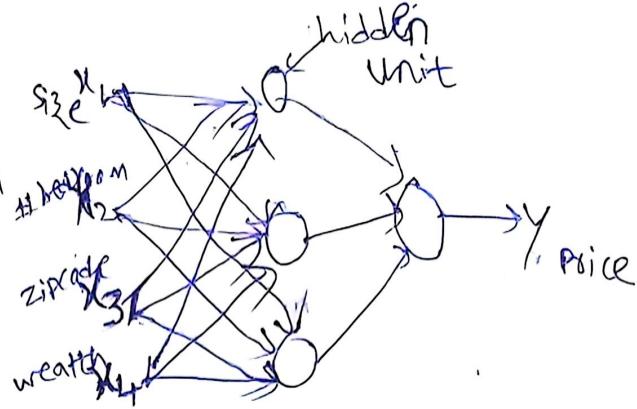


Deep Learning Course

→ Neural network comprises an input layer, and an output layer

→ Deep learning is made up of several hidden layers of neural networks that perform complex operations on massive amounts of structured and unstructured data



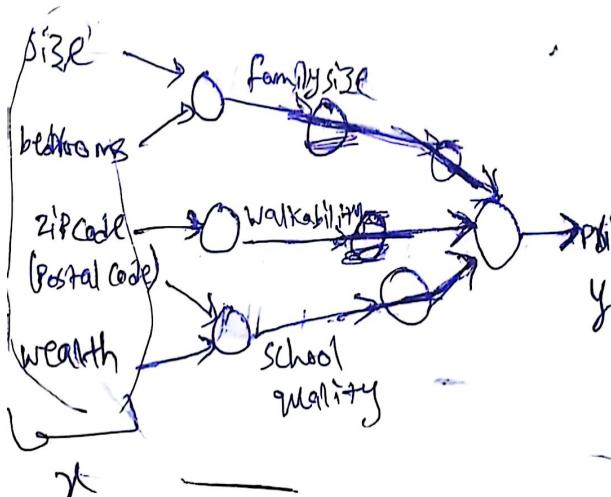
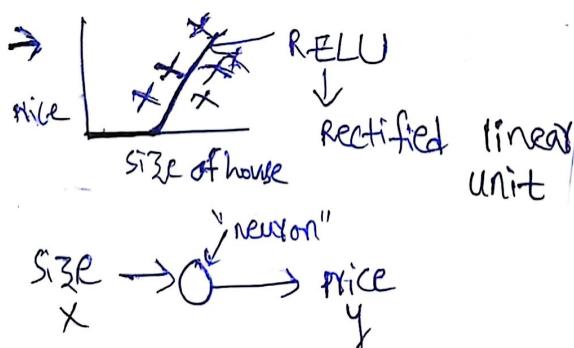
What you'll learn

1. Neural networks and deep learning
2. Improving Deep Neural Networks: hyper parameter tuning, regularization and optimization.
3. Structuring your ML Project
4. Convolution Neural Networks
5. NLP : Building sequence models.

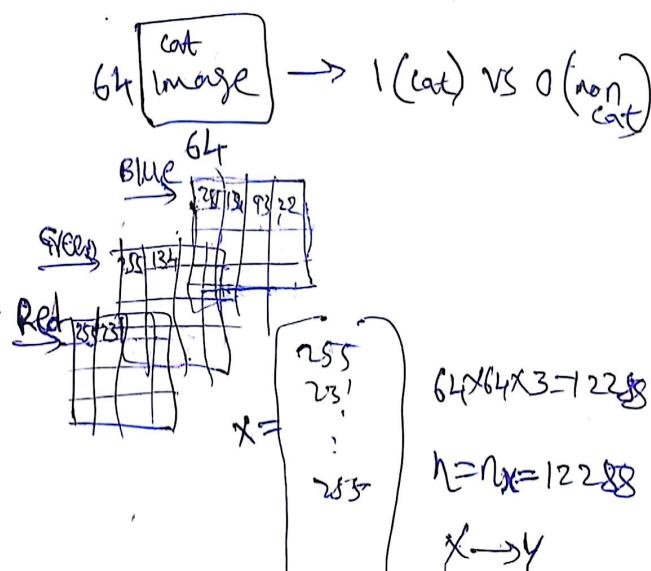
Supervised learning

Input (x)	o/p (y)	standard NN
Home features	Price	Real Estate
Ad/user info	click on add (0/1)	online Advertising
Image	object(1, ..., 100)	CNN Photo tagging
Audio	Text transcript	RNN Speech recognition
English	Chinese	RNN Machine translation
Image, Radar info	Position of other cars	RNN Autonomous driving
		CNN Custom/Hybrid NN

Housing price prediction



Binary classification



Notation 'm' \rightarrow small m is used to represent data

(x_i, y) $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$

input output

m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

M = Train

M_{test} = # test examples

\rightarrow A cost function is a measure of how inaccurate the model is in estimating the connection b/w X and Y

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \begin{matrix} \uparrow \\ n_x \end{matrix}$$

$\leftarrow m \rightarrow$

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

$$X.\text{shape} = (n_x, m)$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

\rightarrow A loss function is a function

that compares the target & predicted o/p values; measures how well the neural network models the training data. When training, we aim to minimize this loss b/w the predicted & target o/p's.

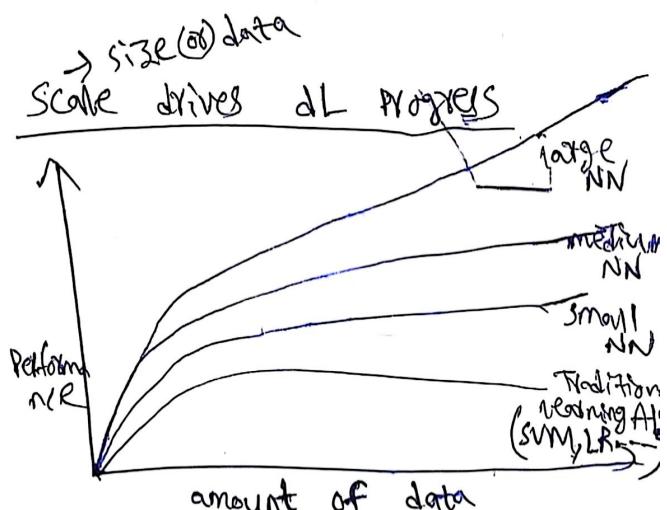
\rightarrow The cost function of a neural network will be the sum of errors in each layer. This is done by finding the error at each layer first and then summing the individual error to get the total error.

