# ·DEEP LEARNING
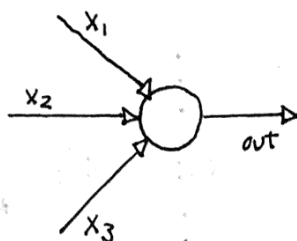
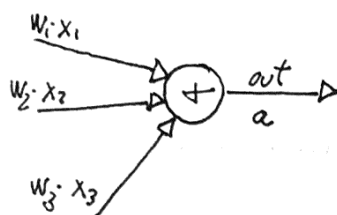## PERCEPTRON
(Neurona Binaria)

$$output = \begin{cases} 0 & if \ W \cdot x + b \leq 0 \\ 1 & if \ W \cdot x + b > 0 \end{cases}$$

$b = bias = -threshold$

$W = weights$

## SIGMOID NEURON

$W_1 \cdot X_1$

$W_2 \cdot X_2$ → out $a$

$W_3 \cdot X_3$

sigmoid function → $\sigma(z) = \dfrac{1}{1 + e^{-z}} \in [0,1]$

$z = \sum_{j} W_j \cdot x_j + b$ (componente a componente)

⇩

$$z = W \cdot x + b$$

## GRADIENT DESCENT

Cost function → $c(w,b) = \dfrac{1}{2n} \sum_{x} \| y(x) - a \|^{2}$ where $\begin{cases} y(x) = test \ vector \ (y=out, x=input) \\ a = neuron \ output \end{cases}$

$\Delta out \approx \sum_{j} \dfrac{\partial out}{\partial w_j} \Delta w_j + \dfrac{\partial out}{\partial b} \Delta b$

$\Delta C \approx \dfrac{\partial C}{\partial v_1} \Delta v_1 + \dfrac{\partial C}{\partial v_2} \Delta v_2$   $(v_1, v_2) \rightarrow direccions$

$v \equiv (w, b)$

$\nabla C = \left( \dfrac{\partial C}{\partial v_1}, \dfrac{\partial C}{\partial v_2} \right)^{T} \longrightarrow \Delta C = \nabla C \cdot \Delta v$

$\Delta v \equiv -\nabla C \cdot \eta$   $\eta = learning \ rate$

SGD → Stocastic Gradient descent

SVM → Support Vector Machine (NOT MACHINE LEARNING ALGORITHM)

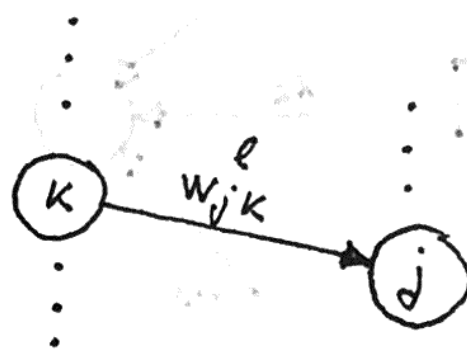Sophisticated algorithm ≤ simple learning + good Training
algorithm                Data

Deep Neural → Networks with many-layers structure (two or more hidden layers)
Networks

Notation → $w_{jk}^{\ell}$ $\begin{cases} \ell \equiv layer \\ K \equiv Neurona \ en \ capa \ (\ell-1) \\ j \equiv Neurona \ en \ capa \ \ell \end{cases}$

$b_j^{\ell} \begin{cases} \ell \to capa \\ j \to Nevron \end{cases} a_j^{\ell}$ [Bias & Activations]



$$a_j^{\ell} = \dagger \left( \sum_K w_{jk}^{\ell} a_K^{\ell-1} + b_j^{\ell} \right) \equiv \dagger \underbrace{\left( w^{\ell} \cdot a^{\ell-1} + b^{\ell} \right)}_{z^{\ell}} \leftarrow \text{MATRIX FORM}$$

↳ weighted input for neurons in layer $\ell$.

