

Neural Networks

\vec{x}

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}$

\vec{x}

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}$

\vec{x}

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}$

$\hat{y} = a$ "activation"

\vec{x}

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}$

$\hat{y} = a$ "activation"

Input Layer
Output Layer
"Visible" Layers

One hidden layer

\vec{x}

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}$

$\hat{y} = a$ "activation"

Se Tom Mitchell

Multiple Networks and the Backpropagation Alg.

• Usado para aprender a NN para decisões não lineares.
• Função Sigmod \rightarrow Output não função dos inputs e dependentes, pois que permitem usar o gradientes descendente.

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

Backpropagation

Algoritmo que usa o gradientes descendente para diminuir o square error da rede neural.

Treinada num mundo de rede neural de múltiplos camadas, logo, precisamos estudar o erro.

$$E(\vec{x}) = \frac{1}{2} \sum_i \sum_j (t_{ij} - o_{ij})^2$$

\rightarrow Se usamos a quadrado no todos os possíveis outputs

\rightarrow busca estudar e encontrar os pesos ideais para todos da rede neural

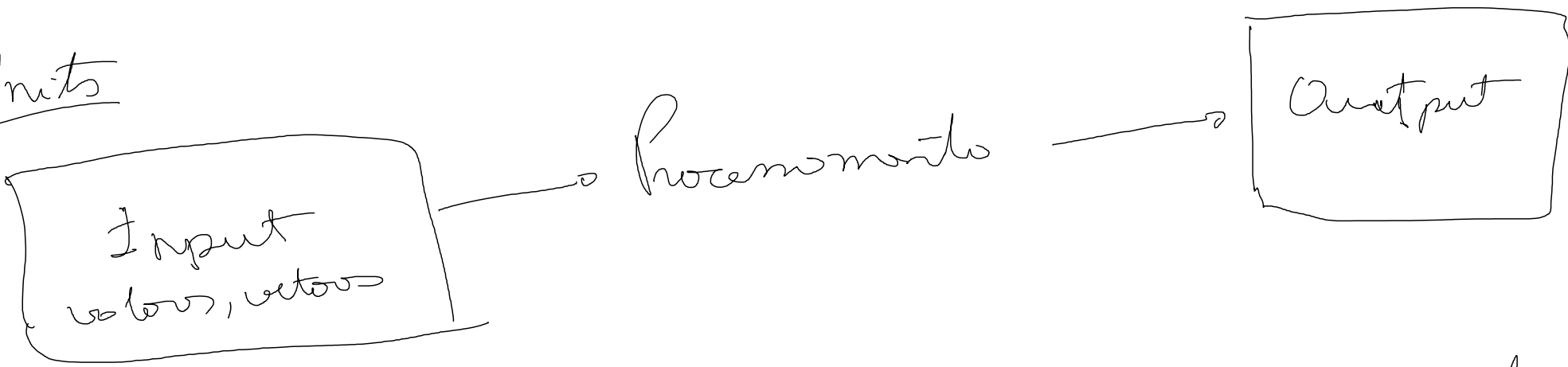
→ Multilayer

Revisão
Genl

* Neural Network Cap7 Stanford

Conjunto de unidades de computação que recebem vetores e retornam valores para a próxima camada.
Deep learning vem da ideia de várias camadas

→ Units



Uma unit se baseia na soma ponderada de inputs junto com um bias de modo que $Z = w \cdot x + b$ isso irá retornar apenas um valor

Após isso, aplicamos uma função de ativação não linear, para evitar que tudo fique linear. Três são comuns: Sigmoid; tanh; ReLU

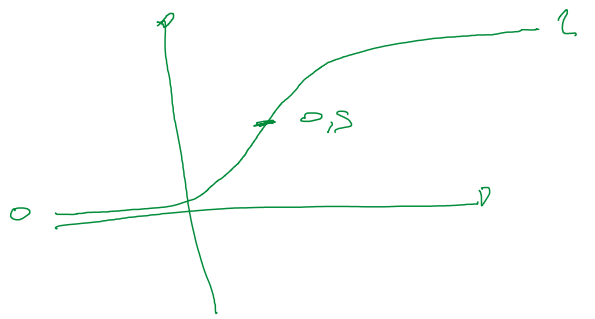
\rightarrow É importante que essas funções sejam deriváveis

\rightarrow Sigmoid irá mapear os resultados entre [0,1]

\rightarrow tanh irá mapear entre -1 e 1

\rightarrow ReLU é o max(0, z), ou seja, z se positivo, zero se negativo

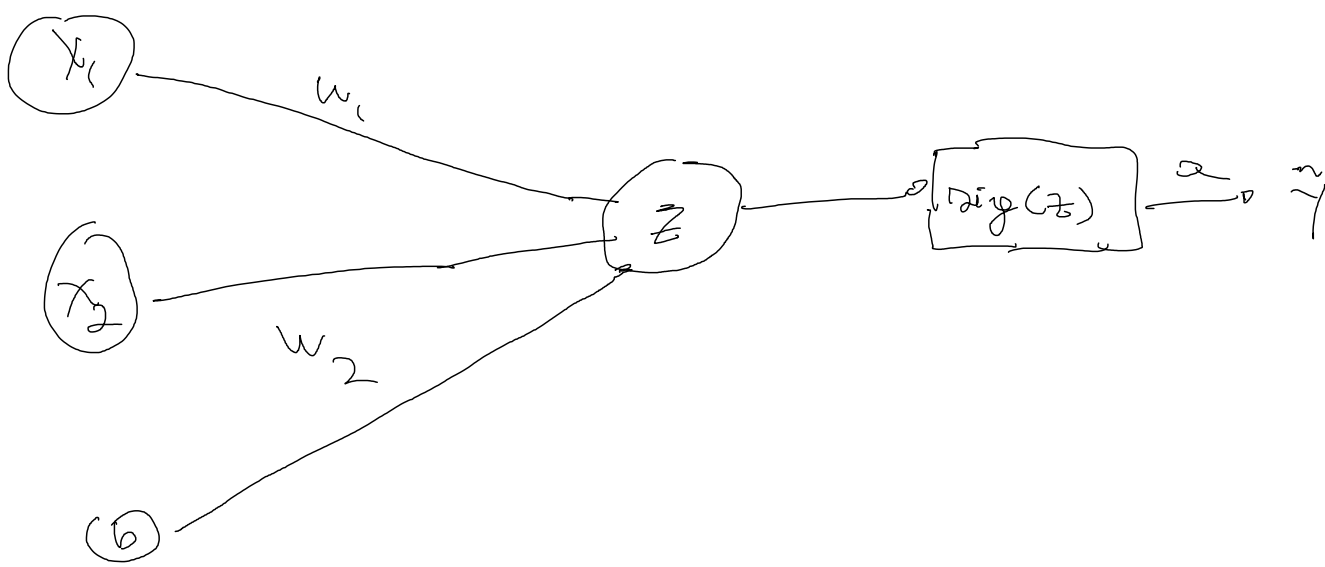
$$\text{Sigmoid} = \frac{1}{1 + e^{-x}}$$



\rightarrow Cada função de ativação tem características próprias que podem ajudar em certos contextos. Por isso ter camadas com diferentes funções de ativação é interessante