Loss function & cost function

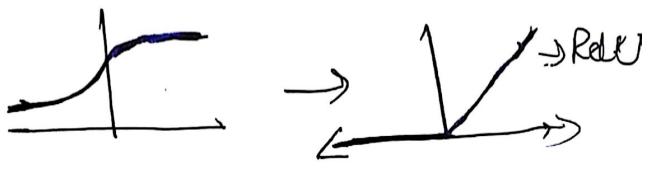
\rightarrow The loss function is to capture the difference between the actual & predicted values for a single record whereas cost functions aggregate the difference for the entire training dataset.

\rightarrow The most commonly used loss functions are mean-squared error and hinge loss.

\rightarrow The cost function measures the model's error on a group of objects, whereas the loss function



\rightarrow Gradient descent converts sigmoid function to ReLU



Logistic Regression!

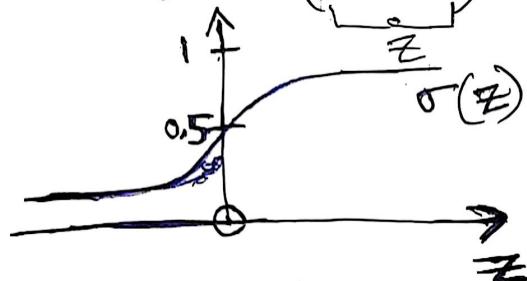
Given X , want $\hat{y} = P(y=1/x)$

$$x \in \mathbb{R}^{n \times 1}$$

Parameters: $w \in \mathbb{R}^{n \times 1}$, $b \in \mathbb{R}$ \rightarrow real number

$$0 \leq \hat{y} \leq 1$$

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



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

If z is large, $\sigma(z) \approx \frac{1}{1+0} = 1$

If z is small (large -ve number)

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\infty} = 0$$

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

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

$$z^{(i)} = w^T x^{(i)} + b$$

Loss (error) function:

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

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

$$\text{If } y=1 : L(\hat{y}, y) = -\log \hat{y}$$

want $\log \hat{y}$ large, for that \hat{y} should be large

$$\text{If } y=0 : L(\hat{y}, y) = -\log(1-\hat{y})$$

want $\log(1-\hat{y})$ large, for that \hat{y} should be small.

$$x_0 = 1, x \in \mathbb{R}^{n \times 1}$$

$$\hat{y} = \sigma(w^T x)$$

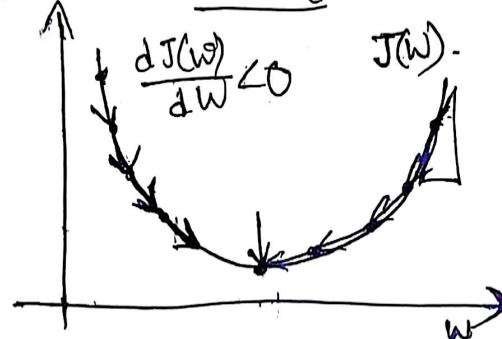
$$w = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad b \leftarrow \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

Cost function:

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

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

Gradient Descent



LR Cost function:

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

$$\text{Given } \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$$

Repeat {

$$W := W - \alpha \frac{d J(W)}{d W}$$

learning rate

}

more derivative examples:

$$f(a) = 2a \quad \frac{d}{da} f(a) = 2a \quad a=2 \quad f(a) = 4$$

$$a=2.001 \quad f(a) \approx 4.004$$

$$f(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2 \quad a=2 \times f(a) = 8$$

$$a=2.001 \quad f(a) \approx 8.012$$

$$f(a) = \log_{(e)}(a) \quad \frac{d}{da} f(a) = \frac{1}{a} \quad a=2 \quad f(a) \approx 0.69315$$

$$a=2.001 \quad f(a) \approx 0.6935$$

$$0.0005 \leftarrow 0.0005$$

$$\frac{d}{da} f(a) = \frac{1}{a}$$

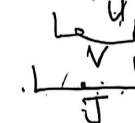
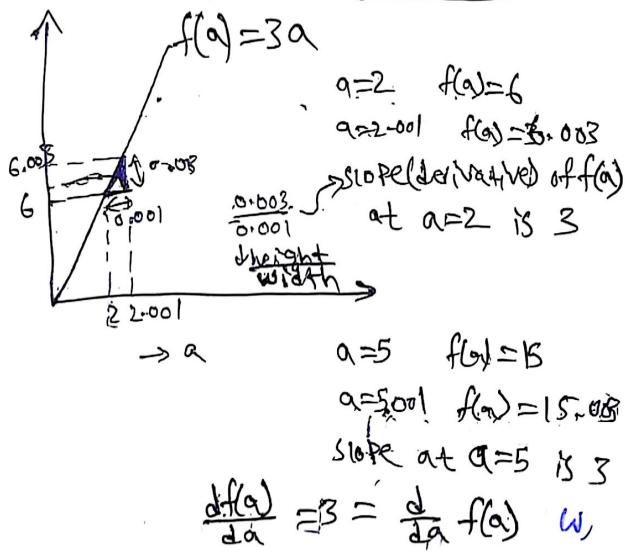
In Logistic Regression

$$J(W, b) = \begin{cases} W := W - \alpha \frac{d J(W, b)}{d W} \\ b := b - \alpha \frac{d J(W, b)}{d b} \end{cases}$$

Computation Graph

$$J(a, b, c) = 3(a+b+c) = 3(5+3+2) = 33$$

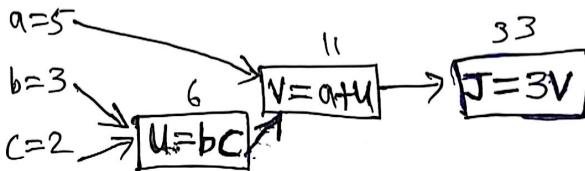
Intuition about derivatives:



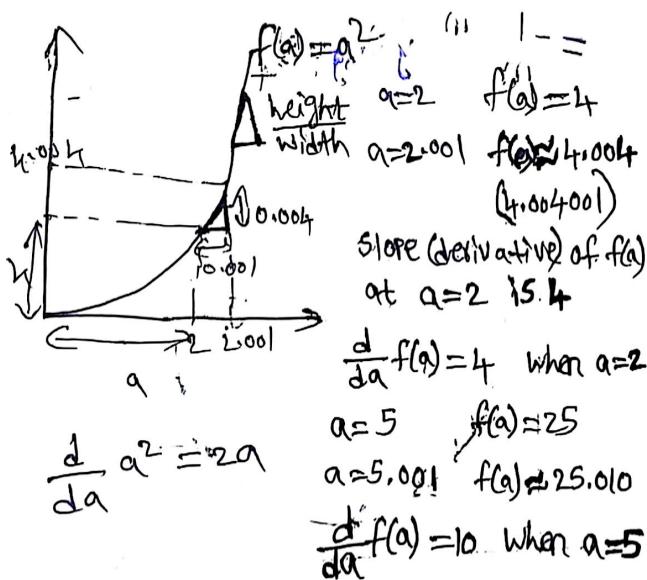
$$U = bC$$

$$V = a+U$$

$$J = 3V$$



Computing derivatives)

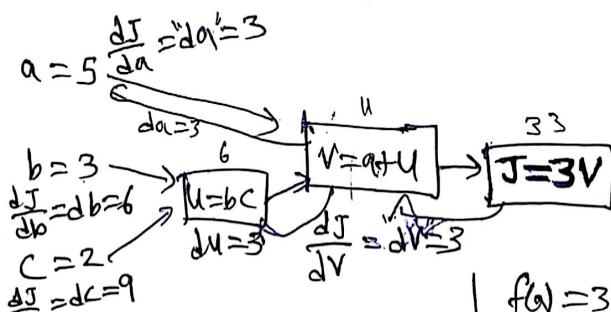


$$\frac{d}{da} a^2 = 2a$$

$$a=5 \quad f(a) = 25$$

$$a=5.001 \quad f(a) \approx 25.001$$

$$\frac{d}{da} f(a) = 10 \text{ when } a=5$$



$$\frac{d}{da} f(a) = 3a$$

$$J = 3V \quad V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$a=5 \rightarrow 5.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{du} = 3 = \underbrace{\frac{dJ}{dv}}_3 \cdot \underbrace{\frac{dv}{du}}_1$$

$$\frac{dJ}{db} = \underbrace{\frac{dJ}{dy}}_3 \cdot \underbrace{\frac{dy}{db}}_2 = 6$$

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9$$

Logistic regression derivatives

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

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

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

$$z = w_1 x_1 + w_2 x_2 + b$$

$$a = \sigma(z)$$

$$L(a, y)$$

$$d\bar{z} = \frac{dL}{dz} = \frac{dL(a, y)}{da}$$

$$da = \frac{dL(a, y)}{da}$$

$$= a - y$$

$$= \frac{y}{a} + \frac{1-y}{1-a}$$

$$= \frac{dL}{da} \cdot \frac{da}{d\bar{z}}$$

$$= a(1-a)$$

By multiplying above two

$$\text{we get } \frac{dL}{d\bar{z}} = a - y = d\bar{z}$$

$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot d\bar{z}$$

$$\frac{\partial L}{\partial w_2} = "dw_2" = x_2 \cdot d\bar{z}$$

$$db = d\bar{z}$$

$$w_1 : w_1 - \alpha dw_1$$

$$w_2 : w_2 - \alpha dw_2$$

$$b : b - \alpha db$$

logistic Regression on m examples

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

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$(x^{(i)}, y^{(i)})$$

$$dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})}_{dw_1^{(i)}}$$

Logistic Regression derivatives

$$J=0; dw_1=0; dw_2=0; db=0$$

for $i=1$ to m :

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J_i = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$d\bar{z}^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 = x_1^{(i)} d\bar{z}^{(i)}$$

$$dw_2 = x_2^{(i)} d\bar{z}^{(i)}$$

$$db = d\bar{z}^{(i)}$$

$$J= m$$

$$|dw_1|=m; |dw_2|=m; |db|=m$$

$$dw_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

What is Vectorization?

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

$$w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

non-vectorized

$$z = 0$$

```
for i in range(n-x):
    z += w[i]*x[i]
```

$$z += b$$

vectorized

$$z = np.dot(w, x) + b$$

$$w^T x$$

① import numpy as np

import time

```
a = np.random.rand(1000000)
b = np.random.rand(1000000)
```

tic = time.time()

c = np.dot(a, b)

toc = time.time()

print(c)

print("Vectorized version! " + str((toc - tic)) + "ms")

c = 0

tic = time.time()

for i in range(1000000):

c += a[i]*b[i]

toc = time.time()

print(c)

print("For loop! " + str(1000*(toc - tic)) + "ms")

OP:

250286.989866

vectorized version: 1.50275230 ms

250286.989866

For loop: 474.29513931 ms

The program is to know which

one is faster vectorized version
or ~~for loop~~ non-vectorized version.

