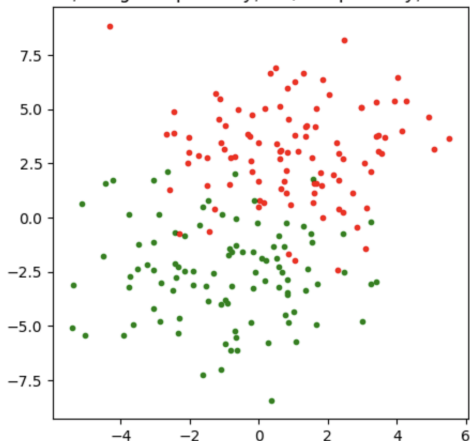


x : Feature vector representing measurements from the equipment (e.g., temperature, pressure, etc.).
 w : Weights indicating the relative importance of each feature in predicting equipment failure.
 y : Binary indicator of failure (1 for failure, 0 for normal operation).

x_i is a green point if $y_i = 0$, red point if $y_i = 1$.



Binary supervised classification

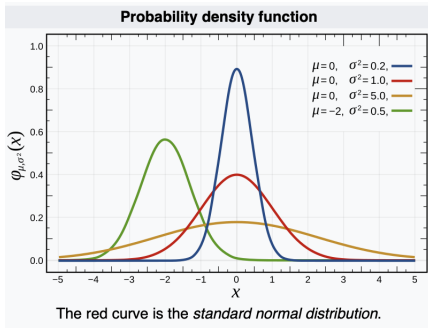
To simplify, we will take centered data in 0, so we only need to find w such that

$$\langle w, x \rangle = w^T x = 0.$$

e.g. for $i = 1, \dots, n$, x_i = random float vectors ($x \in \mathbb{R}^2$) sampled from a “normal” (Gaussian) distribution of mean 0 and variance 1: $\mathcal{N}(\mu, \sigma) = \mathcal{N}(0, 1)$

The form of the associated probability density function is

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Python → `np.random.randn(n)`

- μ : mean of the distribution
- σ : variance of the distribution

Binary supervised classification

Binary supervised classification \rightarrow logistic regression

regression = find a correlation between a binary variable and some observations thanks to an optimization problem

- Decision Trees
- K-Nearest Neighbors (k-NN)
- Probabilistic Models
- Neural Networks

Binary supervised classification

Logistic regression

$$x_1 \rightarrow f(\langle w_1, x_1 \rangle)$$

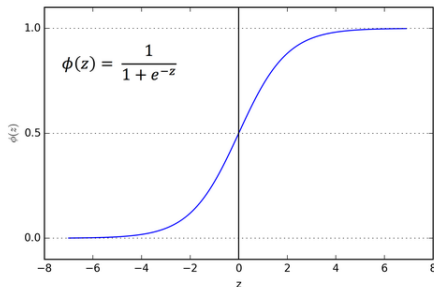
$$\dots \rightarrow \dots$$

$$x_n \rightarrow f(\langle w_n, x_n \rangle)$$

The sigmoid function σ is often used (for f) in logistic regression:

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

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$



Binary supervised classification

Logistic regression

$$x_1 \rightarrow f(\langle w_1, x_1 \rangle)$$

$$\dots \rightarrow \dots$$

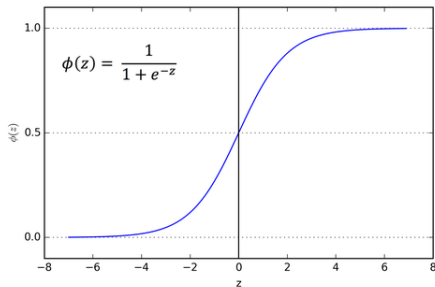
$$x_n \rightarrow f(\langle w_n, x_n \rangle)$$

The sigmoid function σ is often used (for f) in logistic regression:

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

Decision rule:

- if $\sigma(\langle w, x \rangle) > 0.5, y = 1$
- if $\sigma(\langle w, x \rangle) < 0.5, y = 0$
- if $\langle w, x \rangle \gg 0, P(y = 1|x) \simeq 1$
- if $\langle w, x \rangle \ll 0, P(y = 1|x) \simeq 0$



Binary supervised classification

Logistic regression

$$x_1 \rightarrow f(\langle w_1, x_1 \rangle)$$

$$\dots \rightarrow \dots$$

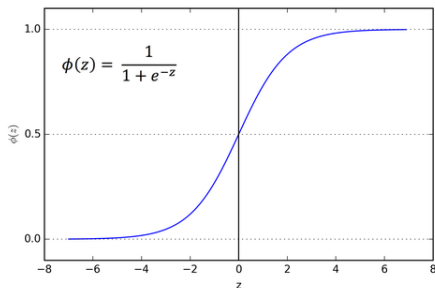
$$x_n \rightarrow f(\langle w_n, x_n \rangle)$$