## BACKPROPAGATION

Objetive → $\left( \dfrac{\partial C}{\partial w}, \dfrac{\partial C}{\partial b} \right)$

Quadratic Cost → $C = \dfrac{1}{2n} \sum_x \| y(x) - a^L(x) \|^2$ $\begin{cases} n = \text{Total number of training examples} \\ L = \text{Number of layers (Final layer)} \\ x = \text{input} \\ y = \text{output (desired)} \end{cases}$

$$C = \frac{1}{n} \sum_x C_x$$

$$C_x = \frac{1}{2} \| y(x) - a^L(x) \|^2 \quad \triangleleft \text{Average}$$

$$C = C(a_L) \quad \triangleleft \text{Depends on activations}$$

$\delta_j^{\ell} \equiv$ Error in $j^{th}$ nevron in $\ell^{th}$ layer → $\delta_j^{\ell} = \dfrac{\partial C}{\partial z_j^{\ell}}$

## 4 EQUATIONS

### 1.- ERROR IN THE OUTPUT LAYER

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \dagger'(z_j^L) = \underbrace{\nabla_a C \cdot \dagger'(z^L)}_{\substack{\text{MATRIX-BASED} \\ \text{FORM}}} = (a^L - y) \cdot \dagger'(z^L)$$

Substituyendo el gradiente del coste respecto de la activación.

2.- ERROR IN LAYER $L$ REFERRED
      TO ERROR IN LAYER $L+1$    $\rightarrow$   $\delta^{\ell} = \left( (w^{\ell+1})^T \delta^{\ell+1} \right) \bullet f'(z^{\ell})$.

· with 1 an 2, we can compute every error of the net from $L$ backwords.

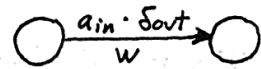3.- RATE OF CHANGE OF THE COST
    WITH RESPECT TO ANY BIAS $\rightarrow$   $\dfrac{\partial C}{\partial b_j^{\ell}} = \delta_j^{\ell}$

4.- RATE OF CHANGE OF THE COST
    WITH RESPECT TO ANY WEIGHT $\rightarrow$ $\dfrac{\partial C}{\partial w_{jk}^{\ell}} = a_k^{\ell-1} \cdot \delta_j^{\ell}$   $\rightarrow$ $\dfrac{\partial C}{\partial w} = a_{in} \cdot \delta_{out}$

"Weights output from low activations
neurons learn slowly (the output neuron
is saturated o111)"

---

PROOFS

1.   $\delta_j^L = \dfrac{\partial C}{\partial z_j^L} = \sum_k \dfrac{\partial C}{\partial a_k^L} \cdot \dfrac{\partial a_k^L}{\partial z_j^L} = \dfrac{\partial C}{\partial a_j^L} \cdot \dfrac{\partial a_j^L}{\partial z_j^L} = \dfrac{\partial C}{\partial a_j^L} f'(z_j^L)$

                                                  $\llcorner a_j^L = f(z_j^L)$

2.   $\delta_j^{\ell} = \dfrac{\partial C}{\partial z_j^{\ell}} = \sum_k \dfrac{\partial C}{\partial z_k^{\ell+1}} \cdot \dfrac{\partial z_k^{\ell+1}}{\partial z_j^{\ell}} = \sum_k \dfrac{\partial z_k^{\ell+1}}{\partial z_j^{\ell}} \cdot \delta_k^{\ell+1} = \sum_k w_{kj}^{\ell+1} \cdot \delta_k^{\ell+1} \cdot f'(z_j^{\ell})$

     $z_k^{\ell+1} = \sum_j w_{kj}^{\ell+1} \cdot a_j^{\ell} + b_k^{\ell+1} = \sum_j w_{kj}^{\ell+1} f(z_j^{\ell}) + b_k^{\ell+1} = w_{kj}^{\ell+1} f'(z_j^{\ell})$

                                                  $\dfrac{\partial}{\partial z_j^{\ell}}$

3.   $\delta_j^{\ell} = \dfrac{\partial C}{\partial z_j^{\ell}} = \dfrac{\partial C}{\partial b_j^{\ell}} \cdot \dfrac{\partial b_j^{\ell}}{\partial z_j^{\ell}} = \dfrac{\partial C}{\partial z_j^{\ell}} \cdot 1$

     $z_k^{\ell} = \sum_j w_{kj}^{\ell} a_j^{\ell-1} + b_k^{\ell} = \sum_j w_{kj}^{\ell} f(z_j^{\ell-1}) + b_k^{\ell} = 1$

