

Logistic Regression

Machine Learning

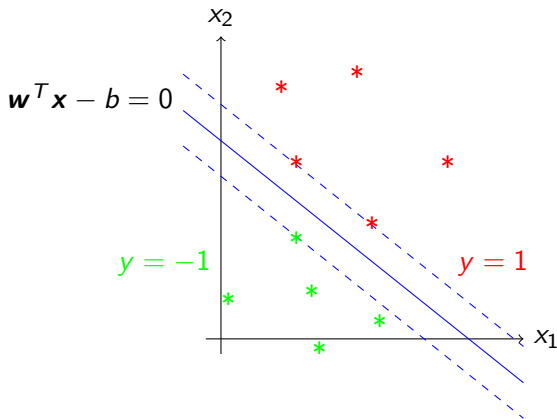
Axel Carlier, Jean-Yves Tournet

2024

Outline

- 1 Logistic regression
- 2 Single-layer perceptron

Reminder: Support Vector Machines



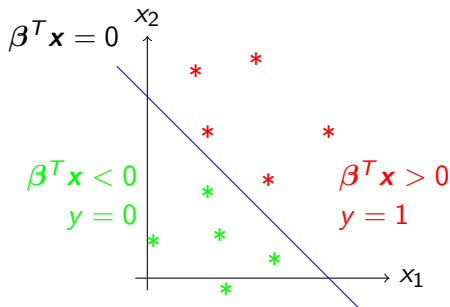
Find an optimal separating hyperplane, such that $|\mathbf{w}^T \mathbf{x}_i - b| \geq 1$.

Logistic regression for binary classification

For linear regression, we again define a linear model:

$$z = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p = \boldsymbol{\beta}^T \mathbf{x}, \text{ with } \mathbf{x} = [1, x_1, \dots, x_p]^T$$

This linear model acts as a **separator** for the 2 classes: $y \in \{0, 1\}$



Back to the Bayesian classifier

Bayesian rule for two classes ω_1 and ω_2

$$\begin{aligned}d^*(\mathbf{x}) = \omega_1 &\Leftrightarrow f(\mathbf{x}|\omega_1)P(\omega_1) \geq f(\mathbf{x}|\omega_2)P(\omega_2) \\ &\Leftrightarrow \frac{f(\mathbf{x}|\omega_1)P(\omega_1)}{f(\mathbf{x}|\omega_2)P(\omega_2)} \geq 1 \\ &\Leftrightarrow \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^T \Sigma_1^{-1}(\mathbf{x} - \mathbf{m}_1)\right] P(\omega_1)}{\exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^T \Sigma_2^{-1}(\mathbf{x} - \mathbf{m}_2)\right] P(\omega_2)} \geq 0\end{aligned}$$

In the Gaussian case with $\Sigma_1 = \Sigma_2 = \Sigma$, one obtains

$$d^*(\mathbf{x}) = \omega_1 \Leftrightarrow \sigma(\beta^T \mathbf{x}) = \frac{1}{1 + e^{-\beta^T \mathbf{x}}} \leq 0.5$$

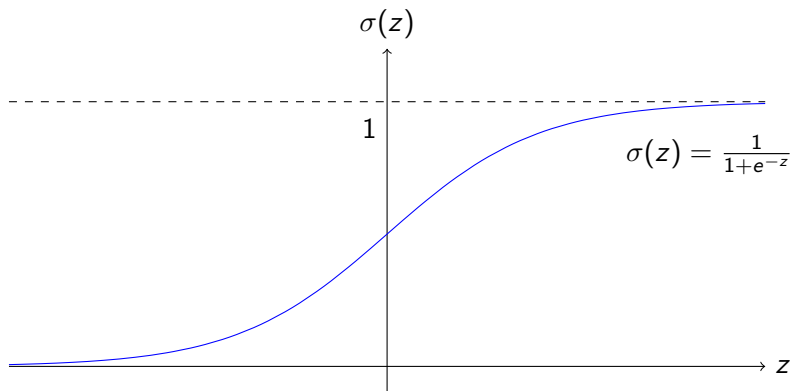
with

$$\beta^T \mathbf{x} = \frac{1}{2} \mathbf{x}^T \Sigma^{-1} (\mathbf{m}_1 - \mathbf{m}_2) + \frac{1}{2} \mathbf{m}_1^T \Sigma^{-1} \mathbf{m}_1 - \frac{1}{2} \mathbf{m}_2^T \Sigma^{-1} \mathbf{m}_2 + \ln \left[\frac{P(\omega_1)}{P(\omega_2)} \right].$$

