

MAC0460 - Introdução ao aprendizado de máquina

Back-propagation 1

Felipe Salvatore

<https://felipessalvatore.github.io/>

Nina S. T. Hirata

<https://www.ime.usp.br/~nina/>

April 23, 2018

IME-USP: Institute of Mathematics and Statistics, University of São Paulo

Revisão: regressão logística

O problema de classificação

Em vários casos a função desconhecida $f : \mathbb{R}^d \rightarrow \mathbb{R}$ que queremos aproximar é uma **distribuição de probabilidade**.

Temos um vetor \mathbf{x} e queremos saber a qual das classes k_1, \dots, k_n ele pertence. Um modo de formular esse problema como um problema de aprendizado supervisionado é coletar um conjunto de dados $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ onde $y_i \in \{k_1, \dots, k_n\}$ e tentar estimar $p(y|\mathbf{x})$ por meio de uma família de modelos $p(y|\mathbf{x}; \theta)$.

Classificação com duas classes

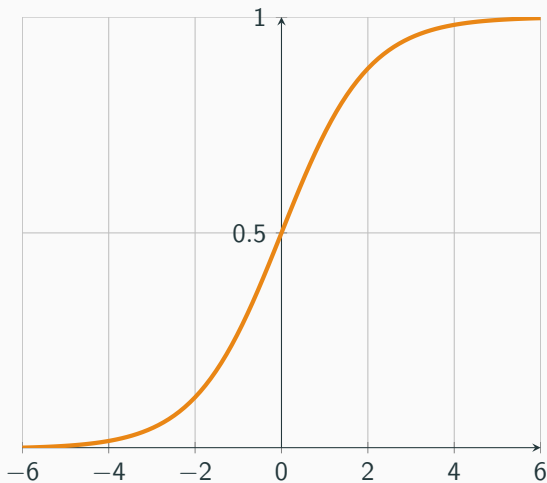
Quando y é uma variável binária definimos o modelo $p(y|\mathbf{x}; \boldsymbol{\theta})$ do seguinte modo:

$$\begin{aligned}\hat{y} &= p(y = 1|\mathbf{x}; \boldsymbol{\theta}) \\ &= h(\mathbf{x}; \boldsymbol{\theta}) \\ &= \sigma(z)\end{aligned}$$

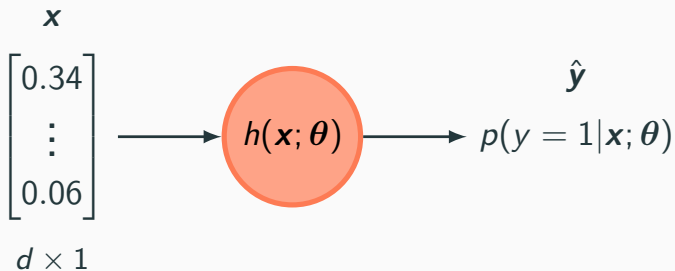
em que

$$z = \mathbf{w}^\top \mathbf{x} + b$$

Revisão: função sigmoide



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



Classificação para várias classes

E quando y é uma variável com n valores definimos $p(y|\mathbf{x}; \theta)$ do seguinte modo:

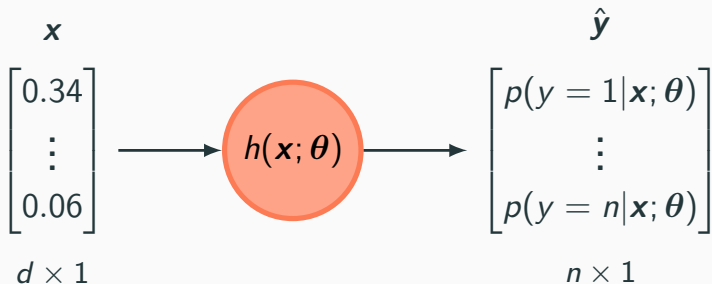
$$\begin{aligned}\hat{y} &= p(y|\mathbf{x}; \theta) \\ &= h(\mathbf{x}; \theta) \\ &= \textit{softmax}(\mathbf{z})\end{aligned}$$

em que

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$$\begin{bmatrix} 3.82 \\ 5.35 \\ 1.44 \\ -1.26 \\ 2.71 \\ 1.98 \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} 0.16115195 \\ 0.74422819 \\ 0.01491471 \\ 0.00100235 \\ 0.05310907 \\ 0.02559374 \end{bmatrix}$$

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$



Os parâmetros θ vão ser adaptados de modo que $p(y|\mathbf{x}; \theta)$ seja a distribuição mais adequada para os dados

$$(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})$$

A função que queremos maximizar é

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}, y \sim p_{data}} \log p(y|\mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{1}{N} \sum_{i=1}^N \log p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})\end{aligned}$$

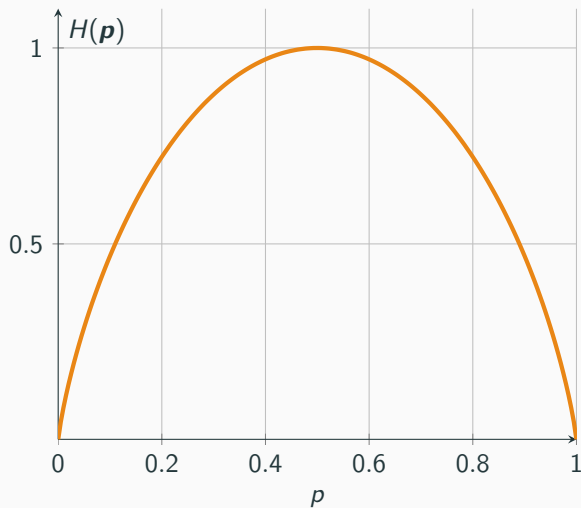
$$\mathbf{p} \quad \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}$$

$$\mathbf{q} \quad \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

$$H(\mathbf{p}) = 0.72 \quad H(\mathbf{q}) = 1$$

$$H(\mathbf{p}) = \sum_i \mathbf{p}_i \log \frac{1}{\mathbf{p}_i}$$

Revisão: entropia



$$\begin{matrix} p \\ \begin{bmatrix} p \\ 1 - p \end{bmatrix} \end{matrix}$$

Revisão: divergência Kullback-Leibler

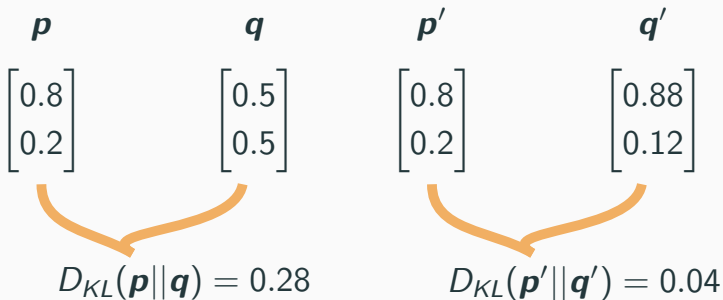


Diagram illustrating the Kullback-Leibler (KL) divergence calculation for two pairs of probability distributions p and q , and p' and q' .

