

Técnicas de optimización numérica

Gradiente

Si $f: \mathbb{R}^n \rightarrow \mathbb{R}$ es un campo escalar entonces el gradiente de f en $\mathbf{x} = (x_1, \dots, x_i, \dots, x_n)$ se define como el campo vectorial $\nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ cuyas componentes son las derivadas parciales del campo escalar, esto es:

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_i}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right)$$

Ejemplo:

1. Dada la función $f(x, y, z) = 2x + 3y^2 - \sin(z)$ su vector gradiente es el siguiente:

$$\nabla f = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right) = (2, 6y, -\cos(z)).$$

Jacobiana

Suponga $\mathbf{F}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ que es una función tal que sus derivadas parciales de primer orden existen en \mathbb{R}^n . Esta función toma un punto $\mathbf{x} \in \mathbb{R}^n$ y devuelve un vector $\mathbf{F}(\mathbf{x}) \in \mathbb{R}^m$. La matriz Jacobiana de \mathbf{F} , denotada por \mathbf{J} , está definida como una matriz de tamaño $m \times n$ cuya (i, j) -ésima entrada es :

$$J_{ij} = \frac{\partial f_i(\mathbf{x})}{\partial x_j}$$

o de forma explícita

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{F}}{\partial x_1} & \dots & \frac{\partial \mathbf{F}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla^T f_1 \\ \vdots \\ \nabla^T f_m \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

donde $\nabla^T f_i(\mathbf{x})$ es la traspuesta del gradiente de la i -ésima componente.

Ejemplo:

Supóngase la función $\mathbf{F} : \mathbb{R}^3 \rightarrow \mathbb{R}^4$, cuyas funciones componentes son:

$$y_1 = x_1$$

$$y_2 = 5x_3$$

$$y_3 = 4x_2^2 - 2x_3$$

$$y_4 = x_3 \operatorname{sen}(x_1)$$

tiene asociada como matriz jacobiana

$$\mathbf{J}_{\mathbf{F}}(x_1, x_2, x_3) = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \\ \frac{\partial y_3}{\partial x_1} & \frac{\partial y_3}{\partial x_2} & \frac{\partial y_3}{\partial x_3} \\ \frac{\partial y_4}{\partial x_1} & \frac{\partial y_4}{\partial x_2} & \frac{\partial y_4}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 5 \\ 0 & 8x_2 & -2 \\ x_3 \cos x_1 & 0 & \operatorname{sen} x_1 \end{bmatrix}.$$

Hessiana

Sea $f: \mathbb{R}^n \rightarrow \mathbb{R}$ un campo escalar cuyas derivadas parciales de segundo orden existen. La **matriz hessiana** de f , denotada por $\operatorname{Hess} f(\mathbf{x})$, $H_f(\mathbf{x})$, es una matriz cuadrada $n \times n$ definida como:

$$\operatorname{Hess} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix}$$

Teorema de Schwarz: El orden de derivación no importa, es decir derivar parcialmente primero con respecto la variable x_1 y después respecto la variable x_2 es lo mismo que derivar parcialmente primero con respecto x_2 y luego respecto x_1 , esto es:

$$\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} = \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_i}$$

La matriz Hessiana de una función f es la matriz jacobiana, J , del gradiente de la misma función:
 $H_f(\mathbf{x}) = J(\nabla f(\mathbf{x}))$

La traza de la matriz Hessiana es equivalente al **operador de Laplace**: $\text{tr}(H_f(\mathbf{x})) = \Delta f(\mathbf{x})$. Esta igualdad se puede demostrar fácilmente, ya que la definición del operador de Laplace es la **divergencia del gradiente** de una función: $\Delta f(\mathbf{x}) = \nabla \cdot (\nabla f(\mathbf{x})) = (\nabla \cdot \nabla)f(\mathbf{x}) = \nabla^2 f(\mathbf{x})$. Por lo tanto su expresión es:

$$\Delta f(\mathbf{x}) = \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} + \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} + \dots + \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2}$$

Ejemplo:

Calcula la hessiana de $f(\mathbf{x}) = x_1^3 - 2x_1x_2 - x_2^6$

1. Calculamos las derivadas de primer orden:

$$f_{x_1}(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial x_1} = \frac{\partial}{\partial x_1}(x_1^3 - 2x_1x_2 - x_2^6) = 3x_1^2 - 2x_2$$

$$f_{x_2}(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial x_2} = \frac{\partial}{\partial x_2}(x_1^3 - 2x_1x_2 - x_2^6) = -2x_1 - 6x_2^5$$

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2} \right) = (3x_1^2 - 2x_2, -2x_1 - 6x_2^5)$$