                                                  $\dfrac{\partial}{\partial z_k^{\ell}}$

**4.**

$$\delta_j^\ell = \frac{\partial C}{\partial z_j^\ell} = \frac{\partial C}{\partial w_{jk}^\ell} \cdot \frac{\partial \dot{w}_{jk}^\ell}{\partial z_j^\ell} =$$

---

## ALGORITHM

1. INPUT. $x \rightarrow$ set $a^\ell$ for input layers

2. FEEDFORWARD $\rightarrow z^\ell = w^\ell \cdot a^{\ell-1} + b^\ell \quad \forall \, \ell \in (2, L)$

3. OUTPUT ERROR $\rightarrow \delta^L = \nabla_a C \cdot f'(z^L)$

4. BACK PROPAGATE ERROR $\rightarrow \delta^\ell = ((w^{\ell+1})^T \cdot \delta^{\ell+1}) \cdot f'(z^\ell) \quad \forall \, \ell \in (L-1, 1)$.

5. OUTPUT $\rightarrow \dfrac{\partial C}{\partial w_{jk}^\ell} = a_k^{\ell-1} \delta_j^\ell \quad ; \quad \dfrac{\partial C}{\partial b_j^\ell} = \delta_j^\ell$

6. GRADIENT DESCENT (update)

$$\begin{cases} w^\ell \rightarrow w^\ell - \dfrac{\eta}{m} \sum_x \delta^{x,\ell} (a^{x,\ell-1})^T \\[4mm] b^\ell \rightarrow b^\ell - \dfrac{\eta}{m} \sum_x \delta^{x,\ell} \end{cases}$$

Se hace con grupos de vectores de test (mini-batch)

## CROSS-ENTROPY COST FUNCTION

$$C = -\frac{1}{n} \sum_x \left[ y \ln(a) + (1-y) \ln(1-a) \right]$$

Avoids the problem of learning.
  slow-down near saturation.
Better than quadratic cost for using
  with sigmoid neurons.

$$\frac{\partial C}{\partial w_j} = \frac{1}{n} \sum_x x_j \left( \sigma(z) - y \right)$$

↑
error, learning depends on error.

$$\frac{\partial C}{\partial b} = \frac{1}{n} \sum_x \left( \sigma(z) - y \right)$$

Many neurons → $C = -\frac{1}{n} \sum_x \sum_y \left[ y_j \ln a_j^L + (1-y_j) \ln(1-a_j^L) \right]$

## SOFTMAX   (Inciso, solución alternativa a slowdown problem)

* New type of output ~~neuron~~ layer in neuronal network.

$$z_j^\ell = \sum_k w_{jk}^L a_k^{L-1} + b_j^L$$

softmax function → $a_j^L = \dfrac{e^{z_j^L}}{\sum_k e^{z_k^L}}$   $\left[ \sum_j a_j^L = 1 \right]$ ← Can be understood as
                                                                    a probability distribution.

Log-likelihood Cost Function → $C \equiv - \ln a_y^L$ ← Desired output for the network

$$\frac{\partial C}{\partial b_j^L} = a_j^L - y_j$$

$$\frac{\partial C}{\partial w_k^\ell} = a_k^{L-1} \left( a_j^L - y_j \right)$$

Works as sigmoids neurons
+ cross-entropy.

- OVERFITTING → What the network learns, no longer generalizes to the test data.

    It is a mayor problem on modern neuronal networks.

Strategy to prevent: EARLY STOPPING

HOLD-OUT METHOD → To use a validation data set. When the classification

accuracy on that set stops, we stop training.

\* Validation data: Training data that helps us to learn good

hyper-parameters.

→ One of the best ways of reducing overfitting is to increase test data size. ↚


- REGULARIZATION → Techniques to reduce over-fitting.

\* WEIGHT DECAY TECHNIQUE / L2 REGULARIZATION

To add a "regularization term" to the cost function.