The sigmoid function σ is often used (for f) in logistic regression:

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

Decision rule:

- if $\sigma(\langle w, x \rangle) > 0.5, y = 1$
- if $\sigma(\langle w, x \rangle) < 0.5, y = 0$
- modeling uncertainty



Binary supervised classification

Logistic regression

$$x_1 \rightarrow f(\langle w_1, x_1 \rangle)$$

$$\dots \rightarrow \dots$$

$$x_n \rightarrow f(\langle w_n, x_n \rangle)$$

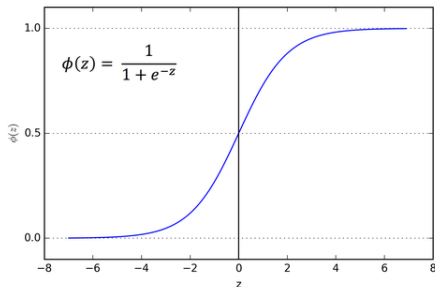
The sigmoid function σ is often used (for f) in logistic regression:

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

- How much more likely is $y = 1$ compared to $y = 0$?
Linear since if the probability to have $y = 1$ increases then the probability of having 0 decreases proportionally.

Likelihood function:

$$\log \left(\frac{P(y=1|x)}{P(y=0|x)} \right) = \langle w, x \rangle$$



Binary supervised classification

Logistic regression

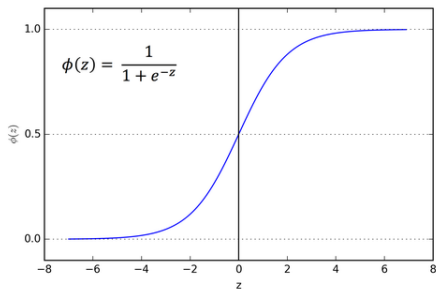
$$x_1 \rightarrow f(\langle w_1, x_1 \rangle)$$

$$\dots \rightarrow \dots$$

$$x_n \rightarrow f(\langle w_n, x_n \rangle)$$

The sigmoid function σ is often used (for f) in logistic regression:

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



- Likelihood function:

$$\log \left(\frac{P(y=1|x)}{P(y=0|x)} \right) = \langle w, x \rangle$$
- $P(y = 1|x) = \sigma(\langle w, x \rangle)$

Likelihood function

$$P(y = 1|x) = \frac{1}{1 + e^{-\langle w, x \rangle}} := \sigma(\langle w, x \rangle)$$

$$P(y = 0|x) = 1 - \frac{1}{1 + e^{-\langle w, x \rangle}} := 1 - \sigma(\langle w, x \rangle)$$

Log-loss function:

$$f(w) = -\frac{1}{n} \sum_{i=1}^n (y_i \log(\sigma(\langle w, x_i \rangle)) + (1 - y_i) \log(1 - \sigma(\langle w, x_i \rangle))) + \lambda \frac{1}{2} \|w\|^2$$

- y_i : true label (0 or 1),
- $\sigma(\langle w, x_i \rangle)$: probability predicted by the model to get $y_i = 1$.

Likelihood optimization

MINIMIZE the log-loss function:

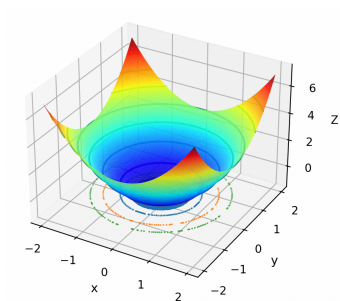
$$f(w) = -\frac{1}{n} \sum_{i=1}^n (y_i \log(\sigma(\langle w, x_i \rangle)) + (1-y_i) \log(1-\sigma(\langle w, x_i \rangle))) + \lambda \frac{1}{2} \|w\|^2$$

→ Gradient descent algorithm !!!

Exercise 1

Write a function `gradf(w)` which computes the gradient of the log-loss function.

Reminders: Gradient methods



Our goal: solving **numerically** the problem

$$\inf_{x \in \mathbb{R}^n} f(x). \quad (P)$$

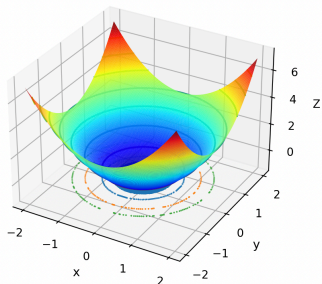
General idea: to compute a sequence $(x_k)_{k \in \mathbb{N}}$ such that

$$f(x_{k+1}) \leq f(x_k), \quad \forall k \in \mathbb{N},$$

the inequality being strict if $\nabla f(x_k) \neq 0$. \rightarrow **Iterative** method.