Neural Network Programming guideline

→ whenever possible, avoid explicit for-loops.

non-vectorized

$$u = Av$$

$$u_i = \sum_j A_{ij} v_j$$

$$u = np.zeros((n, 1))$$

for i --
for j --

$$u[i] += A[i][j] * v[j]$$

vectorized

$$u = np.dot(A, v)$$

Vectors & matrix valued functions

Say you need to apply the exponential operation on every element of a matrix / vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

non-vectorized

$$u = np.zeros((n, 1))$$

for i in range(n):

$$u[i] = math.exp(v[i])$$

vectorized

import numpy as np

$$u = np.exp(v)$$

Vectorizing

logistic

Regression

$$z^{(1)} = w^T x^{(1)} + b$$

$$\alpha^{(1)} = \sigma(z^{(1)})$$

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad (n, m) \text{ shape}$$

$$z^{(2)} = w^T x^{(2)} + b$$

$$\alpha^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = w^T x^{(3)} + b$$

$$\alpha^{(3)} = \sigma(z^{(3)})$$

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = W^T X + \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

1xm

$$= \begin{bmatrix} W^T X^{(1)} + b & W^T X^{(2)} + b & \dots & \dots & W^T X^{(m)} + b \end{bmatrix}$$

1xm

$$W^T X = \begin{bmatrix} \dots \end{bmatrix} \begin{bmatrix} 1 & 1 & \dots & 1 \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \end{bmatrix}$$

$$z = \underbrace{n \cdot \text{dot}(W^T, X)}_{\substack{(1,1) \\ \text{"Broadcasting!"}}} + b$$

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} = \sigma(z)$$

Broadcasting example

$$dz^{(1)} = a^{(1)} - y^{(1)}$$

$$dz^{(2)} = a^{(2)} - y^{(2)} \dots$$

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

1xm

$$A = \begin{bmatrix} a^{(1)} & \dots & a^{(m)} \end{bmatrix} - Y = \begin{bmatrix} y^{(1)} & \dots & y^{(m)} \end{bmatrix}$$

$$dz = A - Y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots \end{bmatrix}$$

Calories from Carbs, Protein, Fats in 100g of different foods:

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

3x4

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} n \cdot \text{sum}(dz)$$

$$dW = \frac{1}{m} \times dZ^T$$

$$= \frac{1}{m} \begin{bmatrix} 1 & 1 & \dots & 1 \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

Calculate % of Calories from Carb, Protein, fat. Can you do this without using explicit for-loop.

Program)

import numpy as np

A = np.array([[[56.0, 0.0, 4.4, 68.0],
[1.2, 104.0, 52.0, 8.0],
[1.8, 135.0, 99.0, 0.9]]])

print(A)

$\text{col} = A.\text{sum}(\text{axis}=0)$

$\text{print}(\text{col})$

O/P:-

[59. 239. 155.4 76.9]

P) Percentage = $100 * \text{A}/\text{col}$

$\text{print}(\text{Percentage})$

O/P:-

[[94.9152 0. - -]
2.0338 - - - -]
3.0508 - - - -]

import numpy as np

a=np.random.randn(5)

$\text{print}(a)$

O/P:-

[0.562052 -0.2969 0.9522 -0.8242
-1.4626]

$\text{print}(a.\text{shape})$

O/P:- (5,)

$\text{print}(a.\text{T})$

O/P:-

O/P is same as above

$\text{print}(\text{np.dot}(a, a.\text{T}))$

O/P:-

4.0657

Broadcasting example

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 204 & 205 & 306 \end{bmatrix}$$

(2,3) (1,3) (m,n) (m,n)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

(m,n) (m,n) ↓ (m,n)

P) $a = \text{np.random.randn}(5,1)$

$\text{print}(a)$

O/P:- [-0.0967
-2.3861 -
-0.3243
-0.962]
[0.544]

P) $\text{print}(a.\text{T})$

[[-0.0967 -2.3861 -
- - - -]]

P) $\text{print}(\text{np.dot}(a, a.\text{T}))$

O/P:- [[0.0093 0.0238
0.2308 : :
0.0313 : :
0.0930 : :]]

General principle:-

$(m,n) + (1,n) \rightsquigarrow (m,n)$

matrix \times $(m,1) \rightsquigarrow (m,n)$

Packages

- numpy is the fundamental package for scientific computing with python.
- h5py is a common package to interact with a dataset that is stored on an h5 file.
- matplotlib is a famous library to plot graphs in python.

Logistic regression cost function.

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

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

If $y=1 : P(y|x) = \hat{y}$

If $y=0 : P(y|x) = 1-\hat{y}$

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

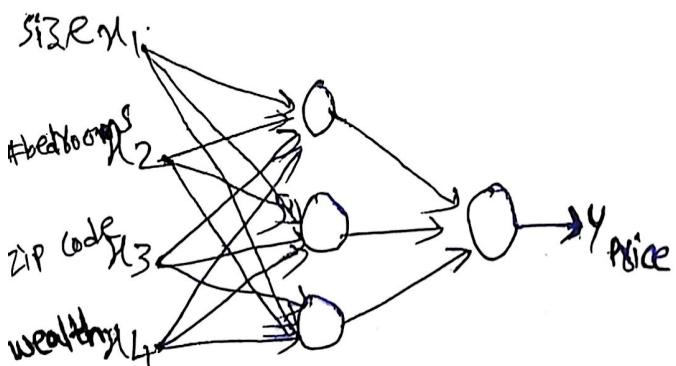
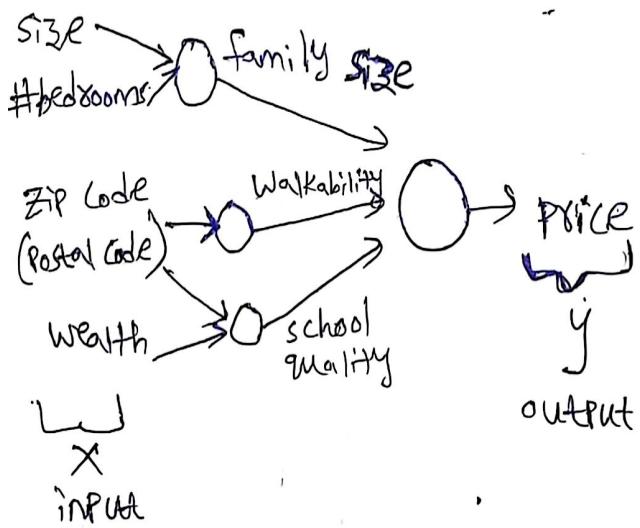
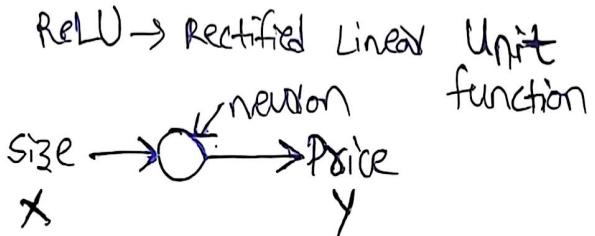
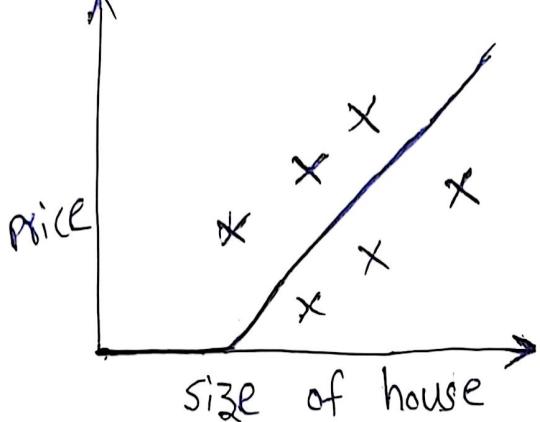
$$\text{If } y=1 : P(y|x) = \hat{y} \quad \cancel{\text{if } y=0}$$

