

Courses in this specialization

Outline of the course

i) Neural Networks and Deep learning

Neural Networks has 2 steps

i) forward propagation ii) Backward propagation

Logistic Regression:

is an algorithm for Binary classification

$$(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$$

$$X = \begin{bmatrix} 1 & | & | & | \\ x_1 & x_2 & x_3 & \dots & x_n \\ 1 & | & | & | \end{bmatrix} \rightarrow n \times m \quad y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix}$$

We would be probably expecting
 $n \times 1 \quad x^T y$

Given x , want $\hat{y} = P(y=1|x)$
 $0 \leq \hat{y} \leq 1$

$x \in \mathbb{R}^{n_x}$ Parameters $w \in \mathbb{R}^{n_w}$, $b \in \mathbb{R}$

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

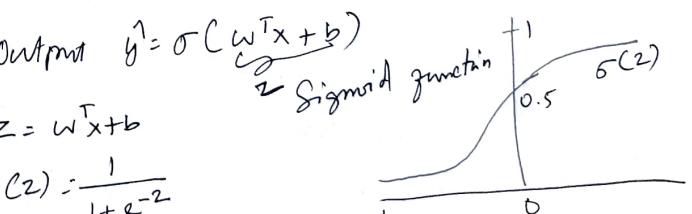
$$z = w^T x + b$$

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

$\approx z$ is large $\sigma(z) \approx \frac{1}{1+0}$
 e^{-z} is very close to 0

So $\sigma(z)$ is close to 1.
Conversely if z is small or

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



Cost function:

To train parameters w and b we need a cost function
 $y = \sigma(w^T x + b)$, where $\sigma(z) = \frac{1}{1+e^{-z}}$

Loss (error) function:

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

It helps to measure the accuracy of our algorithm by comparing \hat{y} (output) of y (actual). We would want the loss function output value to be as small as possible.

$$\begin{aligned} \text{if } y=1 & L(\hat{y}, y) = -y \log \hat{y} \quad \hat{y}=0 \\ & = -\log \hat{y} \Rightarrow -\log 0 = \text{infinity} \\ & \quad \log 0=1 \end{aligned}$$

So basically if $y=1$ you want $\hat{y}=1$

$$\begin{aligned} \text{if } y=0 & L(\hat{y}, y) = -\log(1-\hat{y}) \\ & \quad \hat{y}=1 \Rightarrow \log(0) = \text{infinity} \\ & \quad \hat{y}=0 \Rightarrow \log(1) = 0 \end{aligned}$$

Basically if $y=1$ we try to make \hat{y} large
if $y=0$ we try to make \hat{y} small

Loss function is measured on a single training example.
Whereas Cost function ~~J~~ $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$
 $= \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$

is the cost of your parameters w, b that

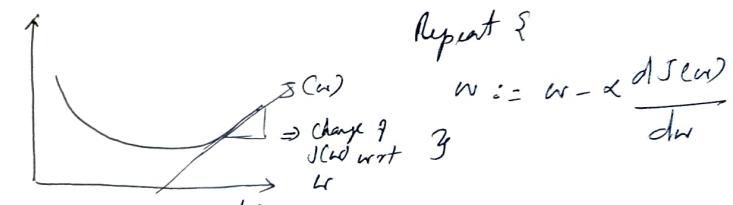
minimizes cost J .

Average of loss function = Cost function

Gradient Descent:

The cost function measures how well your parameter w, b are doing

Gradient descent helps us to find w, b so that we can minimize $J(w, b)$ ~~efficiently~~ efficiently.

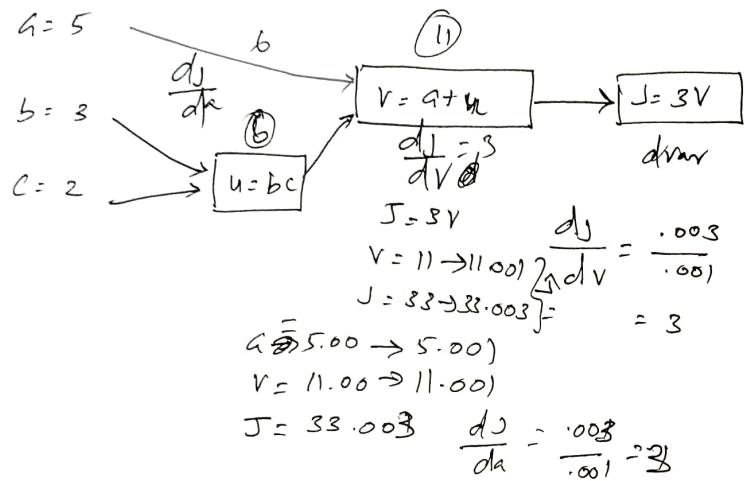


For simplicity we omit b , so that we can deal with one variable. Then continuously minimize w

1) ~~continuously subtract~~

1) Assign w to the difference between w and the derivative of the function $J(w)$ multiplied by a gradient descent constant α .

