | 17 | 1.00 |

- $\pi(x)$ is not a linear function of $x$!
- We rather want to model $\pi(x)$ by a function with a "$\mathcal{S}$" shape $\Rightarrow$ one possibility is the logistic function

## Logistic function

$$
\begin{aligned}
s \quad : \quad & \mathbb{R} \to \, ]0, 1[ \\
& x \to s(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}
\end{aligned}
$$

It is a specific case of the sigmoid function.

# Logistic function

## Logistic function

Shape of the function curve (more general)
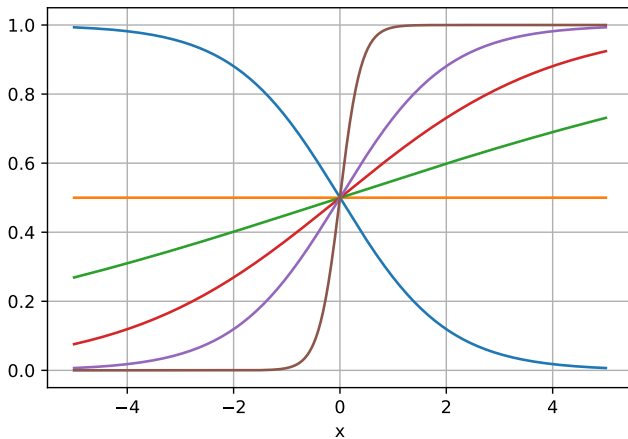
$$s \; : \; \mathbb{R} \to ]0, 1[$$
$$x \to s(x) = \frac{1}{1 + e^{-x\beta}} = \frac{e^{x\beta}}{1 + e^{x\beta}}$$
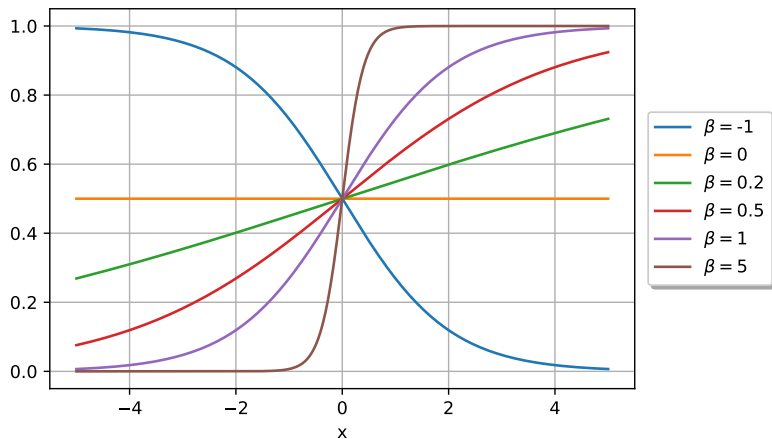
for several $\beta$ values

- $\beta = 0 \Rightarrow$ constant function
- "small" $\beta$ values $\Rightarrow$ wide range of $x$ values for which the function is around 0.5 $\Rightarrow$ difficult discrimination
- "big" $\beta$ values $\Rightarrow$ small range of $x$ values for which the function is around 0.5 $\Rightarrow$ easier discrimination
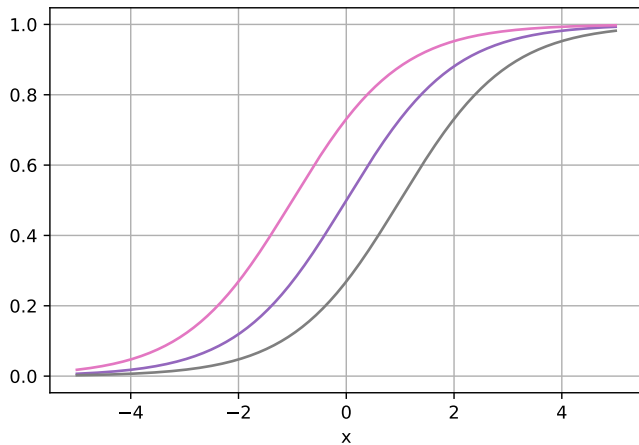- ⚠ It depends on the scale of $x$!

# Logistic model

$$\pi(\mathbf{x}) = \frac{e^{\tilde{\mathbf{x}}\boldsymbol{\beta}}}{1 + e^{\tilde{\mathbf{x}}\boldsymbol{\beta}}} \quad = \quad \frac{e^{\beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d}}{1 + e^{\beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d}}$$

$$\Updownarrow$$

$$logit\Big(\pi(\mathbf{x})\Big) = log\Big(\frac{\pi(\mathbf{x})}{1 - \pi(\mathbf{x})}\Big) \quad = \quad \tilde{\mathbf{x}}\boldsymbol{\beta} = \beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d$$

- Transformation $g$ that creates a linear link between $g\Big(\pi(\mathbf{x})\Big)$ and $\mathbf{x}$: **link function**

**Logit function**

$$logit \quad : \quad ]0, 1[ \rightarrow \mathbb{R}$$
$$p \quad \rightarrow logit(p) = log\Big(\frac{p}{1 - p}\Big)$$

- Let us write out:

$$f_{\boldsymbol{\beta}}(\mathbf{x}_i) = \frac{e^{\tilde{x}_i \boldsymbol{\beta}}}{1 + e^{\tilde{x}_i \boldsymbol{\beta}}} = \mathbb{P}(Y = 1 | X = \tilde{\mathbf{x}}_i)$$

- We want $\boldsymbol{\beta}$ such that $f_{\boldsymbol{\beta}}(\mathbf{x}_i)$ is as closed as possible of $y_i$ for all the training data $\{\mathbf{x}_i, y_i\}_{i=1}^m$ where $\mathbf{x}_i \in \mathbb{R}^d$