$$\text{If } y=0 : P(y|x) = 1-\hat{y}$$

$$\begin{aligned} \log P(y|x) &= \log \hat{y}^y (1-\hat{y})^{(1-y)} \\ &= y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -L(\hat{y}, y) \downarrow \text{decreases.} \end{aligned}$$

Andrew Ng Lectures!

Housing Price Prediction

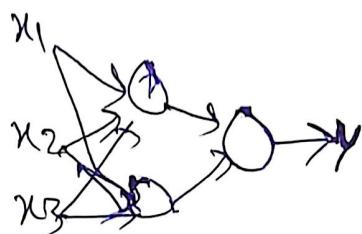


Supervised Learning

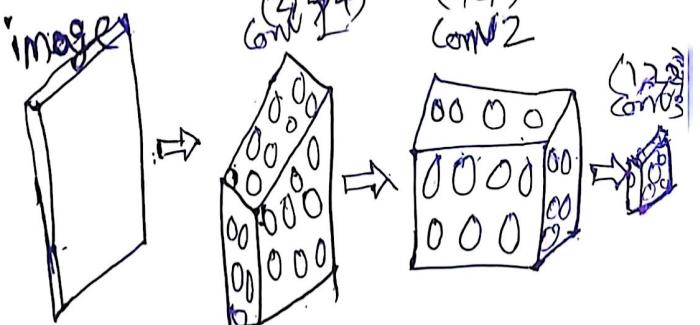
INPUT (X)	OUTPUT (y)	APPLICATION
Home features	Price	Real Estate
Ad, user info	click on ad? (0 or 1)	Online Advertising
Image	Object (1...100)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

Custom Hybrid Neural Network

Neural Network examples



Standard Neural Network
(1, 2, 2)
(2, 4, 1)
(6, 1)



Convolutional Neural Network



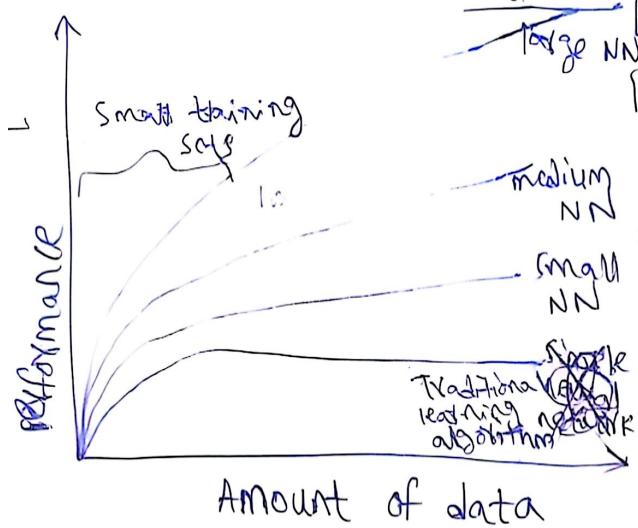
Scale drives deep learning

Course:

Neural networks and Deep

Learning:

Progress:
Large NN



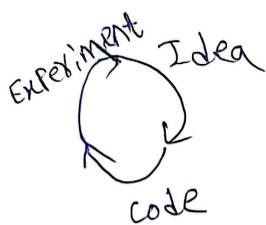
Week 1: Introduction

Week 2: Basics of Neural Network Programming

Week 3: One hidden layer Neural Networks

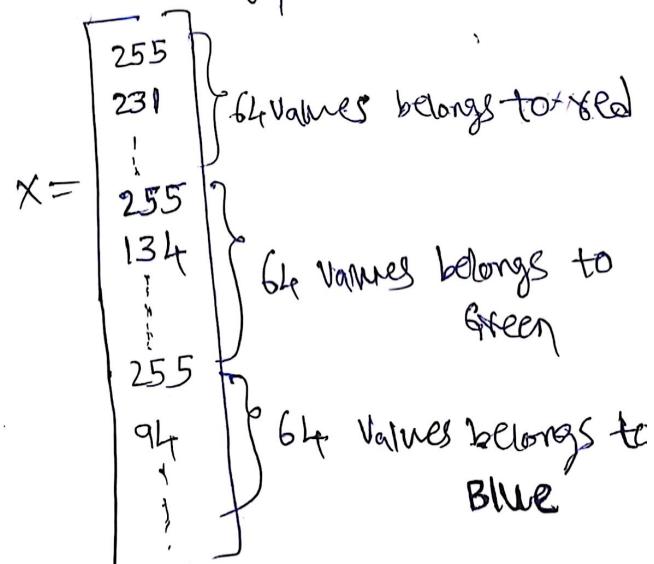
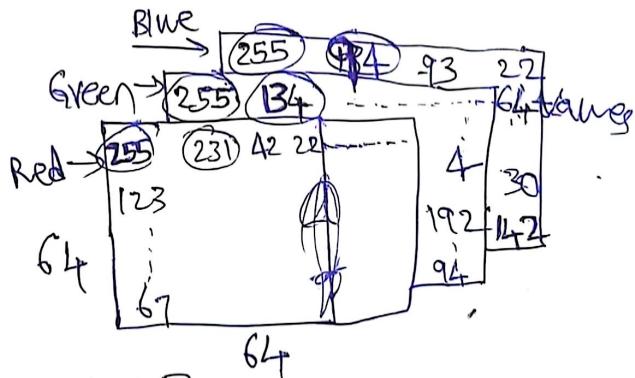
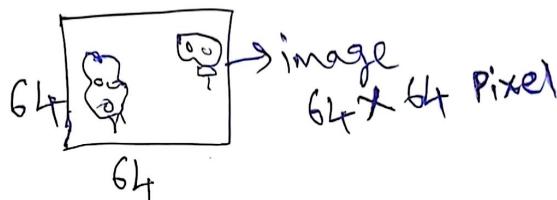
Week 4: Deep Neural Networks

