

Redes Neuronales

Carlos Malanche

21 de mayo de 2019

Las redes neuronales son una de las técnicas del aprendizaje estadístico más prometedoras. Esto, veremos poco más adelante, debido a su poder de *aprender* las funciones de transformación de los vectores independientes a espacios linealmente separables.

1. El antecesor de la red

El concepto de una red neuronal no es algo nuevo, desde hace varias décadas se tenía ya el concepto de interacciones entre elementos no lineales que derivaban en resultados sencillos.

2. Derivando el algoritmo de propagación hacia atrás

Supondremos de manera temporal que la capa ℓ consta de α neuronas, mientras que la capa $\ell + 1$ consta de β neuronas. Con esto, no debiera ser muy complicado notar que el vector saliente de las neuronas de la capa ℓ vive en $\mathbb{R}^{\alpha \times 1}$. Siguiendo esa lógica, el vector saliente de la capa $\ell + 1$ debe vivir en $\mathbb{R}^{\beta \times 1}$. Con ello, y por pura convención, definimos la matriz de pesos que convierte al vector $\mathbf{x}^{(\ell)}$ en el vector $\mathbf{x}^{(\ell+1)}$ como la matriz $\mathbf{W}^{(\ell+1)} \in \mathbb{R}^{\alpha \times \beta}$, y el vector de sesgos $\mathbf{b}^{(\ell+1)} \in \mathbb{R}^{\beta \times 1}$ de modo que

$$\mathbf{x}^{(\ell+1)} = \phi((\mathbf{W}^{(\ell+1)})^T \mathbf{x}^{(\ell)} + \mathbf{b}^{(\ell+1)})$$

Donde la función de activación ϕ ya se considera como una función vectorizada que se aplica a todas las entradas del tensor al que se aplique.

Definiremos dos variables auxiliares que nos ayudarán a calcular el gradiente por capa, primero:

$$\mathbf{z}^{(\ell+1)} := (\mathbf{W}^{(\ell+1)})^T \mathbf{x}^{(\ell)} + \mathbf{b}^{(\ell+1)} \quad (1)$$

Que no es más que meter una variable en medio de las $\mathbf{x}^{(\ell)}$, es decir, $\mathbf{x}^{(\ell)} = \phi(\mathbf{z}^{(\ell)})$

La segunda variable que vamos a introducir es un poco más compleja, y es la derivada de la función de costo respecto a esta misma $\mathbf{z}^{(\ell)}$ para una capa en específico. Vamos a llamar a esta variable *delta*.

$$\delta^{(\ell)} := \frac{\partial \mathcal{L}_n}{\partial \mathbf{z}^{(\ell)}} \quad (2)$$

En donde denoto por \mathcal{L}_n la contribución a la función de costo correspondiente a la medición número n . La variable \mathbf{z} hace la operación conocida como *forward pass*; Dada una $\mathbf{z}^{(\ell)}$, es muy fácil calcular la $\mathbf{z}^{(\ell+1)}$, con lo que estos valores se *propagan* hacia adelante de la red. Lo contrario pasa con la *delta*, que dada una $\delta^{(\ell+1)}$, es sencillo calcular una $\delta^{(\ell)}$. Esperen. Creo que eso aún no ha quedado claro.

Notemos pues que la $\mathbf{z}^{(\ell+1)}$ tiene una dependencia de la $\mathbf{z}^{(\ell)}$. Recordemos que la función de costo es una composición de funciones masomenos del estilo

$$\mathcal{L}_n = \frac{1}{2N} (y - \mathbf{x}^{(L+1)} \circ \mathbf{x}^{(L)} \circ \dots \circ \mathbf{x}^{(0)})^2$$

