

# 1. Classification

$$\mathbb{R}^n \ni X \rightsquigarrow Y \in \mathbb{R} \leftarrow \text{quantitative}$$

$$\mathbb{R}^n \ni X \rightsquigarrow Y \in \{0, 1\} \leftarrow \text{qualitative}$$

- Classifying is a process of assigning observation to categories
- In the classification problem we work with qualitative responses instead of quantitative  
sklearn, datasets, load\_boston  
sklearn, datasets, load\_iris
- Often we want to predict the probability that observation belongs to each of the categories
- Among classification techniques (classifier) we have
  - logistic regression
  - linear discriminant analysis (LDA)
  - quadratic — u — (QDA)
  - K-nearest neighbors
  - tree, random forests
  - SVM (Support Vector Machines)

## 2. Why not linear regression?

- there is a problem with converting a qualitative response variable with more than 2 classes into a qualitative response that can be used in linear regression
- there is a problem with obtaining reasonable estimates of  $P(Y=1|X)$

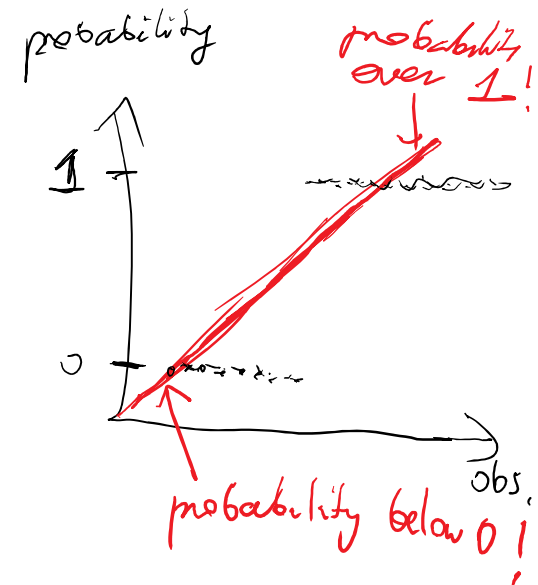
$$Y = \begin{cases} 0 & \text{if cat} \\ 1 & \text{if dog} \\ 2 & \text{if bird} \end{cases}$$

$$Y^{(1)} = \begin{cases} 0 & \text{if not cat} \\ 1 & \text{if cat} \end{cases}$$

$$Y^{(2)} = \begin{cases} 0 & \text{if not dog} \\ 1 & \text{if dog} \end{cases}$$

~~$$Y^{(3)} = \begin{cases} 0 & \text{if not bird} \\ 1 & \text{if bird} \end{cases}$$~~

not necessary



### 3. Logistic Regression - Model

binary classification  $\mathbb{R}^n \ni X \rightsquigarrow Y \in \{0, 1\}$   
 (predictors) (responses) classes/categories

Model

$$p(x) = P(Y=1 | X=x)$$

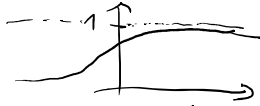
$$= \varphi(\beta_0 + \beta_1 x)$$

$$= \frac{1}{1 + e^{\beta_0 + \beta_1 x}}$$

$$= \frac{e^{\beta_0 + \beta_1 x}}{e^{\beta_0 + \beta_1 x} + 1}$$

$$\varphi: \mathbb{R} \rightarrow [0, 1] \quad \varphi \uparrow$$

$$\varphi(-\infty) = 0, \quad \varphi(\infty) = 1$$



$$\varphi(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

(logistic function/sigmoid)

Hence

$$e^{\beta_0 + \beta_1 x} = p(x) (1 + e^{\beta_0 + \beta_1 x})$$

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

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

$$\log \left( \frac{p(x)}{1 - p(x)} \right) = \beta_0 + \beta_1 x$$

odds  
log odds (logit)

Thus logistic regression has a logit (log odds) that is linear in X,

#### 4. Logistic Regression - fitting

- to fit the model given by  $p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$  we use maximum likelihood method:

$$l(\beta_0, \beta_1) := \prod_{i: y_i=1} p(x_i) \prod_{i: y_i=0} (1 - p(x_i)) \quad \left\{ \begin{array}{l} p(x) = P(Y=1|X) \end{array} \right.$$

- the estimates of  $\beta_0, \beta_1$  are chosen to maximize this likelihood function  $l(\beta_0, \beta_1)$
- intuition: we try to  $\hat{\beta}_0, \hat{\beta}_1$  such that if we put them into  $p(x)$  we get number close to 1 for each observations such that  $Y=1$  and we get number close to 0 for all observations ~~not~~ that  $Y=0$ ,

## 5. Logistic Regression - multiple case

- $$\log \left( \frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

where  $X = (X_1, \dots, X_p)$  - predictors.

- Hence 
$$p(x) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

- We use maximum likelihood method to predict  $\beta_0, \beta_1, \dots, \beta_p$ .

## 6. Logistic Regression - multinomial case

- we can extend two-classes (binary classification) logistic regression to the case of  $K > 2$  classes
- At first, we select a single class to be a baseline, eg. let it be  $K$ th class.

- Then, we replace the model  $p(x) = \frac{e^{\beta_0 + \beta_1 x + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x + \dots + \beta_p x_p}}$  with the model

$$P(Y=j | X=x) = \frac{e^{\beta_{j0} + \beta_{j1} x_1 + \dots + \beta_{jp} x_p}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l0} + \beta_{l1} x_1 + \dots + \beta_{lp} x_p}}$$

$$P(Y=K | X=x) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\beta_{l0} + \beta_{l1} x_1 + \dots + \beta_{lp} x_p}}$$

- Then

$$\log \left( \frac{P(Y=j | X=x)}{P(Y=K | X=x)} \right) = \beta_{j0} + \beta_{j1} x_1 + \dots + \beta_{jp} x_p$$

- Choosing  $K$ th class as the baseline is not important in the sense that we can choose any other class as a baseline

## 7. Logistic Regression - soft-max

- soft-max is an alternative coding for multinomial logistic regression.
- Instead selecting baseline class, we treat all  $K$  classes symmetrically and we assume

$$P(Y=j | X=x) = \frac{e^{\beta_{j0} + \beta_{j1}x_1 + \dots + \beta_{jp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}$$

- Thus, we estimate coefficients for all  $K$  classes and we get

$$\log\left(\frac{P(Y=j | X=x)}{P(Y=j' | X=x)}\right) = (\beta_{j0} - \beta_{j'0}) + (\beta_{j1} - \beta_{j'1})x_1 + \dots + (\beta_{jp} - \beta_{j'p})x_p$$

## 8. Logistic Regression - remarks

- In a linear regression ( $Y = \beta_0 + \beta_1 X$ ) model  $\beta_1$  gives the average change in  $Y$  associated with increasing of  $X$  by one unit

- In a logistic regression model ( $\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 X$ , where  $p(x) = P(Y=1|X=x)$ ) increasing  $X$  by one unit changes the log odds ( $\log\left(\frac{p(x)}{1-p(x)}\right)$ ) by  $\beta_1$ .

- Logistic regression (for binary classification;  $p(x) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$ )

corresponds to neural net with one neuron and activation function  $\varphi(x) = \frac{1}{1 + e^{-x}}$  and cost function

function —  $l(\beta_0, \beta_1)$