Computing derivatives



$$6 \div 3 = 2$$

~~Diagram showing 6 people divided into 3 groups of 2.~~

$\frac{6}{3} = 2$

Symon



The diagram consists of two parts. On the left, there is a large rectangle divided into 6 smaller squares. Inside these squares are 6 stick figures arranged in two rows of three. A horizontal line through the middle separates the two rows. A large bracket above the first three columns indicates they form one group, and another bracket above the last three indicates they form another group. On the right, there is a diagram showing three vertical rectangles, each containing two stick figures. Arrows point from the numbers 2 and 3 written above the rectangles to the corresponding groups of figures in the diagram on the left.

$$Z = W^T x + b \quad \hat{y} = a = \sigma(Z)$$

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

$$\frac{d}{dx} \left(\frac{\ln x}{x} \right) = \frac{1}{x^2} + \frac{1}{1-x}$$

$$\frac{dL(a,y)}{da} = \frac{(1-a) + a}{a(1-a)} - \frac{1}{a(1-a)}$$

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

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

x_1
 w_1
 x_2
 w_2
 b

$z = w_1x_1 + w_2x_2 + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$

"da" = $\frac{\partial L(a, y)}{\partial a}$ "dz" = $\frac{\partial L(a, y)}{\partial z}$
 $= \frac{\partial L}{\partial a} \frac{da}{dz}$ $= -\frac{y}{a} + \frac{1-y}{1-a}$
 $= a - y \rightarrow \text{Sigmoid}$
 $\frac{da}{dz} = a(1-a)$

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

$$\frac{da}{dz} = d \cancel{(c_0 c_2)}$$

$$a = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) \quad (2)$$

$\mu C^{1+\epsilon}$

$$a + a x^{-2} - 1 = 0$$

$$\frac{ds}{dt} = \alpha x^2$$

$$\frac{dx}{dt} = \alpha C(t)x$$

$$\frac{da}{dz} = a + a \cdot \frac{z}{e^{-z}}$$

$$- \partial^2 = a - \lambda^2$$

$$f(x) = g(x)$$

$$\text{unction } f(x) = \frac{g(x)}{h(x)},$$

$$\times h(x)] - [g(x) \times h'(x)]$$

$$\frac{d}{dx} \left[g(x) \times h(x) \right] = g'(x) \times h(x) + g(x) \times h'(x)$$

$$C_1 \frac{h(x)^2}{\left[1 - x^8\right]^{-2}}$$

$$(1+e^{-x})] - \left[\frac{1}{1+e^{-x}} \right] =$$

$$\frac{L}{(1+x^{-2})^2}$$

$$\lambda(1-\lambda) = \frac{1}{4}$$

You need to strongly consider

Vectors and Tensors

It gets tied to "for loops".

Logistic regression (or) function: $y = \sigma(C^T x + b)$ where $\sigma(z) = \frac{1}{1+e^{-z}}$

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

Wish to represent with a neural network!

Shape (num_chan, 3) when 3 is from 3 channels

RGB shape $\rightarrow (3, 3, 3)$

$$\begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

$$\begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

Color image medium

$256 \times 256 \times 256$ or 256×256 64 colors

$$\text{shape}(256, 256) \xrightarrow{\text{64 colors}} \begin{bmatrix} \text{Image} \end{bmatrix}$$

$$64 \times 64 \times 3 = \begin{bmatrix} \text{Image} \end{bmatrix}$$

64 rows

64 columns

64 depth

Then 9

Then $N = 12288$

$\text{shape}(y) = 5$

$A = \sigma(\frac{1}{1+e^{-(C^T x + b)}})$

$$\text{logistic } J = -y \log(a) + (1-y) \log(1-a) = 5$$

$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \text{loss}(A_i)$$

64

Common steps for pre-processing & new dataset we:

- figure out the dimensions and shapes of the problem (train, test, validation, ...)
- reshape the datasets such that each example is now a vector of size (num_pixels * num_pixels * 3, 1)

Standardize the data

$$X = 256 \times 256 \times 256 \quad \text{image value}$$

$$\begin{array}{c} \xrightarrow{R} 64 \leftarrow \text{Height} \rightarrow \text{image} = 64 \text{ px} \\ \xrightarrow{G} 64 \leftarrow \text{Width} \rightarrow \text{image} = 64 \text{ px} \\ \xrightarrow{B} 64 \leftarrow \text{Depth} \rightarrow \text{image} = 64 \text{ px} \end{array}$$

So the image is in color then there is

$$\begin{array}{c} \text{img}_1: \quad \begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \end{bmatrix} \quad \begin{bmatrix} \dots \end{bmatrix} \quad \begin{bmatrix} \dots \end{bmatrix} \quad \dots \\ \text{img}_2: \quad \begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \end{bmatrix} \quad \begin{bmatrix} \dots \end{bmatrix} \quad \begin{bmatrix} \dots \end{bmatrix} \quad \dots \\ \vdots \quad \vdots \quad \vdots \quad \vdots \end{array}$$

So 1 image $\Rightarrow 64 \times 64 \times 3$

Convert the image into $64 \times 64 \times 3$ rows then we have 209 images $\Rightarrow 12288$ rows

we will have

X is the set of all images with 209 rows of 3 matrices of 64×64 flattened vector. Assume 209 images the we will have

$$\Rightarrow X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \\ x^{(4)} \\ x^{(5)} \\ \vdots \\ x^{(209)} \end{bmatrix} \quad \text{we have to set which is first row or 2nd row etc}$$

first row or 2nd row etc

$$\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \\ x^{(4)} \\ x^{(5)} \\ \vdots \\ x^{(209)} \end{bmatrix} \quad \text{but we can do it}$$

$12288 \times 12288 \times 5$ images represented as $A = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \end{bmatrix}$

$$12288 \times 5 \quad \text{matrix}$$

Once you calculate A , remember our $m = 5$

$$A = \frac{1}{1 + e^{-(W^T x + b)}}$$

Now we need to apply the loss function
loss function has to be minimal or zero for our
logistic regression to have worked.