For the first pair:

$$p = \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}, \quad q = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$
$$D_{KL}(p||q) = 0.28$$

For the second pair:

$$p' = \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}, \quad q' = \begin{bmatrix} 0.88 \\ 0.12 \end{bmatrix}$$
$$D_{KL}(p'||q') = 0.04$$

$$D_{KL}(p||q) = \sum_i p_i \log \frac{p_i}{q_i}$$

$$\begin{aligned}CE(\mathbf{p}, \mathbf{q}) &= H(\mathbf{p}) + D_{KL}(\mathbf{p}||\mathbf{q}) \\ &= -\sum_i \mathbf{p}_i \log(\mathbf{q}_i)\end{aligned}$$

$$\arg \min_{\mathbf{q}} CE(\mathbf{p}, \mathbf{q}) = \arg \min_{\mathbf{q}} D_{KL}(\mathbf{p}, \mathbf{q})$$

Entropia cruzada e verossimilhança

Assumindo que \mathbf{y} é one-hot temos que:

$$\begin{aligned} L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}) &= CE(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}) \\ &= - \sum_{k=1}^n \mathbf{y}_k^{(i)} \log p(y = k | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= - \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \end{aligned}$$

Entropia cruzada e verossimilhança

E a função que queremos minimizar é

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{N} \sum_{i=1}^N L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}) \\ &= -\frac{1}{N} \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= -\mathcal{L}(\boldsymbol{\theta}) \end{aligned}$$

$$\arg \max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

- $\hat{\mathbf{y}} = f(\mathbf{x}; \boldsymbol{\theta})$
- $J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})$
- algum algoritmo de otimização (e.g., **SGD**):

$$\boldsymbol{\theta}^{novo} \leftarrow \boldsymbol{\theta}^{velho} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Vamos ver como computar $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ de modo eficiente para uma função arbitrária J .

Grafo de computação (caso escalar)

Considere os seguintes conjuntos de funções:

- $OP_1 = \{\lambda x. -x, \lambda x. x^2, \lambda x. e^x, \lambda x. \log(x), \lambda x. x\}$
- $OP_2 = \{\lambda xy. x + y, \lambda xy. x * y, \lambda xy. \frac{x}{y}\}$
- $OP = OP_1 \cup OP_2$

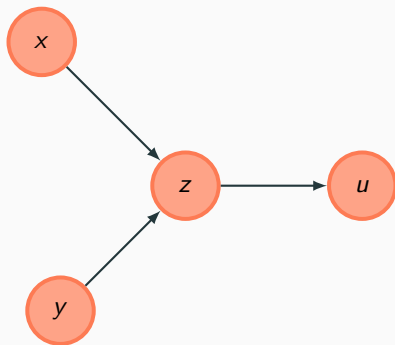
Um grafo de computação definido em OP $\mathcal{G} = (\mathcal{V}, \mathcal{E}_1, \mathcal{E}_2)$ é um grafo acíclico dirigido (DAG) tal que cada elemento $u \in \mathcal{V}$ indica uma variável, se $(x, y) \in \mathcal{E}_1$ então $f(x) = y$ onde $f \in OP_1 \cup \{g(x, \alpha) | \alpha \in \mathbb{R}, g \in OP_2\}$, e se $(x, y) \in \mathcal{E}_2$ então $f(x) = y$ onde $f \in \{g(\alpha, x) | \alpha \in \mathbb{R}, g \in OP_2\}$.

- $Pa(x) = \{y \in \mathcal{V} | (y, x) \in \mathcal{E}_1 \cup \mathcal{E}_2\}$.
- $S(x) = \{y \in \mathcal{V} | (x, y) \in \mathcal{E}_1 \cup \mathcal{E}_2\}$.



- $y = x^2$
- $u = e^y$

Grafo de computação

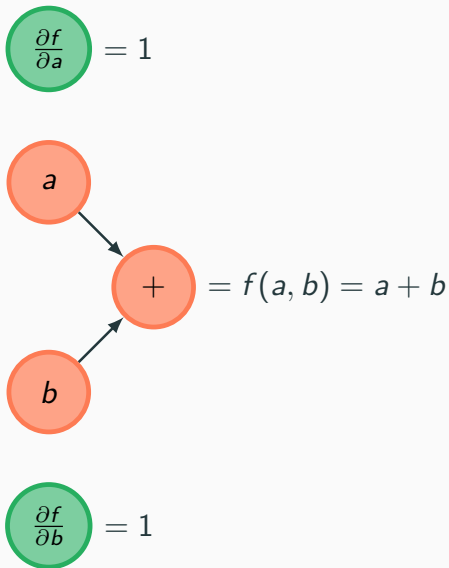


- $z = x + y$
- $u = \log(z)$

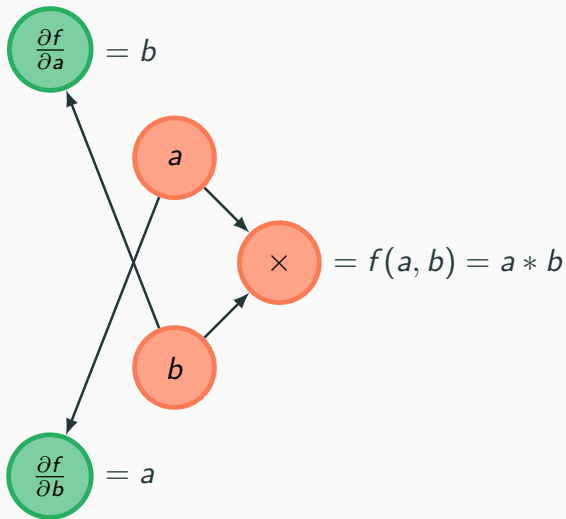
Queremos representar uma função L por um grafo definido em OP pois as derivadas parciais das funções de OP são simples de calcular. E com a **a regra da cadeia** podemos combinar as derivadas das funções locais para obter a derivada parcial de L com respeito a quaisquer parâmetros.

Como todas as funções em OP são diferenciáveis, podemos estender \mathcal{G} em \mathcal{G}' adicionando todas as derivadas parciais dos filhos em relação aos pais junto com as respectivas dependências.

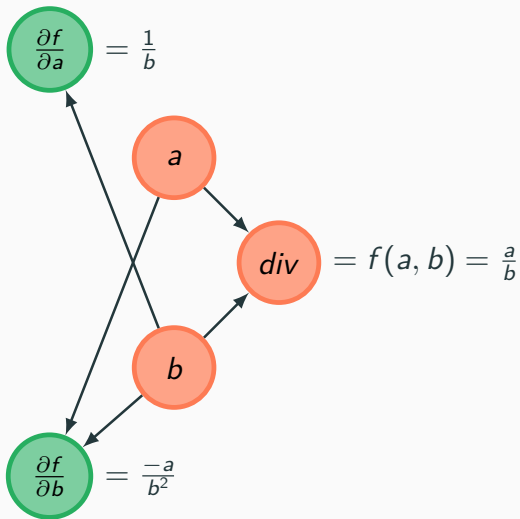
Extendendo o grafo de operações básicas: soma




