

# Machine learning specialization

## Course 1 Machine Learning

WEEK 1

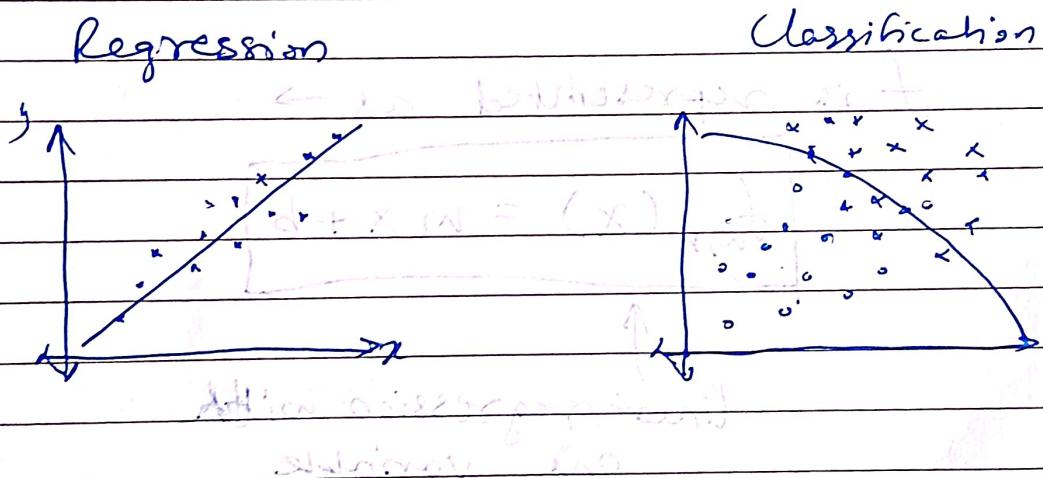
### \* Supervised Machine Learning

algorithms that go from input to output

$$\text{Input } x \rightarrow \text{Output } y$$

→ Regression: predict a number, infinite outputs

→ Classification: Output is limited, predict class



### \* Unsupervised

and words of each word in set

→ finds something in unlabeled data  
not ~~set~~ of words not given in  
(learning) data.

(i) clustering

(ii) anomaly detection

(iii) Dimensionality reduction

$$\{(A), (B)\} \xrightarrow{\text{?}} \{A, B\}$$

notes on linear regression

training set (normal)

↓

learning algorithm

Two types of learning algorithms

$$x \rightarrow f \rightarrow y \text{ or } \hat{y} \text{ "y-hat"}$$

Two types of learning algorithms:  
Stochastic gradient descent  
and Gradient descent (two iterations) ←  
models prediction

(stochastic)

(gradient)

$f$  is represented as →

$$f_{w,b}(x) = w x + b$$

linear regression with  
one variable

brain response

Now we know how to draw the line,

now we need to minimize error. ←

i.e. get  $\hat{y}^{(i)}$  closer to  $y^{(i)}$  for  
all  $(x^{(i)}, y^{(i)})$

$J$  is the cost function. ←

$J$  gives the cost of a line ←

minimizing the function ←

Cost function: Squared error cost function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

$M$  = no. of training examples

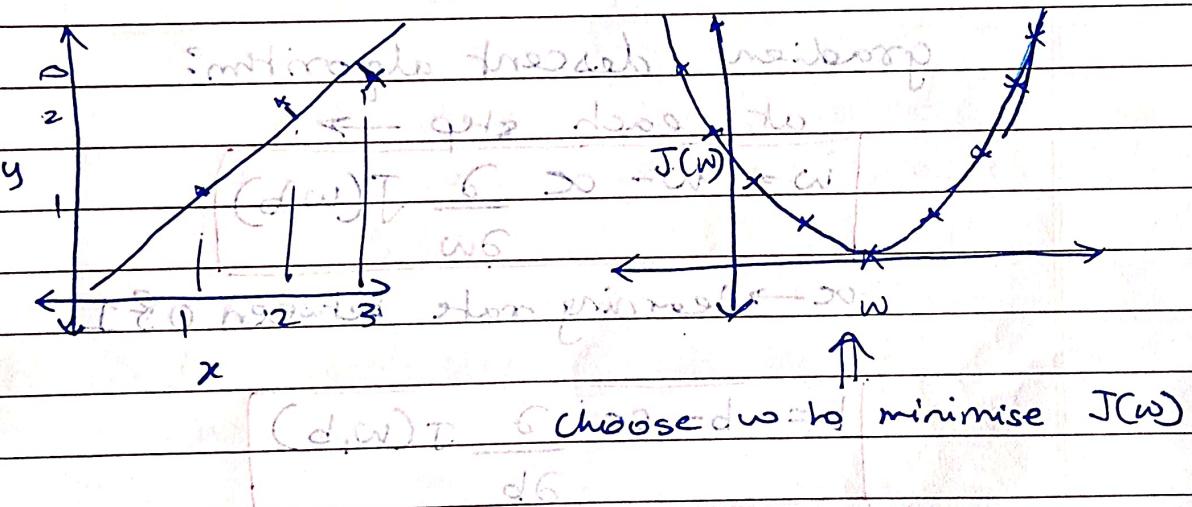
$$y^{(i)} - \hat{y}^{(i)} \rightarrow \text{error}$$

then we square it and sum up all the for the moment we divide by 2M 2 for newer calcs.