Logistic regression for binary classification

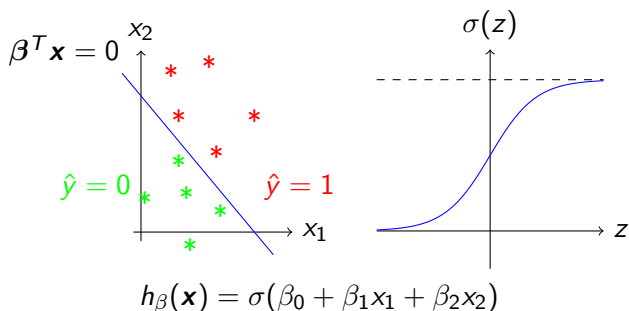
The **sigmoid** or **logistic** function is defined as:

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



Logistic regression for binary classification

2D case where data is linearly separable:



- We predict $\hat{y} = 1$ if $\beta_0 + \beta_1 x_1 + \beta_2 x_2 > 0 \Leftrightarrow \sigma(\beta^T \mathbf{x}) > 0.5$
- We predict $\hat{y} = 0$ if $\beta_0 + \beta_1 x_1 + \beta_2 x_2 \leq 0 \Leftrightarrow \sigma(\beta^T \mathbf{x}) \leq 0.5$

Logistic regression for binary classification

How to estimate β in this framework?

Let $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ be our **training set**, with n examples:

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_p^{(i)} \end{bmatrix} \in \mathbb{R}^p \text{ and } y^{(i)} \in \{0, 1\}, \forall i \in \{1, \dots, n\}$$

and the **model** is defined by the sigmoid function:

$$P(y = 1 | \mathbf{x}; \beta) = \sigma(\beta^T \mathbf{x}) = \frac{1}{1 + e^{-\beta^T \mathbf{x}}}$$

One can define the following classification rule

- Predict $\hat{y} = 1$ if $\sigma(\beta^T \mathbf{x}) > 0.5$
- Predict $\hat{y} = 0$ if $\sigma(\beta^T \mathbf{x}) \leq 0.5$

Logistic regression for binary classification

We define an **objective function** that we want to minimize:

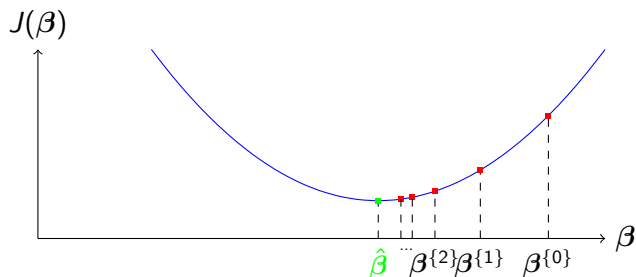
$$J(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \text{loss} \left(h_{\boldsymbol{\beta}}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

The loss function should be small when $h_{\boldsymbol{\beta}}(\mathbf{x}^{(i)})$ is close to $y^{(i)}$ and large otherwise.

We want to estimate:

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} J(\boldsymbol{\beta})$$

Gradient Descent



Algorithm: Gradient descent (\mathcal{D}, α)

Initialize $\beta^{\{0\}} \leftarrow 0, k \leftarrow 0$

WHILE no convergence **DO**

FOR j from 1 to p **DO**

$$\beta_j^{\{k+1\}} \leftarrow \beta_j^{\{k\}} - \alpha \frac{\partial J(\beta^{\{k\}})}{\partial \beta_j}$$

END FOR

$k \leftarrow k + 1$

END WHILE

Logistic regression for binary classification

Back to our problem:

$$\beta^* = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n \text{loss} \left(h_{\beta}(\mathbf{x}^{(i)}), y^{(i)} \right)$$

One could use for example the squared error:

$$\text{loss} \left(h_{\beta}(\mathbf{x}^{(i)}), y^{(i)} \right) = \frac{1}{2} \left(h_{\beta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 = \frac{1}{2} \left(\frac{1}{1 + e^{-\beta^T \mathbf{x}^{(i)}}} - y^{(i)} \right)^2$$

Unfortunately, this function is **not convex** which means the minimum that can be found with gradient descent is not necessary global!