$$W^T x = [w_1 \dots w_{12248}] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ \vdots \\ x_{12248} \end{bmatrix}$$

$$\Rightarrow W^T x = \underbrace{\begin{bmatrix} n^1 & n^2 & n^3 & n^4 & n^5 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}}_{5 \text{ images}} + \underbrace{W \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{n^{28 \times 1}} + \underbrace{b \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{5 \text{ images}}$$

$$W^T x + b = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \end{bmatrix} + b$$

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

$$dw = \cancel{dW^T x} \Rightarrow [A - Y]^T x$$

Suppose $Y = \text{constant value}$
 $A = \text{constant value}$
When you change w by dw then $A - Y$ is constant
 $dw = A - Y$

That is initially $Y = \text{constant value}$
 $A = \text{constant value}$

You change $A - Y = \cancel{dA}$ by dw then

$A - Y = F + \text{the change in } w$

$$\Rightarrow A = \begin{bmatrix} s_0 & s_1 & s_2 & s_3 & s_4 \end{bmatrix}$$

The probability that
 ~~$P(A|X) = 1$~~
The chance the $Y=1$ provided the probability
of A is 1 then the chance that $Y=0$
 $P(A|X) = 1 - A$

$$Y=1 \quad P(A|X)=1$$

$$Y=0 \quad P(A|X)=1-A$$

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

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

$$Y=1 \quad P(A|X)=1$$

Combining the two equations

$$\cancel{P(A|X)} + (1-A) \log(1-A)$$

$$Y \log A + (1-A) \log(1-A)$$

$$J(w, b) = [y_0 - s_0 y_1 - s_1 y_2 - s_2 y_3 - s_3 y_4 - s_4 y_5]$$

Cost $J(w, b) = -\frac{1}{m} \sum_{i=1}^m J(w, b)$ Get it from
we have to minimize the cost with \cancel{A} that
depends on w, x, b

y is the actual value

$$dw = \mathbf{A}^{-1} \cdot (\mathbf{A} - \mathbf{Y})^T$$

forward & backward propagation

$$\mathbf{x} = \begin{bmatrix} x_{1,1} & x_{1,2} & x_{1,3} & x_{1,4} & x_{1,5} \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ x_{2,1} & x_{2,2} & x_{2,3} & x_{2,4} & x_{2,5} \\ x_{3,1} & x_{3,2} & x_{3,3} & x_{3,4} & x_{3,5} \end{bmatrix}$$

$$\text{Given } \Delta_{\text{err}} = [0 \ 1 \ 0.5 \ 0.75 \ 1]^T$$

~~\rightarrow~~

$$\text{All gradient of } \text{err} \text{ will be } 0$$

~~\rightarrow~~

Shape of \mathbf{w} and \mathbf{dw} will be the same. Thus will be repeat above for so many neurons

$$\begin{bmatrix} dw_1 \\ dw_2 \\ dw_3 \\ dw_4 \\ dw_5 \end{bmatrix}$$

$$\begin{bmatrix} dw_1 \\ dw_2 \\ dw_3 \\ dw_4 \\ dw_5 \end{bmatrix}$$

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m \left[\frac{\partial J}{\partial w_1} + \frac{\partial J}{\partial w_2} + \dots + \frac{\partial J}{\partial w_5} \right]$$

$$\frac{dJ}{dw} = \frac{dJ}{dA} \times \frac{dA}{dz} \times \frac{dz}{dw}$$

$$\frac{dA}{dA} = \frac{dA}{dw}$$

$$= \frac{dJ}{dw}$$

which known

$$w = w - \alpha dw$$

loop:

Initialize parameters

$$b = \sum_{i=1}^m (y_i - f(x_i))$$

find for new w, b



When A increase by da
Change of $J(A, b)$ with w

$$\frac{dJ(w, b)}{dw} = \frac{d}{dw} \left(\frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2 \right)$$

When Z increases by dz

A increase. The rate of change of A w.r.t Z is

$$\frac{dA}{dA} = \frac{d}{dA} \left(\frac{1}{1+e^{-x}} \right)$$

$$\Rightarrow \frac{d}{dA} \left(\frac{1}{1+e^{-x}} \right)$$

$$\frac{dA}{dz} = \frac{d}{dz} \left(\frac{1}{1+e^{-x}} \right)$$

$$\frac{dA}{dw} = \frac{d}{dw} \left(\frac{1}{1+e^{-x}} \right)$$

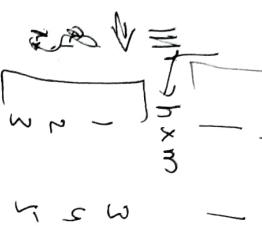
$$= dw^T x + b$$

Neural Network

Activation Function

Similar to sigmoid function or sigmoid activation function

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_m \end{bmatrix}$$



$$z = WX + b$$

$$m = 3 \times 6$$

$$\begin{aligned} W^D &= \begin{bmatrix} w_1^{(1)} \\ w_2^{(1)} \\ \vdots \\ w_3^{(1)} \end{bmatrix} \Rightarrow \begin{bmatrix} \theta_1^{(1)(1)} & \theta_1^{(1)(2)} & \theta_1^{(1)(3)} \\ \theta_2^{(1)(1)} & \theta_2^{(1)(2)} & \theta_2^{(1)(3)} \\ \vdots & \vdots & \vdots \\ \theta_3^{(1)(1)} & \theta_3^{(1)(2)} & \theta_3^{(1)(3)} \end{bmatrix} \\ A^D &= \sigma(Z^D) \end{aligned}$$