Alguns são facilmente diferenciáveis, etc

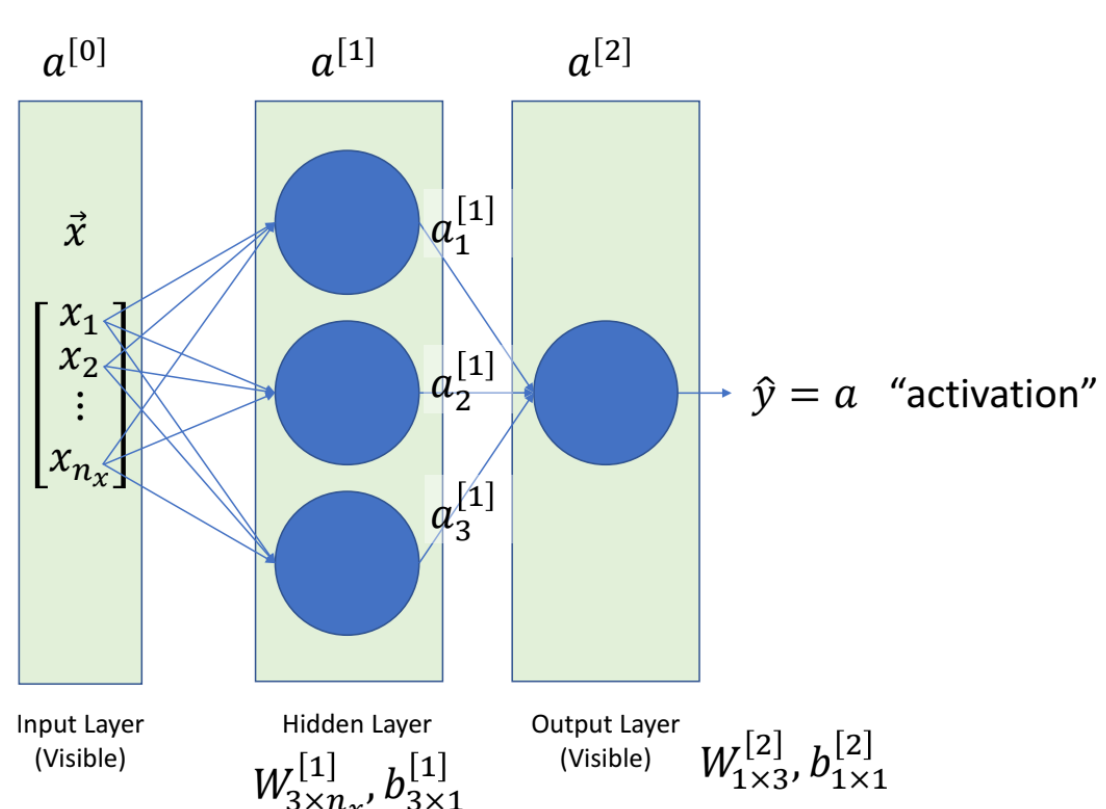
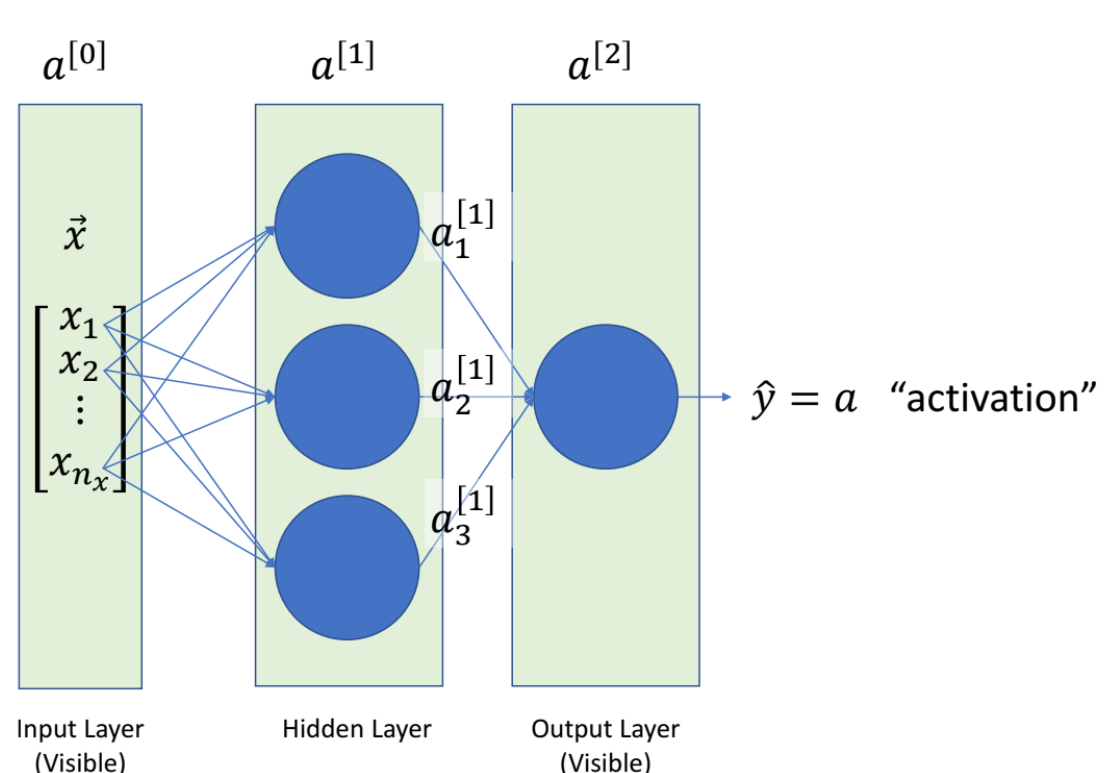
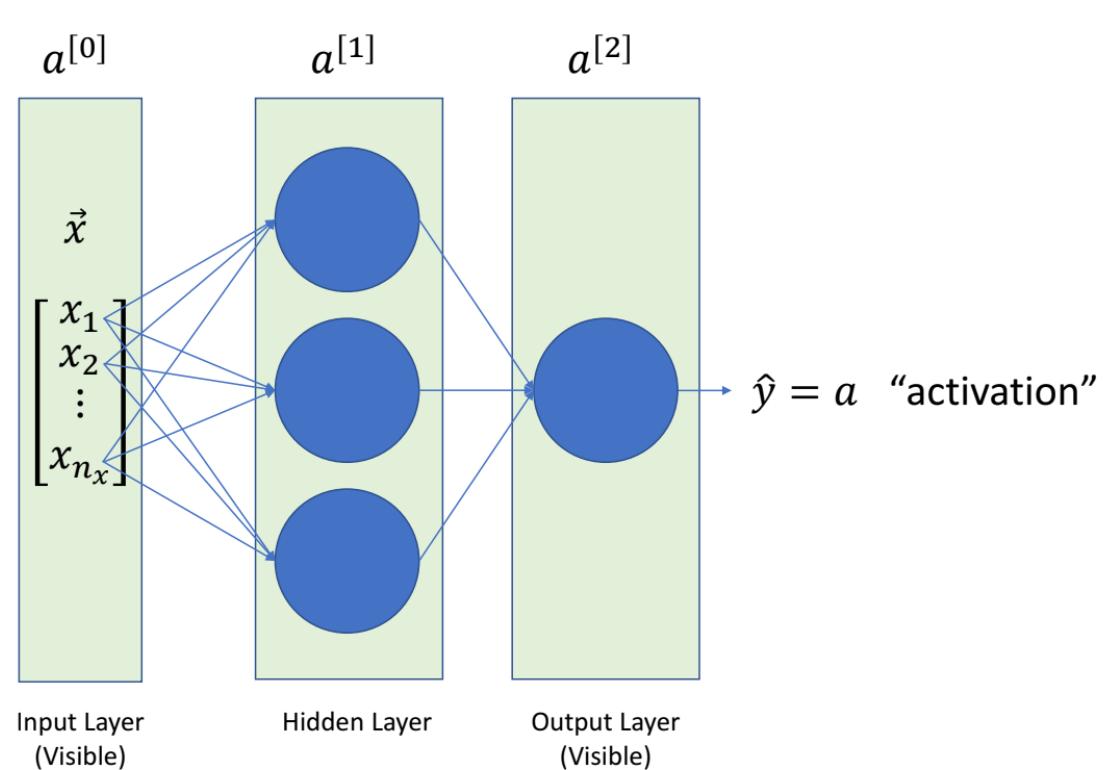
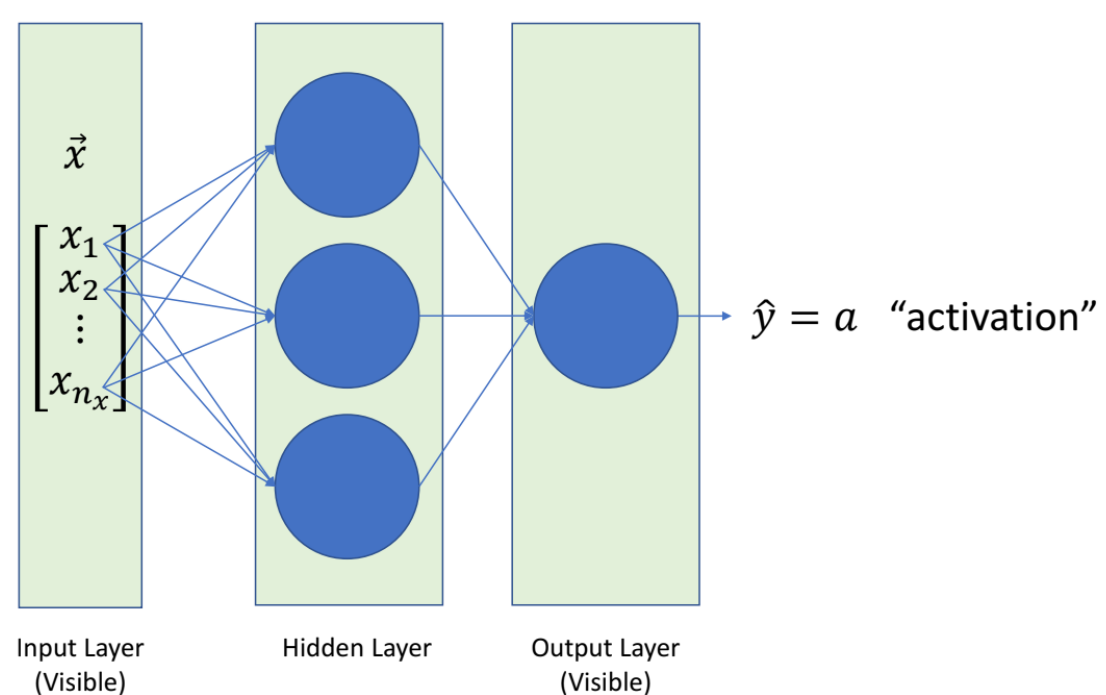
Computational Graph