Extendendo o grafo de operações básicas: multiplicação

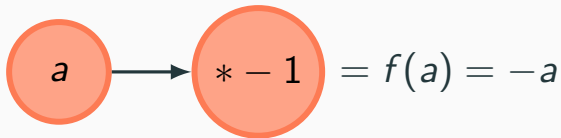


Extendendo o grafo de operações básicas: divisão

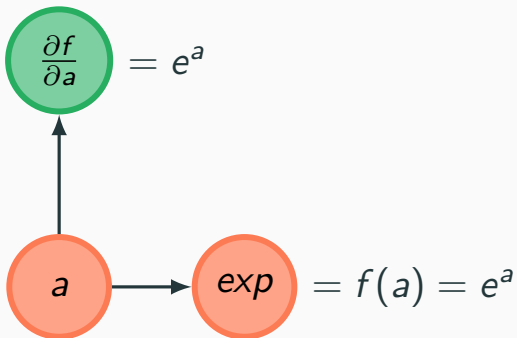


Extendendo o grafo de operações básicas: negativo

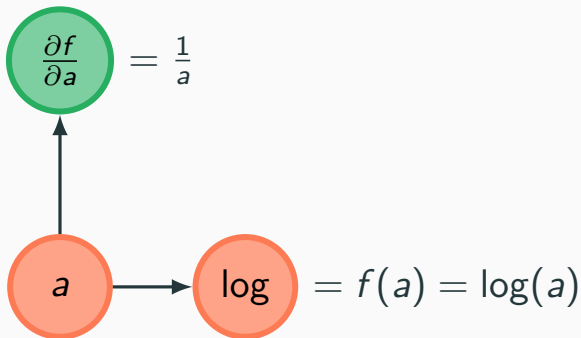

$$\frac{\partial f}{\partial a} = -1$$


$$a \rightarrow * - 1 = f(a) = -a$$

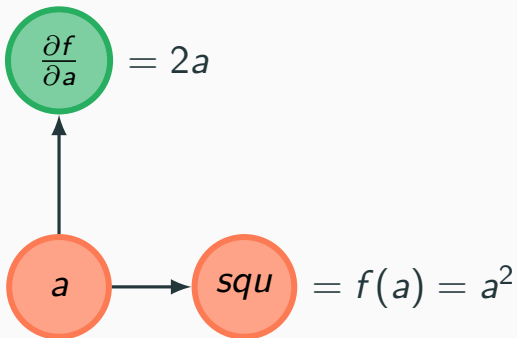
Extendendo o grafo de operações básicas: exponenciação



Extendendo o grafo de operações básicas: logarítimo



Extendendo o grafo de operações básicas: ao quadrado



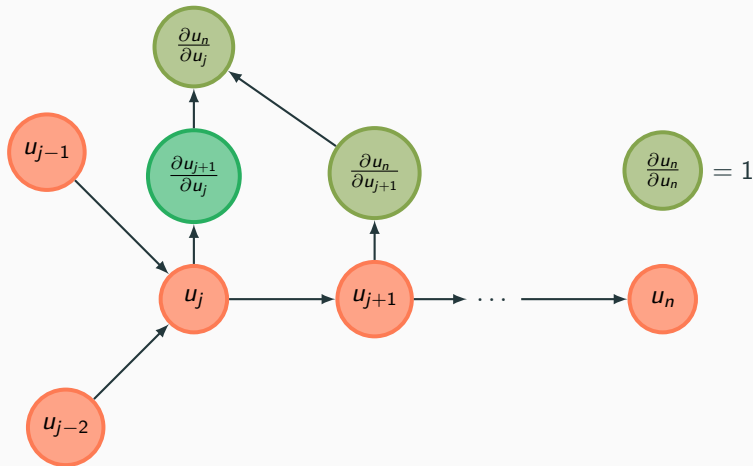
Regra da cadeia

- $f : \mathbb{R} \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}$.
- $y = g(x)$
- $u = f(g(x)) = f(y)$



$$\frac{\partial u}{\partial x} = \frac{\partial u}{\partial y} \frac{\partial y}{\partial x}$$

Aplicando a regra da cadeia



- $$\frac{\partial u_n}{\partial u_j} = \frac{\partial u_n}{\partial u_{j+1}} \frac{\partial u_{j+1}}{\partial u_j}$$

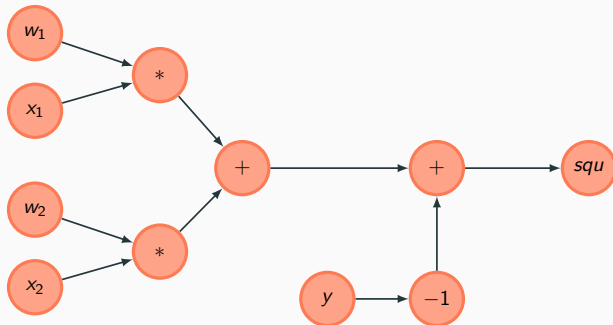
Exemplo 1: regressão linear

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{y}_i) \\ &= \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \end{aligned}$$

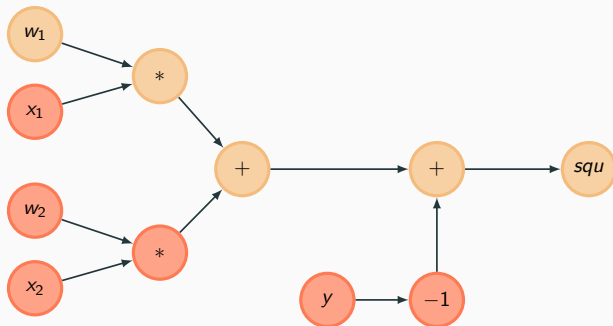
- $\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