Another idea: maximum likelihood estimation

Assuming that

$$P(Y_i = 1 | \mathbf{x}_i; \boldsymbol{\beta}) = \sigma(\boldsymbol{\beta}^T \mathbf{x}_i) = \frac{1}{1 + e^{-\boldsymbol{\beta}^T \mathbf{x}_i}} = h_{\boldsymbol{\beta}}(\mathbf{x}_i)$$

and

$$P(Y_i = 0 | \mathbf{x}_i; \boldsymbol{\beta}) = 1 - P(Y_i = 1 | \mathbf{x}_i; \boldsymbol{\beta}) = 1 - h_{\boldsymbol{\beta}}(\mathbf{x}_i),$$

the **likelihood** of the vector y_1, \dots, y_n can be written

$$L(y_1, \dots, y_n; \boldsymbol{\beta}) = \prod_{i=1}^n P(Y_i = y_i | \mathbf{x}_i; \boldsymbol{\beta}) = \prod_{i=1}^n [h_{\boldsymbol{\beta}}(\mathbf{x}_i)]^{y_i} [1 - h_{\boldsymbol{\beta}}(\mathbf{x}_i)]^{1-y_i}.$$

Thus the **negative log-likelihood** is

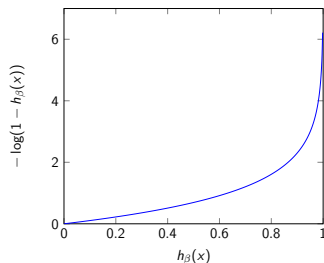
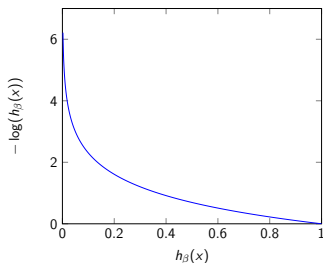
$$-\ln [L(y_1, \dots, y_n; \boldsymbol{\beta})] = \sum_{i=1}^n -y_i \ln [h_{\boldsymbol{\beta}}(\mathbf{x}_i)] - (1 - y_i) \ln [1 - h_{\boldsymbol{\beta}}(\mathbf{x}_i)],$$

We would like to maximize the likelihood of observing the training samples, i.e. minimize the negative log-likelihood.

Logistic regression for binary classification

We introduce the **logistic loss function (or binary cross entropy)** defined as:

$$\begin{aligned} \text{loss}[h_{\beta}(\mathbf{x}), y] &= \begin{cases} -\ln[h_{\beta}(\mathbf{x})] & \text{if } y = 1 \\ -\ln[1 - h_{\beta}(\mathbf{x})] & \text{if } y = 0 \end{cases} \\ &= -y \ln[h_{\beta}(\mathbf{x})] - (1 - y) \ln[1 - h_{\beta}(\mathbf{x})] \end{aligned}$$



This loss function is small when $y = 1$ and $h_{\beta}(\mathbf{x})$ is large and when $y = 0$ and $h_{\beta}(\mathbf{x})$ is small, i.e., when there are few classification errors!

Logistic regression for binary classification

The objective function $J(\beta)$ can be written:

$$J(\beta) = \frac{1}{n} \sum_{i=1}^n -y^{(i)} \ln(h_{\beta}(\mathbf{x}^{(i)})) - (1 - y^{(i)}) \ln(1 - h_{\beta}(\mathbf{x}^{(i)})).$$

This function is **CONVEX** !

→ We can use gradient descent to find its minimum.

Logistic regression for binary classification

Gradient descent requires partial derivatives to be computed:

$$\begin{aligned}\frac{\partial J}{\partial \beta_j} &= \frac{\partial}{\partial \beta_j} \frac{1}{n} \sum_{i=1}^n \left\{ -y^{(i)} \ln(h_\beta(\mathbf{x}^{(i)})) - (1 - y^{(i)}) \ln(1 - h_\beta(\mathbf{x}^{(i)})) \right\} \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ -y^{(i)} \frac{\partial}{\partial \beta_j} \ln(h_\beta(\mathbf{x}^{(i)})) - (1 - y^{(i)}) \frac{\partial}{\partial \beta_j} \ln(1 - h_\beta(\mathbf{x}^{(i)})) \right\} \\ &= \frac{1}{n} \sum_{i=1}^n \left[h_\beta(\mathbf{x}^{(i)}) - y^{(i)} \right] \mathbf{x}_j^{(i)} \quad (\text{all computations made})\end{aligned}\tag{1}$$

Comments on the previous slide

In order to compute $\frac{\partial}{\partial \beta_j} \ln(h_\beta(\mathbf{x}))$, we define $z = \beta^T \mathbf{x}$, $u = \sigma(z)$ and $v = \ln(u)$

$$\begin{aligned} \frac{\partial}{\partial \beta_j} \ln[h_\beta(\mathbf{x})] &= \frac{\partial}{\partial \beta_j} \ln(\sigma(\beta^T \mathbf{x})) \\ &= \frac{\partial v}{\partial u} \frac{\partial u}{\partial z} \frac{\partial z}{\partial \beta_j} && \text{chain-rule} \\ &= \frac{1}{\sigma(z)} \sigma(z)(1 - \sigma(z))x_j && \text{since } \sigma'(z) = \sigma(z)[1 - \sigma(z)] \\ &= (1 - h_\beta(\mathbf{x}))x_j \end{aligned}$$

Similarly, we prove $\frac{\partial}{\partial \beta_j} \ln(1 - h_\beta(\mathbf{x})) = -h_\beta(\mathbf{x})x_j$

Logistic regression for multinomial classification

How to proceed with $k > 2$ **classes** ?

