

INTRODUCTION

Artificial Intelligence (AI)

Neural Network

Neuron function

↓
Rectified Linear Unit (ReLU)

Supervised Learning

Real Estate / Standard
Online Advertising / CNN
Photo tagging / CNN
Speech Recognition / RNN
Machine Translation / RNN
Autonomous driving / Hybrid

Output

Input

Data

Structured Data - Computer

Unstructured Data - Human

Scale drives deep learning process

[- Data - labelled

- Size of NN

- Computation

- Algorithms - Sigmoid → ReLU

parameter changes slowly
gradient → extreme

process



Neural Network & Deep Learning 4
Improving Deep Neural Networks 3
Structuring Machine Learning projects 2
Convolution Neural Networks
Sequence Models

NEURAL NETWORK BASICS

Binary Classification

Logistic Regression

Image → RGB Channels

$\begin{matrix} r \\ g \\ b \end{matrix}$

pixel intensity value → Feature Value

Setup

Given X

Output \hat{y}

$$\hat{y} = P(y=1/x)$$

$$= \sigma(w^T X + b)$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

Parameter

$$w \in R^{n_x}, b \in R$$

Sigmoid function

! Easy

$$b \rightarrow x_0 = 1$$

$$\text{Now } X \sim R^{n_x+1}$$

$$\hat{y} = \sigma(W^T X)$$

$$X \in R^{n_x}$$

$$y \in \{0, 1\}$$

m training example

$$\text{example} - (X^{(i)}, y^{(i)})$$

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & \dots \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{matrix} \downarrow \\ n_x \\ 1 \end{matrix}$$

$$X \in R^{n_x \times m}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots \end{bmatrix} \in R^{1 \times m}$$

Cost function: i^{th} - q^{th} example

Loss (error) function
(single training example)

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

Logistic → non convex

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Cost function
(training example set)

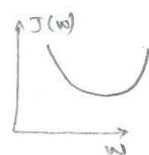
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m \text{loss function}$$

Gradient Descent

Cost function

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



Repeat { learning rate
partial derivative
 $w := w - \alpha \frac{dJ(w)}{dw}$ }
- α dw

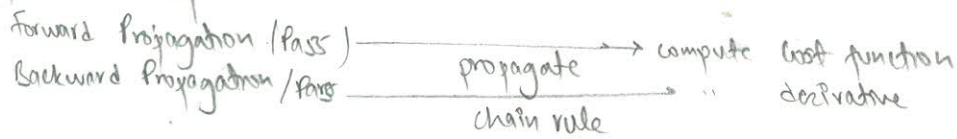
$$y=1 \quad P(y/x) = \hat{y}$$

$$y=0 \quad P(y/x) = 1 - \hat{y}$$

$$P(y/x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\max \log P(y/x)$$

Computation Graph



Logistic Regression

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

$$\hat{y} = \sigma(z)$$

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

- i) Single Example
- ii) Training set

$$w' = w - \alpha \frac{dw}{dw}$$

Vectorization

for loop → vectorization algebra

broadcasting - matrix manipulation

forward propagation
backward propagation

np.dot()

np.exp()

axis

reshape

CPU parallelization

GPU

SIMD - Single instruction multiple data

Jupyter & Python Notebook

numpy - rank 1 array
! vector

Neural Network

[layer]

(individual training example)

Layers

Input hidden output

activation

$a^{[L]}$

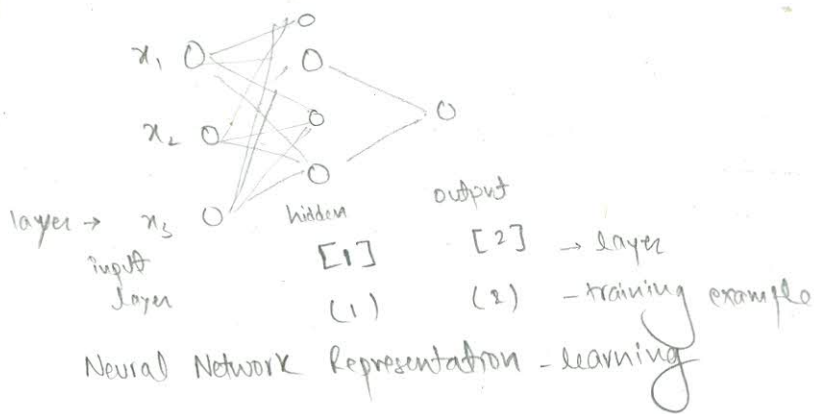
$a^{[layer]}$
in size

no of layer = 1 input layer

dimension of
 $a^{[i]}$ (no. of nodes)
 = (no. of input nodes)

Shallow - neural - networks

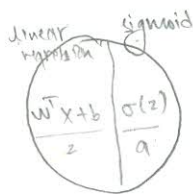
Computational graph



no. of layer - 1 input layer
hidden layer + output layer
 \hookrightarrow parameter associated

dimension
 $w^{[2]} = \text{no. of nodes} \times \text{no. of input nodes / features}$
 $b^{[2]} = \text{no. of nodes} \times 1$

$a^{[0]}$ $a^{[1]}$ $a^{[2]}$ - activation layer
 $z \rightarrow \text{node}$



$$\frac{\sigma(w^T \cdot x + b)}{z}$$

A activation

- 1) Single Training Example
- 2) Entire data \rightarrow across multiple example

n - features
m - training example

Activation functions

Sigmoid

non-linear function

- 1) hidden layer
- 2) output activation function

② tanh - hyperbolic -1 to 1
mean = 0 centering of data

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

③ ReLU - Rectified linear unit
④ Leaky ReLU

$$\max(0, z)$$

$$\max(0.2, z)$$

Choice - Design
- Initialization
- Activation function

Non-linear activation function

linear activation function
activation identity function

composition of two linear function is linear function

Derivative of activation function \rightarrow Backward propagation

Sigmoid $g(z) = \frac{1}{1+e^{-z}}$ $g'(z) = g(z)(1-g(z))$ - easy to compute

Tanh $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ $= 1 - (\tanh(z))^2$

ReLU step function

Random initialization \rightarrow small
Symmetric gradient problem
" breaking "

Deep-neural Network

Network structure

Shallow } Neural Network
Deep }

- no. of hidden layer

- hyper parameter

Forward Propagation

Backward Propagation

$$a^{[0]} \rightarrow a^{[L]} = w^{[L]} \cdot a^{[L-1]} + b^{[L]} \quad (\text{current layer } \times 1)$$

$$da^{[L]} \rightarrow da^{[L-1]} = \left[\frac{dw^{[L]}}{dz^{[L-1]}} \cdot da^{[L]} + db^{[L]} \right] \cdot \left(\frac{da^{[L-1]}}{dz^{[L-1]}} \right)$$

(current layer) (previous layer) (previous layer x no. of example)

L = no. of layer

$n^{[L]}$ = # neuron in layer L

$a^{[L]}$ - activation in layer L

$$a^{[L]} = g^{[L]}(z^{[L]})$$

Why deep learning

low level features

basic unit

complex pattern

Parameter, Hyperparameters

w
 b

learning rate

no. of hidden layer

hidden unit

activation function

iteration

Momentum

mini batch

regularisation parameter



Learning procedure

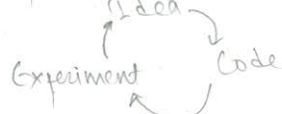
1) over simplified analogy neuro scientist \rightarrow human brain

2) Circuit Theory {logic gate} - function

complexity $O(\log n)$

exponentially large

Empirical priors



$X \rightarrow Y$ mapping

/NN