- $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

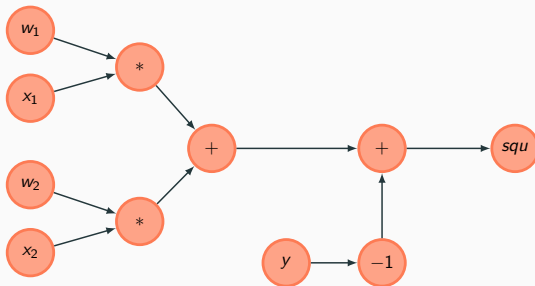
Grafo de $L(\hat{y}, y)$



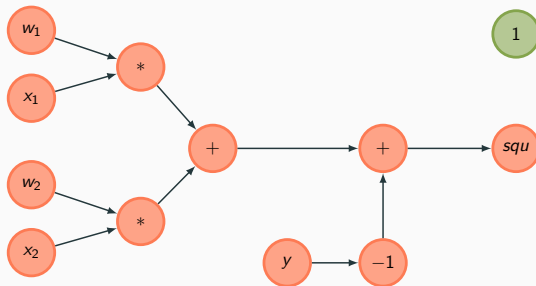
Caminho de w_1



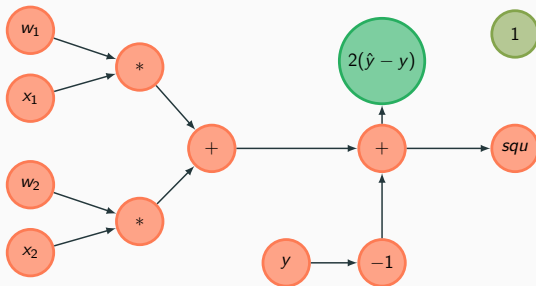
Derivade de L em relação a w_1



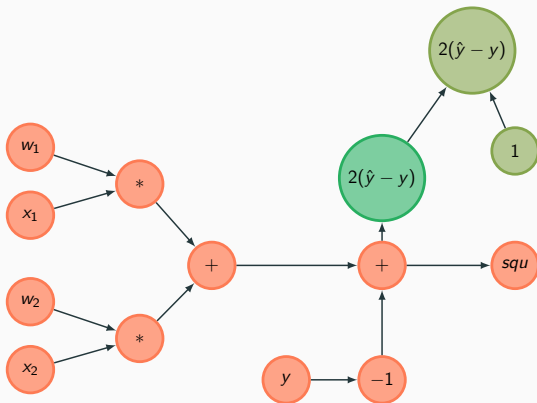
Derivade de L em relação a w_1



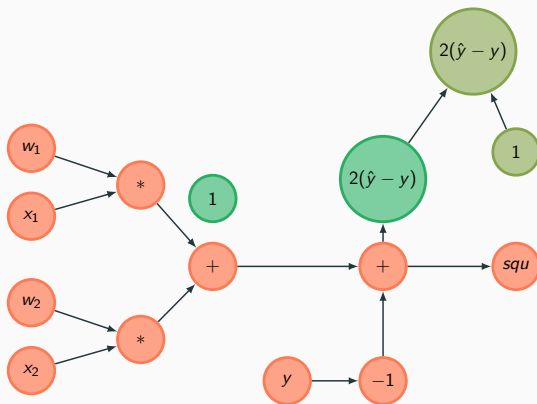
Derivade de L em relação a w_1



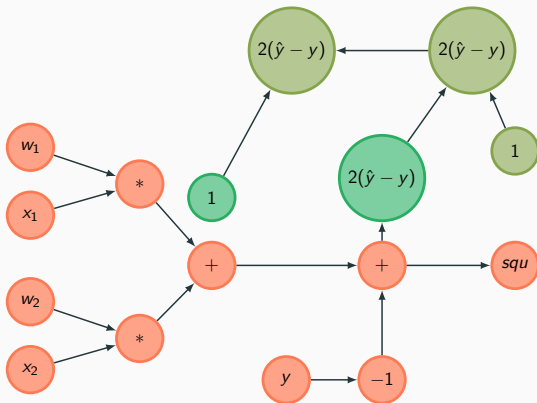
Derivade de L em relação a w_1



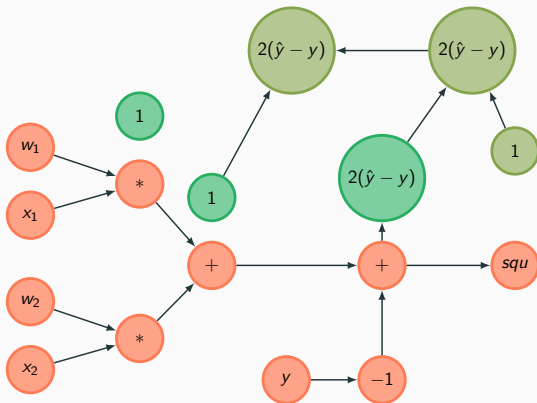
Derivade de L em relação a w_1



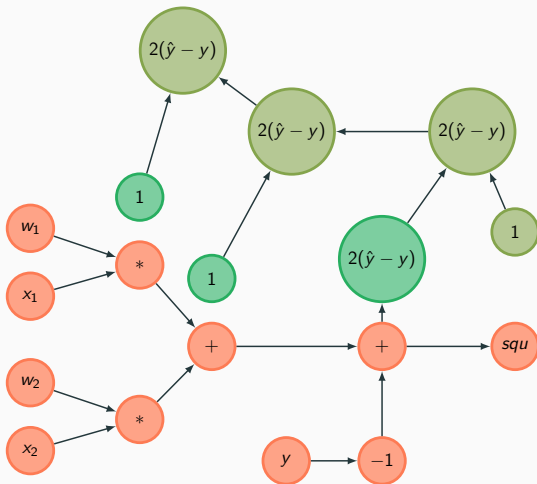
Derivade de L em relação a w_1



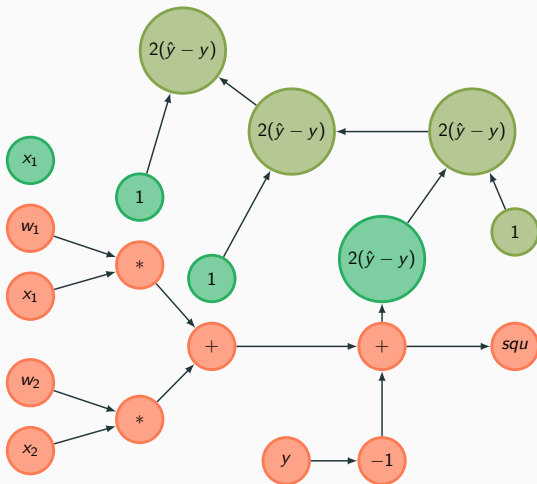
Derivade de L em relação a w_1



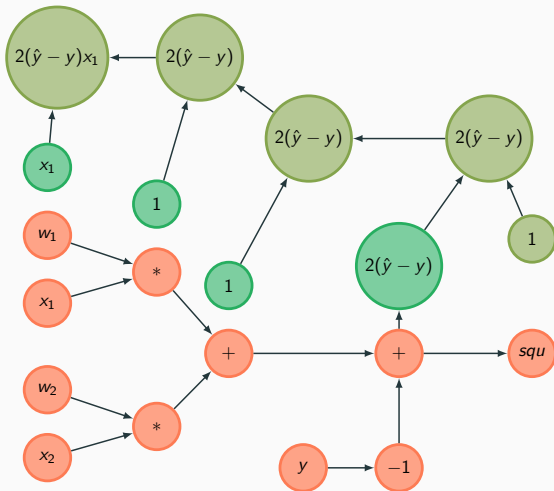
Derivade de L em relação a w_1



Derivade de L em relação a w_1



Derivade de L em relação a w_1

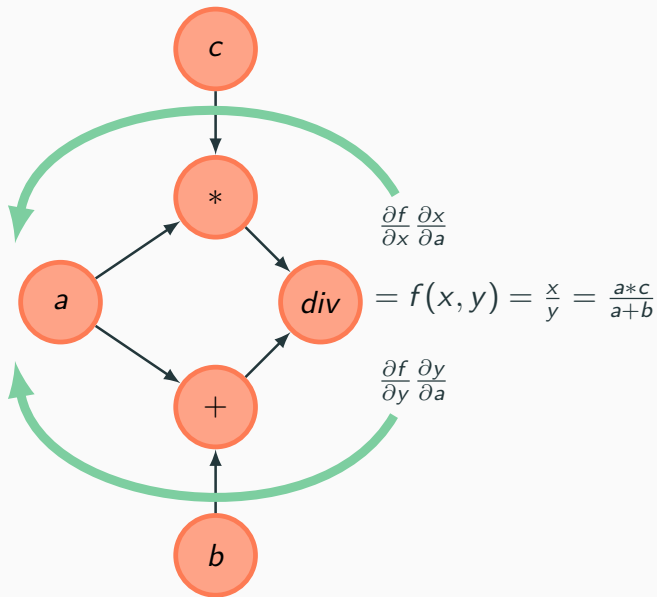


Regra da cadeia para várias variáveis

- $z = f(x, y)$
- $x = f_1(a)$.
- $y = f_2(a)$

$$\frac{\partial z}{\partial a} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial a} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial a}$$

Exemplo



Exemplo 2: regressão logística

$$\hat{\mathbf{y}} = \textit{softmax}(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$L(\mathbf{y}, \hat{\mathbf{y}}) = CE(\mathbf{y}, \hat{\mathbf{y}})$$

$$L(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i \mathbf{y}_i \log \left(\frac{\exp(\sum_k \mathbf{W}_{i,k} \mathbf{x}_k + \mathbf{b}_i)}{\sum_j \exp(\sum_k \mathbf{W}_{j,k} \mathbf{x}_k + \mathbf{b}_j)} \right)$$

Simplificação

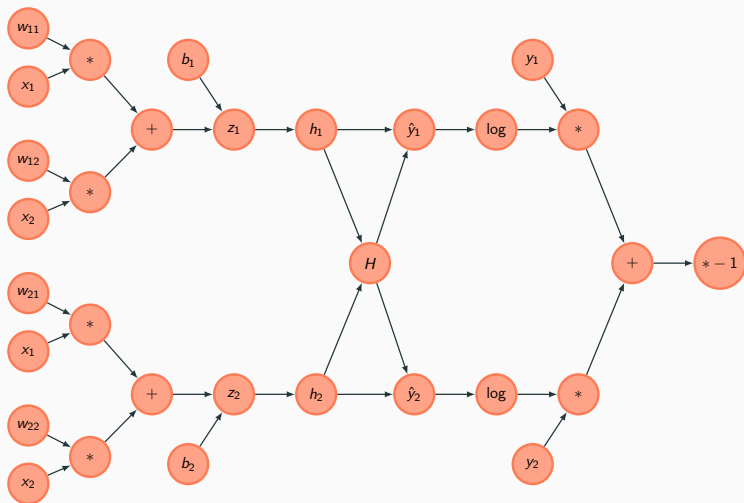
$$\bullet \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$\bullet \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} \exp(z_1) \\ \exp(z_2) \end{bmatrix}$$

