



deeplearning.ai

Multi-class classification

Trying a softmax classifier

Understanding softmax

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ -3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

softmax라는 이름은 hard max의 반대라는 의미에서 가져온 것인데 hard max는 벡터 z 를 보고 가장 큰 element의 자리만 1로 하고 나머지는 0

"hard max"

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

반면 softmax는 z 에서 이 확률 값으로 가는 more gentle mapping이다.

$$a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.742 \\ 0.042 \\ 0.052 \\ 0.114 \end{bmatrix}$$

Softmax regression generalizes logistic regression to C classes.

If $C=2$, softmax reduces to logistic regression.

$$a^{[L]} = \begin{bmatrix} 0.742 \\ 0.154 \end{bmatrix}$$

→ redundant computation

이제 softmax output layer로 neural network를 어떻게 학습하는지 알아보자. 그러려면 학습에 사용할 loss function을 정의해야된다.

Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{ - cat}$$

$y_2 = 1$
 $y_1 = y_3 = y_4 = 0$

$$a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \leftarrow \text{not good classifier,}$$

이 때 loss function은 아래와 같이 정의된다.

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j$$

이게 의미하는 것은 learning algorithm이 $\mathcal{L}(\hat{y}, y)$ 를 작게 만들려고 하고(왜냐하면 gradient descent가 트레이닝 데이터셋 상에서 loss를 작게하려고 하기 때문) 결과적으로 y_{hat_2} 즉, ground truth가 1인 y_2 에 해당하는 확률을 최대한 크게 만드는 것이다.

maximum likelihood estimation의 형태임

위에서 본 것은 single training example에 대한 loss이다. entire training set에 대한 cost function J의 경우는 어떻게?

$$J(W^{[L]}, b^{[L]}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

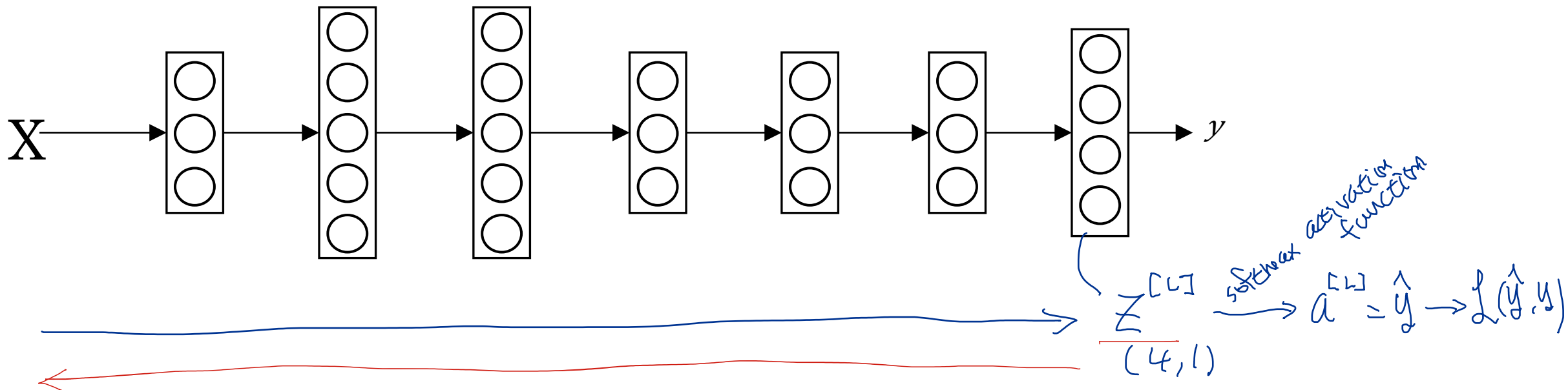
small
make \hat{y}_2 big.

$$= -y_2 \log \hat{y}_2$$
$$= -\log \hat{y}_2$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$
$$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix} (4, m)$$

$$\hat{Y} = [\hat{y}^{(1)}, \dots, \hat{y}^{(m)}]$$
$$= \begin{bmatrix} 0.3 & \dots \\ 0.2 & \dots \\ 0.1 & \dots \\ 0.4 & \dots \end{bmatrix} (4, m)$$

Summary of softmax classifier



Backprop : $\frac{dz^{[L]}_{(4,1)}}{dz^{[L]}_{(4,1)}} = \hat{y} - y$

$\frac{\partial J}{\partial z^{[L]}}$

(partial derivative of the cost function with respect to $z^{[L]}$)

gradient descent를 위한 back propagation 과정에서 key step은 last layer에서 dz 에 대한 derivative를 구하는 것이다.