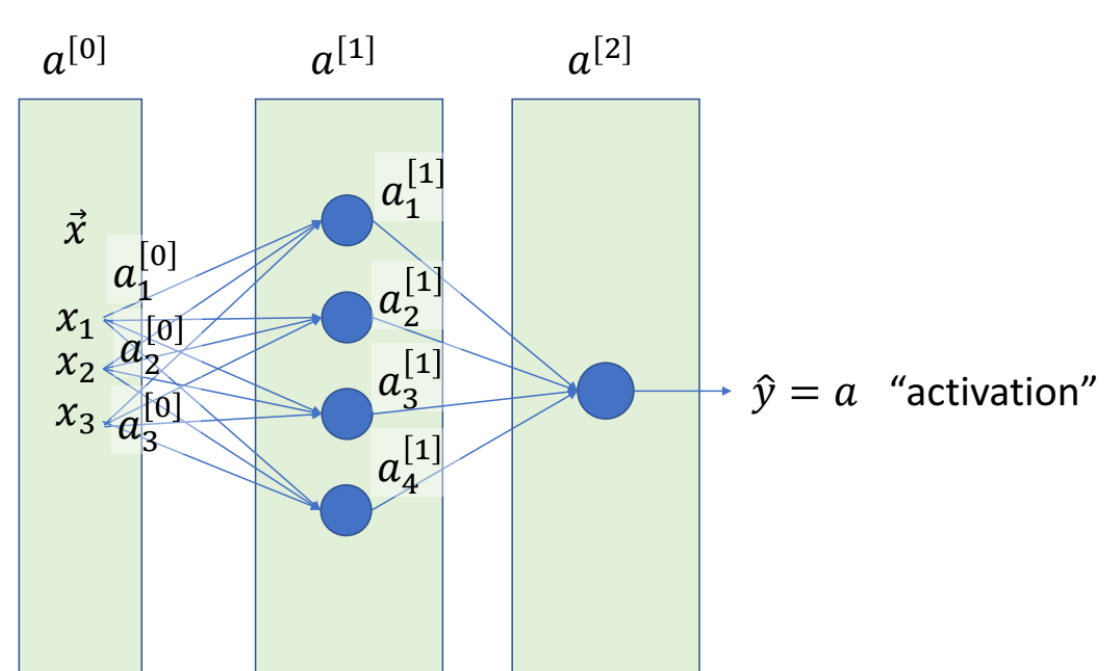
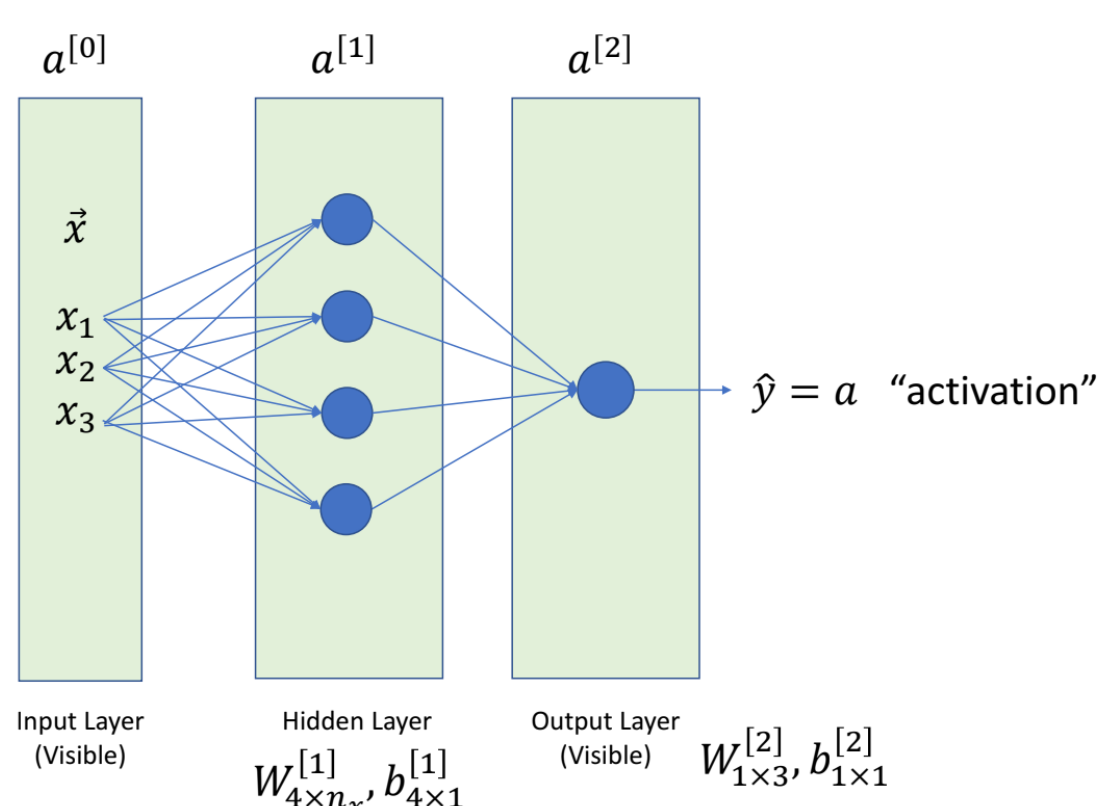
\vec{x}

$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}$

$\hat{y} = a$ "activation"



Computing an NN's output



Feedforward Propagation NN

Multilayer Perceptron

↳ Basicamente contém três partes

- Input
- Hidden layers = hidden perceptrons
- Output layers

↳ Cada hidden layer usa uma função de ativação não-linear e aplica uma função de ativação nos dados

↳ Cada unit de cada layer tem como input todos os dados da rede

↳ Cada layer tem como entrada um vetor x . A matriz W é o conjunto de todos os W_x

↳ Por isso a entrada é a matriz W , porque são todos os W_x

↳ Bias se torna um vetor.