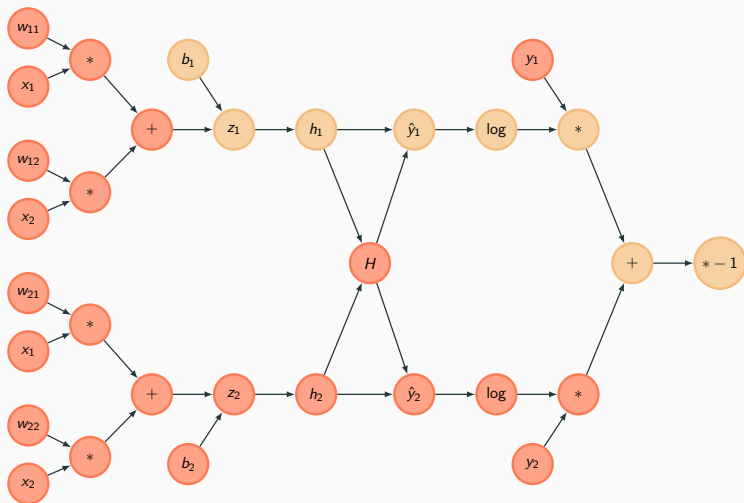
$$\bullet H = h_1 + h_2$$

$$\bullet \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} = \begin{bmatrix} \frac{h_1}{H} \\ \frac{h_2}{H} \end{bmatrix}$$

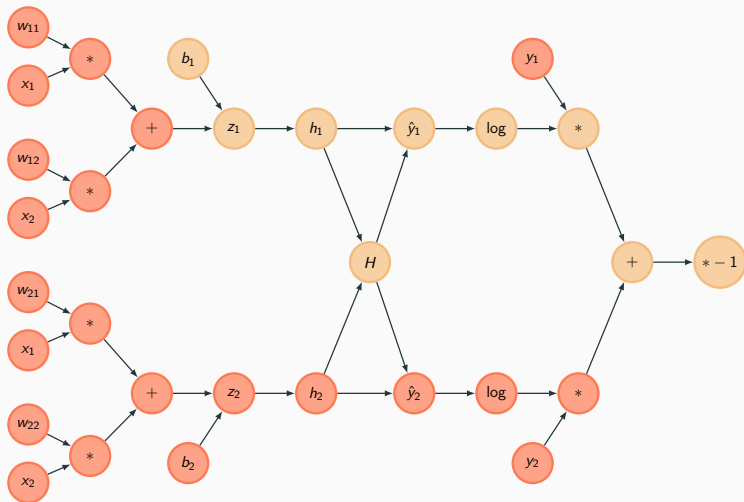
Grafo de $L(\hat{y}, y)$



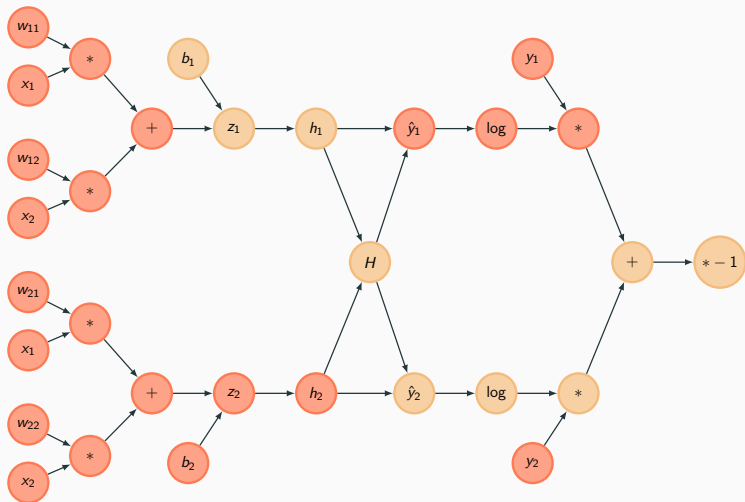
Caminho de b_1 : 1



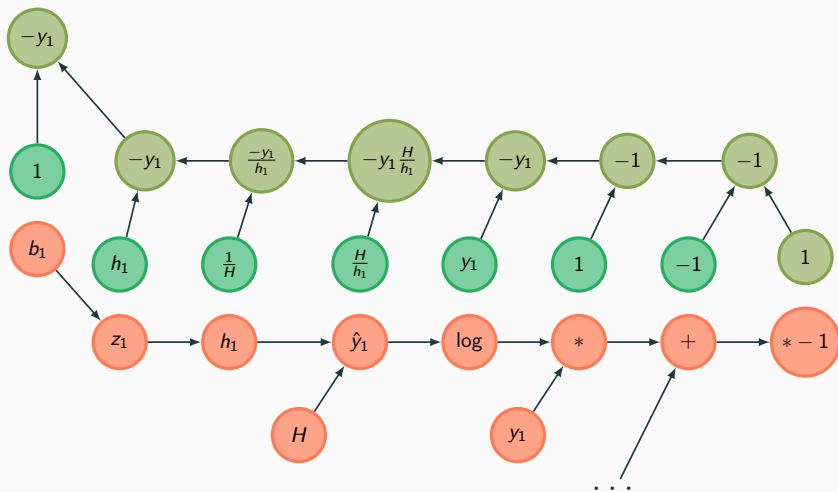
Caminho de b_1 : 2



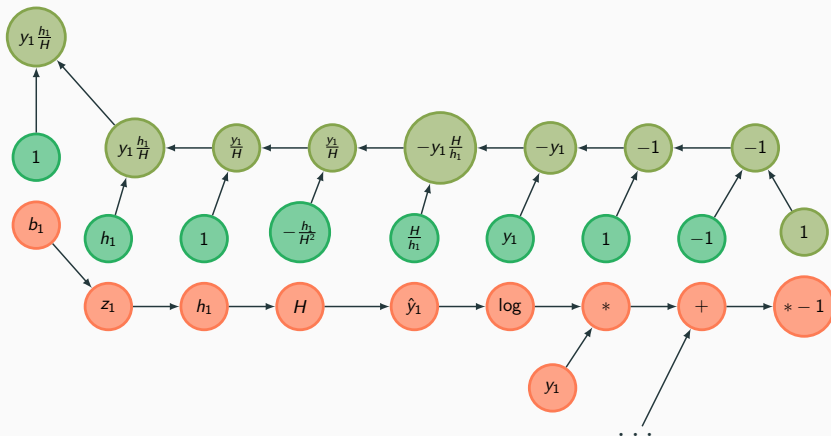
Caminho de b_1 : 3



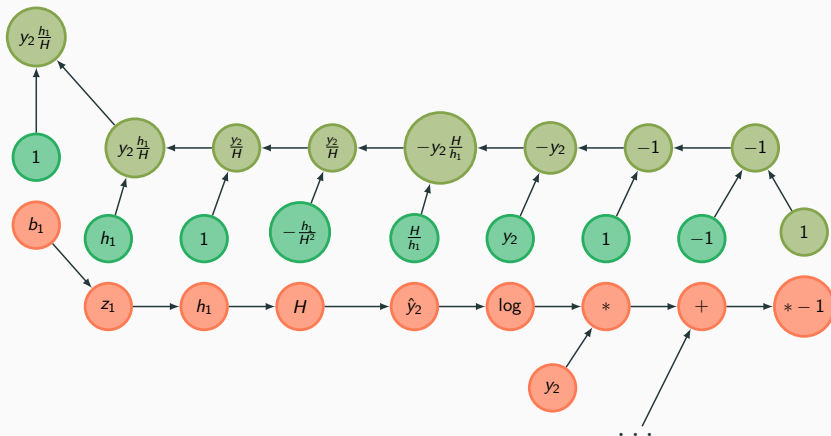
Derivada parcial de L com respeito a b_1 : 1