2. Calculamos las derivadas de segundo orden

$$f_{x_1x_1}(\mathbf{x}) = \frac{\partial}{\partial x_1} \left(\frac{\partial f(\mathbf{x})}{\partial x_1} \right) = \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} = \frac{\partial}{\partial x_1}(3x_1^2 - 2x_2) = 6x_1$$

$$f_{x_1x_2}(\mathbf{x}) = \frac{\partial}{\partial x_2} \left(\frac{\partial f(\mathbf{x})}{\partial x_1} \right) = \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} = \frac{\partial}{\partial x_2}(3x_1^2 - 2x_2) = -2$$

$$f_{x_2x_1}(\mathbf{x}) = \frac{\partial}{\partial x_1} \left(\frac{\partial f(\mathbf{x})}{\partial x_2} \right) = \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} = \frac{\partial}{\partial x_1}(-2x_1 - 6x_2^5) = -2$$

$$f_{x_2x_2}(\mathbf{x}) = \frac{\partial}{\partial x_2} \left(\frac{\partial f(\mathbf{x})}{\partial x_2} \right) = \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} = \frac{\partial}{\partial x_2}(-2x_1 - 6x_2^5) = -30x_2^4$$

$$H_f(\mathbf{x}) = J(\nabla^T f(\mathbf{x})) = J \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \end{pmatrix} = \begin{pmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} \end{pmatrix}$$

$$H_f(\mathbf{x}) = \begin{pmatrix} 6x_1^2 & -2 \\ -2 & -30x_2^4 \end{pmatrix}$$

Hessiana Orlada

La matriz hessiana orlada es una variante de la matriz hessiana utilizada en problemas de optimización condicionada.

Dada la función $f: \mathbb{R}^n \rightarrow \mathbb{R}$ y las condiciones o restricciones son: $g_1(\mathbf{x}) = c_1$, $g_2(\mathbf{x}) = c_2$, ..., $g_m(\mathbf{x}) = c_m$, es decir, $\mathbf{g}(\mathbf{x}) = \mathbf{c}$, donde $\mathbf{g}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, entonces la matriz hessiana orlada de la función Lagrangiana, $\Lambda(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T [\mathbf{g}(\mathbf{x}) - \mathbf{c}]$ asociada al problema de extremos condicionados es: $Hess \Lambda(\mathbf{v}) = J_v(\nabla^T \Lambda(\mathbf{v}))$, de tamaño $(n + m) \times (n + m)$, donde $\mathbf{v} = (\mathbf{x}, \boldsymbol{\lambda})$ es un vector de tamaño $1 \times (n + m)$

$$Hess \Lambda(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} J_x(\nabla_x^T \Lambda) & \nabla^T \mathbf{g}(\mathbf{x}) \\ \nabla \mathbf{g}(\mathbf{x}) & 0 \end{bmatrix}$$

$$Hess \Lambda(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} \frac{\partial^2 \Lambda}{\partial x_1^2} & \dots & \frac{\partial^2 \Lambda}{\partial x_1 \partial x_n} & \frac{\partial g_1}{\partial x_1} & \dots & \frac{\partial g_m}{\partial x_1} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \Lambda}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 \Lambda}{\partial x_n^2} & \frac{\partial g_1}{\partial x_n} & \dots & \frac{\partial g_m}{\partial x_n} \\ \frac{\partial g_1}{\partial x_1} & \dots & \frac{\partial g_1}{\partial x_n} & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_m}{\partial x_1} & \dots & \frac{\partial g_m}{\partial x_n} & 0 & \dots & 0 \end{bmatrix}$$

Si hay m condiciones, el bloque de ceros en la esquina inferior derecha es de tamaño $m \times m$ y hay m filas y m columnas bordeando por abajo y por la derecha.

Ejemplo: Calcular la hessiana orlada de la siguiente función de Lagrange, $\Lambda(\mathbf{x}, \boldsymbol{\lambda})$, con función objetivo $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ y restricciones $\mathbf{g}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$.

Función de costo:

$$f(\mathbf{x}) = f(x_1, x_2)$$

Sujeto a las restricciones:

$$g_1(x_1, x_2) = c_1$$

$$g_2(x_1, x_2) = c_2$$

donde c_1 y c_2 son constantes reales

$$\Lambda(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \lambda_1(g_1(\mathbf{x}) - c_1) + \lambda_2(g_2(\mathbf{x}) - c_2)$$

Derivadas de primer orden

$$\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1} = \frac{\partial f(\mathbf{x})}{\partial x_1} + \lambda_1 \frac{\partial g_1(\mathbf{x})}{\partial x_1} + \lambda_2 \frac{\partial g_2(\mathbf{x})}{\partial x_1}$$

$$\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_2} = \frac{\partial f(\mathbf{x})}{\partial x_2} + \lambda_1 \frac{\partial g_1(\mathbf{x})}{\partial x_2} + \lambda_2 \frac{\partial g_2(\mathbf{x})}{\partial x_2}$$

$$\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} = g_1(\mathbf{x}) - c_1$$

$$\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2} = g_2(\mathbf{x}) - c_2$$

$$\nabla \Lambda(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} \frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1} \\ \frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_2} \\ \frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} \\ \frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2} \end{bmatrix} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} + \lambda_1 \frac{\partial g_1(\mathbf{x})}{\partial x_1} + \lambda_2 \frac{\partial g_2(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} + \lambda_1 \frac{\partial g_1(\mathbf{x})}{\partial x_2} + \lambda_2 \frac{\partial g_2(\mathbf{x})}{\partial x_2} \\ g_1(\mathbf{x}) - c_1 \\ g_2(\mathbf{x}) - c_2 \end{bmatrix}$$

