

Course 1: Neural Networks and Deep Learning

Content:

Week 1: introduction to deep learning

Week 2: Neural Networks Basics

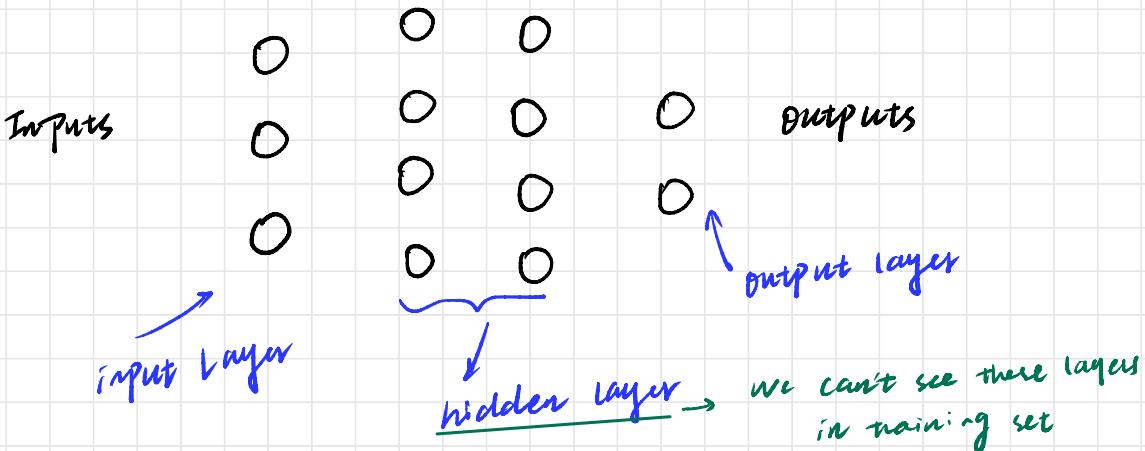
Week 3: Shallow Neural Networks

Week 4: Deep Neural Networks

Coursera

Deep learning Specialization





Structured vs. unstructured data

things that has a defined meaning. e.g. Price, age

e.g. Pixel, raw audio, text

RNN: rectified linear unit

CNN: convolutional neural network \rightarrow useful in computer vision

RNN: Recurrent neural network \rightarrow useful in speech recognition

Standard NN: useful in structured data

Sigmoid: $\hat{y} = P(y=1|x)$ $y \in [0,1]$ prob.

$$\hat{y} = \text{sigmoid}(w^T x + b) = \sigma(w^T x + b)$$

$$\text{sigmoid}(z) = 1/(1+e^{-z}), \quad \sigma'(z) = \sigma(z)(1-\sigma(z))$$

loss function:

$$\text{square root error: } L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

↓ lead us to optimization prob. which is non convex

$$\text{we use: } L(\hat{y}^{(i)}, y^{(i)}) = - (y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

① if $y^{(i)} = 1$, $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$

$\log(\hat{y}^{(i)})$ & $\hat{y}^{(i)}$ close to 1.

② if $y^{(i)} = 0$, $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1-\hat{y}^{(i)})$

$\log(1-\hat{y}^{(i)})$ & $\hat{y}^{(i)}$ close to 0.

cost function:

↓

Average of loss function of entire training set

Goal: find w^*, b s.t minimize overall cost

$$\begin{aligned} J(w, b) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})] \end{aligned}$$

* loss func: measures how well the model is doing on single

training example

* cost func: measures how well parameters w^*, b are doing on entire training set

Gradient Descent

Predict w & b that minimizes cost func

initialize w & b to 0, 0 or to random value

↙
logistic regression

$$w := w - \alpha \underbrace{\frac{\partial J(w, b)}{\partial w}}_{\text{d}w}, \quad b := b - \alpha \underbrace{\frac{\partial J(w, b)}{\partial b}}_{\text{d}b}$$

In logistic regression:

$$\hat{y}^{(i)} = \delta(w^T x^{(i)} + b) \quad \text{where } \delta(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$$

$$\text{single: } z = w^T x + b, \quad \hat{y} = a = \delta(z)$$

$$L(a, y) = -y \log(a) + (1-y) \log(1-a)$$

$$m : \quad z^{(i)} = w^T x^{(i)} + b, \quad \hat{y}^{(i)} = a^{(i)} = \delta(z^{(i)})$$

$$J(w, b) = \frac{1}{m} \sum_i \hat{y}^{(i)}, y^{(i)}$$

$$= -\frac{1}{m} \sum_i [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

$$\text{d}w_i = \frac{1}{m} \sum_i x_i^{(i)} (a^{(i)} - y^{(i)})$$

$$\text{d}b = \frac{1}{m} \sum_i (a^{(i)} - y^{(i)})$$

Vectorization

in coding → avoid for loops.
in python

$$x : [N_x, m]$$

$$\text{matrix } Y : [N_y, m]$$

$$\text{Compute: } [z_1, z_2, \dots, z_m] = \underbrace{w^T x + [b, \dots, b, b]}_{\text{some constants}}$$

Python:

$$(1, m) \quad z = \text{np. dot}(w.T, x) + b$$

$$(1, m) \quad A = \text{sigmoid}(z)$$

$$(1, m) \quad dA = A - Y$$

$$(N_x, 1) \quad dw = \frac{1}{m} * \text{np. dot}(X, dA.T) \quad | \quad \text{vec.} \rightarrow \text{bias computation}$$
$$db = \frac{1}{m} * \text{np. sum}(dA)$$

$$w = w - \text{alpha} * dw \quad | \quad \text{update parameters}$$

$$b = b - \text{alpha} * db$$

↓
Coding part refer to lab assignments.

