## I/ The model

Let consider a dataset with $\underline{y} = (y_1, \ldots, y_m)'$ and $\underline{x} = (\underline{x}_1', \ldots, \underline{x}_m')'$ when $y_i \in \{0, 1\}$ and $\underline{x}_i = (x_{i1}, \ldots, x_{id})' \in \mathbb{R}^d$.

Let also denote $\mathcal{D} = \{(\underline{x}_1, y_1), \ldots, (\underline{x}_m, y_m)\}$ the dataset composed with each instance with its label.

Ex

| y | $X_1$ | ..... | $X_d$ |
|---|---|---|---|
| yes | 20 | ---- | 5 |
| yes | 25 | ---- | 4 |
| no | 14 | ---- | 0 |

The goal is to predict $y_i$ based on $\underline{x}_i$

$$\boxed{y_i \leadsto B(\pi(\underline{x}_i))}$$ $y_i$ is assumed to follow a Bernoulli distribution (success or fail) where the probability of success depends on $\underline{x}_i$ (the covariates).

$$P(y_i \mid \underline{x}_i) = \begin{cases} \pi(\underline{x}_i) & \text{if } y_i = 1 \\ 1 - \pi(x_i) & \text{if } y_i = 0 \end{cases}$$ or equivalently $P(y_i \mid \underline{x}_i) = \pi(\underline{x}_i)^{y_i} (1 - \pi(\underline{x}_i))^{1-y_i}$

Let also notice that $E[y_i \mid \underline{x}_i] = \pi(\underline{x}_i)$

Logistic regression assumes that $$\boxed{\begin{aligned} \text{logit}(\pi(\underline{x}_i)) &:= \log\left(\frac{\pi(\underline{x}_i)}{1 - \pi(\underline{x}_i)}\right) \\ &= \beta_0 + \beta_1 x_{i1} + \cdots + \beta_d x_{id} \\ &= \beta' \tilde{\underline{x}}_i \end{aligned}}$$

$$\text{logit}^{-1}, \quad \pi(\underline{x}_i) = \text{logit}^{-1}(\beta' \tilde{\underline{x}}_i)$$
$$= \sigma(\beta' \tilde{\underline{x}}_i) \iff$$

with $\sigma : z \longmapsto \dfrac{1}{1 + \exp(-z)}$

with $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{pmatrix}$ the vector of parameters to estimate

and $\tilde{\underline{x}}_i = (1, \underline{x}_i')' = \begin{pmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{pmatrix}$

Case $d = 1$



data points

logistic regression curve $\pi(\underline{x}_i)$

The statistical inference issue consist in estimating $\beta$ (the vector of parameters) based on the dataset $\mathcal{D}$.

Strategy :   → Build an estimator $\hat{\beta}$ of $\beta$ by maximum likelihood

⚠ $\hat{\beta}$ is a random variable since it depend on the data $y_i$ which is random

→ Study the properties of $\hat{\beta}$ : asymptotically normal/Gaussian thus allows to compute confidence interfalls and to make tests on the parameters.

# II / Parameters estimation

Let consider $p(\underline{y} | \underline{x}) = p(y_1, ..., y_m | \underline{x}_1, ....., \underline{x}_m)$ it is the probability of the responses given the covariates (conditionnal likehood) ~~our model~~

Since the data $(\underline{x}_1, y_1), ..., (\underline{x}_m, y_m)$ are assumed to be independent we can write :
$$\underbrace{p(\underline{y} | \underline{x})}_{: \mathcal{L}(\beta)} = \prod_{i=1}^{m} p(y_i | \underline{x}_i) = \prod_{i=1}^{m} \pi(\underline{x}_i)^{y_i} (1 - \pi(\underline{x}_i))^{1-y_i}$$

The likelihood which is a function of $\beta$ $\left(\text{since } \pi(\underline{x}_i) \text{ is a function of } \beta\right)$

We often consider the log-likelihood $\ell(\beta) = \log \mathcal{L}(\beta)$.