$$= \frac{\partial g_1(\mathbf{x})}{\partial \mathbf{x}_2}$$

$$\frac{\partial}{\partial \lambda_2} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \mathbf{x}_2} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2 \partial \mathbf{x}_2} = \frac{\partial}{\partial \lambda_2} \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_2} \right) + \frac{\partial}{\partial \lambda_2} \left(\lambda_1 \frac{\partial g_1(\mathbf{x})}{\partial \mathbf{x}_2} \right) + \frac{\partial}{\partial \lambda_2} \left(\lambda_2 \frac{\partial g_2(\mathbf{x})}{\partial \mathbf{x}_2} \right)$$

$$= \frac{\partial g_2(\mathbf{x})}{\partial \mathbf{x}_2}$$

$$\frac{\partial}{\partial x_1} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1 \partial \lambda_1} = \frac{\partial}{\partial x_1} (g_1(\mathbf{x}) - c_1) = \frac{\partial g_1(\mathbf{x})}{\partial x_1}$$

$$\frac{\partial}{\partial x_2} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_2 \partial \lambda_1} = \frac{\partial}{\partial x_2} (g_1(\mathbf{x}) - c_1) = \frac{\partial g_1(\mathbf{x})}{\partial x_2}$$

$$\frac{\partial}{\partial \lambda_1} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1^2} = \frac{\partial}{\partial \lambda_1} (g_1(\mathbf{x}) - c_1) = 0$$

$$\frac{\partial}{\partial \lambda_2} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2 \partial \lambda_1} = \frac{\partial}{\partial \lambda_2} (g_1(\mathbf{x}) - c_1) = 0$$

$$\frac{\partial}{\partial x_1} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1 \partial \lambda_2} = \frac{\partial}{\partial x_1} (g_2(\mathbf{x}) - c_2) = \frac{\partial g_2(\mathbf{x})}{\partial x_1}$$

$$\frac{\partial}{\partial x_2} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_2 \partial \lambda_2} = \frac{\partial}{\partial x_2} (g_2(\mathbf{x}) - c_2) = \frac{\partial g_2(\mathbf{x})}{\partial x_2}$$

$$\frac{\partial}{\partial \lambda_1} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1 \partial \lambda_2} = \frac{\partial}{\partial \lambda_1} (g_2(\mathbf{x}) - c_2) = 0$$

$$\frac{\partial}{\partial \lambda_2} \left(\frac{\partial \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2} \right) = \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2^2} = \frac{\partial}{\partial \lambda_2} (g_2(\mathbf{x}) - c_2) = 0$$

$$Hess \Lambda(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \mathbf{x}_1^2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_2 \partial \mathbf{x}_1} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1 \partial \mathbf{x}_1} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2 \partial \mathbf{x}_1} \\ \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1 \partial x_2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \mathbf{x}_2^2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1 \partial x_2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2 \partial x_2} \\ \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \mathbf{x}_1 \partial \lambda_1} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \mathbf{x}_2 \partial \lambda_1} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1^2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2 \partial \lambda_1} \\ \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_1 \partial \lambda_2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_2 \partial \lambda_2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_1 \partial \lambda_2} & \frac{\partial^2 \Lambda(\mathbf{x}, \boldsymbol{\lambda})}{\partial \lambda_2^2} \end{bmatrix}$$

$$Hess \Lambda(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x}_1^2} + \lambda_1 \frac{\partial^2 g_1(\mathbf{x})}{\partial \mathbf{x}_1^2} + \lambda_2 \frac{\partial^2 g_2(\mathbf{x})}{\partial \mathbf{x}_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial \mathbf{x}_1} + \lambda_1 \frac{\partial^2 g_1(\mathbf{x})}{\partial x_2 \partial \mathbf{x}_1} + \lambda_2 \frac{\partial^2 g_2(\mathbf{x})}{\partial x_2 \partial \mathbf{x}_1} & \frac{\partial g_1(\mathbf{x})}{\partial \mathbf{x}_1} & \frac{\partial g_2(\mathbf{x})}{\partial \mathbf{x}_1} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} + \lambda_1 \frac{\partial^2 g_1(\mathbf{x})}{\partial x_1 \partial x_2} + \lambda_2 \frac{\partial^2 g_2(\mathbf{x})}{\partial x_1 \partial x_2} & \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x}_2^2} + \lambda_1 \frac{\partial^2 g_1(\mathbf{x})}{\partial \mathbf{x}_2^2} + \lambda_2 \frac{\partial^2 g_2(\mathbf{x})}{\partial \mathbf{x}_2^2} & \frac{\partial g_1(\mathbf{x})}{\partial x_2} & \frac{\partial g_2(\mathbf{x})}{\partial x_2} \\ & \frac{\partial g_1(\mathbf{x})}{\partial x_1} & 0 & 0 \\ & \frac{\partial g_2(\mathbf{x})}{\partial x_1} & \frac{\partial g_2(\mathbf{x})}{\partial x_2} & 0 \end{bmatrix}$$

Gradient descent optimization algorithms

1. Adadelta

`trainAdadelta(params, lr=1.0, rho=0.9, eps=1e-06, weight_decay=0)`

Implements Adadelta algorithm.

input : γ (lr), θ_0 (params), $f(\theta)$ (objective), ρ (decay), λ (weight decay)
initialize : $v_0 \leftarrow 0$ (square avg), $u_0 \leftarrow 0$ (accumulate variables)

for $t = 1$ **to** ... **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$v_t \leftarrow v_{t-1} \rho + g_t^2 (1 - \rho)$

$\Delta x_t \leftarrow \frac{\sqrt{u_{t-1} + \epsilon}}{\sqrt{v_t + \epsilon}} g_t$

$u_t \leftarrow u_{t-1} \rho + \Delta x_t^2 (1 - \rho)$

$\theta_t \leftarrow \theta_{t-1} - \gamma \Delta x_t$

return θ_t

For further details regarding the algorithm we refer to [ADADELTA: An Adaptive Learning Rate Method](#).