↳ A partir disso, você consegue sempre calcular os dados

↳ A hidden layer usa os dados como output: $h = f(Wx + b)$, sendo

"f" a função de ativação

↳ Basicamente vai aplicar p/ todos os dados $W_1x_1 + b_1, W_2x_2 + b_2$ e etc

Primeira codificação para uma feedforward

for i in 1..n:
 $z^{(i)} = W^{(i)} a^{(i-1)} + b^{(i)}$
 ↳ resultado da camada anterior
 ↳ bias e weight da camada
 $a^{(i)} = g^{(i)}(z^{(i)})$
 ↳ função de ativação
 $\hat{y} = a^{(n)}$
 ↳ finaliza a rede neural

*Treinando Redes Neurais

O objetivo é ajustar os pesos "W" e o bias "b" para obtermos um \hat{y} mais próximo do valor real "y"

- Steps:
- Loss function: Calcula a erro entre \hat{y} e y
 - Gradient Descent: Encontra os parâmetros que minimizam a loss function
 - Backpropagation: Atribui os derivadas parciais p/ o gradiente

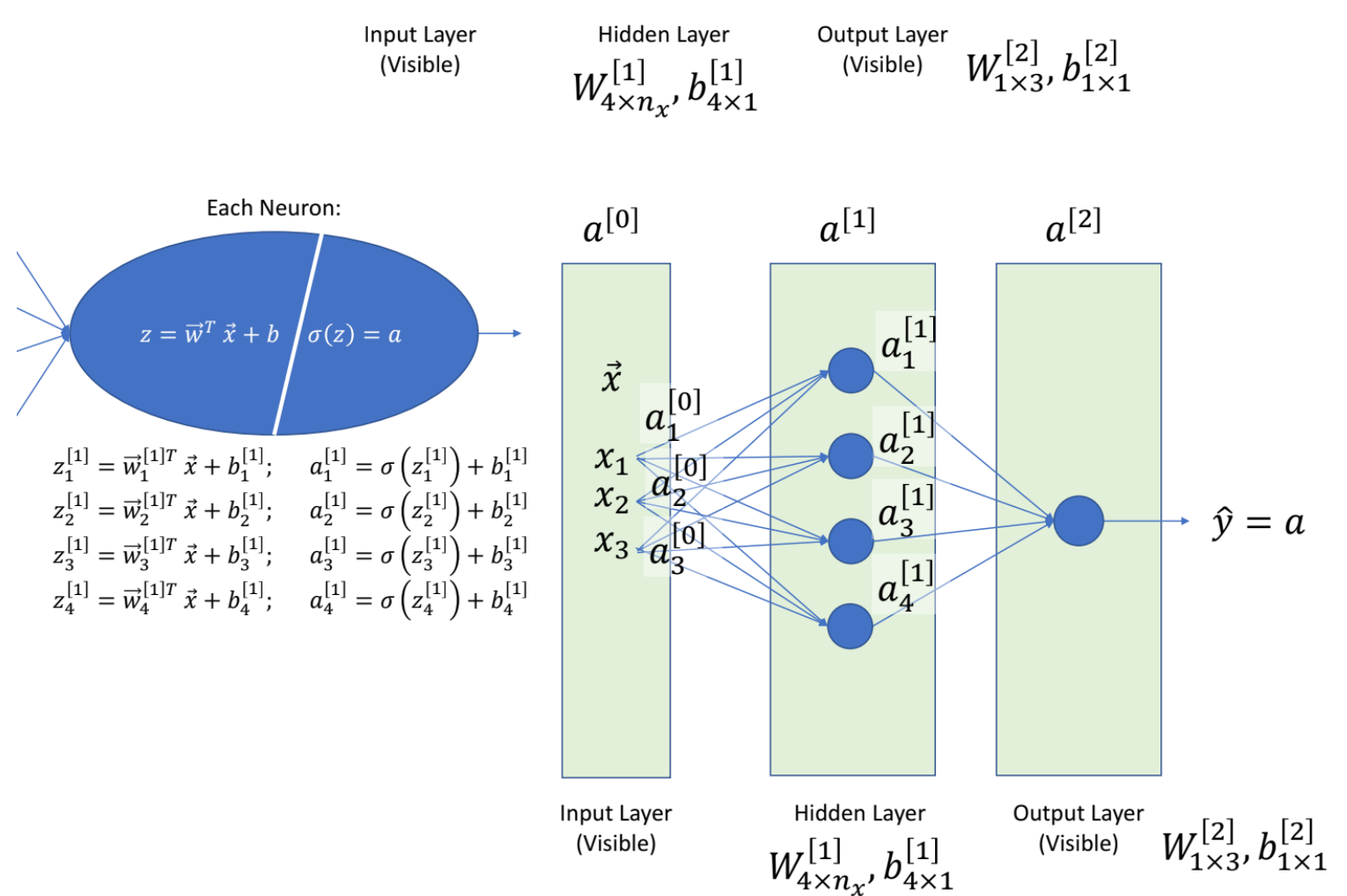
*Andrew Ng

↳ Vectorization

Forma de tirar os for loops da codificação, ajudando a reduzir o tempo em grandes datasets

A ideia é transformar os dados em operações de matriz, que são realizadas mais rapidamente até em GPUs. Todos os dados são transformados em vetores/matrizes

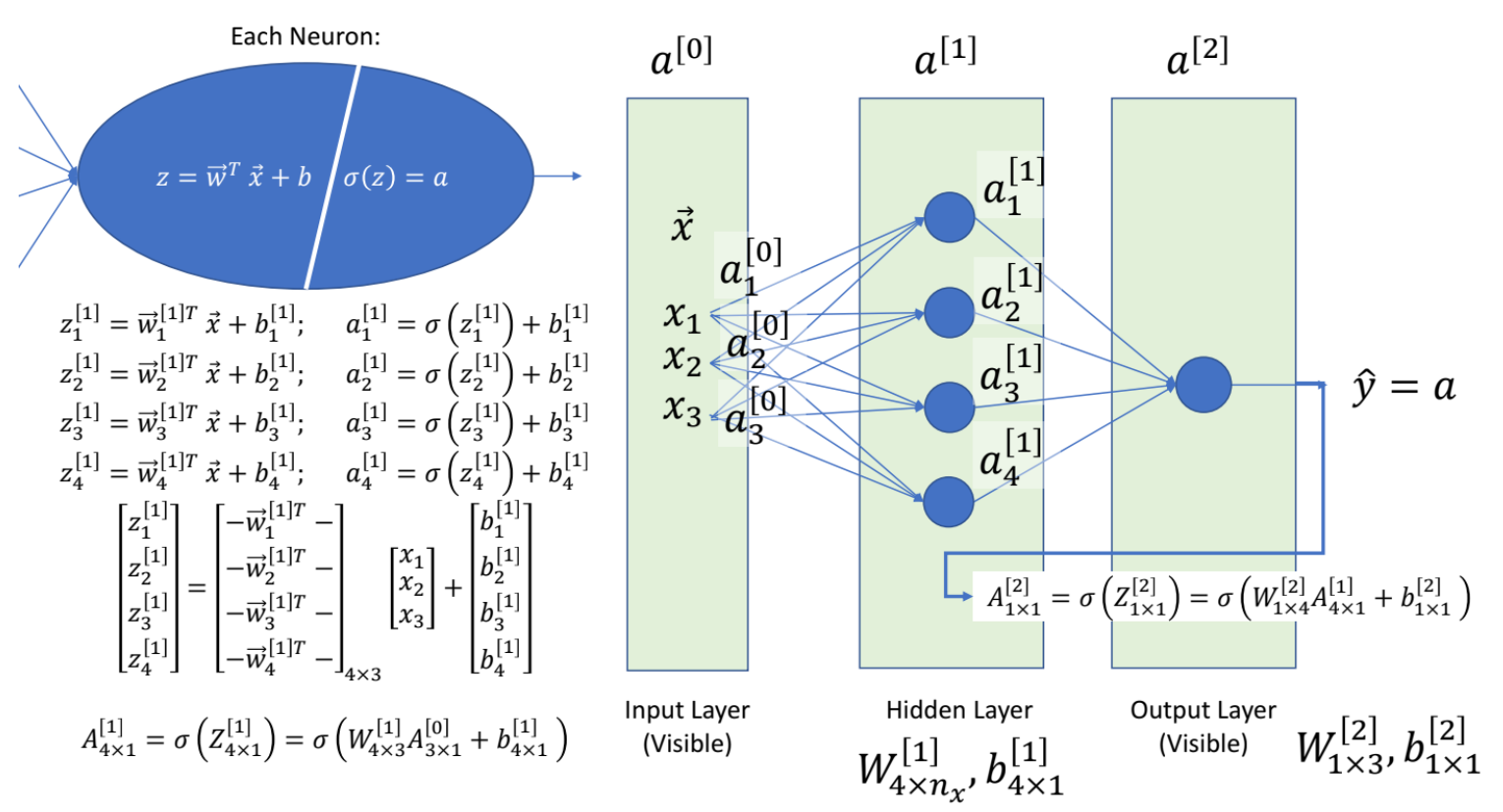
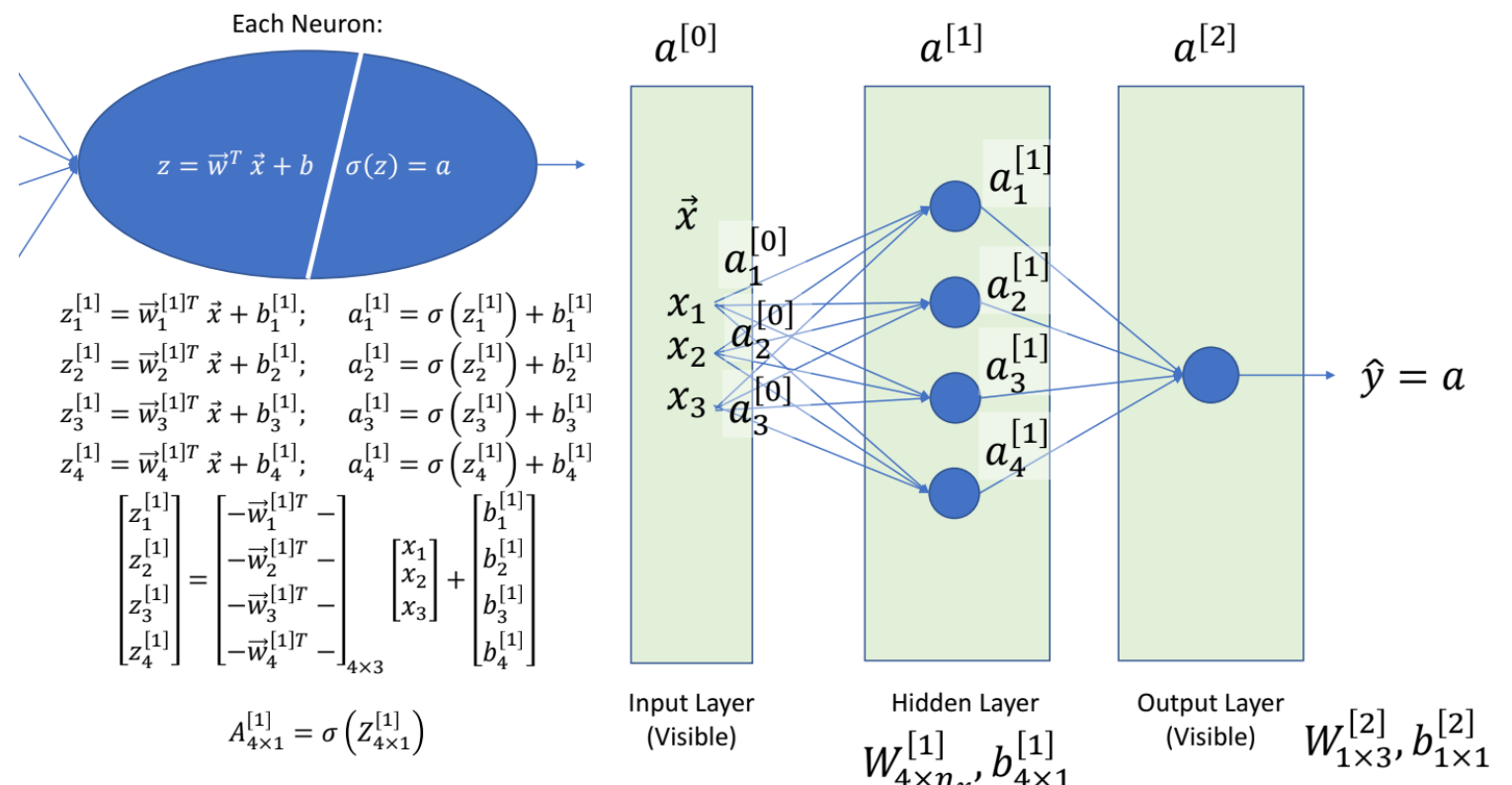
↳ Isso é feito em GPU, o CPU também tem