How to compute x_{k+1} ?

Reminders: Gradient methods



Our goal: solving **numerically** the problem

$$\inf_{x \in \mathbb{R}^n} f(x). \quad (P)$$

General idea: to compute a sequence $(x_k)_{k \in \mathbb{N}}$ such that

$$f(x_{k+1}) \leq f(x_k), \quad \forall k \in \mathbb{N},$$

the inequality being strict if $\nabla f(x_k) \neq 0$. \rightarrow **Iterative** method.

How to compute x_{k+1} ?

Gradient methods

Main idea of gradient methods.

Let $x_k \in \mathbb{R}^n$. Let d_k be a descent direction at x_k . Let $\alpha > 0$. Then

$$f(x_k + \alpha d_k) = f(x_k) + \underbrace{\alpha \langle \nabla f(x_k), d_k \rangle}_{<0} + o(\alpha).$$

Therefore, if α is small enough,

$$f(x_k + \alpha d_k) < f(x_k).$$

We can set

$$x_{k+1} = x_k + \alpha d_k.$$

Gradient methods

Definition 1

Let $x \in \mathbb{R}^n$ and let $d \in \mathbb{R}^n$. The vector d is called **descent direction** if

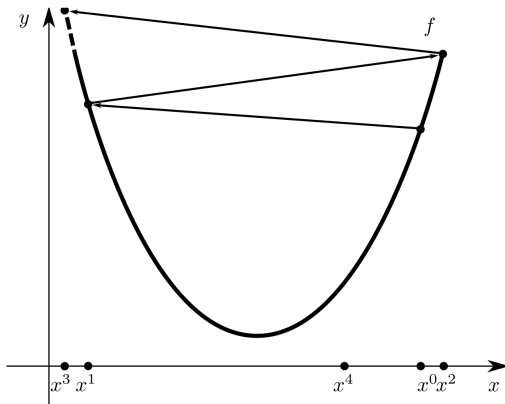
$$\langle \nabla f(x), d \rangle < 0.$$

Remark. If $\nabla f(x) \neq 0$, then $d = -\nabla f(x)$ is a descent direction. Indeed,

$$\langle \nabla f(x), d \rangle = -\langle \nabla f(x), \nabla f(x) \rangle = -\|\nabla f(x)\|^2 < 0.$$

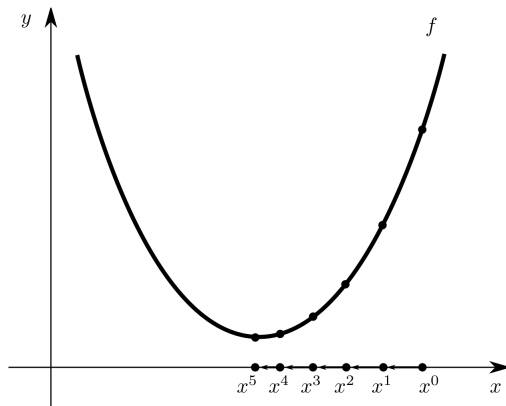
Gradient methods

On the choice of α_k .



Gradient methods

On the choice of α_k .



Gradient methods

Definition 3

Let $0 < c_1 < c_2 < 1$. We say that $\alpha > 0$ satisfies **Wolfe's rule** if

$$\phi_k(\alpha) < \phi_k(0) + c_1 \phi'_k(0)\alpha \quad \text{and} \quad \phi'_k(\alpha) \geq c_2 \phi'_k(0).$$

Armijo condition implies that the function must decrease.

The condition $\phi'_k(\alpha) \geq c_2 \phi'_k(0)$ implies that the directional derivative must increase sufficiently to approach the local minimum.

Newton's method

$$\alpha d_k = x_{k+1} - x_k.$$

$$\begin{aligned} 0 = F(x_{k+1}) &\simeq F(x_k) + \langle \nabla F(x_k), x_{k+1} - x_k \rangle. \\ &= F(x_k) + DF(x_k)(x_{k+1} - x_k). \end{aligned}$$

$$0 = F(x_k) + DF(x_k)(x_{k+1} - x_k),$$

$$-F(x_k) = DF(x_k)(x_{k+1} - x_k),$$

$$-DF(x_k)^{-1}F(x_k) = x_{k+1} - x_k$$

$$x_{k+1} = x_k - D^2f(x_k)^{-1}\nabla f(x_k).$$