(can also be handwritten as  $\text{f}(\text{f}(x))$  or  $\text{f}^2(x)$ )

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$f_w(x) \rightarrow J(w)$$



# INTUITION: Wie ist es möglich intuitiv zu denken?

Now speak French at few times

$$\text{Model: } f_{w,b}(x) = w_1x + b$$

: spw n: northmonegami, Yosra)

Parameters:  $T, w, b$  do ->  $w$  = weight

$$\text{Cost function} = J(\vec{w}; b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2$$

Goal : minimize  $J(w, b)$

## \* Gradient Descent

→ You can apply it to any function, to minimize that function. It's not just losses, rather not  $J(w, b)$ .

### Outline:

Start with some  $J(w, b)$  and do me

Keep changing  $w, b$  to reduce  $J(w, b)$

Until we settle at or near a minimum.

(a) there can be  $> 1$  minimum (b)

gradient descent algorithm:

at each step →

$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$\alpha \rightarrow$  learning rate between 0 & 1

$$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

Keep updating  $w$  &  $b$  simultaneously until  $w$  &  $b$  don't change much.

correct implementation in code:

$$\text{temp\_}w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$\text{temp\_}b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

$$w = \text{temp\_}w$$

$$b = \text{temp\_}b$$

Lecture 4

$$\frac{\partial}{\partial w} J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2 / x^{(i)}$$

$$\frac{\partial}{\partial w} J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

↑ notice how after derivative,

the square over the bracket  
is gone.

At each step we can (say),

$$w = w - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$b = b - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

This gradient descent is also called "batch" gradient descent. At each step, or gradient descent makes use of all training examples.

## Week 2

\* Linear Regression with multiple variables.

$$(y - (\vec{w} \cdot \vec{x}))^2 \leq 1.6 = (\vec{w}, \vec{x})^T \vec{w}$$

Col1	Col2	Col3	Col4	Output
$x_1$	$x_2$	$x_3$	$x_n$	$y$
$(\vec{w} \cdot \vec{x}) + b$	$\leq 1.6$	$= (\vec{w}, \vec{x})^T \vec{w}$		

$x_j$  = j<sup>th</sup> feature

n = number of features

$\vec{x}^{(i)}$  = features of i<sup>th</sup> example

$$\vec{x}^{(2)} = [1, 4, 16, 3, 2, 40]$$

This is a vector

$x_j^{(i)}$  = value of feature j of i<sup>th</sup> example.

$$x_3^{(2)} = 2$$

scalar

Previously:  $f_{w,b}(x) = w \cdot x + b$

Now:

$$f_{w,b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$$\vec{w} = [w_1, w_2, w_3, \dots, w_n]$$

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

b is a number

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

dot product  
of vectors

Code implementation of  $f(\vec{w}, \vec{x})$  + error analysis

$$\text{Let } \vec{w} = [1.0 \ 2.5 \ -3.3] \quad b = 4 \\ \vec{x} = [10 \ 20 \ 30]$$

We can initialize these matrices or vectors in python using NumPy as:

~~using below code we can initialize~~  $w = np.array([1.0, 2.5, -3.3])$

~~using below code we can initialize~~  $b = 4$

~~using below code we can initialize~~  $x = np.array([10, 20, 30])$

Now,  $f_{\vec{w}, b}(\vec{x}) = w_1x_1 + w_2x_2 + w_3x_3 + b$

so (WITHOUT VECTORISATION) python code

~~defining f to find out if i.e.~~  $f = w[0]*x[0] +$

~~w[1]\*x[1] + w[2]\*x[2] + b~~

~~write new line for all values~~

~~else or if i.e.~~  $f_{\vec{w}, b}(\vec{x}) = \left( \sum_{j=1}^n w_j x_j \right) + b$

WITHOUT VECTORISATION:

$$f = 0 \quad \text{(initialising 0)}$$

~~for j in range(0, n)~~

$$f = f + w[j]*x[j]$$

~~really slow~~

Now using Vectorisation

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

One line  
+ much faster

Code:

$$[0.8 \ 0.5 \ 0.1] = 5.5$$

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

np.dot() uses vectorisation in the computer's hardware to run multiple calculations in parallel. Thus, significantly reducing compute time.