$$\ell(\beta) = \sum_{i=1}^{m} \left[ y_i \log \pi(\underline{x}_i) + (1-y_i) \log (1 - \pi(\underline{x}_i)) \right]$$

$$\ell(\beta) = \sum_{i=1}^{m} \left[ y_i \log \left( \frac{1}{1 + \exp(-\beta' \underline{\tilde{x}}_i)} \right) + (1-y_i) \log \left( \frac{\exp(-\beta' \underline{\tilde{x}}_i)}{1 + \exp(-\beta' \underline{\tilde{x}}_i)} \right) \right]$$

$$\ell(\beta) = \sum_{i=1}^{m} y_i \beta' \underline{\tilde{x}}_i - \log (1 + \exp(\beta' \underline{\tilde{x}}_i))$$

$\ell : \mathbb{R}^{d+1} \to \mathbb{R}$.   We want to maximise $\ell$ with respect to $\beta$ thus we will compute the gradiant.

$$\frac{\partial \ell(\beta)}{\partial \beta} = \sum_{i=1}^{m} y_i \tilde{x}_i - \tilde{x}_i \frac{\exp(\beta'\tilde{x}_i)}{1+\exp(\beta'\tilde{x}_i)} = \sum_{i=1}^{m} \tilde{x}_i (y_i - \pi(\tilde{x}_i))$$

$$\frac{\partial \ell(\beta)}{\partial \beta} \in \mathbb{R}^{d+1}.$$ Solving $\frac{\partial \ell(\beta)}{\partial \beta} = 0$ has not closed-form

thus it is needed to use an iterative algorithm such as Newton-Raphson or gradient descent.

Newton Raphson
$$\beta^{(k+1)} = \beta^{(n)} - \left(\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta'}\right)^{-1}\Bigg|_{\beta=\beta^{(n)}} \times \left(\frac{\partial \ell(\beta)}{\partial \beta}\right)\Bigg|_{\beta=\beta^{(n)}}$$

$$\underbrace{\hspace{3cm}}_{\text{Hessian matrix}}$$

We have
$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta'} = -\tilde{X}' V \tilde{X}$$

where $\tilde{X}$ is the matrix with $m \times (d+1)$ composed with $\tilde{x}_i$ in rows

and $V$ is the diagonal matrix of $\pi(\tilde{x}_i)(1-\pi(\tilde{x}_i))$

$$V = \begin{pmatrix} \pi(\tilde{x}_1)(1-\pi(\tilde{x}_1)) & & & 0 \\ & \ddots & & \\ & & \pi(\tilde{x}_i)(1-\pi(\tilde{x}_i)) & \\ 0 & & & \\ & & & \pi(\tilde{x}_m)(1-\pi(\tilde{x}_m)) \end{pmatrix}$$

Let notice that $\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta'}$ is definite negative if $\tilde{X}$ is of maximal

rank $\left(\text{ie. } d+1 \text{ if } d+1 \leq m\right)$.

## III / Properties of $\hat{\beta}$

Since $\hat{\beta}$ is the maximum likelihood estimator it is asymptotically unbiased and follows asymptotically a normal distribution

with asymptotic covariance matrix $\hat{V}(\hat{\beta}) = \left[-\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta'}\Big|_{\beta=\hat{\beta}}\right]^{-1}$

which is the inverse of the Fisher information matrix.

Thus $\qquad \hat{\beta} \underset{\text{as } m \to +\infty}{\approx} \mathcal{N}_{d+1}\left(\beta, \left(\tilde{X}'\hat{V}\tilde{X}\right)^{-1}\right)$

$\Rightarrow$ Thus possible to make similar computations as in the linear regression. Since $m$ assumed to be large we only rely on ~~Student (t)~~ ~~distri~~ <u>normal distribution</u> and not anymore on student (t)

## IV / Tests, confidence intervals and model choice

### A/ Test on $\beta_j$

Let consider testing if variable $j$ has an effect or not :

$$H_0 : \beta_j = 0 \qquad \text{against} \qquad H_1 : \beta_j \neq 0$$

There are three possibilities :