$$Z^D = W^D X + b$$

$$X = \begin{bmatrix} 1 & 3 & 6 & 1 & 1 & 1 \\ 2 & 4 & 7 & 2 & 2 & 2 \\ 3 & 5 & 8 & 3 & 3 & 3 \end{bmatrix}$$

$$Z^D = \begin{bmatrix} 1 \times w_1^{(1)(1)} + 2 \times w_1^{(1)(2)} + 3 \times w_1^{(1)(3)} \\ 1 \times w_2^{(1)(1)} + 2 \times w_2^{(1)(2)} + 3 \times w_2^{(1)(3)} \\ \vdots \\ 1 \times w_3^{(1)(1)} + 2 \times w_3^{(1)(2)} + 3 \times w_3^{(1)(3)} \end{bmatrix} = 2 \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_3 \end{bmatrix}$$

$$A^D = \sigma(Z^D)$$

$$X = 3 \times 4$$

$$Z^D = \begin{bmatrix} 1 & 3 & 6 & 1 & 1 & 1 \\ 2 & 4 & 7 & 2 & 2 & 2 \\ 3 & 5 & 8 & 3 & 3 & 3 \end{bmatrix}$$

$$A^D = \sigma(Z^D)$$

$$A^D = \begin{bmatrix} 0.5 & 0.7 & 0.9 & 0.1 \\ 0.6 & 0.8 & 0.9 & 0.2 \\ 0.7 & 0.9 & 0.9 & 0.3 \end{bmatrix}$$

Given x :

$$Z^D = W^D x + b$$

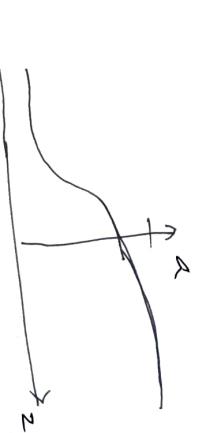
$$A^D = \sigma(Z^D)$$

$$Z^D = W^D x + b$$

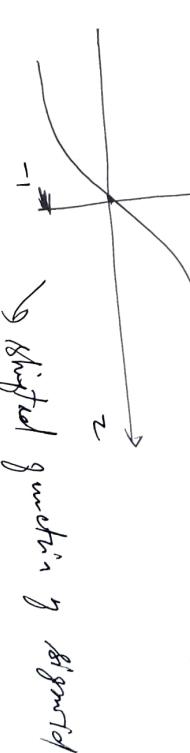
$$A^D = \sigma(Z^D)$$

"S" could be a non-linear function, which could be sigmoid function

$$a = \frac{1}{1+e^{-x}}$$

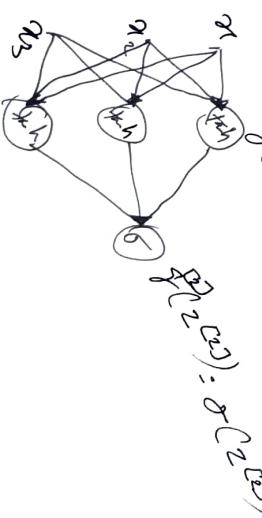


$$a = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



$$a = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

The function is the output layer when the function is binary $\hat{y}^{(0)}$: $\tanh(\hat{y}^{(0)})$



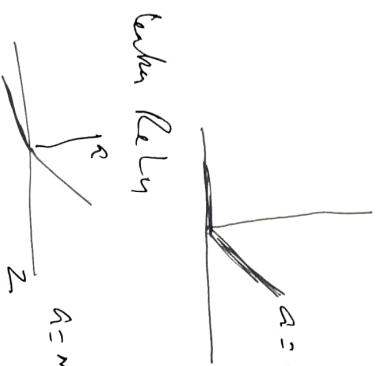
$$\hat{y}^{(0)} = \sigma(Z^{(0)})$$

With my high + low Z values the slope of

the function becomes very smaller.

When the slope of the function becomes very small, it will show down the gradient descent.

The situation is Rectified linear Unit ReLu



$$a^{[L]} = \max(0.012, z)$$

when ReLu

$$\frac{da}{dz} = \max(0.012, z)$$

Forward propagation for layer 1
Input to L1

→ output $a^{[L-1]}, d_w^{[L]}, d_b^{[L]}$

$$d_z^{[L]} = da^{[L]} \cdot g^{[L]'}(z)$$

$$d_w^{[L]} = d_z^{[L]} \cdot a^{[L-1]}$$

$$d_b^{[L]} = d_z^{[L]} \cdot b^{[L-1]}$$

$$da^{[L-1]} = w^{[L]} \cdot d_z^{[L]}$$

$$\frac{da}{dz} = \frac{1}{1+e^{-z}}$$

Note of change in $a^{[L-1]}$ with

$$\frac{da}{dz}$$

Let's change $a^{[L-1]}$ with respect to z = $\frac{da}{dz}$

$$= \frac{d}{dz} (\sigma(z))$$

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

$$dw = dw \cdot x$$

Forward propagation for layer 1
Input to L1

$$Output a^{[L]} = \text{tanh}(z^{[L]})$$

$$z^{[L]} = h^{[L]} \cdot a^{[L-1]} + b^{[L]}$$

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

What are hyperparameters?