Derivada parcial de L com respeito a b_1 : 2



Derivada parcial de L com respeito a b_1 : 3



Derivada parcial de L com respeito a b_1

$$\frac{\partial L}{\partial b_1} = -y_1 + y_1 \frac{h_1}{H} + y_2 \frac{h_1}{H}$$

Derivada parcial de L com respeito a b_1

$$\begin{aligned}\frac{\partial L}{\partial b_1} &= -y_1 + y_1 \frac{h_1}{H} + y_2 \frac{h_1}{H} \\ &= y_1 \left(\frac{h_1}{H} - 1 \right) + y_2 \left(\frac{h_1}{H} - 0 \right)\end{aligned}$$

Derivada parcial de L com respeito a b_1

$$\begin{aligned}\frac{\partial L}{\partial b_1} &= -y_1 + y_1 \frac{h_1}{H} + y_2 \frac{h_1}{H} \\ &= y_1 \left(\frac{h_1}{H} - 1 \right) + y_2 \left(\frac{h_1}{H} - 0 \right) \\ &= y_1 (\hat{y}_1 - 1) + y_2 (\hat{y}_1 - 0)\end{aligned}$$

Derivada parcial de L com respeito a b_1

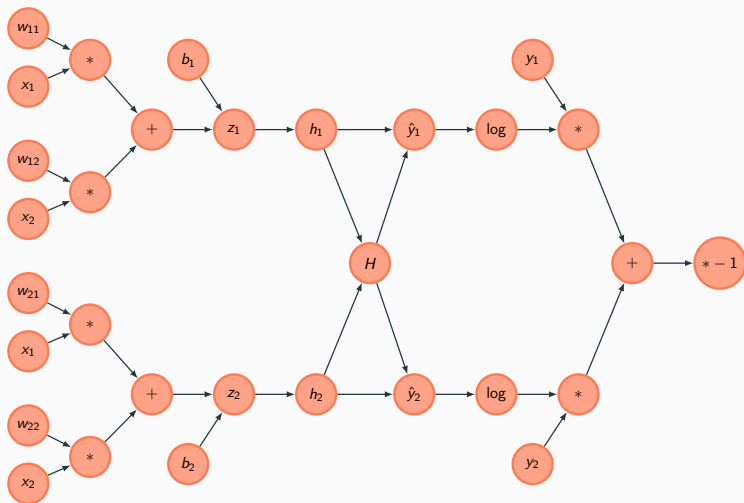
$$\begin{aligned}\frac{\partial L}{\partial b_1} &= -y_1 + y_1 \frac{h_1}{H} + y_2 \frac{h_1}{H} \\ &= y_1 \left(\frac{h_1}{H} - 1 \right) + y_2 \left(\frac{h_1}{H} - 0 \right) \\ &= y_1 (\hat{y}_1 - 1) + y_2 (\hat{y}_1 - 0) \\ &= \hat{y}_1 - y_1 \quad (\text{quando } y \text{ é um vetor one-hot})\end{aligned}$$

Exemplo

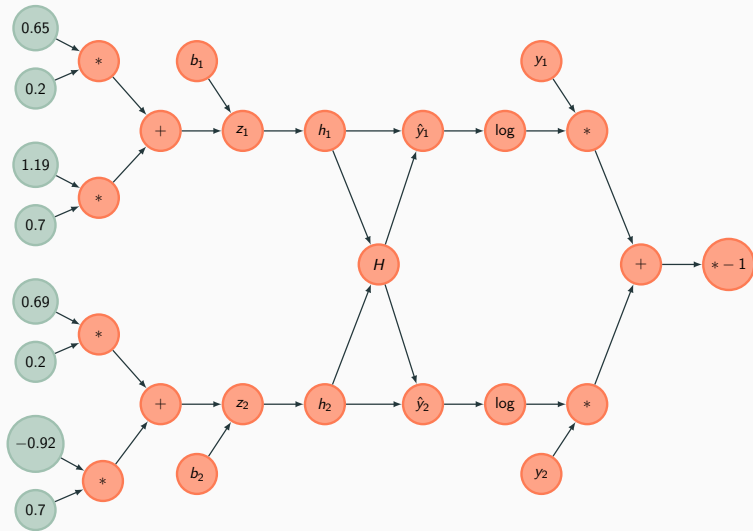
$$\mathbf{W} = \begin{bmatrix} 0.65 & 1.19 \\ 0.69 & -0.92 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 0.2 \\ 0.7 \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