2. Adagrad

`trainAdagrad(params, lr=0.01, lr_decay=0, weight_decay=0, initial_accumulator_value=0, eps=1e-10)`

Implements Adagrad algorithm.

input : γ (lr), θ_0 (params), $f(\theta)$ (objective), λ (weight decay),
 τ (initial accumulator value), η (lr decay)

initialize : $state_sum_0 \leftarrow 0$

for $t = 1$ **to** \dots **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

$\tilde{\gamma} \leftarrow \gamma / (1 + (t - 1)\eta)$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$state_sum_t \leftarrow state_sum_{t-1} + g_t^2$

$\theta_t \leftarrow \theta_{t-1} - \tilde{\gamma} \frac{g_t}{\sqrt{state_sum_t} + \epsilon}$

return θ_t

For further details regarding the algorithm we refer to [Adaptive Subgradient Methods for Online Learning and Stochastic Optimization](#).

3. Adam

`trainAdam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False, maximize=False)`

Implements Adam algorithm.

input : γ (lr), β_1, β_2 (betas), θ_0 (params), $f(\theta)$ (objective)

λ (weight decay), *amsgrad*, *maximize*

initialize : $m_0 \leftarrow 0$ (first moment), $v_0 \leftarrow 0$ (second moment), $\widehat{v}_0^{max} \leftarrow 0$

for $t = 1$ **to** ... **do**

if *maximize* :

$g_t \leftarrow -\nabla_{\theta} f_t(\theta_{t-1})$

else

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$

$\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

if *amsgrad*

$\widehat{v}_t^{max} \leftarrow \max(\widehat{v}_t^{max}, \widehat{v}_t)$

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t^{max}} + \epsilon)$

else

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

return θ_t

For further details regarding the algorithm we refer to [Adam: A Method for Stochastic Optimization](#).

4. AdamW

`trainAdamW(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08, weight_decay=0.01, amsgrad=False, maximize=False)`

Implements AdamW algorithm.

input : $\gamma(\text{lr})$, $\beta_1, \beta_2(\text{betas})$, $\theta_0(\text{params})$, $f(\theta)(\text{objective})$, ϵ (epsilon)
 $\lambda(\text{weight decay})$, *amsgrad*, *maximize*

initialize : $m_0 \leftarrow 0$ (first moment), $v_0 \leftarrow 0$ (second moment), $\widehat{v}_0^{\text{max}} \leftarrow 0$

for $t = 1$ **to** ... **do**

if *maximize* :

$g_t \leftarrow -\nabla_{\theta} f_t(\theta_{t-1})$

else

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

$\theta_t \leftarrow \theta_{t-1} - \gamma \lambda \theta_{t-1}$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$

$\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

if *amsgrad*

$\widehat{v}_t^{\text{max}} \leftarrow \max(\widehat{v}_t^{\text{max}}, \widehat{v}_t)$

$\theta_t \leftarrow \theta_t - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t^{\text{max}}} + \epsilon)$

else

$\theta_t \leftarrow \theta_t - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

return θ_t

For further details regarding the algorithm we refer to [Decoupled Weight Decay Regularization](#).

5. Adamax

`trainAdamax(params, lr=0.002, betas=(0.9, 0.999), eps=1e-08, weight_decay=0)`

Implements Adamax algorithm (a variant of Adam based on infinity norm).

input : γ (lr), β_1, β_2 (betas), θ_0 (params), $f(\theta)$ (objective), λ (weight decay)
 ϵ (epsilon)

initialize : $m_0 \leftarrow 0$ (first moment), $u_0 \leftarrow 0$ (infinity norm)

for $t = 1$ **to** ... **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$u_t \leftarrow \max(\beta_2 u_{t-1}, |g_t| + \epsilon)$

$\theta_t \leftarrow \theta_{t-1} - \frac{\gamma m_t}{(1 - \beta_1^t) u_t}$

return θ_t

For further details regarding the algorithm we refer to [Adam: A Method for Stochastic Optimization](#).

6. Lion

The Lion optimizer is a stochastic-gradient-descent method that uses the sign operator to control the magnitude of the update, unlike other adaptive optimizers such as Adam that rely on second order moments. This makes Lion more memory-efficient as it only keeps track of the momentum. According to the authors (see reference), its performance gain over Adam grows with the batch size. Because the update of Lion is produced through the sign operation, resulting in a larger norm, a suitable learning rate for Lion is typically 3-10x smaller than that for AdamW. The weight decay for Lion should be in turn 3-10x larger than that for AdamW to maintain a similar strength ($\text{lr} * \text{wd}$).