Parameters $\mu^{[L]}, \mu^{[B]}, \mu^{[L_2]}, b^{[L_2]}, \mu^{[L_3]}, b^{[L_3]} \dots$

\Rightarrow Hyperparameters: learning rate &

of iterations for gradient descent,

of hidden layers

of hidden units

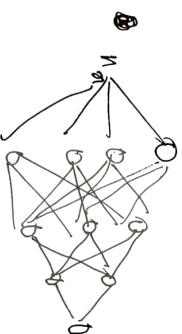
choice of activation functions,

know we chose the values for $n^{[L_i]}$ in a way

called layers

$\mu^{[L]} \rightarrow \text{hyperdim} [n_{[L]}, 4, 3, 2, 1]$

$$\begin{matrix} 1 \\ \alpha^{[L]} \end{matrix} \quad \alpha^{[L]} = \alpha^{[L_1]} \alpha^{[L_2]} \dots \alpha^{[L_L]}$$



A, I

$$X \Rightarrow Z_1 = w^1x + b^1 \rightarrow A = \sigma(Z_1) \rightarrow Z_2 = w^2A + b^2$$

$$\frac{dA^2}{dz^2} \quad A^2 = \sigma(Cz^2)$$

$$\frac{dA^2}{dz^2} \quad \frac{dL}{da^2} \quad L(A^2, y^1)$$

$$\frac{dL}{dw^1} \quad \frac{dL}{db^1}$$

Want

Image

Num-pix, Num-pix, 3
example $4, 4, 3$

$$\left[\begin{matrix} 4 \times 4 & 4 \times 4 & 4 \times 4 \\ \begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 1 \\ 3 & 4 & 1 & 2 \\ 4 & 1 & 2 & 3 \end{bmatrix}, \begin{bmatrix} \quad \end{bmatrix}, \begin{bmatrix} \quad \end{bmatrix} \end{matrix} \right]$$

$$\left[\begin{matrix} 2 \times 2, 2 \times 2, 2 \times 2 \end{matrix} \right]$$

$$\left[\begin{matrix} 4 \times 4, 4 \times 4, 4 \times 4 \\ \begin{bmatrix} \quad \end{bmatrix}, \begin{bmatrix} \quad \end{bmatrix}, \begin{bmatrix} \quad \end{bmatrix} \end{matrix} \right]$$

Here $m = 3$

Bias / Variance:

High variance: look at dev set performance.
 Solution: open data, regularization
 C make more appropriate NN architecture

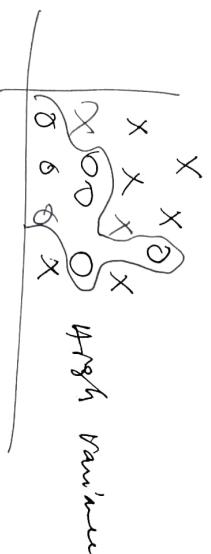
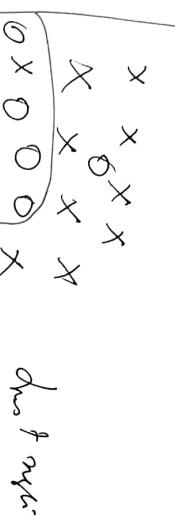
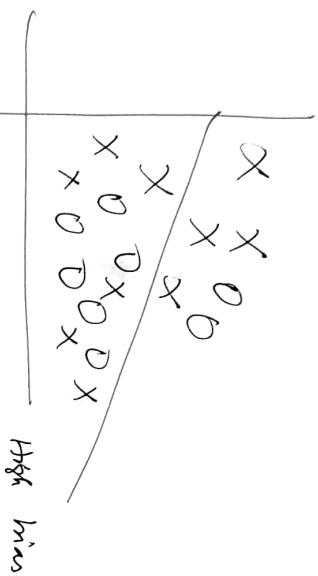
Normal Network Regularization:

$$\frac{\lambda}{2m} \sum_{k=1}^K \|W^{(k)}\|^2$$

$$\|W^{(k)}\|^2 + \|W^{(k+1)}\|^2 + \dots$$

$$\|W^{(K)}\|^2 = \sum_{i=1}^n \sum_{j=1}^m (W_{ij}^{(K)})^2$$

feature norm



High bias? look at the training set

Solution \rightarrow bigger network, hidden units

more features
 (NN architecture that is better trained)

Convolutional Neural Networks:

Edge detection example:

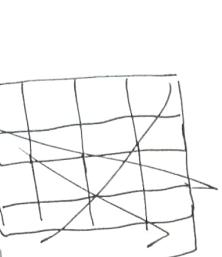
A 6×6 convolved with 3×3 matrix produces a 4×4

To select a 2×2 slice at the upper left corner of a matrix $(5, 5, 3)$, you would do

$$A_{\text{slice_prior}} = A_{\text{prior}}[0:2, 0:2, :]$$

height 0-2
width 0-2

Filters or channels A^{11}



m

$= 1$

n

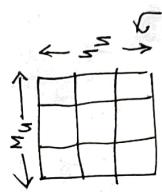
$= 1$

Crochital larger backward pars:..

Compting d'A:

Computing dA:
This is the formula for computing dA
with respect to the area from a certain filter
with given training samples.

$$dA = \sum_{k=0}^n \sum_{\mu=0}^{n_k} w_c \times dZ_{k\mu}$$



This is an
important bridge
that needs to be
crossed.

