

Week3 - Neural Network

Freitag, 17. Juli 2020 14:21

Standard Notation:

Standard notations for Deep Learning

This document has the purpose of discussing a new standard for deep learning mathematical notations.

1 Neural Networks Notations.

General comments:

· superscript (i) will denote the i^{th} training example while superscript [l] will denote the l^{th} layer

Sizes:

· m : number of examples in the dataset

· n_x : input size

· n_y : output size (or number of classes)

· $n_h^{[l]}$: number of hidden units of the l^{th} layer

In a for loop, it is possible to denote $n_x = n_h^{[0]}$ and $n_y = n_h^{[number\ of\ layers + 1]}$.

· L : number of layers in the network.

Objects:

· $X \in \mathbb{R}^{n_x \times m}$ is the input matrix

· $x^{(i)} \in \mathbb{R}^{n_x}$ is the i^{th} example represented as a column vector

· $Y \in \mathbb{R}^{n_y \times m}$ is the label matrix

· $y^{(i)} \in \mathbb{R}^{n_y}$ is the output label for the i^{th} example

· $W^{[l]} \in \mathbb{R}^{number\ of\ units\ in\ next\ layer \times number\ of\ units\ in\ the\ previous\ layer}$ is the weight matrix, superscript [l] indicates the layer

· $b^{[l]} \in \mathbb{R}^{number\ of\ units\ in\ next\ layer}$ is the bias vector in the l^{th} layer

· $\hat{y} \in \mathbb{R}^{n_y}$ is the predicted output vector. It can also be denoted $a^{[L]}$ where L is the number of layers in the network.

Common forward propagation equation examples:

$a = g^{[l]}(W_x x^{(i)} + b_1) = g^{[l]}(z_1)$ where $g^{[l]}$ denotes the l^{th} layer activation function

$\hat{y}^{(i)} = softmax(W_h h + b_2)$

· General Activation Formula: $a_j^{[l]} = g^{[l]}(\sum_k w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]}) = g^{[l]}(z_j^{[l]})$

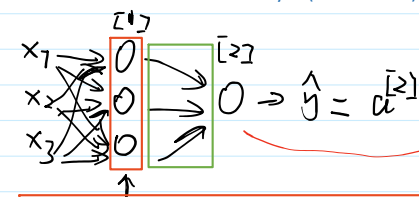
· $J(x, W, b, y)$ or $J(\hat{y}, y)$ denote the cost function.

Examples of cost function:

· $J_{CE}(\hat{y}, y) = -\sum_{i=0}^m y^{(i)} \log \hat{y}^{(i)}$

· $J_1(\hat{y}, y) = \sum_{i=0}^m |y^{(i)} - \hat{y}^{(i)}|$

Structure of NN with one hidden layer (shallow NN)



$$w^{[1]} = z^{[1]} = w^{[1]} x + b^{[1]}$$

$$a^{[1]} = G(z^{[1]})$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix}$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

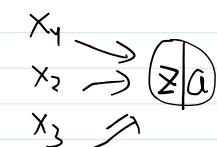
$$a^{[2]} = G(z^{[2]})$$

$$\hat{y} \rightarrow L(a^{[2]}, y)$$

$$\Rightarrow \begin{matrix} w^{[1]} & 3 \times 3 \\ b^{[1]} & 3 \times 1 \\ 3 \text{ node} & * & 3 \text{ inputs} \end{matrix}$$

Unit number of current layer * unit number of previous layer

What happens in one node / Mathematical representation of one node:



$$z = w^T x + b$$

$$a = G(z)$$

Vectorized:

$$\begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} -b_1^{[1]} \\ b_2^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \end{bmatrix}$$

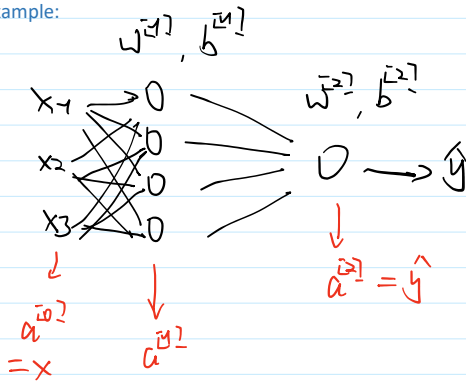
$a^{[l]}$ z -node in layer l -layer

\Rightarrow happens at each node

$$u = u(z)$$

$$\underbrace{\begin{bmatrix} w_1^{[1]} \\ w_2^{[1]} \\ \vdots \\ w_4^{[1]} \end{bmatrix}}_{\substack{w^{[1]} \\ 4 \times 3}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{\substack{x \\ 3 \times 1}} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}}_{\substack{b^{[1]} \\ 4 \times 1}} = \underbrace{\begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{bmatrix}}_{\substack{z^{[1]} \\ 4 \times 1}} \xrightarrow{\text{activation}} a^{[1]} \quad 4 \times 1$$