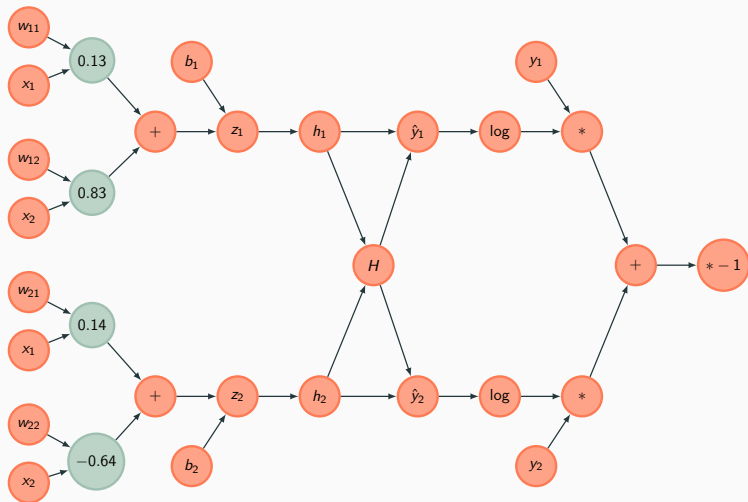
Forward



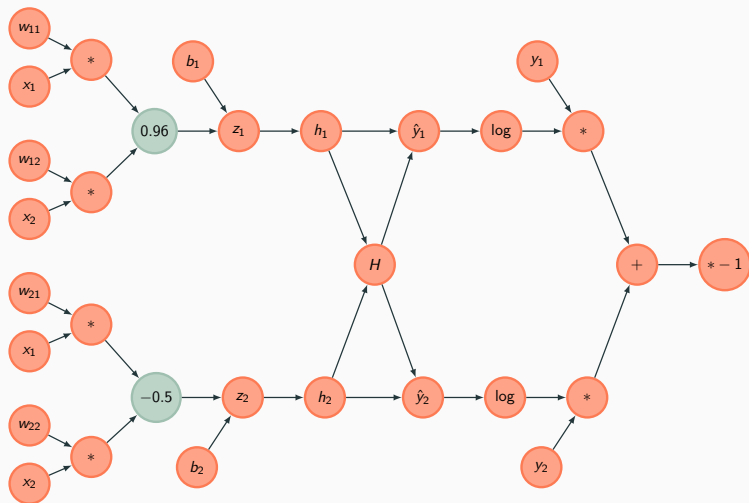
Forward



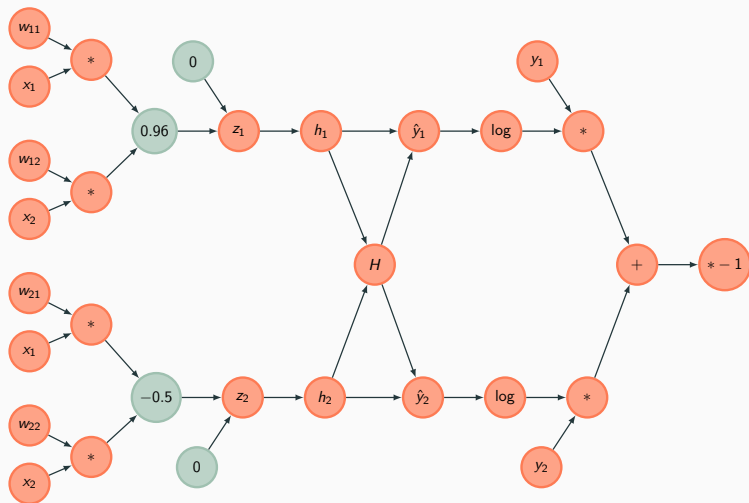
Forward



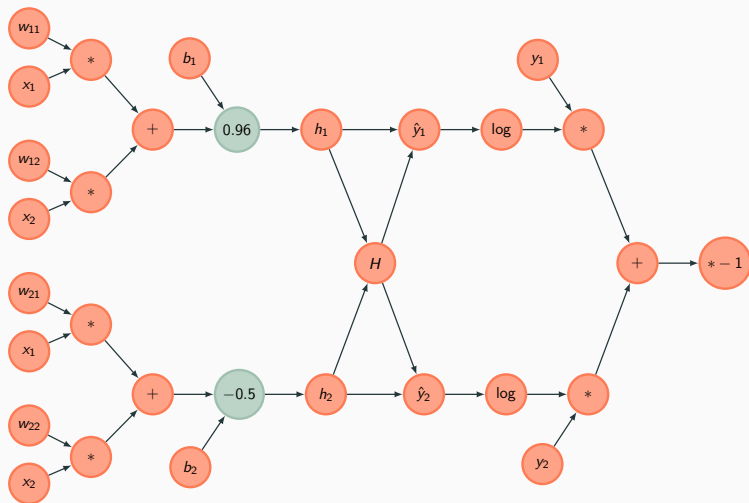
Forward



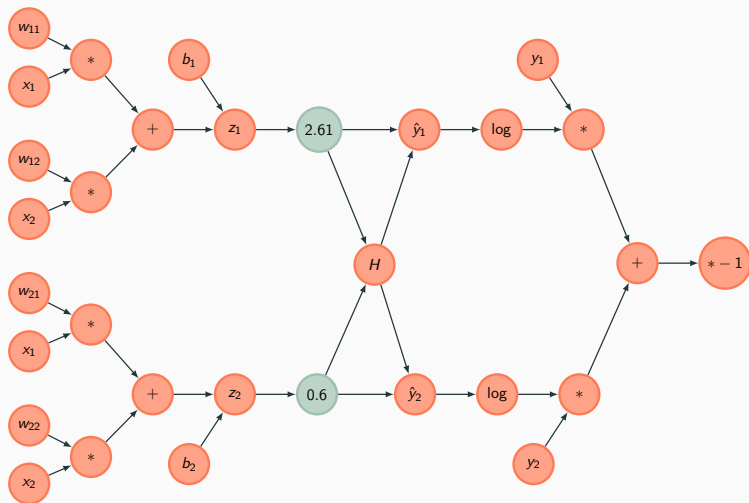
Forward



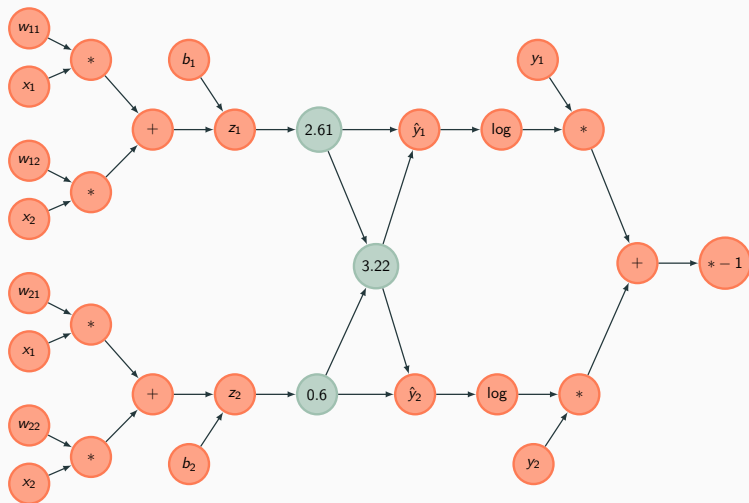
Forward



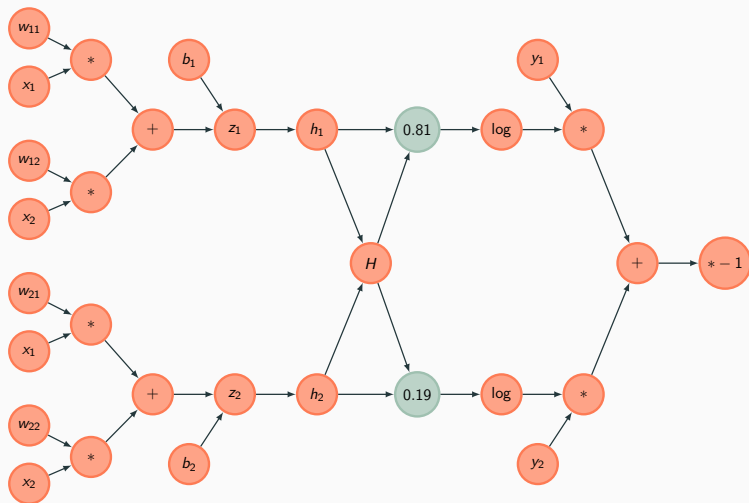
Forward



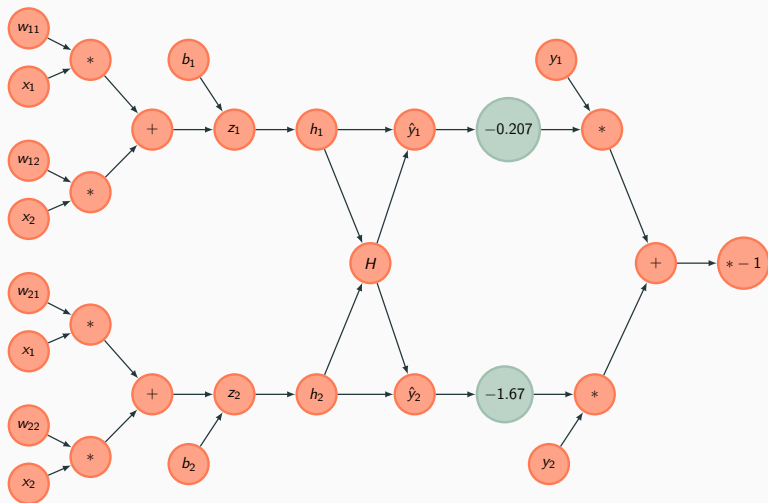
Forward



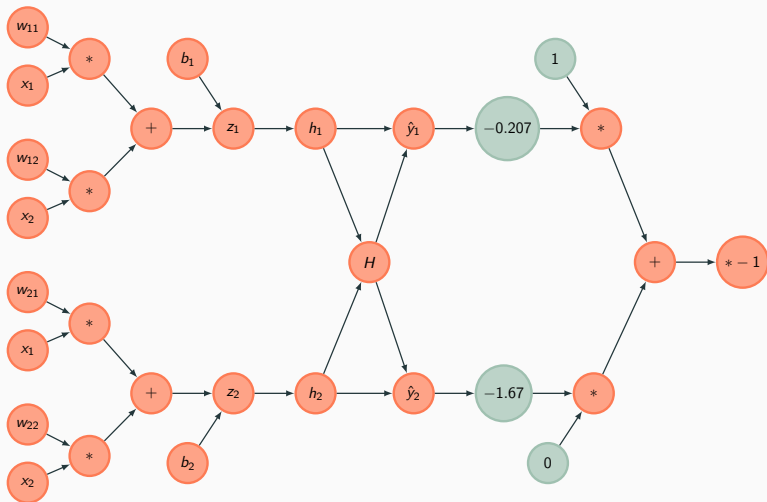
Forward



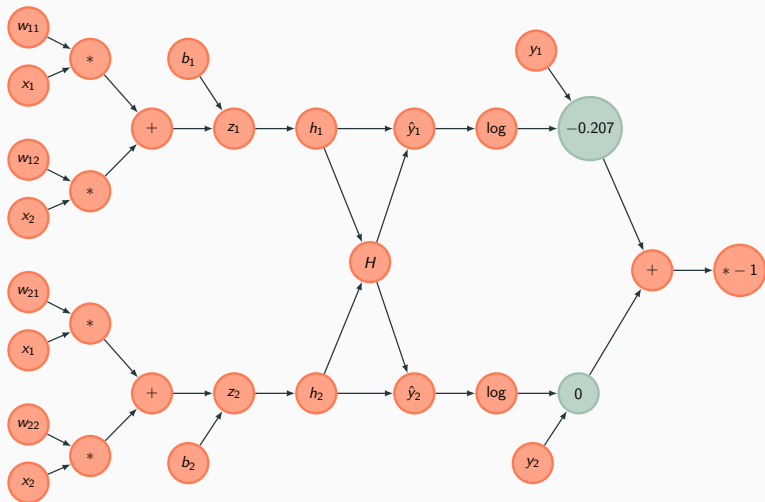
Forward



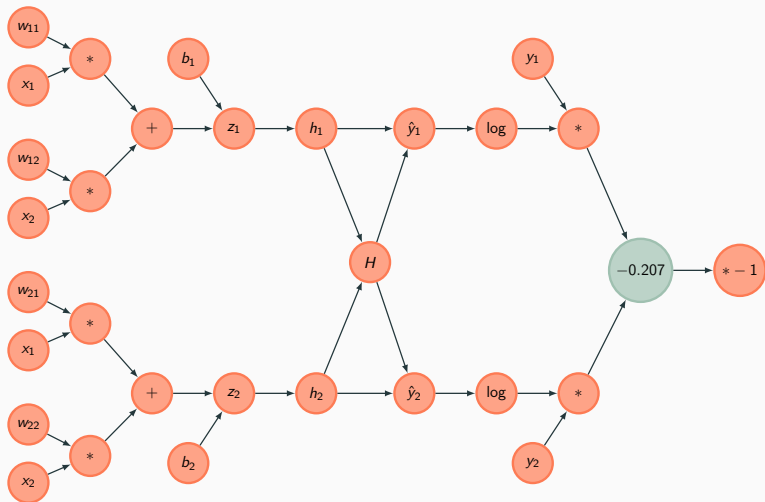
Forward



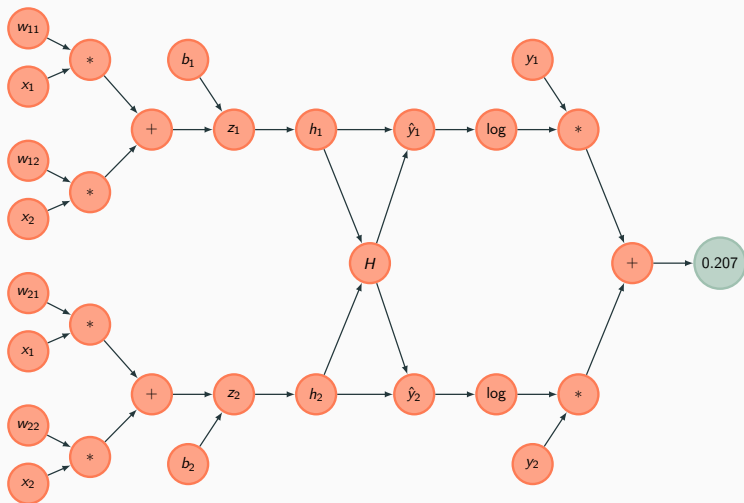
Forward



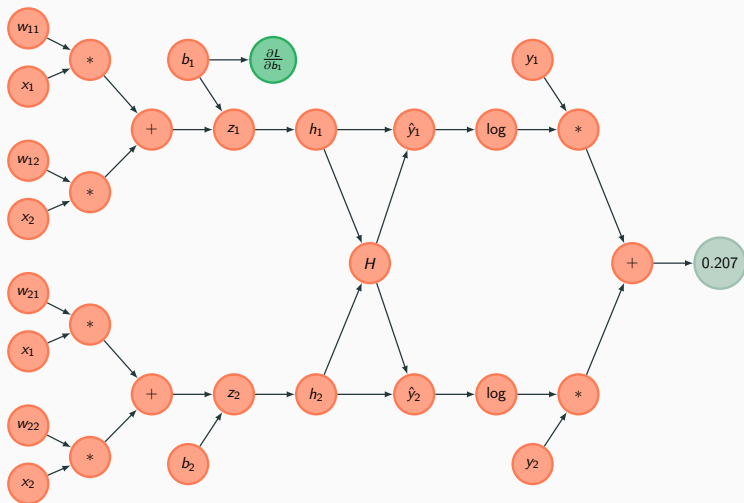
Forward



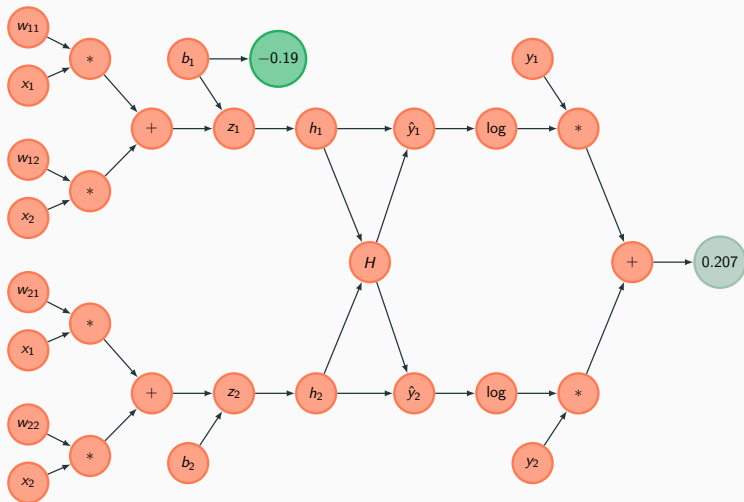
Forward



Forward



Backward



Algoritmo de back-propagation (caso escalar)

Algorithm 1 Back-propagation (scalar case)

- 1: **Require:** Computational graph $\mathcal{G} = (\{u_1, \dots, u_n\}, \mathcal{E}_1, \mathcal{E}_2)$, where u_n is a leaf node.
 - 2: Initialize *grad_table*, a data structure that will store the derivatives that have been computed (at the end $\text{grad_table}[u_i] = \frac{\partial u_n}{\partial u_i}$).
 - 3: $\text{grad_table}[u_n] \leftarrow 1$
 - 4: **for** $j = n - 1$ down to 1 **do**
 - 5: $\text{grad_table}[u_j] \leftarrow \sum_{u_i \in S(u_j)} \text{grad_table}[u_i] \frac{\partial u_i}{\partial u_j}$
 - 6: **end for**
 - 7: **return** *grad_table*
-



I. Goodfellow, Y. Bengio, and A. Courville.

Deep Learning.

MIT Press, 2017.