Algorithm 2 Lion Optimizer (ours)

given $\beta_1, \beta_2, \lambda, \eta, f$
initialize $\theta_0, m_0 \leftarrow 0$
while θ_t not converged **do**
 $g_t \leftarrow \nabla_{\theta} f(\theta_{t-1})$
 update model parameters
 $c_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$
 $\theta_t \leftarrow \theta_{t-1} - \eta_t (\text{sign}(c_t) + \lambda \theta_{t-1})$
 update EMA of g_t
 $m_t \leftarrow \beta_2 m_{t-1} + (1 - \beta_2) g_t$
end while
return θ_t

7. NAdam

trainNAdam(params, lr=0.002, betas=(0.9, 0.999), eps=1e-08, weight_decay=0, momentum_decay=0.004)

Implements NAdam algorithm.

input : γ_t (lr), β_1, β_2 (betas), θ_0 (params), $f(\theta)$ (objective)
 λ (weight decay), ψ (momentum decay)

initialize : $m_0 \leftarrow 0$ (first moment), $v_0 \leftarrow 0$ (second moment)

for $t = 1$ **to** ... **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$\mu_t \leftarrow \beta_1 \left(1 - \frac{1}{2} 0.96^{t\psi}\right)$

$\mu_{t+1} \leftarrow \beta_1 \left(1 - \frac{1}{2} 0.96^{(t+1)\psi}\right)$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$\widehat{m}_t \leftarrow \mu_{t+1} m_t / \left(1 - \prod_{i=1}^{t+1} \mu_i\right)$

$+ (1 - \mu_t) g_t / \left(1 - \prod_{i=1}^t \mu_i\right)$

$\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

return θ_t

For further details regarding the algorithm we refer to [Incorporating Nesterov Momentum into Adam](#).

8. RAdam

`trainRAdam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08, weight_decay=0)`

Implements RAdam algorithm.

input : γ (lr), β_1, β_2 (betas), θ_0 (params), $f(\theta)$ (objective), λ (weightdecay), ϵ (epsilon)

initialize : $m_0 \leftarrow 0$ (first moment), $v_0 \leftarrow 0$ (second moment),
 $\rho_\infty \leftarrow 2/(1 - \beta_2) - 1$

for $t = 1$ **to** ... **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$

$\rho_t \leftarrow \rho_\infty - 2t\beta_2^t / (1 - \beta_2^t)$

if $\rho_t > 5$

$l_t \leftarrow \frac{\sqrt{(1 - \beta_2^t)}}{\sqrt{v_t} + \epsilon}$

$r_t \leftarrow \sqrt{\frac{(\rho_t - 4)(\rho_t - 2)\rho_\infty}{(\rho_\infty - 4)(\rho_\infty - 2)\rho_t}}$

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t r_t l_t$

else

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t$

return θ_t

For further details regarding the algorithm we refer to [On the variance of the adaptive learning rate and beyond](#).

9. RMSprop

`trainRMSprop(params, lr=0.01, alpha=0.99, eps=1e-08, weight_decay=0, momentum=0, centered=False)`

Implements RMSprop algorithm.

input : α (alpha), γ (lr), θ_0 (params), $f(\theta)$ (objective)
 λ (weight decay), μ (momentum), *centered*

initialize : $v_0 \leftarrow 0$ (square average), $\mathbf{b}_0 \leftarrow 0$ (buffer), $g_0^{ave} \leftarrow 0$

for $t = 1$ **to** ... **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

if $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$v_t \leftarrow \alpha v_{t-1} + (1 - \alpha) g_t^2$

$\tilde{v}_t \leftarrow v_t$

if *centered*

$g_t^{ave} \leftarrow g_{t-1}^{ave} \alpha + (1 - \alpha) g_t$

$\tilde{v}_t \leftarrow \tilde{v}_t - (g_t^{ave})^2$

if $\mu > 0$

$\mathbf{b}_t \leftarrow \mu \mathbf{b}_{t-1} + g_t / (\sqrt{\tilde{v}_t} + \epsilon)$

$\theta_t \leftarrow \theta_{t-1} - \gamma \mathbf{b}_t$

else

$\theta_t \leftarrow \theta_{t-1} - \gamma g_t / (\sqrt{\tilde{v}_t} + \epsilon)$

return θ_t

For further details regarding the algorithm we refer to [lecture notes](#) by G. Hinton. and centered version [Generating Sequences With Recurrent Neural Networks](#).

10. SMORMS3

SMORMS3 (Squared Mean Over Root Mean Squared Cubed) is an optimization algorithm used in machine learning and deep learning to optimize the training of neural networks. It is a variant of the RMSProp algorithm and was introduced in 2017 by Daniel Fortunato, et al.

SMORMS3 modifies the way the moving average of the squared gradients is calculated in RMSProp. Instead of taking the simple average of the squared gradients, SMORMS3 takes the cube root of the moving average of the cube of the squared gradients. This modification helps to normalize the scale of the moving average, which can prevent the learning rate from decreasing too quickly.

Overall, SMORMS3 is a powerful optimization algorithm that can help accelerate the training of deep neural networks and improve their performance, particularly in situations where the gradients have a high variance.