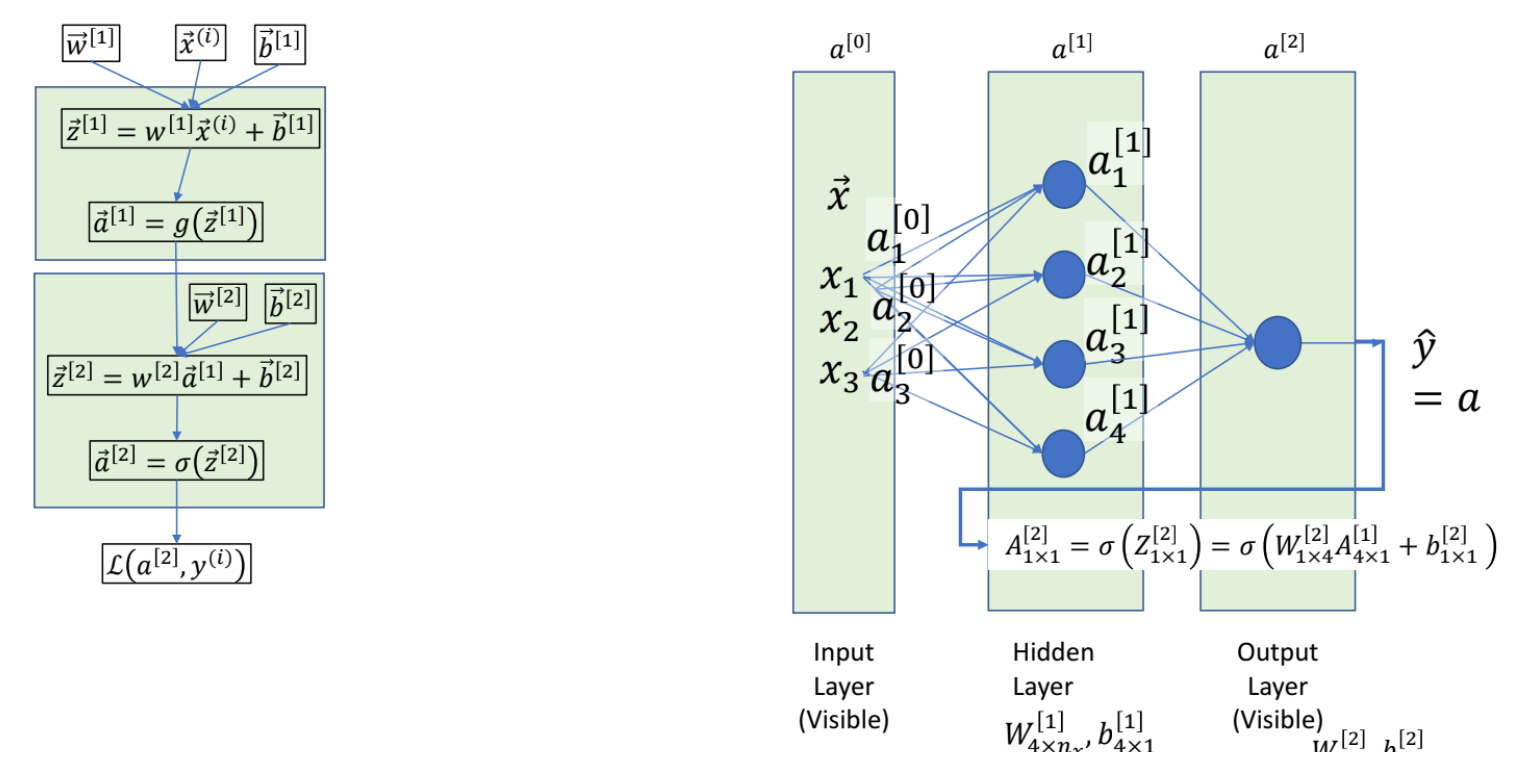
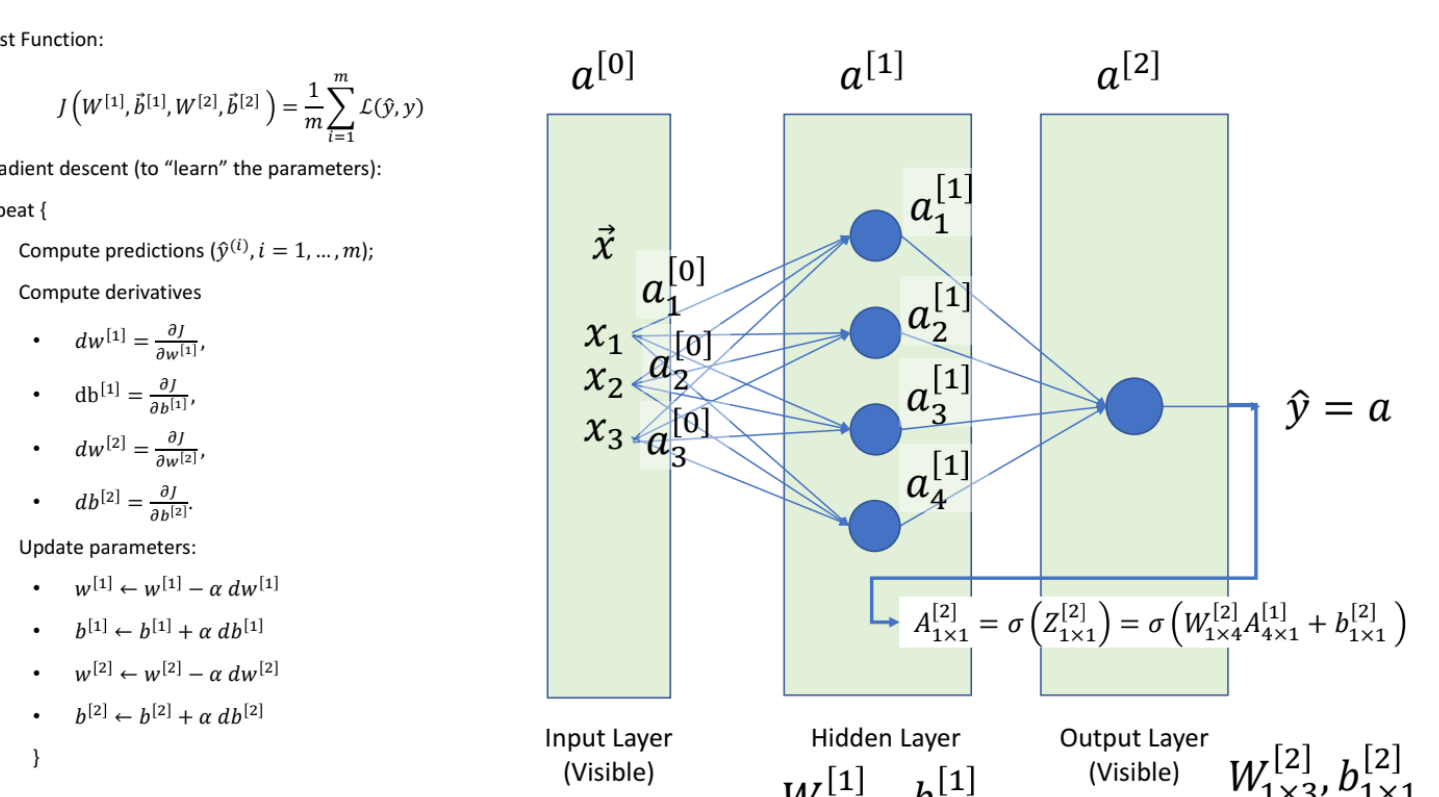
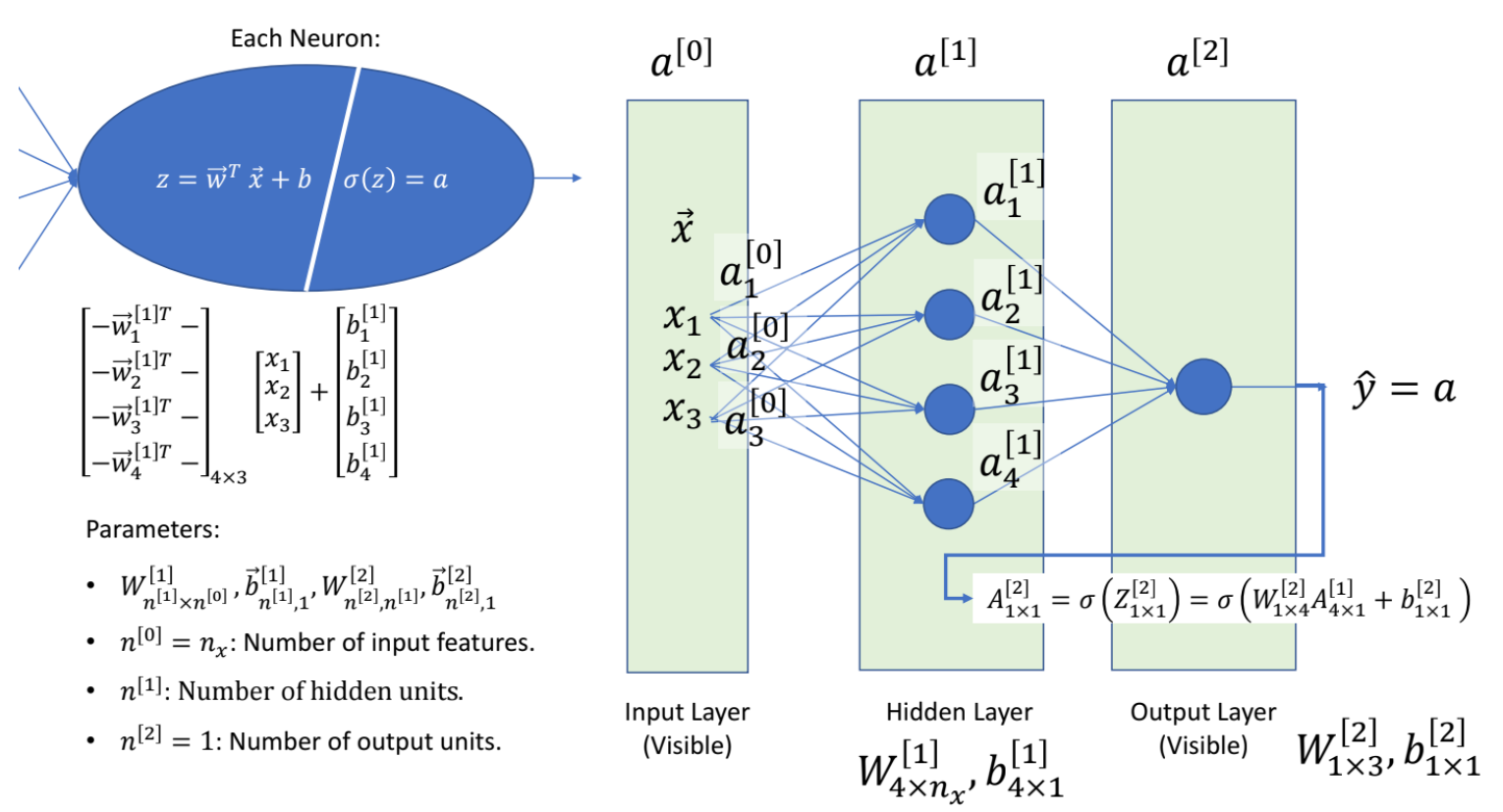
↳ CPU x GPU → Não tudo x GPU
alguns algoritmos de paralelismo que o numpy usa