Single Examples

) shape

$X \in (N_x, 1)$ input with N_x features.

$w^{(1)} \in (N_1, N_x)$ matrix of 1st hidden layer

$b^{(1)} \in (N_1, 1)$ bias term of hidden layer

$z^{(1)} \in (N_1, 1)$ result of $z^{(1)} = w^{(1)}x + b$

$a^{(1)} \in (N_1, 1)$ result of $a^{(1)} = \text{sigmoid}(z^{(1)})$, act. of hid. layer

$w^{(2)} \in (1, N_1)$, $b^{(2)} \in (1, 1)$, $z^{(2)} \in (1, 1)$, $a^{(2)} \in (1, 1)$

Multiple Examples

$X \rightarrow (N_x, m)$, $z^{(1)} \rightarrow (N_1, m)$, $A^{(1)} \rightarrow (N_1, m)$

Activation function

$a = \frac{1}{1 + e^{-z}}$ \rightarrow derivative

* Sigmoid $a = 1 / (1 + e^{-z})$ output value bit 0,1

tanh $a = (e^z - e^{-z}) / (e^z + e^{-z})$ works better than sigmoid

$\begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$ \downarrow $1 - a^2$ for hidden unit w.r.t the mean of its output closer to 0.

* ReLU $a = \max(0, z)$ most widely used

Leaky ReLU $a = \max(0, 0.01z, z)$ improved ReLU

$\begin{cases} 0.01z & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$ \downarrow this can be a parameter

gradient descent

x	(n_x, m)			variables
$z^{(1)}$	(n_1, m)	$z^{(2)}$	$(n_2, 1)$	
$A^{(1)}$	(n_1, m)	$A^{(2)}$	$(n_2, 1)$	parameters
$w^{(1)}$	(n_1, n_x)	$w^{(2)}$	(n_2, n_1)	
$b^{(1)}$	$(n_1, 1)$	$b^{(2)}$	$(n_2, 1)$	

Deep N.N.

$n^{(0)}$ # of neurons in input layer

$n^{(l)}$. . . in l th layer

$w^{(l)}$ weights of l th layer shape $(n^{(l)}, n^{(l-1)})$

$b^{(l)}$ bias term of l th layer shape $(n^{(l)}, 1)$

$z^{(l)}$ affine result of the layer shape $(n^{(l)}, m)$, $z^{(l)} = w^{(l)}A^{(l-1)} + b^{(l)}$

$g^{(l)}$ activation fn of l th layer

$A^{(l)}$ activation output of l th layer shape $(n^{(l)}, m)$

$$A^{(l)} = g^{(l)}(z^{(l)})$$

Summary of Dimensions:

$w^{[l]}$	$(n^{[l]}, n^{[l-1]})$	$dw^{[l]}$	$(n^{[l]}, n^{[l-1]})$
$b^{[l]}$	$(1^{[l]}, 1)$	$db^{[l]}$	$(n^{[l]}, 1)$
$z^{[l]}$	$(n^{[l]}, m)$	$dz^{[l]}$	$(n^{[l]}, m)$
$A^{[l]}$	$(n^{[l]}, m)$	$dA^{[l]}$	$(n^{[l]}, m)$

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 $n_y = n_h^{[\text{number of layers} + 1]}$.

· 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

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

· $b^{[l]} \in \mathbb{R}^{\text{number of units in next layer}}$ 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

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

· 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:

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

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

2 Deep Learning representations

For representations:

- nodes represent inputs, activations or outputs
- edges represent weights or biases

Here are several examples of Standard deep learning representations

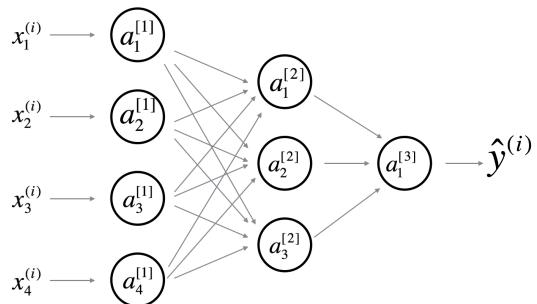


Figure 1: Comprehensive Network: representation commonly used for Neural Networks. For better aesthetic, we omitted the details on the parameters ($w_{ij}^{[l]}$ and $b_i^{[l]}$ etc...) that should appear on the edges

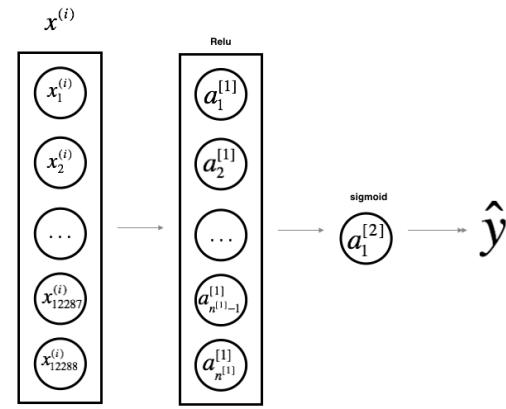


Figure 2: Simplified Network: a simpler representation of a two layer neural network, both are equivalent.