Algorithm 4 Algorithm of SMORMS3

Require: lr : Learning rate

Require: θ_0 : Initial parameter

Require: $f(\theta)$: Loss function with parameter θ

Require: $g(\theta)$: Gradient of loss function with parameter θ

Require: ε : Small constant

$t \leftarrow 0$

$m_0 \leftarrow 0$

$v_0 \leftarrow 0$

$s_0 \leftarrow 1$

while Exit condition is not satisfied **do**

$t \leftarrow t + 1$

$g(\theta_t) \leftarrow \nabla_{\theta_t} f(\theta_t)$

$s_t \leftarrow 1 + (1 - x_{t-1}) \cdot s_{t-1}$

$\rho_t \leftarrow \frac{1}{s_t + 1}$

$m_t \leftarrow (1 - \rho_t) \cdot m_{t-1} + \rho_t \cdot g(\theta_t)$

$v_t \leftarrow (1 - \rho_t) \cdot v_{t-1} + \rho_t \cdot g^2(\theta_t)$

$x_t \leftarrow \frac{m_t^2}{v_t + \varepsilon}$

$\theta_t \leftarrow \theta_{t-1} - \frac{\min\{lr, x_t\}}{\sqrt{v_t + \varepsilon}} \cdot g(\theta_t)$

end while

11. Rprop

`trainRprop(params, lr=0.01, etas=(0.5, 1.2), step_sizes=(1e-06, 50))`

Implements the resilient backpropagation algorithm.

input : $\theta_0 \in \mathbf{R}^d$ (params), $f(\theta)$ (objective),
 $\eta_{+/-}$ (etaplus, etaminus), $\Gamma_{max/min}$ (step sizes)
initialize : $g_{prev}^0 \leftarrow 0$, $\eta_0 \leftarrow lr$ (learning rate)

for $t = 1$ **to** \dots **do**
 $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$
 for $i = 0, 1, \dots, d - 1$ **do**
 if $g_{prev}^i g_t^i > 0$
 $\eta_t^i \leftarrow \min(\eta_{t-1}^i \eta_+, \Gamma_{max})$
 else if $g_{prev}^i g_t^i < 0$
 $\eta_t^i \leftarrow \max(\eta_{t-1}^i \eta_-, \Gamma_{min})$
 $g_t^i \leftarrow 0$
 else
 $\eta_t^i \leftarrow \eta_{t-1}^i$
 $\theta_t \leftarrow \theta_{t-1} - \eta_t \text{sign}(g_t)$
 $g_{prev} \leftarrow g_t$

return θ_t

For further details regarding the algorithm we refer to the paper [A Direct Adaptive Method for Faster Backpropagation Learning: The RPROP Algorithm](#).

12. SGD

`trainSGD(params, lr=<required parameter>, momentum=0, dampening=0, weight_decay=0, nesterov=False, maximize=False)`

Implements stochastic gradient descent (optionally with momentum).

input : γ (lr), θ_0 (params), $f(\theta)$ (objective), λ (weight decay),
 μ (momentum), τ (dampening), *nesterov*, *maximize*

```
for  $t = 1$  to ... do
     $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ 
    if  $\lambda \neq 0$ 
         $g_t \leftarrow g_t + \lambda \theta_{t-1}$ 
    if  $\mu \neq 0$ 
        if  $t > 1$ 
             $\mathbf{b}_t \leftarrow \mu \mathbf{b}_{t-1} + (1 - \tau) g_t$ 
        else
             $\mathbf{b}_t \leftarrow g_t$ 
        if nesterov
             $g_t \leftarrow g_t + \mu \mathbf{b}_t$ 
        else
             $g_t \leftarrow \mathbf{b}_t$ 
    if maximize
         $\theta_t \leftarrow \theta_{t-1} + \gamma g_t$ 
    else
         $\theta_t \leftarrow \theta_{t-1} - \gamma g_t$ 
return  $\theta_t$ 
```

Nesterov momentum is based on the formula from [On the importance of initialization and momentum in deep learning](#).

13. SGD with U-Clip

SGD with U-Clip, a novel method for achieving on-average unbiased stochastic gradient clipping.

Algorithm 1 SGD with U-Clip

Input: $f(\theta, \cdot)$ (loss function), S (dataset), $\theta_1 \in \mathbb{R}^p$ (initial parameters), $\eta \in \mathbb{R}_{>0}$ (learning rate), $\gamma \in \mathbb{R}_{>0}$ (clipping region), $T \in \mathbb{N}$ (number of timesteps).

Return: θ_{T+1} (optimized parameters).

$\Delta_1 \leftarrow 0$

for $t = 1, \dots, T$ **do**

 Sample $Z_t \sim \text{Unif } S$

$g_t \leftarrow \nabla_{\theta} f(\theta_t, Z_t)$

$u_t \leftarrow \text{clip}(g_t + \Delta_t, \gamma)$

$\Delta_{t+1} \leftarrow g_t + \Delta_t - u_t$

$\theta_{t+1} \leftarrow \theta_t - \eta u_t$

end for

Proposition 2.4 (Carry bound: component-wise clipping). *In this result, all variables are real scalars. Consider the following update rules with $\Delta_1 = 0$*

$$u_t = \text{clip}(g_t + \Delta_t, \gamma) \quad \text{and} \quad \Delta_{t+1} = g_t + \Delta_t - u_t$$

where $\text{clip}(x, \gamma) = \min\{\max\{x, -\gamma\}, \gamma\}$ projects x onto the interval $[-\gamma, \gamma]$.

