

Introducción a "Statistical Machine Learning V2"
Club de Ciencia de Datos La Paz

Introducción a Redes Neuronales

Perceptrón Multicapa

Rafael Villca Poggian

Abril de 2019



Introducción

- ¿Qué es una Red Neuronal?
- ¿Cómo se modela matemáticamente?
- Componentes del modelo
- ¿Cómo ajustar los parámetros del modelo?

Grafo Computacional

- ¿Por qué usarlo?
- Interpretación
- Ejemplo
- Regresión Logística

Funciones de Activación

- Definición y Uso
- Sigmoide
- Relu

Perceptrón Multicapa

- Modelo
- Forward & Backprop
- Entrenamiento

Introducción

¿Qué es una Red Neuronal?



Conocido formalmente como "Perceptrón Multicapa", es un modelo matemático que nació bajo la inspiración de modelar la codificación de patrones en la estructura neuronal, como respuesta a impulsos externos.

Introducción

¿Cómo se modela matemáticamente?



- ▶ Se puede entender una red neuronal como una regresión logística con capas intermedias, las cuales la convierten en un aproximador universal de funciones, y un selector inteligente de variables.
- ▶ Si antes vimos que una Regresión Logística se modelaba como una composición de funciones derivables, la red neuronal mantiene esta misma estructura, sólo un poco más complicada y con más componentes.



- ▶ \mathbf{W}_l : Matriz de parámetros de transición entre la capa l de la red.
- ▶ \mathbf{b}_l : Vector de Sesgos o bias.
- ▶ \mathbf{A}_l : Activación en de las neuronas visibles en la capa l
- ▶ $\mathbf{g}_l(\mathbf{x})$: Función de Activación para la capa l .

Introducción

¿Cómo ajustar los parámetros del modelo?



- ▶ Debemos dar una pasada por los datos para obtener las activaciones en cada capa de la red neuronal (Forward Propagation).
- ▶ Calculamos la derivada de todos los pasos en reversa (Backpropagation).
- ▶ Usamos los parámetros ya calculados al navegar por el grafo para evaluar los gradientes al recorrer el grafo de manera inversa.

Introducción

¿Cómo ajustar los parámetros del modelo?



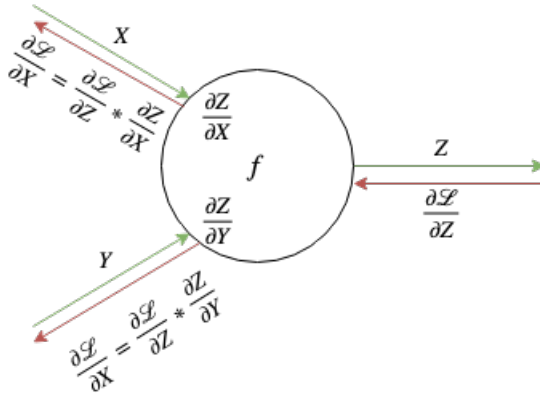
- ▶ Una vez evaluados los gradientes actualizamos los valores de los parámetros bajo la regla del gradient descent para dar un paso hacia los valores óptimos.
- ▶ Para poder almacenar esta información de manera óptima utilizamos un grafo de computación (DAG).

Grafo Computacional

¿Por qué usarlo?

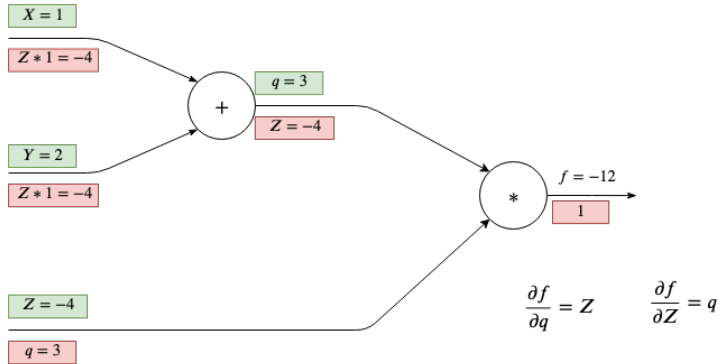


- ▶ Podemos descomponer una función compleja en sus distintas operaciones componentes básicas.
- ▶ De esta forma es mucho más fácil computar derivadas de funciones compuestas.
- ▶ Es gráficamente interpretativo.