- Regularization Cross-Entropy →
$$C = \frac{-1}{n} \sum_{xj} \left[ y_j \ln(a_j^L) + (1-y_j) \ln(1-a_j^L) \right] + \underbrace{\frac{\lambda}{2n} \sum_w w^2}$$

$\lambda > 0$ (Regularization Parameter)

- Regularized Quadratic Cost →
$$C = \frac{1}{2n} \sum_x \| y - a^L \|^2 + \frac{\lambda}{2n} \sum_w w^2$$

- In general →
$$C = C_0 + \frac{\lambda}{2n} \sum_w w^2$$

Compromise between finding small weights and minimizing the
original cost function.

$$\begin{cases} \lambda \downarrow : \text{Minimize original cost function} \\ \lambda \uparrow : \text{Find small weights.} \end{cases}$$

$$\left( \frac{\partial C}{\partial w} = \frac{\partial C_0}{\partial w} + \frac{\lambda}{n} w \quad ; \quad \frac{\partial C}{\partial b} = \frac{\partial C_0}{\partial b} \right) \twoheadleftarrow \text{To Apply SGD}$$

$$w \to \left(1 - \frac{\eta \lambda}{n}\right) w - \eta \frac{\partial C_0}{\partial w} \qquad ; \quad b \to b - \eta \frac{\partial C_0}{\partial b}$$

Weight decay

· OTHER TECHNIQUES

- $L1$ REGULARIZATION $\Rightarrow C = C_0 + \frac{\lambda}{n} \sum_w |w|$

  (It also penalize large weights)

  $$\frac{\partial C}{\partial w} = \frac{\partial C_0}{\partial w} + \frac{\lambda}{n} sgn(w)$$

  $$w' \rightarrow w - \frac{\eta \lambda}{n} sgn(w) - \eta \frac{\partial C_0}{\partial w}$$

- DROPOUT $\Rightarrow$ Modify the network itself

  $\rightarrow$ 1. Randomly (and temporarily) deleting half the hidden neurons.

  2. Apply the backpropagation algorithm

  3. Back to 1.

- ARTIFIAL EXPANDING THE TRAINING DATA.

  Modify the training data to obtain a new one (i.e. rotating the set of images a little amount). It can be applied to many fields like noise recognition (i.e. adding noise or changing velocity of speech).

WEIGHT INITIALIZATION

We initialized the weights with a Gaussian distribution with ($\bar{x} = 0$ & $\sigma = 1$)

($\sigma =$ standard deviation)

⇩

New Approach $\rightarrow \bar{x} = 0$, $\sigma = \frac{1}{\sqrt{n_{in}}}$  ($n_{in} =$ number of input weights.)

## HOW TO CHOOSE NEURAL NETWORK HYPER PARAMETERS?

→ BROAD STRATEGY

- Learning Rate: The only parameter that can be tunned using the cost, because it represents the step in SGD.

- Number of epochs ⟹ <u>Early stopping</u> (automatic)

  ↓

  When the classification accuracy has not improved in quite some time.

- Learning Rate Schedule → Start with a big learning rate and when accuracy is getting worse, reduce it by a factor (2, 10, ...).

- Regularization Parameter → Start with $\lambda = 0$ and tunned $\eta$.

  Put $\lambda = 1.0$ and tunned it in factors of 10.

  Then tunned again $\eta$.

- Minibatch size → Compromise value that maximizes the speed of learning (most rapid improvement in performance).

→ AUTOMATED TECHNIQUES → Grid Search.

  READ → "Practical Bayesian optimization of machine learning algorithm" Jasper Snoek, Hug Larochelle & Ryan Adams.

## OTHER TECHNIQUES

- Hessian Technique → $C(w + \Delta w) \approx C(w) + \nabla C \cdot \Delta w + \frac{1}{2} \Delta w^T H \Delta w$

  $$\Delta w = -H^{-1} \nabla C$$

  $H$ = Hessian Matrix

  In practice → $\Delta w = -\eta H^{-1} \nabla C$

  $$w' \to w + (-\eta H^{-1} \nabla C)$$

  Muy difícil de llevar a la práctica.

# · MOMEMTUM-BASED GRADIENT DESCENT

Introduce velocity $\rightarrow$ v (one per variable)

$$v \rightarrow v' = \mu v - \rho \nabla C \quad (\mu \equiv friction) \quad \mu = 1 \rightarrow \text{NO FRICTION}$$

$$w \rightarrow w' = w + v'$$

momentum coefficient.

## OTHER MODELS OF NEURONS

· tanh neurons $\rightarrow$ $a = \tanh(w \cdot x + b)$

· Rectified linear neurons $\rightarrow$ $a = max(0, wx + b)$