14. Gradient Descent

trainGD(params, lr=<required parameter>)

implements gradient descent

inputs: γ (lr = 0.001: learning rate), θ_0 (params), $f(\theta)$ (objective)

for $t = 1$ **to** ... **do**

$g_t \leftarrow \nabla f(\theta_{t-1})$

$\theta_t \leftarrow \theta_{t-1} - \gamma g_t$

return θ_t

15. Gradient Descent with momentum

`trainGDM(params, lr=<required parameter>, momentum=<required parameter>)`

implements gradient descent with momentum

```
inputs:  $\gamma$  ( $lr = 0.001$ : learning rate),  $\theta_0$ (params),  $f(\theta)$  (objective),  $\mu$  (momentum = 0.90)
 $g_0 \leftarrow \nabla f(\theta_0)$ 
 $\Delta\theta_0 \leftarrow \gamma g_0$ 
for  $t = 1$  to ... do
     $\Delta\theta_t \leftarrow \mu \Delta\theta_{t-1} - (1 - \mu) \gamma g_{t-1}$ 
     $\theta_t \leftarrow \theta_{t-1} + \Delta\theta_t$ 
     $g_t \leftarrow \nabla f(\theta_t)$ 
return  $\theta_t$ 
```

16. Variable Learning Rate Gradient Descent

`trainGDX(params, lr=<required parameter>, momentum=<required parameter>, lr_dec
=<required parameter>, lr_inc=<required parameter>, max_per_inc=<required parameter>)`

inputs: γ (lr : learning rate), θ_0 (params), $f(\theta)$ (objective), μ (momentum = 0.90), γ_{dec} (lr_{dec}
= 0.70: Ratio to decrease learning rate), γ_{inc} (lr_{inc} = 1.05: Ratio to increase learning rate),
 max_per_inc (maximum performance increase = 1.04)

```
 $g_0 \leftarrow \nabla f(\theta_0)$ 
 $\Delta\theta_0 \leftarrow \gamma g_0$ 
for  $t = 1$  to ... do
     $\Delta\theta_t \leftarrow \mu \Delta\theta_{t-1} - (1 - \mu) \gamma g_{t-1}$  gradiente actual
     $\theta_t \leftarrow \theta_{t-1} + \Delta\theta_t$ 
    if  $\frac{f(\theta_t)}{f(\theta_{t-1})} > max\_perf\_inc$ 
         $\gamma \leftarrow \gamma \gamma_{dec}$ 
         $g_t \leftarrow g_{t-1}$ 
         $\Delta\theta_t \leftarrow \gamma g_t$ 
    else
        if  $f(\theta_t) < f(\theta_{t-1})$ 
             $\gamma = \gamma \gamma_{inc}$ 
        end
         $g_t \leftarrow \nabla f(\theta_t)$ 
    end
return  $\theta_t$ 
```

17. Levenberg-Marquardt

`trainLM(params, lr=<required parameter>, lr_dec=<required parameter>, lr_inc=<required parameter>, lr_max=<required parameter>)`

inputs: γ ($lr = 0.001$: Initial damping factor), θ_0 (params), $f(\theta)$ (objective), γ_{dec} ($lr_dec = 0.1$: decrease damping factor), γ_{inc} ($lr_inc = 10.0$: increase damping factor), γ_{max} ($lr_max = 10^{10}$: máximo damping factor)

```
for  $t = 1$  to ... do
    while ( $\gamma \leq \gamma_{max}$ )
        //  $(J'_f(\theta_{t-1}) J_f(\theta_{t-1}) + \gamma I) \Delta\theta_t = -J'_F(\theta_{t-1}) f(\theta_{t-1})$ 
         $\Delta\theta_t = -(J'_f(\theta_{t-1}) J_f(\theta_{t-1}) + \gamma I)^{-1} J'_f(\theta_{t-1}) f(\theta_{t-1})$ 
         $\theta_t = \theta_{t-1} + \Delta\theta_t$ 
        if  $f(\theta_t) < f(\theta_{t-1})$ 
             $\gamma = \gamma \gamma_{dec}$ 
            break // terminates the execution of a for while loop.
        end
         $\gamma = \gamma \gamma_{inc}$ 
    end
end
```

Validación cruzada

Los procesos iterativos por los cuales se efectúa el ajuste de un modelo, en base a datos experimentales, generalmente necesitan un criterio que les indique en qué momento detenerse. El criterio que se elija para detener el proceso de ajuste o entrenamiento debe tomar en cuenta algunas consideraciones, por ejemplo, la calidad de la salida del modelo, la cantidad en que se filtra el ruido experimental, la presencia de un sobreajuste de los datos, etc.

La validación cruzada es una técnica que indica el momento oportuno para detener el entrenamiento de una RNA, de modo que indica la capacidad de generalización de la RNA durante su proceso de entrenamiento.

La validación cruzada consiste en dividir los datos experimentales de entrada ($X = [x_{i,p}]$ para $i=1,\dots,n$ y $p=1,\dots,q$) y datos de salida deseados $T = [t_{j,p}]$ para $j=1,\dots,m$ y $p=1,\dots,q$; en dos conjuntos. Uno de ellos es el *conjunto de entrenamiento* y el otro es llamado *conjunto de validación*. El conjunto de entrenamiento es utilizado para ajustar los parámetros de la red durante el proceso de entrenamiento, a su vez, el conjunto de validación se utiliza para medir el rendimiento de la red mientras es entrenada de forma iterativa.