We define the **softmax** function:

$$P(y = i | \mathbf{x}, \boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}_i^T \mathbf{x}}}{\sum_{j=1}^k e^{\boldsymbol{\beta}_j^T \mathbf{x}}}$$

For a sample from class j , we define a one-hot vector \mathbf{y} of dimension k such that

$$y_i = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{otherwise} \end{cases} \quad (2)$$

Logistic regression for multinomial classification

For each sample $\mathbf{x}^{(i)}$, we can predict its probability to belong to each of the k classes.

$$h_{\beta}(\mathbf{x}^{(i)}) = \begin{bmatrix} p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \beta) \\ p(y^{(i)} = 2 | \mathbf{x}^{(i)}; \beta) \\ \vdots \\ p(y^{(i)} = k | \mathbf{x}^{(i)}; \beta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\beta_j^T \mathbf{x}^{(i)}}} \begin{bmatrix} e^{\beta_1^T \mathbf{x}^{(i)}} \\ e^{\beta_2^T \mathbf{x}^{(i)}} \\ \vdots \\ e^{\beta_k^T \mathbf{x}^{(i)}} \end{bmatrix}$$

Logistic regression for multinomial classification

The objective function can be written as:

$$J(\beta) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k \mathbb{I}(y^{(i)} = j) \ln \left[\frac{e^{\beta_j^T \mathbf{x}^{(i)}}}{\sum_{l=1}^k e^{\beta_l^T \mathbf{x}^{(i)}}} \right]. \quad (3)$$

Its gradient is:

$$\nabla_{\beta_j} J(\beta) = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \left[\mathbb{I}(y^{(i)} = j) - P(y^{(i)} = j | \mathbf{x}^{(i)}; \beta) \right].$$

Outline

- 1 Logistic regression
- 2 Single-layer perceptron

Artificial Neuron

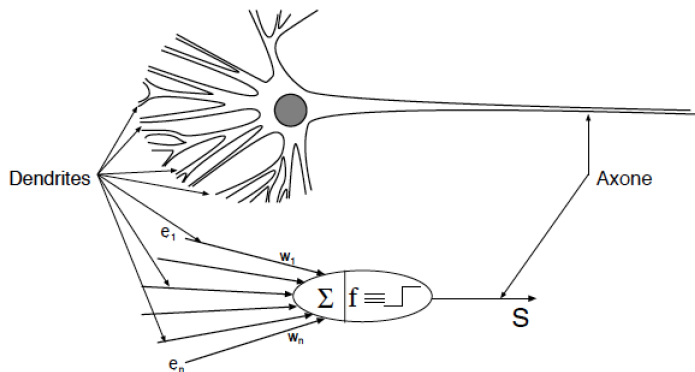
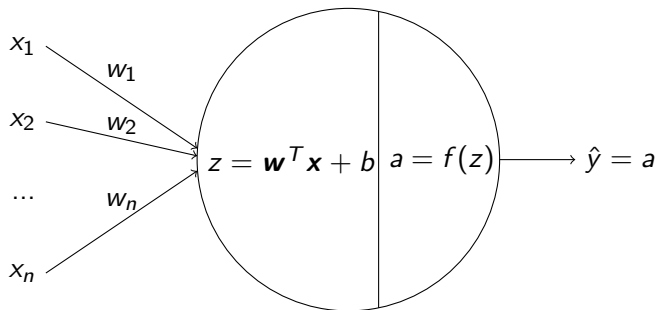


Figure: Structure of biological and artificial neurones

Artificial Neuron

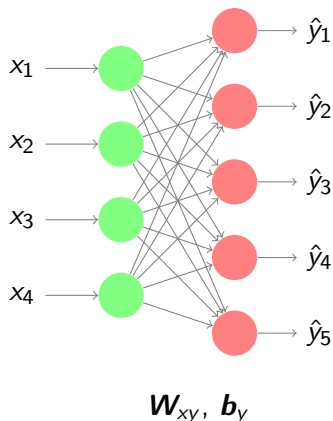


An artificial neuron is defined as a triplet of synaptic weights \mathbf{w} , bias b and activation function f .