## Logistic model

- Model likelihood:

$$L(\boldsymbol{\beta}, y) = \prod_{i=1}^{m} \mathbf{P}(Y = y_i | X = \mathbf{x}_i) = \prod_{i=1}^{m} f_{\boldsymbol{\beta}}(\mathbf{x}_i)^{y_i} \times (1 - f_{\boldsymbol{\beta}}(\mathbf{x}_i))^{1-y_i}$$

## Logistic model

- Model likelihood:

$$L(\boldsymbol{\beta}, y) = \prod_{i=1}^{m} \mathbf{P}(Y = y_i | X = \mathbf{x}_i) = \prod_{i=1}^{m} f_{\boldsymbol{\beta}}(\mathbf{x}_i)^{y_i} \times (1 - f_{\boldsymbol{\beta}}(\mathbf{x}_i))^{1-y_i}$$

- Log-likelihood

$$\mathcal{L}(\boldsymbol{\beta}, y) = \sum_{i=1}^{m} \Big[ y_i log\Big(f_{\boldsymbol{\beta}}(\mathbf{x}_i)\Big) + (1 - y_i) log\Big(1 - f_{\boldsymbol{\beta}}(\mathbf{x}_i)\Big) \Big]$$

- Model likelihood:

$$L(\boldsymbol{\beta}, y) = \prod_{i=1}^{m} \mathbf{P}(Y = y_i | X = \mathbf{x}_i) = \prod_{i=1}^{m} f_{\boldsymbol{\beta}}(\mathbf{x}_i)^{y_i} \times (1 - f_{\boldsymbol{\beta}}(\mathbf{x}_i))^{1-y_i}$$

- Log-likelihood

$$\mathcal{L}(\boldsymbol{\beta}, y) = \sum_{i=1}^{m} \left[ y_i log\left( f_{\boldsymbol{\beta}}(\mathbf{x}_i) \right) + (1 - y_i) log\left( 1 - f_{\boldsymbol{\beta}}(\mathbf{x}_i) \right) \right]$$

- Maximisation of (log-)likelihood
  $\Rightarrow$ cancel the derivatives of $L(\boldsymbol{\beta}, y)$
  - **no explicit analytical solution**
  - there is a need for optimisation iterative algorithms

- **Objective**: find the best $\beta$ to minimise the global (quadratic) cost:

$$\underset{\beta}{\mathrm{argmin}}\Big(\mathbf{J}(\beta)\Big)$$

where

$$\mathbf{J}(\boldsymbol{\beta}) = -\frac{1}{m}\sum_{i=1}^{m}\Big[y_i log\Big(f_{\boldsymbol{\beta}}(\mathbf{x}_i)\Big) + (1 - y_i)log\Big(1 - f_{\boldsymbol{\beta}}(\mathbf{x}_i)\Big)\Big]$$

# Principle

- Objective: find the minimum of the cost-function
- Principle: iterative algorithm
  1. initialisation: $\beta^{(0)}$
  2. at each step $k$, change the $\boldsymbol{\beta}^{(k-1)}$ value to decrease $\mathbf{J}(\boldsymbol{\beta}^{(k)})$
  3. stop the process when a minimum is reached

## Iteration $k$ of the gradient descent algorithm

For the $\beta_j$ parameter

$$\beta_j^{(k)} \quad := \quad \beta_j^{(k-1)} - \alpha \frac{\partial}{\partial \beta_j} \mathbf{J}(\boldsymbol{\beta}^{(k-1)})$$

where:

- $\frac{\partial}{\partial \beta_j}$:
- $\alpha$:

- Objective: find the minimum of the cost-function
- Principle: iterative algorithm
  1. initialisation: $\beta^{(0)}$
  2. at each step $k$, change the $\boldsymbol{\beta}^{(k-1)}$ value to decrease $\mathbf{J}(\boldsymbol{\beta}^{(k)})$
  3. stop the process when a minimum is reached

## Iteration $k$ of the gradient descent algorithm

For the $\beta_j$ parameter

$$\beta_j^{(k)} \quad := \quad \beta_j^{(k-1)} - \alpha \frac{\partial}{\partial \beta_j} \mathbf{J}(\boldsymbol{\beta}^{(k-1)})$$

where:

- $\frac{\partial}{\partial \beta_j}$: partial derivatives
- $\alpha$: learning rate (hyperparameter to be determined by the user)

- Iteration $k$ of the gradient descent algorithm?

**Iteration $k$ of the gradient descent algorithm - logistic regression**

$$\beta_j^{(k)} \quad := \quad \beta_j^{(k-1)} - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{\boldsymbol{\beta}^{(k-1)}}(\mathbf{x}_i) - y_i \right) x_{ij}$$

Note: $\forall i, x_{i0} = 1$

# For the logistic regression

- Iteration $k$ of the gradient descent algorithm?
- Same formula than for the linear regression...

**Iteration $k$ of the gradient descent algorithm - logistic regression**

$$\beta_j^{(k)} \quad := \quad \beta_j^{(k-1)} - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{\boldsymbol{\beta}^{(k-1)}}(\mathbf{x}_i) - y_i \right) x_{ij}$$

Note: $\forall i, x_{i0} = 1$

**Optimisation**

There are more advanced versions of the gradient descent algorithm:

- conjugated gradient, BFGS[1], (limited-memory) L-BFGS
    - no need to fix $\alpha$ (it is done automatically by the algorithm)
    - faster (it requires less iterations) than the classical gradient descent algorithm
    - but more complex to implement (these variants are avaible in most of the software/programming languages)

---

[1]Broyden-Fletcher-Goldfarb-Shanno