Cuando se construye una gráfica de la función error a partir del conjunto de validación respecto del número de iteraciones, Figura 2.7, se observa un comportamiento irregular, es decir, el error presenta muchos picos. La evolución del error de validación, aunque irregular, empieza a disminuir hasta alcanzar un mínimo para después aumentar. Por lo tanto, el error de validación funciona como criterio para determinar el momento adecuado en que debe detenerse el proceso de entrenamiento.

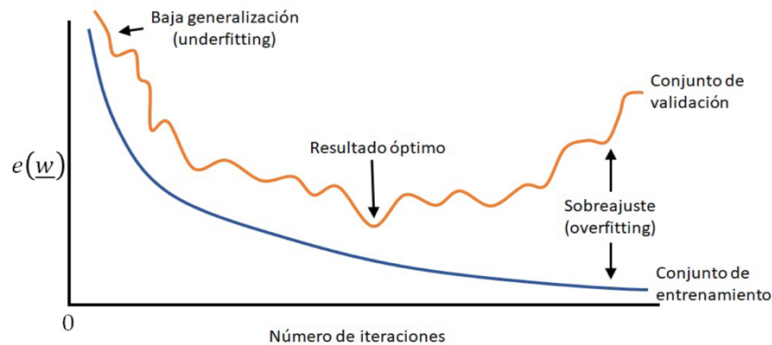


Figura 2.7: Comportamiento de los errores del conjunto de entrenamiento y validación durante el ajuste de parámetros de la RNA.

Durante el proceso de entrenamiento es posible identificar tres etapas: 1) la etapa de sobreajuste (overfitting), Figura 2.8a, es decir, que la red no tiene la capacidad de distinguir entre la señal y el ruido que contienen los datos de entrada, el entrenamiento debe detenerse antes de llegar a este punto; 2) la etapa de baja generalización (underfitting), Figura 2.8c, se presenta cuando los parámetros no se ajustan adecuadamente y la red no puede predecir el comportamiento real del sistema, generalmente ocurre cuando el proceso de entrenamiento es detenido antes de tiempo; 3) resultado óptimo, Figura 2.8b, es donde la red presenta la mejor capacidad de generalización.

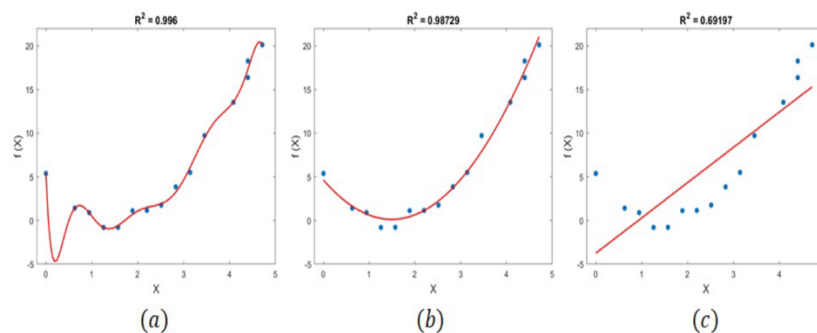


Figura 2.8: Distintos tipos de ajuste en el entrenamiento de una RNA: a) sobreajuste, b) ajuste óptimo, c) bajo-ajuste.

En el entrenamiento usando los distintos algoritmos de optimización, se detiene cuando ocurre cualquiera de estas condiciones:

- 1) Se alcanza el número máximo de épocas (repeticiones).
- 2) Se ha excedido la cantidad máxima de tiempo.
- 3) El rendimiento se ha minimizado hasta la meta de la función objetivo (***goal*** = 10^{-8}).
- 4) El gradiente de rendimiento cae por debajo de ***mingrad*** = 10^{-10}
- 5) El rendimiento de la validación ha aumentado más de (max_fail = 6) veces desde la última vez que disminuyó (al usar la validación).

La idea de construir un modelo de Machine Learning no es sólo que lo haga bien con los conjuntos de entrenamiento y validación. La idea es que además lo haga bien con datos que nunca antes haya visto, pues de esta manera veremos qué tan robusto es y cómo se comportará con nuevos datos. Esto se conoce como generalización.

Para ver la capacidad de generalización del modelo seleccionado debemos usar un tercer conjunto de datos, que se conoce como el *conjunto de prueba*. Este conjunto de prueba es un conjunto de datos diferente al conjunto de entrenamiento y validación pero proveniente de la misma distribución. Con este tercer conjunto podemos finalmente poner a prueba el modelo para ver qué tan bien lo hace con datos nunca antes vistos, es decir su capacidad de generalización.

Idealmente para entrenar diferentes modelos de Machine Learning y poder escoger el mejor de todos deberíamos tener tres conjuntos: **entrenamiento, validación y prueba**

- El ***conjunto de entrenamiento*** se usa para entrenar como tal cada modelo (es decir, para aprender los parámetros)
- El ***conjunto de validación*** lo usamos para poder escoger el valor más adecuado de los hiperparámetros de cada modelo y para elegir al final el mejor modelo de todos los que hayamos entrenado
- El ***conjunto de prueba*** lo usamos para medir la capacidad de generalización del modelo seleccionado.

Es importante además tener en cuenta que los tres conjuntos deben ser diferentes pero a la vez deben provenir de la misma distribución.