The activation a of a neuron answering to an input \mathbf{x} is:

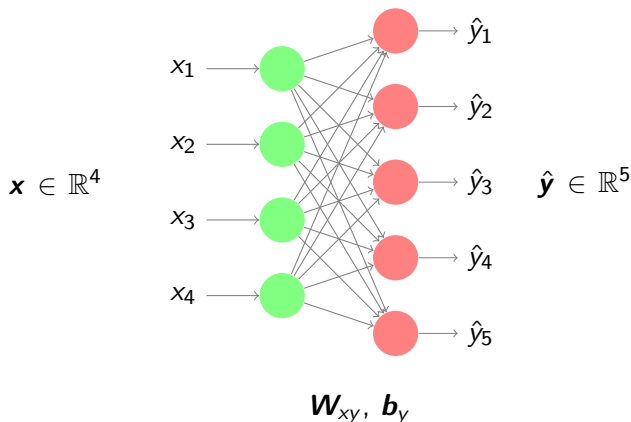
$$a = f(\mathbf{w}^T \mathbf{x} + b)$$

Neural network: Single-layer perceptron



A single-layer perceptron is in fact a little bit more general and can be composed of several neurons assembled in a single layer ($p = 4$ is the data dimension, and $K = 5$ is the number of classes).

Neural network: Single-layer perceptron

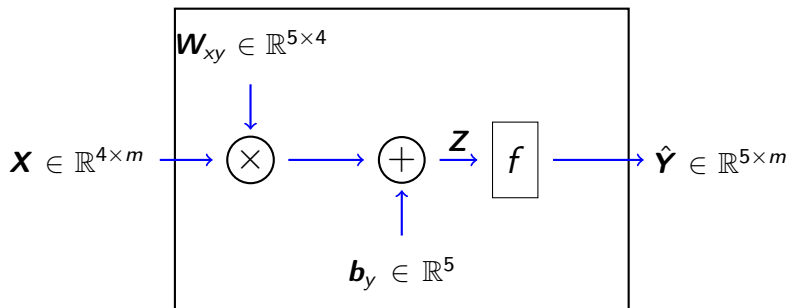


The perceptron computation is implemented as a matrix product:

$$\hat{\mathbf{y}} = f(\mathbf{W}_{xy}\mathbf{x} + \mathbf{b}_y)$$

where f is an activation function

Neural network: Single-layer perceptron

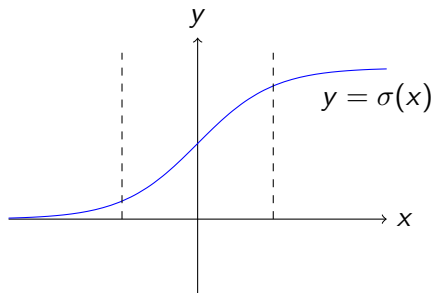


Another representation of the single-layer perceptron for prediction and classification.

Single-layer perceptron

Activation functions, denoted as f , can often be decomposed in three subdomains, as for sigmoid functions below.

- a non-activated part, below a certain threshold;
- a transition part, in the threshold neighbourhood;
- an activated part, above the threshold.



Single-layer perceptron

Algorithm 1 Single-layer perceptron training

- 1 Initialization of weights $\mathbf{W}_{xy}^{\{0\}}$ and biases $\mathbf{b}_y^{\{0\}}$.
- 2 Presentation of a set of inputs $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ and corresponding outputs $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}$
- 3 Model prediction and objective function computation:

$$\hat{\mathbf{y}}^{(i)\{k\}} = f(\mathbf{W}_{xy}^{\{k\}} \mathbf{x}^{(i)} + \mathbf{b}_y^{\{k\}}) \text{ et } J = \sum_{i=1}^n \text{loss}(\hat{\mathbf{y}}^{(i)\{k\}}, \mathbf{y}^{(i)})$$

- 4 Parameter updates

$$\mathbf{W}_{xy}^{\{k+1\}} = \mathbf{W}_{xy}^{\{k\}} - \alpha \frac{\partial J}{\partial \mathbf{W}_{xy}}, \text{ and } \mathbf{b}_y^{\{k+1\}} = \mathbf{b}_y^{\{k\}} - \alpha \frac{\partial J}{\partial \mathbf{b}_y}$$

where α is the learning rate ($0 \leq \alpha \leq 1$).

- 5 Loop to 2 until convergence (i.e., $\hat{\mathbf{y}}^{(i)\{k\}} \approx \mathbf{y}^{(i)}$).

Visualization

<https://playground.tensorflow.org/>

