

# Logistic Regression

Nathanaël Carraz Rakotonirina

Mathématiques Informatique et Statistique Appliquées (MISA)  
Université d'Antananarivo

It is a classification model  $p(y|\mathbf{x}; \boldsymbol{\theta})$ .

- ▶  $\mathbf{x} \in \mathbb{R}^D$  : input
- ▶  $y \in \{1, \dots, C\}$  : class label
- ▶  $\boldsymbol{\theta}$  : parameters

If  $C = 2$ , it is called **binary logistic regression** and if  $C > 2$ , it is known as **multiclass logistic regression**.

# Binary logistic regression

Since we want to predict  $y \in 0, 1$  given some inputs  $\mathbf{x}$ , the model is of the form

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y; f(\mathbf{x}; \boldsymbol{\theta}))$$

where  $f(\mathbf{x}; \boldsymbol{\theta})$  is a function giving the parameter of the distribution hence must satisfy  $0 \leq f(\mathbf{x}; \boldsymbol{\theta}) \leq 1$ . To allow  $f$  to be any function, we use:

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y; \sigma(f(\mathbf{x}; \boldsymbol{\theta})))$$

where  $\sigma$  is the **sigmoid** (S-shaped) or **logistic** function:

$$\begin{aligned} \sigma : \mathbb{R} &\rightarrow [0, 1] \\ z &\mapsto \sigma(z) = \frac{1}{1 + e^{-z}} \end{aligned}$$

$z$  is called the logit or the pre-activation.

For logistic regression, we choose a linear function  $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}^\top \mathbf{x} + b$ . The model has the form

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y; \sigma(\mathbf{w}^\top \mathbf{x} + b))$$

This means

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \sigma(\mathbf{w}^\top \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)}}$$

# Decision boundary

During prediction, we have

$$\hat{y} = \begin{cases} 1 & \text{if } p(y = 1|\mathbf{x}; \boldsymbol{\theta}) > 0.5 \\ 0 & \text{if } p(y = 1|\mathbf{x}; \boldsymbol{\theta}) < 0.5 \end{cases}$$

which is the same as

$$\hat{y} = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} + b > 0 \\ 0 & \text{if } \mathbf{w}^\top \mathbf{x} + b < 0 \end{cases}$$

The **decision boundary** is  $\mathbf{w}^\top \mathbf{x} + b = 0$ . It is a linear hyperplane with normal vector  $\mathbf{w}$  and an offset  $b$  from the origin. It separates the space into 2 half-spaces.

The data is said to be **linearly separable** when the examples can be perfectly separated by the linear hyperplane.

# Maximum likelihood estimation

We note  $\mu_i = \sigma(z_i) = \sigma(\mathbf{w}^\top \mathbf{x}_i)$ . The negative log likelihood is

$$\begin{aligned}\text{NLL}(\mathbf{w}) &= - \sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = - \sum_{i=1}^N \log \text{Ber}(y_i; \mu_i) \\ &= - \sum_{i=1}^N \log [\mu_i^{y_i} + (1 - \mu_i)^{1-y_i}] \\ &= - \sum_{i=1}^N [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)] \\ &= \sum_{i=1}^N \mathbb{H}(y_i, \mu_i)\end{aligned}$$

where  $\mathbb{H}$  is the **binary cross entropy**. This objective is convex and can be minimized using gradient-based methods.

# Multinomial logistic regression

It is a classification model of the form

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y; f(\mathbf{x}; \boldsymbol{\theta}))$$

We note  $\mu = f(\mathbf{x}; \boldsymbol{\theta})$  (here  $f : \mathbb{R}^D \rightarrow \mathbb{R}^C$ ) . It must satisfy  $0 \leq \mu_i \leq 1$  and  $\sum_{i=1}^C \mu_i = 1$ . To allow  $f$  to be any function, we pass it to the **softmax** function

$$\mathcal{S} : \mathbb{R}^C \rightarrow [0, 1]^C$$
$$z \mapsto \mathcal{S}(z) = \left[ \frac{e^{z_1}}{\sum_{i=1}^C e^{z_i}}, \dots, \frac{e^{z_C}}{\sum_{i=1}^C e^{z_i}} \right]$$

You might want to use the log-sum-exp trick to avoid numerical overflow when computing the softmax.

We use a linear function  $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\mathbf{x} + \mathbf{b}$  where  $\mathbf{W}$  is a  $C \times D$  matrix and  $\mathbf{b}$  is a  $C$  dimensional vector. The model is of the form

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y; \mathcal{S}(\mathbf{W}\mathbf{x} + \mathbf{b}))$$

If we note  $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$  the  $C$  dimensional vector of logits, we have

$$p(y = c|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{z_c}}{\sum_{i=1}^C e^{z_i}}$$



# Maximum likelihood estimation

We keep  $\boldsymbol{\mu} =$ . The negative log likelihood is

$$\begin{aligned}\text{NLL}(\mathbf{w}) &= -\log \prod_{i=1}^N \prod_{c=1}^C \mu_{ij}^{y_{ic}} \\ &= -\sum_{i=1}^N \sum_{i=c}^C y_{ic} \log \mu_{ic} \\ &= \sum_{i=1}^N \mathbb{H}(y_i, \mu_i)\end{aligned}$$

where  $\mu_{ic} = p(y_i = c | \mathbf{x}_i; \boldsymbol{\theta}) = (\mathcal{S}(\mathbf{W}\mathbf{x}_i + \mathbf{b}))_c$  and  $y_{ic} = \mathbb{I}(y_i = c)$ . This objective is also convex and can be minimized using gradient descent.

# Do you want more ?

Of course you do !

## Explore further

- ▶ Robust logistic regression
- ▶ Bayesian logistic regression
- ▶ Multilabel classification
- ▶ Hierarchical classification