Isso é o que fazemos no numpy, evitando de usar for loops e sim recorrendo
a alguns vetores mas performáticos

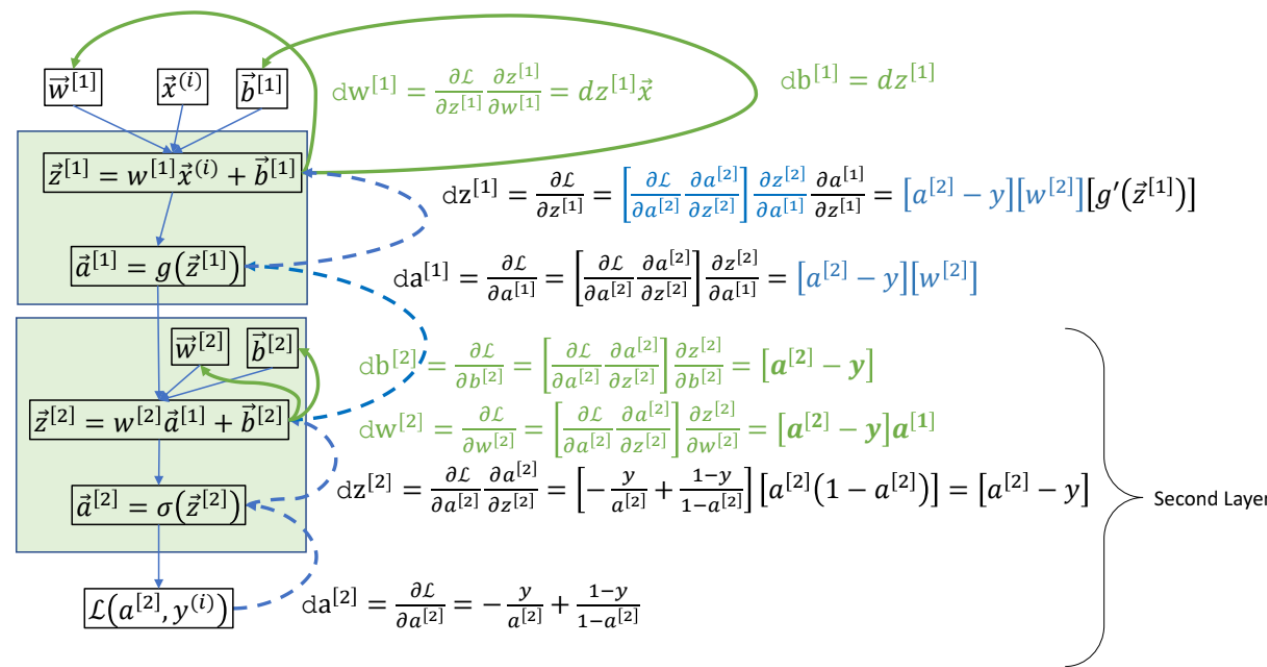
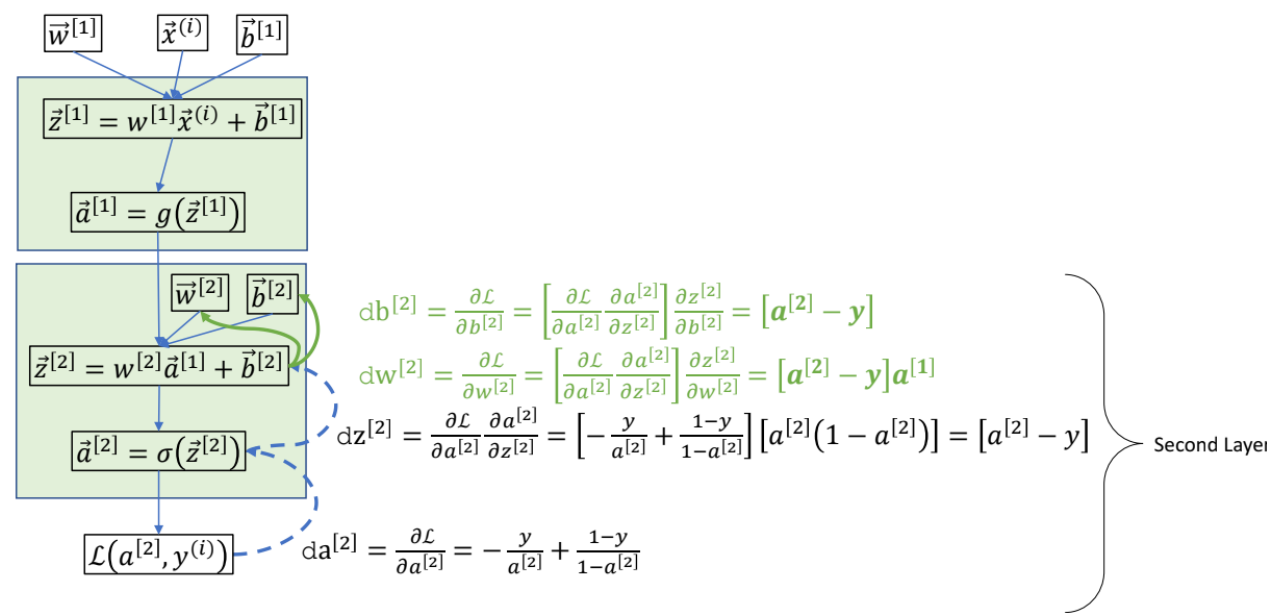
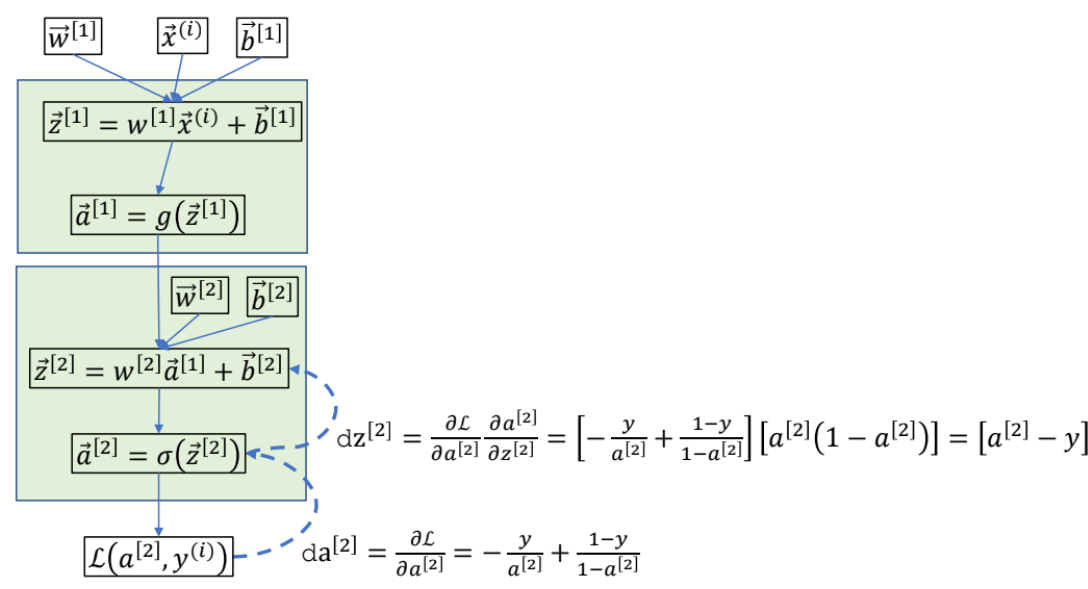
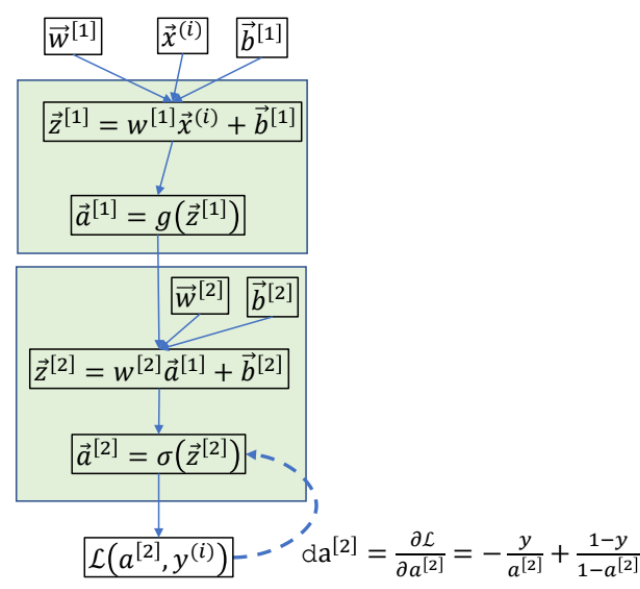
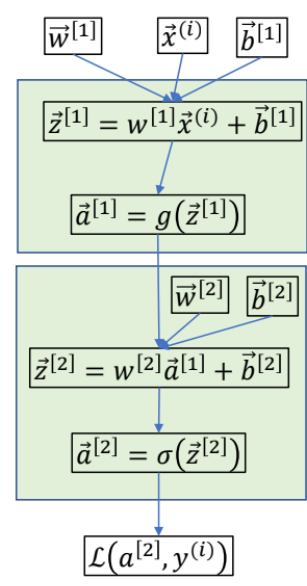
* NN Computing



Gradient descent for NNs (one hidden layer)