Object Detection:

Anchor Boxes:

YOLO algorithm:

Triplet losses:

$$J(A, P, N) =$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2$$

+ λ

$$X - \quad z = \text{W} \cdot \text{tanh} \quad y = \sigma(z^2)$$

l1

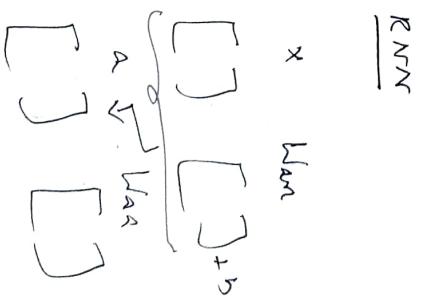
Back propagation

$$x = \begin{bmatrix} 1 & 5 & 2 \\ 3 & 4 & 8 \\ 9 & 8 & 5 \end{bmatrix} \quad w = \begin{bmatrix} 4 & 5 & 6 \\ 8 & 4 & 10 \\ 2 & 4 & 3 \end{bmatrix}$$

$$z = h(x + b)$$

$$h = \sigma(z^2)$$

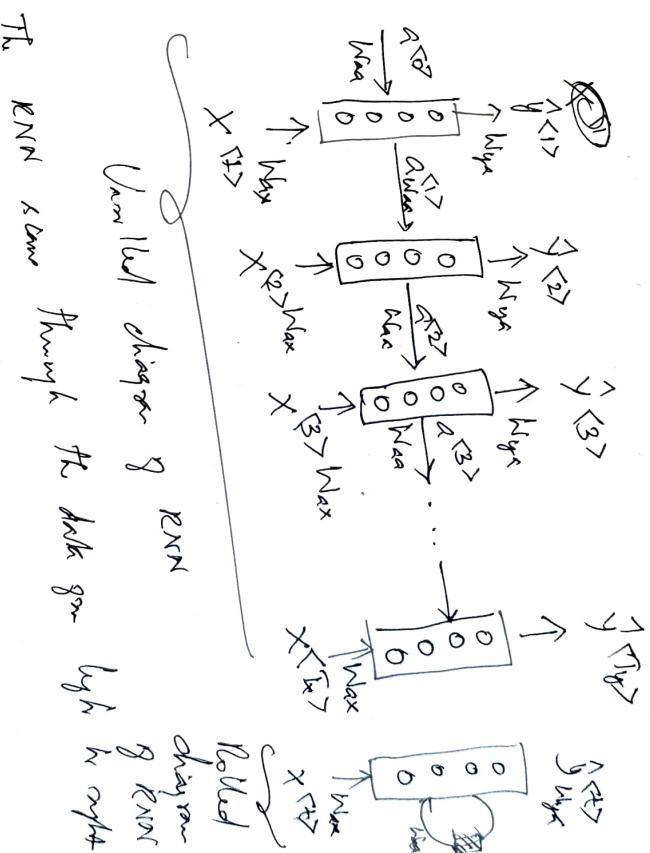
Vanishing gradients in RNN



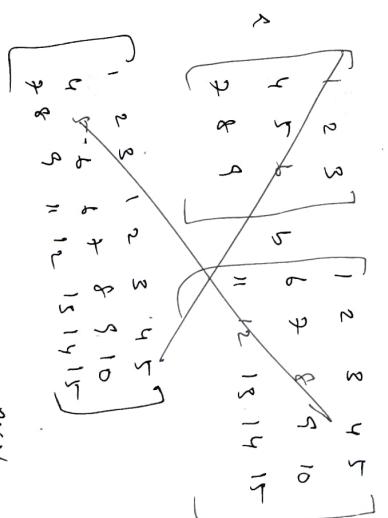
which is

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} + b$$

RNN:
Recurrent Neural Networks



The RNN can store through the data you high in right



RNN

Van

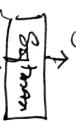
Van

Van

Gated Recurrent Unit (GRU)

RNN unit

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$



↓
Hidden layer of RNN unit

GRU Simplified

$C = \text{memory cell}$

$$C^{(t)} = a^{(t)}$$

$$\tilde{C}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c) \Rightarrow \text{Candidate cell update } C^{(t)}$$

$$C^{(t)} = \sigma(C^{(t-1)}, x^{(t)}) + b_u$$

↳ "update"

$$C^{(t)} = \underbrace{T_u * \tilde{C}^{(t)}}_{\text{Update rule}} + (1 - T_u) * C^{(t-1)}$$

↳ T_u or the forget gate is 1 then

$$\text{but } C^{(t)} = \tilde{C}^{(t)}$$

$$\text{else } C^{(t)} = C^{(t-1)}$$

$$T_u = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$$

$$C^{(t)} = T_u * \tilde{C}^{(t)} + T_f * C^{(t-1)}$$

$$T_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$$

GRU

System $\rightarrow \hat{y}^{(t)}$



$x^{(t)}$

* $C^{(t)}$ can be a vector

Long Short Term Memory - LSTM ::

LSTM-Unit

Arne Hochreiter & Schmidhuber

long short term memory ::

$$\tilde{C}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c)$$

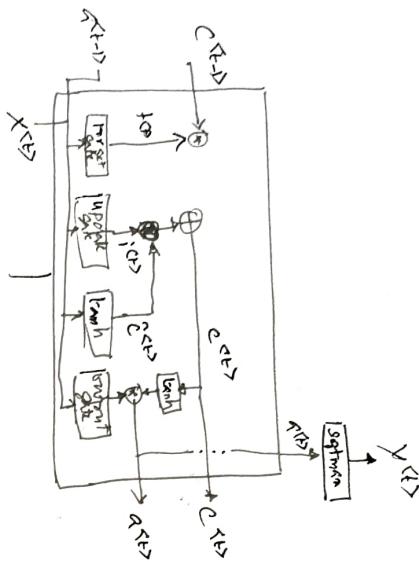
$$C^{(t)} = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$$

$$T_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$$

$$\text{Output Gate } T_o = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

$$C^{(t)} = T_u * \tilde{C}^{(t)} + T_f * C^{(t-1)}$$

$$T_o = T_o * C^{(t)}$$



Embedding:

MILK

Word 2 Vec

Skip 6 lines:
Content
C ("comes,") → T

卷之二

```

graph TD
    O_c[O_c] --> E[E]
    E --> c_c[c_c]
    c_c --> embedding[embedding]
    embedding --> O[O]
    O --> y[y]
    
```

The diagram illustrates a hierarchical embedding process. It starts with an input O_c , which is processed by a function E to produce a feature vector c_c . This vector is then passed through an "embedding" layer, represented by a box containing a circular arrow, to produce a final output O . The output O is then mapped to a label y .

VCC, 1908 Section 2
the right extent
and c

Sigmoid: $P(C|I_c) = \frac{e^{aT_{IC}}}{\sum_{j=1}^{10000} e^{aT_{Ij}}}$

Lancet word

$$p(t|c) = \frac{e^{t \theta_c}}{\sum_{j=1}^{10000} e^{\theta_j t}}$$

GRU is a simpler model, computationally faster
LSTM more powerful & flexible

"*O*" no the parameter associated with output "t".

Wm. French Mr. Lyman

$$\ell(\hat{y}, y) = - \sum_{i=1}^{1000} y_i \log \hat{y}_i$$

The target will be a medium rectangular

✓

How to sample content?
Sample uniformly random from a taxonomy
as oppose

Negative Sampling:

I want a glass of orange juice to go along with my meal

Orange: We are going to predict
Juice is it a content target?

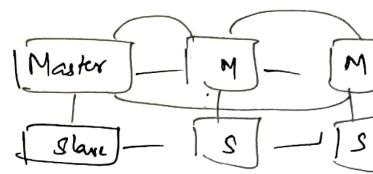
Content word Target?

Orange juice 1

Orange buy 0

(logistic regression model)

$$P(y=1 | c, t) = \sigma(\theta_t^T \phi_c)$$



Cosine Similarity:

To measure how similar two words are, we need to measure the degree of similarity between two embedding vectors for the two words. Given two vectors u and v , cosine similarity is defined as follows:

$$\text{Cosine Similarity}(u, v) = \frac{u \cdot v}{\|u\|_2 \|v\|_2} = \cos(\theta)$$

GloVe Word Vectors:

↪ Global Vectors for Word Representation

I want a glass of orange juice to go along with my meal

$$\partial A \subseteq P$$

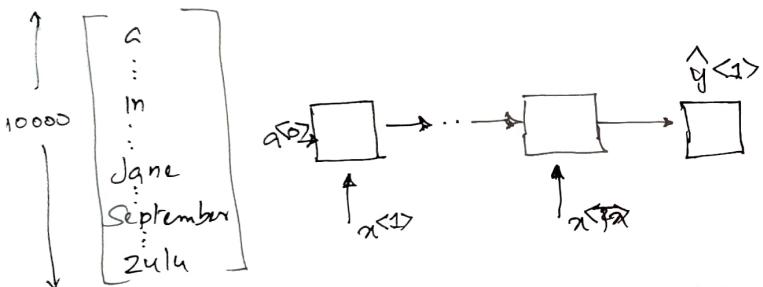
$$P + B = R$$

Sequence to Sequence Models:

Beam Search:

Jane visit Africa in september

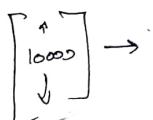
→ Jane is visiting Africa in september



Beam search "B" is beam width. When $b=3$ beam search will consider 3 possibilities at a given time.

$$b = 3$$

⇒



Bleu Score:

Bilingual Evaluation Understudy

Bleu score on bigrams

→ pairs of words appearing next to each other.

Attention Model:

Content is sum of the different features at different time steps weighted by attention weights.

$$\sum_{t'} \alpha^{<1, t'} = 1 \quad C^{<1>} = \sum_{t'} \alpha^{<1, t'} z^{<t'>} \quad \text{Sum of all attention weights should be 1}$$

$$\text{Rep } \alpha^{<t, t'} = \text{amount of attention given} \quad \text{z}^{<t'>} \quad \text{should stay low}$$

$$\alpha^{<t, t'} = \frac{\exp(e^{<t, t'}))}{\sum_{t'=1}^m \exp(e^{<t, t'}))}$$

$$b = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \quad c = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad d = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

$$b+c = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$a = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \quad c = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 \\ 2 \end{bmatrix} - \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

-2

4

$$\overrightarrow{r}, \overrightarrow{s} \quad \overrightarrow{r}, \overrightarrow{s}$$

10

$$\frac{125}{161}$$

$$|\vec{r}| = 5 \quad |s| = \sqrt{|b|}$$

$$\frac{2}{\sqrt{12}} = \frac{1}{6}$$

$$\overrightarrow{r}, \overrightarrow{s} = 10$$

$$\frac{182}{255} * \begin{bmatrix} 3 \\ -4 \\ 0 \end{bmatrix} = \frac{b}{5}$$

$$\begin{aligned} &= \frac{-4}{5} \\ &= 194 \end{aligned}$$

$$a = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

$$a+b \Rightarrow |a+b| = \sqrt{25+25b}$$

$$a+b \Rightarrow |a+b| = \sqrt{25+144}$$

$$|a| = \sqrt{25} \quad |b| = \sqrt{144}$$

$$a = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ -2 \\ 1 \end{bmatrix}$$