- <u>likelihood ratio test</u> :

$$LRT = 2 \log \frac{\max_{\beta} \mathcal{L}(\beta)}{\max_{\substack{\beta \\ \text{s.c } \beta_j = 0}} \mathcal{L}(\beta)} = 2 \log \frac{\max_{\beta} \mathcal{L}_{H_1}(\beta)}{\max_{\beta} \mathcal{L}_{H_0}(\beta)}$$

$$= 2 \left[ \underbrace{\ell(\hat{\beta})}_{\substack{\text{maximum} \\ \text{log-likelihood}}} - \underbrace{\ell(\hat{\beta}_{H_0})}_{\substack{\text{maximum log-likelihood} \\ \text{under } H_0}} \right] = D_0 - D_1 \quad \text{with} \quad \begin{array}{l} D_0 = -2\,\ell(\hat{\beta}_{H_0}) \\ \text{and} \\ D_1 = -2\,\ell(\hat{\beta}) \end{array}$$

Under $H_0$ : $\quad LRT \rightsquigarrow \chi_1^2$

- <u>Wald test</u> $\qquad$ We know that $\quad \dfrac{\hat{\beta}_j - \beta_j}{\hat{\sigma}(\hat{\beta}_j)} \approx \mathcal{N}(0,1)$
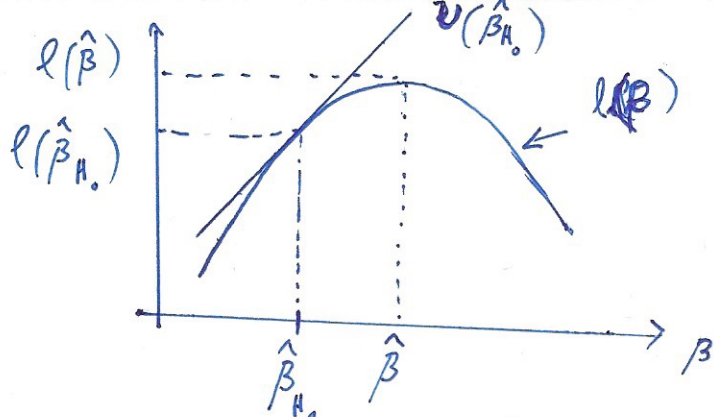
Thus under $H_0$ : $\quad \dfrac{\hat{\beta}_j}{\hat{\sigma}(\hat{\beta}_j)} \approx \mathcal{N}(0,1) \Rightarrow \dfrac{\hat{\beta}_j^2}{\hat{\sigma}(\hat{\beta}_j)^2} \rightsquigarrow \chi_1^2$

- <u>Score test</u> $\qquad U(\hat{\beta}_{H_0})' \hat{V}(\hat{\beta}_{H_0}) U(\hat{\beta}_{H_0}) \longrightarrow \chi_1^2$

where $\hat{V}(\hat{\beta}_{H_0})$ is the inverse of the Fisher information matrix and $U(\hat{\beta}_{H_0})$ the vector of partial derivative of the log-likelihood both estimated under $H_0$.

$\ell(\hat{\beta})$  
$\ell(\hat{\beta}_{H_o})$  
$U(\hat{\beta}_{H_o})$  
$\ell(\beta)$  
$\hat{\beta}_{H_o}$   $\hat{\beta}$   $\beta$

$q \times (k-1)$

## B/ Test on a linear combination of the parameter $C\beta = 0$

We can still use the test defined in previous section extended to the multivariate framework.

Let first notice that $C\hat{\beta} \approx \mathcal{N}\left(C\beta, \; C(\tilde{X}'\hat{V}\tilde{X})^{-1}C'\right)$

Thus under $H_o$ : $\quad C\hat{\beta} \approx \mathcal{N}\left(0, \; C(\tilde{X}'\hat{V}\tilde{X})^{-1}C'\right)$

And consequently $\quad (C\hat{\beta})'\left(C(\tilde{X}'\hat{V}\tilde{X})^{-1}C'\right)^{-1}(C\hat{\beta}) \approx \chi_q^2$

Thus possible to use ~~test similar to~~ Wald test.