- Data
- Computation
- Algorithms



Binary Classification

→ Logistic Regression is an algorithm for binary classification.



$$64 \times 64 \times 3 = 12288$$

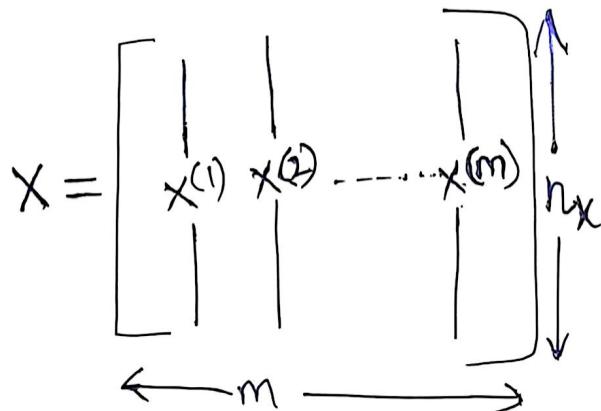
→ PCY one 64×64 pixels image $M = M_{\text{train}}$ $M_{\text{test}} = M_{\text{test}}$
 there are three matrices (Red, Green, Blue) (RGB), each of size 64×64 .

~~Red~~ Per one image

$$\text{Red} \rightarrow 64 \times 64$$

$$\text{Green} \rightarrow 64 \times 64$$

$$\text{Blue} \rightarrow 64 \times 64$$



→ The resultant feature vector contains $64 \times 64 \times 3$ values

→ The unfold of the matrices happens like, first row of Red matrices happens and then second row and so on like that until 64th row of Red matrices and then Green matrices unfolds in the same way and then Blue matrices.

$$64 \times 64 \times 3 = 12288$$

→ The feature vector is represented

as $n_x = 12288$

dimension of input feature vector

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

$$X.\text{shape} = (n_x, m)$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

Problem: Cat(1) or non-Cat(0)
Logistic Regression:

Given X , want $\hat{y} = P(y=1 | x)$
 $X \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$, $0 \leq \hat{y} \leq 1$

OUTPUT $\hat{y} = w^T x + b$
 ↳ Linear regression

It shows the probability of image, either cat or non-cat

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

$\hat{z} = 1$

If z is large, \hat{y} approaches 1
 If z is small, \hat{y} approaches 0

If z is large (z)
 If z is large, \hat{y} approaches 1
 If z is small, \hat{y} approaches 0

Notation

(x, y) , $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$
 Single training example

m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

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

\rightarrow If z is large, $\sigma(z) \approx \frac{1}{1+0} = 1$

\rightarrow If z is large negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\infty} \approx 0$$

$$x_0 = 1, x \in \mathbb{R}^{N_x+1}$$

$$\hat{y} = \sigma(w^T x), \text{ if } \hat{y}$$

$$w = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{N_x} \end{bmatrix} \quad \left\{ \begin{array}{l} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{N_x} \end{array} \right\} \left\{ \begin{array}{l} b \\ w \end{array} \right\}$$

This notation is
not important

Logistic Regression Cost Function:

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

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$,
Want $\hat{y}^{(i)} \approx y^{(i)}$

for m training examples

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

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

$$z^{(i)} = w^T x^{(i)} + b$$

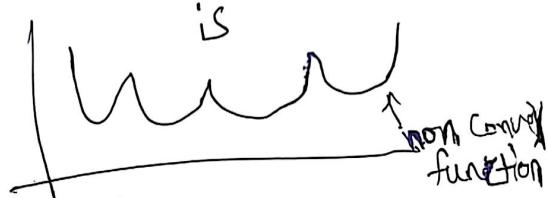
For i th training example, it is
represented as

Loss (Error) function:

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

↳ square error

The graph of square error



So, square error curve is not possible to find global minimum or global maximum.

\rightarrow So, in logistic regression we use

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

$$\text{If } y=1: L(\hat{y}, y) = -\log \hat{y}$$

want $\log \hat{y}$ large, want \hat{y} large

$$\text{If } y=0: L(\hat{y}, y) = -\log(1-\hat{y})$$

want $\log(1-\hat{y})$ large, want \hat{y} small

Cost Function:

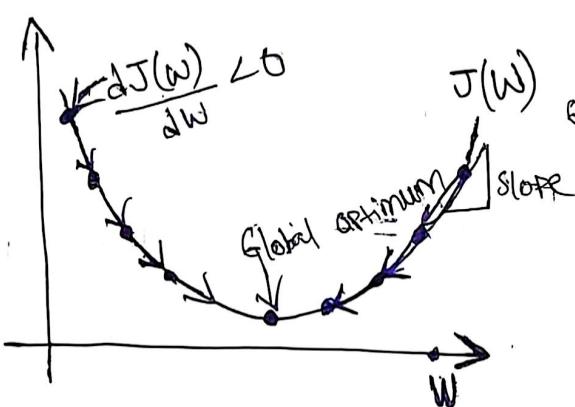
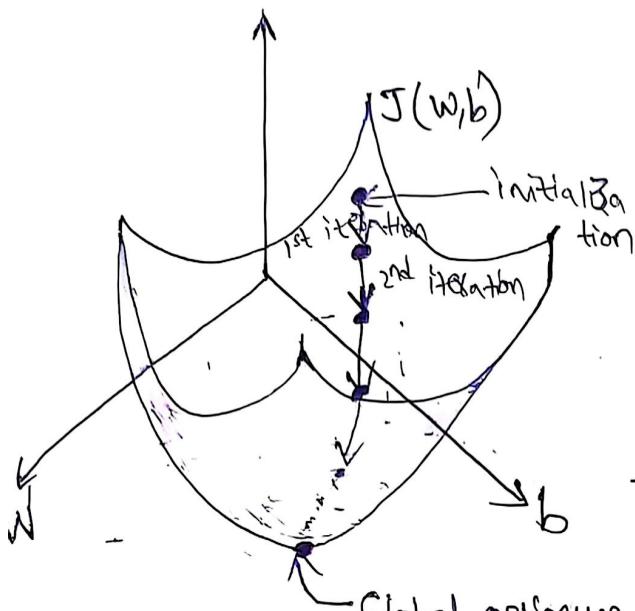
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