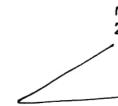
$$a = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad c = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$a = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \quad c = \begin{bmatrix} -3 \\ 1 \\ -2 \end{bmatrix}$$

$$\begin{aligned} &\underline{+} \quad \underline{+} \quad \underline{-} \\ \frac{3}{2} &+ \frac{-3}{-2} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

11:58

12:59



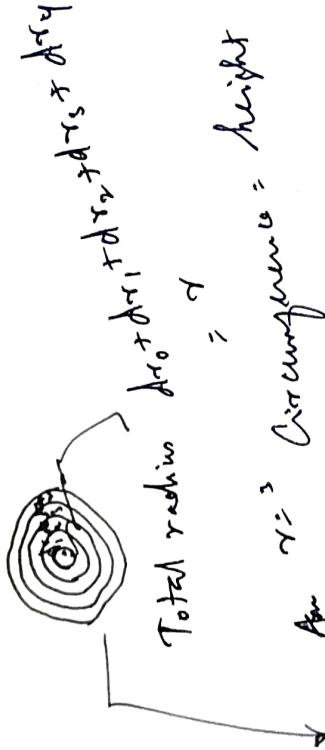
$$\frac{3}{2}$$

2

Calculus:

Arc of a circle πr^2

Diameter
in medium grow
one ring by next



$$r \approx \text{Circumference} = \text{height of the strip} = 2\pi dr$$

$$\bullet 2\pi(0.0)(0.1)$$

$$+ 2\pi(0.1)(0.1)$$

$$+ 2\pi(0.2)(0.1)$$

$$+ \dots$$

$$\approx 2\pi(2.3)(0.1)$$



What is a derivative?

$$(r + dr)^2$$

$$\begin{aligned}
 & (r + dr)^2 = (r + dr)(r + dr) \\
 & = r^2 + 2rdr + dr^2 \\
 & \cancel{r^2} + \cancel{2rdr} + dr^2 \\
 & \Rightarrow r^2 + 2rdr
 \end{aligned}$$

$$r^3 + 3r^2 dr$$

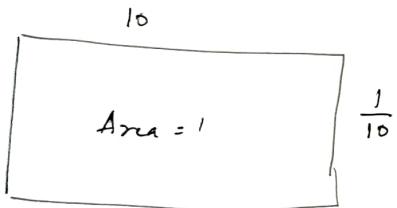
X

$$(x + dx) \times d\left(\frac{1}{x}\right) + \left(\frac{1}{x} - d\left(\frac{1}{x}\right)\right) \times dx$$

$$-x d\left(\frac{1}{x}\right) + dx d\left(\frac{1}{x}\right) + \frac{1}{x} dx - d\left(\frac{1}{x}\right) dx$$

$$-x d\left(\frac{1}{x}\right) - 2dx d\left(\frac{1}{x}\right) + \frac{1}{x} dx$$

What is the derivative of $\frac{1}{x}$



$$\frac{1}{10} - x = \frac{1}{12}$$

$$\frac{1}{10} = \frac{1}{12} + x$$

$$x = \frac{1}{10} - \frac{1}{12}$$

$$\frac{6 - 5}{60} = \frac{1}{60}$$

$$\begin{aligned} \frac{1}{10} - x &= \frac{1}{11} \\ x &= \frac{1}{10} - \frac{1}{11} = \frac{11 - 10}{110} \\ x &= \frac{1}{110} \end{aligned}$$

$$\begin{array}{c} 2 \\ | \\ 10, 12 \\ | \\ 5, 6 \\ | \\ 6 \end{array}$$

$$\begin{array}{c} 1 \\ | \\ 10 \\ | \\ 60 \\ | \\ 5 \\ | \\ 12 \\ | \\ 6 \end{array}$$

$$\frac{6 - 1}{60} = \frac{5}{60} = \frac{1}{12}$$

$$\begin{array}{c} 1 \\ | \\ 10 \\ | \\ 110 \\ | \\ 11 \\ | \\ 10 \end{array}$$

$$\begin{array}{c} 10 \\ | \\ 10, 110 \\ | \\ 11 \\ | \\ 10 \end{array}$$

$$\frac{11 - 10}{110} = \frac{1}{110}$$

What is mean:

Average of the dataset

$$x = \{x_1, x_2, x_3, \dots, x_n\}$$

$$M = \frac{1}{n} \sum_{i=1}^n x_i$$

What is ~~mean~~: standard deviation

The average distance of points or data from the mean.

$$\text{Ex: } \{1, 2, 3, 4, 5\}$$

$$M = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

So standard deviation σ would be: $\{\sqrt{-2, -1, 0, 1, 2}\}$

$$\text{The values have to be absolute} \\ \Rightarrow \sqrt{|-2| + |-1| + |0| + |1| + |2|} = \sqrt{\frac{6}{5}} = 1.2$$

On an average they are all 1.2 away from the mean

↳ What is normalizing?

$$\frac{x - M}{\sigma} \Rightarrow \boxed{z}$$

$$A^{-N}$$

\bullet $\sim \sim \sim \sim \sim$