Example:



Given input $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

$$z^{[1]} = \overset{4 \times 1}{w^{[1]}} \overset{3 \times 1}{x} + \overset{4 \times 1}{b^{[1]}}$$

$$a^{[1]} = G(z^{[1]})$$

$$z^{[2]} = \overset{1 \times 1}{w^{[2]}} \overset{4 \times 1}{a^{[1]}} + \overset{1 \times 1}{b^{[2]}}$$

$$a^{[2]} = \hat{y} = G(z^{[2]})$$

For m training examples

$$\begin{aligned} x^{(1)} &\rightarrow a^{[1]}(1) \\ x^{(2)} &\rightarrow a^{[1]}(2) \\ x^{(3)} &\rightarrow a^{[1]}(3) \\ x^{(m)} &\rightarrow a^{[1]}(m) \end{aligned}$$

$$X = \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & 1 \end{bmatrix} \quad (n \times m)$$

$$z^{[1]} = \begin{bmatrix} z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \end{bmatrix}$$

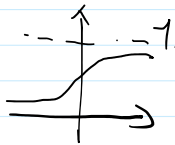
$$a^{[1]} = \begin{bmatrix} a^{[1]}(1) & \dots & a^{[1]}(m) \end{bmatrix} \quad \text{units}$$

m examples

Different Activation Functions:

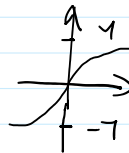
1. Sigmoid

$$\frac{1}{1 + e^{-z}}$$



2. Tanh

$$g = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

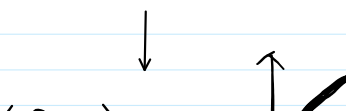


Mean of tanh == 0 and the dataset has usually also mean of 0 (standardized), increases the process; However the output layer has to be sigmoid because (0,1)

Yet, one problem for tanh(z):

When $z \rightarrow \text{infinite}$, gradient really small which slow down gradient descent.

3. ReLU



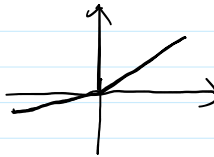
3. ReLU

$$a = \max(0, z)$$



4. Leaky ReLU

$$a = \max(0.01z, z)$$



Default setting:

For output layer: sigmoid

For hidden layer: ReLU

Gradient of different activation function:

Sigmoid:

$$g'(z) = g(z)(1-g(z)) = a(1-a)$$

Tanh:

$$g'(z) = 1 - (\tanh(z))^2 = 1 - a^2$$

ReLU

$$g'(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$$

Leaky ReLU:

$$g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases}$$

Implementation of 2 layer NN:

$$\text{Parameters} = \begin{matrix} (n^{[1]}, n^{[0]}) & (n^{[1]}, n^{[0]}) & (n^{[2]}, n^{[1]}) & (n^{[2]}, n^{[1]}) \\ w^{[1]} & b^{[1]} & w^{[2]} & b^{[2]} \end{matrix} \quad \left. \begin{array}{l} n^{[0]} = n^{[0]} \text{ input} \\ n^{[1]} - \text{number of hidden units} \\ n^{[2]} - \text{number of output units} \end{array} \right\}$$

$$\text{cost: } J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}, y)$$

Repeat i iterations:

$$\text{computer prediction } \hat{y} = A^{[2]}$$

$$d_z^{[2]} = A^{[2]} - y \rightarrow (n^{[2]}, m)$$

$$dw^{[2]} = \frac{1}{m} A^{[1]} d_z^{[2]T} \rightarrow (n^{[2]}, n^{[1]})$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(d_z^{[2]}, \text{axis}=1, \text{keepdim=True}) \rightarrow (n^{[2]}, 1)$$

$$d_z^{[1]} = \underbrace{w^{[2]T} d_z^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[1]}, m)} \rightarrow (n^{[1]}, m)$$

$$dW^{(L)} = \frac{1}{m} \times dZ^{(L)T} \rightarrow (\eta^{(L)}, \gamma)$$

How to implement NN

Take out from the programming assignment:

Random Initilization

The general guidance to build a NN:

1. Define the neural network structure (# of input units, # of hidden units, etc).
2. Initialize the model's parameters
3. Loop:
 - Implement forward propagation
 - Compute cost
 - Implement backward propagation to get the gradients
 - Update parameters (gradient descent)

After all these, the model can be used to predict.

Usually build helper functions to compute and then merge them into one function / `nn_model()`. Once you've built `nn_model()` and learnt the right parameters, you can make predictions on new data calling the `nn_model()`.

During the implementation, i personally think paying attention to the demensions of each matrix is extremly helpful.

Last but not least, the design of hidden layer size was observed with different numbers of units in the hidden layer: 1,2,3,4,5,20,50. The result shows that more than 5 units in hidden layer doesnt helps to increase the accuracy, because the model is overfitted to the training set, which wont perform well on the other datasets.