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

Gradient Descent :

Want to find w, b , that

minimize $J(w, b)$



Repeat {

$$w := w - \alpha \frac{dJ(w)}{dw}$$

$$w := w - \alpha \frac{dJ(w)}{dw}$$

~~$\frac{dJ(w)}{dw} = ?$~~

$J(w, b) \rightarrow w$ updates as

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$\rightarrow b: \text{Updates as } b := b - \alpha \frac{dJ(w, b)}{db}$$

$$\frac{\partial J(w, b)}{\partial w} \Rightarrow dw \xrightarrow{\text{in code it is }} dw$$

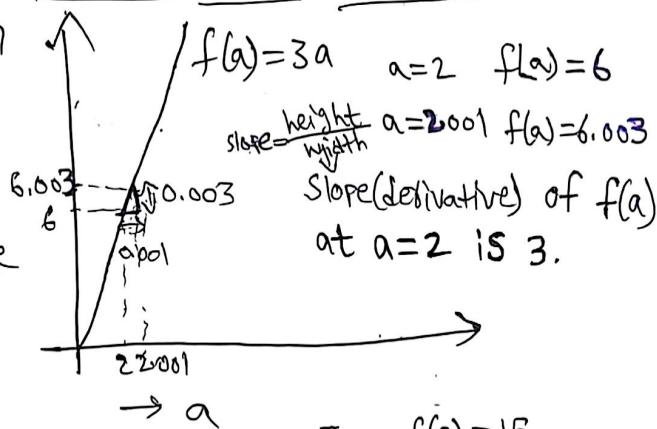
This shows the slope of $J(w, b)$ in the direction of w .

$$\frac{\partial J(w, b)}{\partial b} \Rightarrow db \xrightarrow{\text{in code it is shown as }} db$$

This shows the slope of $J(w, b)$ in the direction of b .

15/01/2024

Intuition about derivatives:



$$a=5 \quad f(a)=15$$

$$a=5.001 \quad f(a)=15.003$$

Slope at $a=5$ is also 3

$$\frac{d.f(a)}{da} = 3 = \frac{d}{da} f(a)$$

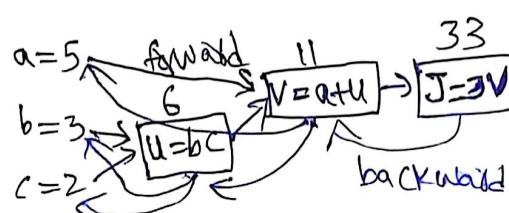
Computation graph

$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \cdot 2) = 33$$

$$U = bc$$

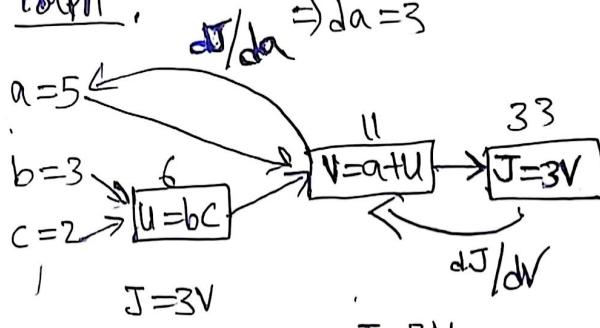
$$V = a + U$$

$$J = 3V$$



Derivatives with a Computation

Graph:



$$J = 3V$$

$$\frac{dJ}{dV} = 3$$

$$J = 3V \\ V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

Increase in 'V' increases 'J' (three times)

$$\frac{dJ}{da} = 3 = \frac{dJ}{dV} \cdot \frac{dV}{da}$$

$$a = 5 \xrightarrow{\text{increased}} 5.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$= 3 \cdot (1)$$

$$= 3$$

$$\frac{dJ}{da} = 3$$

$$\frac{dV}{da} = \frac{d(a+u)}{da}$$

$$= (1)$$

Final Output Variable

↳ d variable

→ When we are implementing it in software, we can give this as "dJ/dw1x1"

→ In the code we can write it as "dL/dw1x1"

$$\frac{dJ}{dV} \Rightarrow dV = 3, \frac{dJ}{da} \Rightarrow da = 3$$

$$\frac{dJ}{du} = 3$$

$$= \frac{dJ}{dV} \cdot \frac{dV}{du} \\ = (3) (1) \\ = 3.$$

$$u = 6 \xrightarrow{\text{increased}} 6.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} \\ 3 \quad 2$$

$$b = 3 \rightarrow 3.001 \\ u = bc = 6 \rightarrow 6.002 \\ c = 2$$

$$V = 11.002 \\ J = 33.006$$

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} \\ 3 \quad 3$$

$$= 9$$

18/01/2024

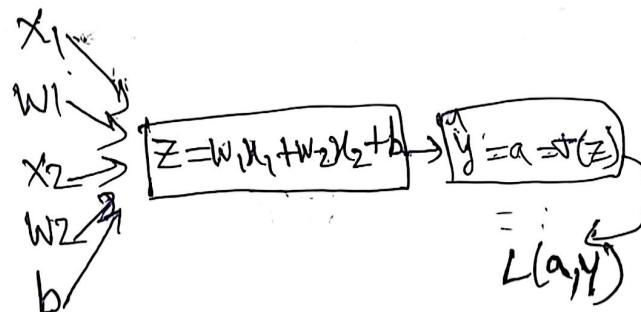
Logistic Regression Gradient Descent

Recap

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

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

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



$$\frac{dL}{da} = \frac{dL(a, y)}{da} = -\left(\frac{d}{da}(\sigma(y \log a + (1-y) \log(1-a)))\right)$$