$\mathbf{w}^1_{1 \times 2} \mathbf{w}^1_{1 \times 1}$



Six key equations

- $d\mathbf{z}^2 = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^2} \frac{\partial \mathbf{a}^2}{\partial \mathbf{z}^2} = \left[-\frac{\mathbf{y}}{\mathbf{a}^2} + \frac{1-\mathbf{y}}{1-\mathbf{a}^2} \right] [\mathbf{a}^2(1-\mathbf{a}^2)] = [\mathbf{a}^2 - \mathbf{y}]$
- $d\mathbf{w}^2 = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^2} \frac{\partial \mathbf{a}^2}{\partial \mathbf{w}^2} = \left[\frac{\partial \mathcal{L}}{\partial \mathbf{a}^2} \frac{\partial \mathbf{a}^2}{\partial \mathbf{z}^2} \right] \frac{\partial \mathbf{z}^2}{\partial \mathbf{w}^2} = [\mathbf{a}^2 - \mathbf{y}] \mathbf{a}^1$
- $d\mathbf{b}^2 = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^2} \frac{\partial \mathbf{a}^2}{\partial \mathbf{b}^2} = \left[\frac{\partial \mathcal{L}}{\partial \mathbf{a}^2} \frac{\partial \mathbf{a}^2}{\partial \mathbf{z}^2} \right] \frac{\partial \mathbf{z}^2}{\partial \mathbf{b}^2} = [\mathbf{a}^2 - \mathbf{y}]$
- $d\mathbf{z}^1 = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^1} \frac{\partial \mathbf{a}^1}{\partial \mathbf{z}^1} = \left[\frac{\partial \mathcal{L}}{\partial \mathbf{a}^2} \frac{\partial \mathbf{a}^2}{\partial \mathbf{z}^2} \right] \frac{\partial \mathbf{z}^2}{\partial \mathbf{a}^1} = [\mathbf{a}^2 - \mathbf{y}] [\mathbf{g}'(\mathbf{z}^1)]$
- $d\mathbf{w}^1 = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^1} \frac{\partial \mathbf{z}^1}{\partial \mathbf{w}^1} = d\mathbf{z}^1 \mathbf{x} = [\mathbf{a}^2 - \mathbf{y}] [\mathbf{w}^2] [\mathbf{g}'(\mathbf{z}^1)] \mathbf{x}$
- $d\mathbf{b}^1 = d\mathbf{z}^1 = [\mathbf{a}^2 - \mathbf{y}] [\mathbf{w}^2] [\mathbf{g}'(\mathbf{z}^1)]$

For m samples

$$\begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \hat{g}^{(l)}(1) & \hat{g}^{(l)}(2) & \dots & \hat{g}^{(l)}(m) & \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{n^{(l)} \times m} = \begin{bmatrix} - & \hat{w}^{(l)T} & - \\ - & \hat{w}^{(l)T} & - \\ \vdots & \vdots & \vdots \\ - & \hat{w}^{(l)T} & - \end{bmatrix}_{n^{(l)} \times n^{(l)-1}} \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \hat{g}^{(l-1)}(1) & \hat{g}^{(l-1)}(2) & \dots & \hat{g}^{(l-1)}(m) & \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{n^{(l-1)} \times m} + \begin{bmatrix} b^{(l)} \\ \vdots \\ b^{(l)} \end{bmatrix}_{n^{(l)} \times 1}$$

$$Z^{(l)} = W^{(l)T}X + B$$

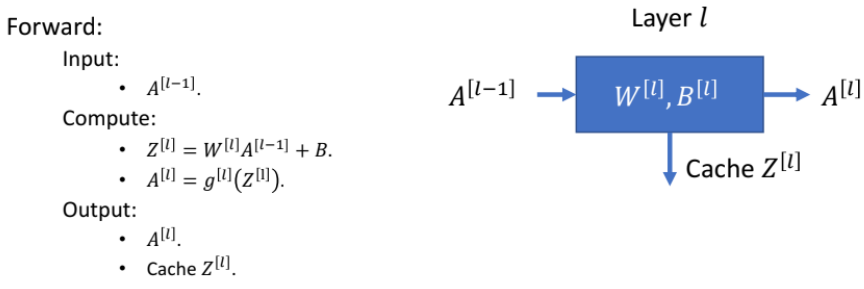
$$A^{(l)} = g^{(l)}(Z) = g^{(l)}(W^{(l)T}X + B)$$

Deep Neural Networks

DL Building Blocks

Forward:

- Input:
 - $A^{(l-1)}$
- Compute:
 - $Z^{(l)} = W^{(l)}A^{(l-1)} + B$
 - $A^{(l)} = g^{(l)}(Z^{(l)})$
- Output:
 - $A^{(l)}$
 - Cache $Z^{(l)}$



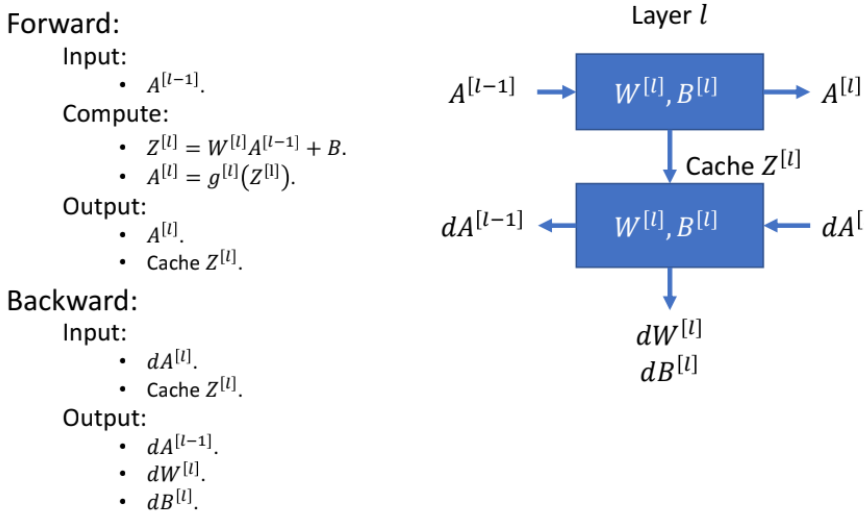
DL Building Blocks

Forward:

- Input:
 - $A^{(l-1)}$
- Compute:
 - $Z^{(l)} = W^{(l)}A^{(l-1)} + B$
 - $A^{(l)} = g^{(l)}(Z^{(l)})$
- Output:
 - $A^{(l)}$
 - Cache $Z^{(l)}$

Backward:

- Input:
 - $dA^{(l)}$
 - Cache $Z^{(l)}$
- Output:
 - $dA^{(l-1)}$
 - $dW^{(l)}$
 - $dB^{(l)}$



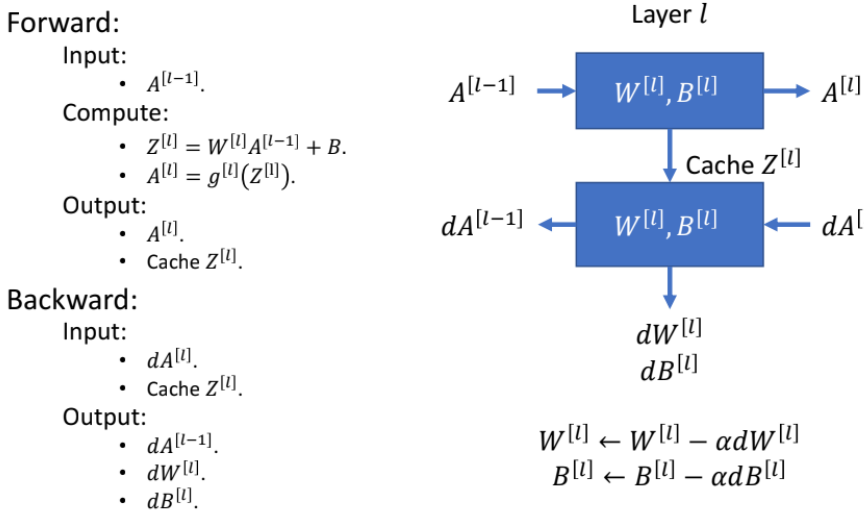
DL Building Blocks

Forward:

- Input:
 - $A^{(l-1)}$
- Compute:
 - $Z^{(l)} = W^{(l)}A^{(l-1)} + B$
 - $A^{(l)} = g^{(l)}(Z^{(l)})$
- Output:
 - $A^{(l)}$
 - Cache $Z^{(l)}$

Backward:

- Input:
 - $dA^{(l)}$
 - Cache $Z^{(l)}$
- Output:
 - $dA^{(l-1)}$
 - $dW^{(l)}$
 - $dB^{(l)}$



Training adjustment

- Parameters:
 - $W^{(l)}, B^{(l)}$
- Hyperparameters:
 - Learning rate α .
 - Number of iterations.
 - Number of hidden units $n^{[1]}, n^{[2]}, \dots$
 - Choice of activation function.

Thank you

$$\begin{aligned} dz &= \frac{dA}{dL} \\ dw &= \frac{dw}{dz} = a(-\phi) \times \\ db &= \phi \end{aligned}$$