De aquí se sigue que si queremos derivar la función de costo respecto a una $\mathbf{z}^{(\ell)}$, podemos utilizar la regla de la cadena para hacer aparecer la primer función inmediata que a su vez depende de la misma $\mathbf{z}^{(\ell)}$:

$$\delta^{(\ell)} := \frac{\partial \mathcal{L}_n}{\partial \mathbf{z}^{(\ell)}} = \frac{\partial \mathbf{z}^{(\ell+1)}}{\partial \mathbf{z}^{(\ell)}} \frac{\partial \mathcal{L}_n}{\partial \mathbf{z}^{(\ell+1)}} = \frac{\partial \mathbf{z}^{(\ell+1)}}{\partial \mathbf{z}^{(\ell)}} \delta^{(\ell+1)} \quad (3)$$

Ahora sí se ha vuelto más evidente la dependencia. La derivada que nos falta calcular no es tan complicada realmente

$$\frac{\partial \mathbf{z}^{(\ell+1)}}{\partial \mathbf{z}^{(\ell)}} = \frac{\partial}{\partial \mathbf{z}^{(\ell)}} ((\mathbf{W}^{(\ell+1)})^T \phi(\mathbf{z}^{(\ell)}) + \mathbf{b}^{(\ell+1)}) = \text{diag}(\phi'(\mathbf{z}^{(\ell)})) \mathbf{W}^{(\ell+1)}$$

Con estas dos variables que hemos definido, calcular los gradientes se vuelve *casi* trivial. Primer, el vector de sesgos

$$\frac{\partial \mathcal{L}_n}{\partial \mathbf{b}^{(\ell)}} = \frac{\partial \mathbf{z}^{(\ell)}}{\partial \mathbf{b}^{(\ell)}} \frac{\partial \mathcal{L}_n}{\partial \mathbf{z}^{(\ell)}} = \mathbf{I} \boldsymbol{\delta}^{(\ell)} = \boldsymbol{\delta}^{(\ell)} \quad (4)$$

La derivada de la matriz de pesos la haremos entrada por entrada, pues la vida se nos puede complicar mucho si lo intentamos hacer matricialmente (aunque la expresión final será matricial).

$$\frac{\partial \mathcal{L}_n}{\partial \mathbf{W}_{i,j}^{(\ell+1)}} = \underbrace{\frac{\partial \mathbf{z}^{(\ell+1)}}{\partial \mathbf{W}_{i,j}^{(\ell+1)}} \frac{\partial \mathcal{L}_n}{\partial \mathbf{z}^{(\ell+1)}}}_{\text{producto interno}} = \sum_k \left(\frac{\partial}{\partial \mathbf{W}_{i,j}^{(\ell+1)}} \sum_s \mathbf{W}_{s,k}^{(\ell+1)} \mathbf{x}_s^{(\ell)} \right) \boldsymbol{\delta}_k^{(\ell+1)} = \sum_k \mathbf{x}_i^{(\ell)} \boldsymbol{\delta}_k^{(\ell+1)} \delta_{jk} = \mathbf{x}_i^{(\ell)} \boldsymbol{\delta}_j^{(\ell+1)}$$

Que escrito de forma matricial

$$\frac{\partial \mathcal{L}_n}{\partial \mathbf{W}^{(\ell+1)}} = \mathbf{x}^{(\ell)} (\boldsymbol{\delta}^{(\ell+1)})^T \quad (5)$$

Con lo que nuestro descenso de gradiente para cada capa queda escrito como:

$${}_{it+1} \mathbf{b}^{(\ell+1)} = {}_{it} \mathbf{b}^{(\ell+1)} - \mu ({}_{it} \boldsymbol{\delta}) \quad (6)$$

$${}_{it+1} \mathbf{W}^{(\ell+1)} = {}_{it} \mathbf{W}^{(\ell+1)} - \mu \mathbf{x}^{(\ell)} ({}_{it} \boldsymbol{\delta}^{(\ell+1)})^T \quad (7)$$

Donde ya abusé un poco de la notación para colocar el índice de iteración como subíndice que precede.

3. Regularizador de las redes neuronales