$$= -\left(\frac{y}{a} - \frac{1-y}{1-a}\right)$$

$$\frac{dL}{dz} = \frac{dL(a, y)}{dz} = \frac{da}{da} \cdot \frac{da}{dz}$$

$$\boxed{\frac{dL}{dz} = a - y = \frac{dL}{dz}}$$

$$\Rightarrow \frac{dL}{da} \cdot \frac{da}{dz} \Rightarrow -\left(\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot \frac{da}{dz} \Rightarrow -\left(\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot (a(1-a))$$

$$\Rightarrow \frac{da}{dz} = \frac{d}{dz} \sigma(z) \quad \left[\because a = \sigma(z) = \frac{1}{1+e^{-z}} \right] \Rightarrow -\left(y(1-a) + a(1-y)\right)$$

$$= \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) \Rightarrow -\left(y(1-a) + a(1-y)\right)$$

~~$$\therefore \frac{d}{dx} \left(\frac{1}{A} \right) = -\frac{1}{A^2} \frac{dA}{dx}$$~~

$$= \frac{-1}{(1+e^{-z})^2} \cdot \frac{d}{dz}(1+e^{-z})$$

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

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

Add 1 to numerator & subtract 1 to numerator

$$= \frac{1+e^{-z}-1}{(1+e^{-z})^2} \quad \left[\text{split numerator} \right]$$

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

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

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

$$= a(1-a)$$

$$dz = a - y = \frac{dL}{dz}$$

At final step to adjust weights and 'b'

$$\frac{\partial L}{\partial w_1} = \Delta w_1 = x_1 \cdot dz$$

$$\Delta w_2 = x_2 \cdot dz$$

$$db = dz$$

$$w_1 := w_1 - \alpha \Delta w_1$$

$$w_2 := w_2 - \alpha \Delta w_2$$

$$b := b - \alpha db$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$= \left(\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot (a(1-a)) \cdot \frac{\partial z}{\partial w_1}$$

$$\left(\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot (a(1-a)) \cdot \frac{\partial z}{\partial w_1} = (a-y) \cdot \frac{\partial z}{\partial w_1}$$

$$\frac{\partial z}{\partial w_1} = \frac{\partial}{\partial w_1} \left(\frac{w_1 x_1 + w_2 x_2 + b}{1+e^{-z}} \right) = x_1 \cdot dz$$

$$= x_1 \cdot dz$$

$$\frac{\partial L}{\partial w_1} = x_1 \cdot dz$$

Gradient Descent on m

examples:

Logistic Regression on m

examples:

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

$$\rightarrow a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_j} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_j} L(a^{(i)}, y^{(i)})}_{\partial w_j^{(i)} - (x^{(i)}, y^{(i)})}$$

lets initialize

$$J = 0, \partial w_1 = 0, \partial w_2 = 0, \partial b = 0$$

for $i = 1$ to m

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += - \left[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)}) \right]$$

$$\partial z^{(i)} = a^{(i)} - y^{(i)}$$

$$\partial w_1 += x_1^{(i)} \cdot \partial z^{(i)}$$

$$\partial w_2 += x_2^{(i)} \cdot \partial z^{(i)}$$

$$\underbrace{\partial w_{\text{features}}}_{\text{features}} += x_{\text{features}}^{(i)} \cdot \partial z^{(i)}$$

$$\partial b += \partial z^{(i)}$$

so, we are calculating averages
we need to divide each
with 'm'

$$J / m, \partial w_1 / m;$$

$$\partial w_2 / m; \partial b / m$$

$$\partial w_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \partial w_1$$

$$w_2 := w_2 - \alpha \partial w_2$$

$$b := b - \alpha \partial b$$

Vectorization

→ it is used to get rid of extra for loops in the code

what is vectorization?

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

$$w = \begin{bmatrix} \vdots \\ i \\ \vdots \\ \vdots \end{bmatrix}$$

$$w \in \mathbb{R}^{n_w}$$

$$x = \begin{bmatrix} \vdots \\ j \\ \vdots \\ \vdots \end{bmatrix}$$

$$x \in \mathbb{R}^{n_x}$$

non-vectorial implementation

$$z = 0$$

for i in range ($n \rightarrow x$):

$$z += w[i] * x[i]$$

~~z += b~~

$$z += b$$

vectorial implementation

$$z = \underbrace{n_p \cdot \text{dot}(w, x)}_{w^T x} + b$$

w^T x

*
Note }

Whenever possible, avoid
explicit for-loops.

```
import numpy as np
```

```
import time
```

```
a = np.random.rand(1000000)
```

```
b = np.random.rand(1000000)
```

```
tic = time.time()
```

```
c = np.dot(a, b)
```

```
toc = time.time()
```

```
print(c)
```

```
print("Vectorized version: " + str(1000 * (toc - tic)))
```

```
c = 0
```

```
tic = time.time()
```

```
for i in range(1000000):
```

```
    c += a[i] * b[i]
```

```
toc = time.time()
```

```
print(c)
```

```
print("For loop: " + str(1000 * (toc - tic)))
```

MR) 249660.39691457726

vectorized version: 1.52587890625 ms

249660.39691457726

for loop: 1200.25777-- ms

Ex:-

$$U = A \cdot V$$

$$U_i = \sum_j A_{ij} V_j$$

Non-vectorized implementation

$$U = np.zeros((10, 1))$$

for i --- ←

for j --- ←

$$U[i] += A[i][j] * V[j]$$

Vectorized implementation:

$$U = np.dot(A, V)$$

Vectors and matrix valued functions

Say

→ In above code vectorized

method takes less time to

compute compared to non-vectorized

method.