

Tarea 5: Problema 1

Victor Manuel Gómez Espinosa

30 de noviembre de 2020

1 a)

Recordando de la Tarea 1, haciendo backpropagation hacia la capa de salida $a_2 = \hat{y}$, con función softmax $\hat{y}_j = \text{softmax}(z_{2,j})$, donde $z_2 = U'a_1$ (a_1 es la capa oculta), tenemos lo siguiente: $dz = \frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_2} = (\hat{y} - y)$.

Ahora si para este caso sabemos que $z_{2,j} = u'_{w_j} v_{w_i}$, y su derivada parcial es $\frac{\partial z_{2,j}}{\partial v_{w_i}} = u'_{w_j}$

o $\frac{\partial z_2}{\partial v_{w_i}} = U'$, entonces $\frac{\partial L}{\partial v_{w_i}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_2} \frac{\partial z_2}{\partial v_{w_i}} = U (\hat{y} - y)$.

1 b)

De forma similar, su derivada parcial es $\frac{\partial z_{2,j}}{\partial u_{w_j}} = v_{w_i}$ o $\frac{\partial z_2}{\partial U} = v_{w_i}$, entonces

$$\frac{\partial L}{\partial U} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_2} \frac{\partial z_2}{\partial U} = v_{w_i} (\hat{y} - y)'$$

1 c)

Si ahora se tiene la siguiente función de costo:

$$L(v_{w_i}, u_{w_j}) = -\ln(\sigma(u'_{w_j} v_{w_i})) - \sum_{k=1}^K \ln(\sigma(-u'_{w_k} v_{w_i})), \text{ o con un cambio de variable:}$$

$L(g) = -\ln(g) - \sum_{k=1}^K \ln(g)$, $g(f) = \sigma(f)$, $f(v_{w_i}, u_{w_j}) = u'_{w_j} v_{w_i}$ para las palabras en el contexto, o $f(v_{w_i}, u_{w_r}) = -u'_{w_r} v_{w_i}$ si no lo están, y sus respectivas derivadas son:

$$L' = -\frac{g'}{g} - \sum \frac{g'}{g}, \quad g' = g(1-g)f' \quad \text{donde} \quad f' = \frac{\partial f}{\partial v_{w_i}} = u'_{w_j} \quad \text{o} \quad f' = \frac{\partial f}{\partial u_{w_j}}, \text{ en esta ultima}$$

hay dos posibles casos, donde la palabra esta en el contexto o no lo está, por lo tanto

$$f' = \frac{\partial f}{\partial u_{w_j=j}} = v_{w_i} \quad \text{o} \quad f' = \frac{\partial f}{\partial u_{w_j \neq j}} = 0$$

$$\text{Entonces: } L' = -(1 - \sigma(f))f' - \sum (1 - \sigma(f))f',$$

Para la parte derecha hay dos opciones, si las palabras están en el contexto o no lo estan y sabiendo que $1 - \sigma(-u'_{w_k} v_{w_i}) = \sigma(u'_{w_k} v_{w_i})$:

$$L' = \begin{cases} (\sigma(u'_{w_j} v_{w_i}) - 1)f' + \sum_{k=1}^k (\sigma(u'_{w_k} v_{w_i}) - 1)f' & , si, k = j \\ (\sigma(u'_{w_k} v_{w_i}))f' + \sum_{k=1}^k (\sigma(u'_{w_k} v_{w_i}))f' & , si, k \neq j \end{cases}$$

Y, por lo tanto:

$$\frac{\partial L}{\partial v_{w_i}} = \begin{cases} \sum_{k=1}^{k-1} (\sigma(u'_{w_k} v_{w_i}) - 1)u'_{w_j} & , si, k = j \\ \sum_{k=1}^{k-1} (\sigma(u'_{w_k} v_{w_i}))u'_{w_j} & , si, k \neq j \end{cases}$$

y

$$\frac{\partial L}{\partial u'_{w_j}} = \begin{cases} (\sigma(u'_{w_j} v_{w_i}) - 1)v_{w_i} & , si, k = j \\ (\sigma(u'_{w_j} v_{w_i}))v_{w_i} & , si, k \neq j \end{cases}$$

1 d)

Esta ultima función de costo es mucho mas eficiente ya que reduce el problema multinomial (softmax) a uno binario. Por ejemplo, si se tiene un vocabulario de 10,000 palabras, cada que se llame a la función softmax se tendría que hacer la sumatoria sobre ese tamaño enorme del vocabulario, al contrario, en este ultimo caso, se reduce a un problema binario, es decir, de regresión logística, con una ventana de palabras fuera de contexto de mucho menor tamaño.