The likelihood ratio test is defined by

$$D_o - D_1 = 2\left(\ell(\hat{\beta}) - \ell(\hat{\beta}_{H_o})\right) \underset{H_o}{\approx} \chi_q^2$$

Where $\ell(\hat{\beta}_{H_o}) = \underset{\beta}{max} \; \ell(\beta)$
$\qquad\qquad\qquad\qquad$ s.s. $C\beta = 0$

And the score test

$$U(\hat{\beta}_{H_o})'\hat{V}(\hat{\beta}_{H_o})U(\hat{\beta}_{H_o}) \longrightarrow \chi_q^2$$

## C/ Model choice

$\qquad$ BIC $= -2\ell(\hat{\beta}) + (d+1)\log n \qquad$ where $d+1$ is the number of estimated parameters

$\qquad$ AIC $= -2\ell(\hat{\beta}) + (d+1)$

AIC and BIC can be used to put models into competitions
the goal is to find the model ( subset of variables ) minimizing
the criterium. This can be done for instance by using a stepwise
approach. (forward, backward, or forward-backward)

BIC tends to select model of lower dimension than AIC

By default the step function of R uses AIC.

## D/ Confidence interval

Since $\hat{\beta}_j \underset{appro}{\sim} \mathcal{N}(\beta_j, \hat{\sigma}(\hat{\beta}_j.))$

One can deduce a $(1-\alpha)$ confidence interval by the

formula $\hat{\beta}_j. \pm z_{\alpha/2} \, \hat{\sigma}(\hat{\beta}_j.)$ with $z_{\alpha/2}$ the $\alpha/2$ upper quantile

of the normal distribution.

## V Case of categorical features

Let assume that a variable $x_j$ is categorical, with $x_j. \in \{1, ... J\}$

Thus a binary coding can be used

| $x_j.$ | $X_{j,blue}$ | $X_{j,red}$ | $X_{j,green}$ |
|--------|--------------|-------------|---------------|
| blue   | 1            | 0           | 0             |
| red    | 0            | 1           | 0             |
| green  | 0            | 0           | 1             |

Reference level not used in design matrix

By default in R the first
level of the variable is used
as reference level ( the column )
is not used to fit
the model

Testing if variable $j$ has an effect consist in testing

$$H_0 : \beta_{red} = \beta_{green} = 0 \quad vs \quad H_1 : \beta_{red} \neq 0 \text{ or } \beta_{green} \neq 0.$$

(possible : see V B )

$$Odds(\underline{x}_i) = \frac{\pi(\underline{x}_i)}{1-\pi(\underline{x}_i)} = \exp(\beta'\tilde{\underline{x}}_i)$$

$$Odds\text{-ratio}(\underline{x}_i, \underline{x}_{i'}) = \frac{odds(\underline{x}_i)}{odds(\underline{x}_{i'})} = \exp(\beta'(\tilde{\underline{x}}_i - \tilde{\underline{x}}_{i'}))$$

If $i$ and $i'$ differ for only one variable $x_{ij} \neq x_{i'j}$ then

$$Odds\text{-ratio}(\underline{x}_i, \underline{x}_{i'}) = \exp(\beta_j(x_{ij} - x_{i'j})) \quad \text{and this variable}$$

is categorical

$$\widehat{Odds\text{-ratio}}(\underline{x}_i, \underline{x}_{i'}) = \exp(\hat{\beta}_j(x_{ij} - x_{i'j}))$$

Then it is possible to derive a confidence interval on
$\widehat{Odds\text{-ratio}}(\underline{x}_i, \underline{x}_{i'})$ :

$$\left[\exp\left(\left(\hat{\beta}_j \pm z_{\alpha/2}\,\hat{\sigma}(\hat{\beta}_j)\right)(x_{ij}-x_{i'j})\right)\right]$$

If the variable is categorical $\exp(\beta_j)$ gives the odds ratio between the considered level and the reference level _____

Remark can add interaction in the model by creating new variables for instance $x_{new} = x_{i1} \times x_{i2}$

## VI Multi-class logistic regression  $y_i \in \{1, ..., K\}$

$$P(y_i = k \mid x_i) = \frac{\exp(\beta_k'\tilde{x}_i)}{1+\sum_{h=1}^{K-1}\exp(\beta_h'\tilde{x}_i)} \quad \text{for } k \in \{1,...,K-1\}$$

$$\text{where } \beta_k = \begin{pmatrix}\beta_{k0}\\\beta_{k1}\\\vdots\\\beta_{kd}\end{pmatrix}$$

$$P(y_i = K \mid x_i) = \frac{1}{1+\sum_{h=1}^{K-1}\exp(\beta_h'\tilde{x}_i)}$$

$\Rightarrow$ Maximum likelihood ...
Can need regularization if too many parameters to estimate