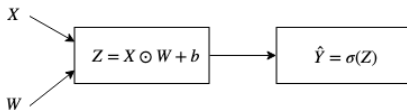
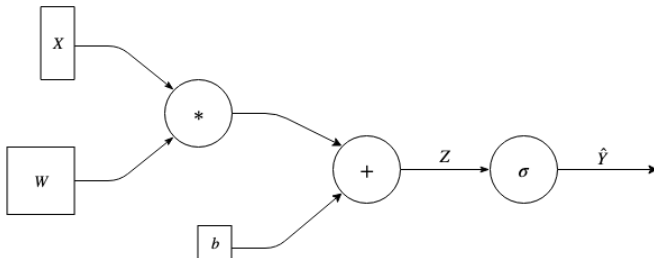




$$f(X, Y, Z) = (X + Y) * Z = q * Z$$



$$\frac{\partial q}{\partial X} = 1 \quad \frac{\partial f}{\partial X} = \frac{\partial f}{\partial q} * \frac{\partial q}{\partial X} \quad , \quad \frac{\partial q}{\partial Y} = 1 \quad \frac{\partial f}{\partial Y} = \frac{\partial f}{\partial q} * \frac{\partial q}{\partial Y}$$



Funciones de Activación

Definición y Uso



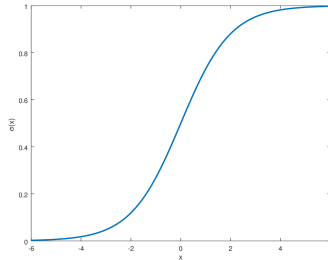
- ▶ Es una función no lineal que se utiliza para aumentar la capacidad de ajuste a los datos de la red neuronal.
- ▶ Si no se usara, la red sería simplemente un clasificador lineal.
- ▶ Es una función

$$g_l(z_l)$$

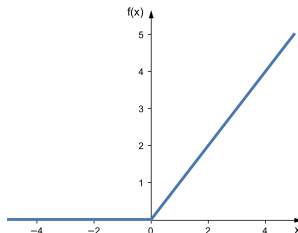
en cada capa oculta.

Funciones de Activación

Sigmoide



- ▶ Utilizada bajo inspiración biológica.
- ▶ Costosa de computar.
- ▶ Los gradientes causan problemas de desvanecimiento.



- ▶ Igualmente no lineal.
- ▶ Puede causar desvanecimiento de gradientes en el peor de los casos.
- ▶ Fácil de computar.



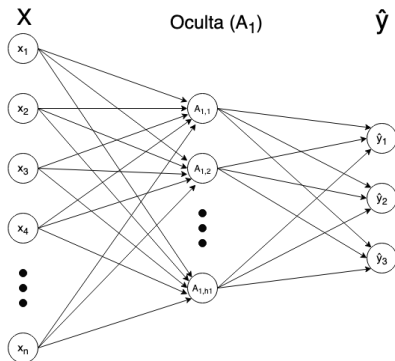
Rectified Linear Unit

$$Relu(Z) = \max(0, Z) \forall Z_j, j : 0, 1, \dots, h_l$$

Dónde h_l es el número de elementos en Z

Derivada

$$\frac{\partial Relu}{\partial Z} = \begin{cases} 1, & \text{si } Z_j > 0 \\ 0, & \text{si } Z_j \leq 0 \end{cases}$$



$$Z_1 = X \odot W_1 + b_1$$

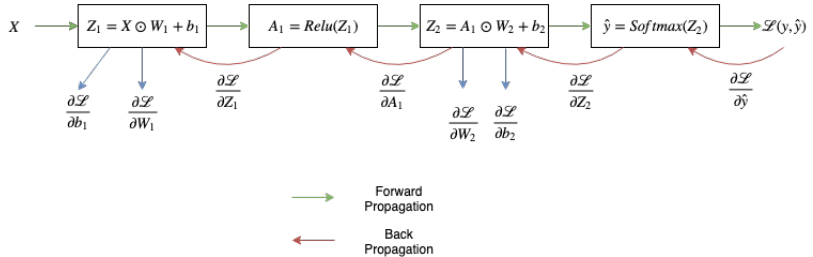
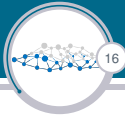
$$A_1 = X \odot Relu(Z_1)$$

$$Z_2 = A_1 \odot W_2 + b_2$$

$$\hat{Y} = softmax(Z_2)$$

Perceptrón Multicapa

Forward & Backprop





Algorithm 1 Neural Network Training

```
1: procedure GRADIENT DESCENT( $X, Y, \text{layer\_shapes}, n\_layers, \alpha$ )
2:    $W \leftarrow [\dots], b \leftarrow [\dots]$ 
3:   for  $l = 1$  to  $n\_layers$  do
4:      $W[l] \leftarrow \text{random}(\text{layer\_shapes}[l-1], \text{layer\_shapes}[l])$ 
5:      $b[l] \leftarrow \text{zeros}(\text{layer\_shapes}[l])$ 
6:   for  $i = 1$  to  $iters$  do
7:      $\hat{Y}, \text{cache} \leftarrow \text{forward\_propagation}(X, W, b)$ 
8:      $\mathcal{L} \leftarrow \text{cross\_entropy}(\hat{Y}, Y)$ 
9:      $\nabla W, \nabla b \leftarrow \text{backpropagation}(X, W, b, \text{cache})$ 
10:     $W \leftarrow W - \alpha * \frac{\partial \mathcal{L}}{\partial W}$ 
11:     $b \leftarrow b - \alpha * \frac{\partial \mathcal{L}}{\partial b}$ 
12:   return predicción
```