### \* Gradient Descent: Implementation

$$\vec{w}' = (\alpha w_1, \alpha w_2, \dots, \alpha w_3, \dots, \alpha w_{16})$$

$$\vec{d} = (d_1, d_2, \dots, d_{16}) \leftarrow \text{derivatives}$$

$$w = np.array([0.5, 1.3, \dots, 3.4])$$

$$d = np.array([0.3, 0.2, \dots, 0.1])$$

compute  $w_j = w_j - 0.1 * d_j$  for  $j = 1 \dots 16$

W/O vectorisation:

(for  $j$  in range(0, 16):

$$w[j] = w[j] - 0.1 * d[j]$$

$$d[j] = ?$$

plus  
create

With vectorisation:

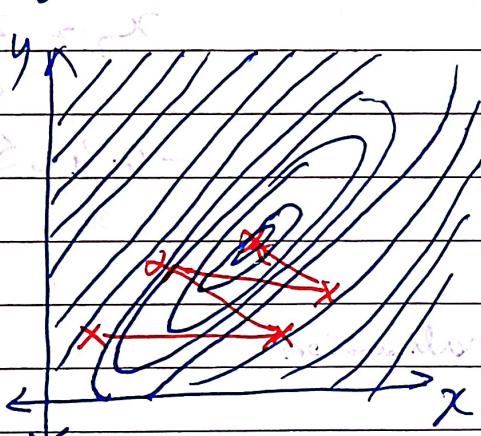
initial priors also somewhat makes no sense

$$w = w - 0.1 * d \quad \leftarrow \text{descreases for all values of } w$$

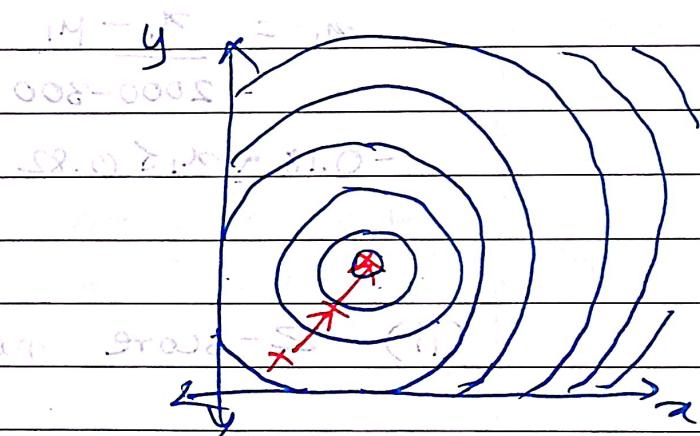
(gradient descent and update parallelly)

# Feature Scaling:

If difference between parameters is too too high, gradient descent could take a lot of time



Skinny Contour plot  
Inefficient for GD



Round Contour plot  
Efficient for GD

hence we normalize incoming inputs  
between 0 to 1 so that it's more efficient.

for eg)  $300 \leq x_1 \leq 2000$

$$x_1' = \frac{x_1}{2000}$$

feature scaling

$$0.15 \leq x_1' \leq 1$$

$$x_2' = \frac{x_2}{5}$$

$$0 \leq x_2' \leq 1$$

we can also feature scale using other methods.

(i) Mean normalisation:

$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$\mu_1$  = average

$\mu_2 = 2.3$  (assumption)

foot note 2:  $\approx 600$  (assumed) of constraint 41

if  $x_1$  has a value below  $\mu_1$  then bring

$$x_1 = x_1 - \mu_1$$

$$x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

$$2000 - 300$$

$$\approx 5 - 0$$

$$-0.18 \leq x_1 \leq 0.82$$

$$-0.46 \leq x_2 \leq 0.54$$

(ii) Z-score normalisation:

why we do?  $\sigma$ : standard deviation

foot note 3

(i) not position

$\sigma_1$ : SD of 1<sup>st</sup>

$\sigma_2$ : SD of 2<sup>nd</sup>

z-score formula:  $x_1 - \mu_1$  or  $x_2 - \mu_2$

constraint  $x_1 = \mu_1$  or 1 of 0  $x_2 = \mu_2$

$$\sigma_1$$

$$\text{the same } \sigma_2$$

$$-0.67 \leq z \leq 3.1$$

$$0.005 > \text{range} (-1.6 \leq x_2 \leq 1.9)$$

$$\downarrow$$

$$\frac{\sigma_1}{\sigma_2} = 5$$

$$1 \geq x_1 \geq 0$$

$$\downarrow$$

$$\frac{\sigma_2}{\sigma_1} = 15$$

$$\downarrow$$

$$1 \geq x_2 \geq -1.0$$

envelope  
rule

# In feature scaling, aim for about

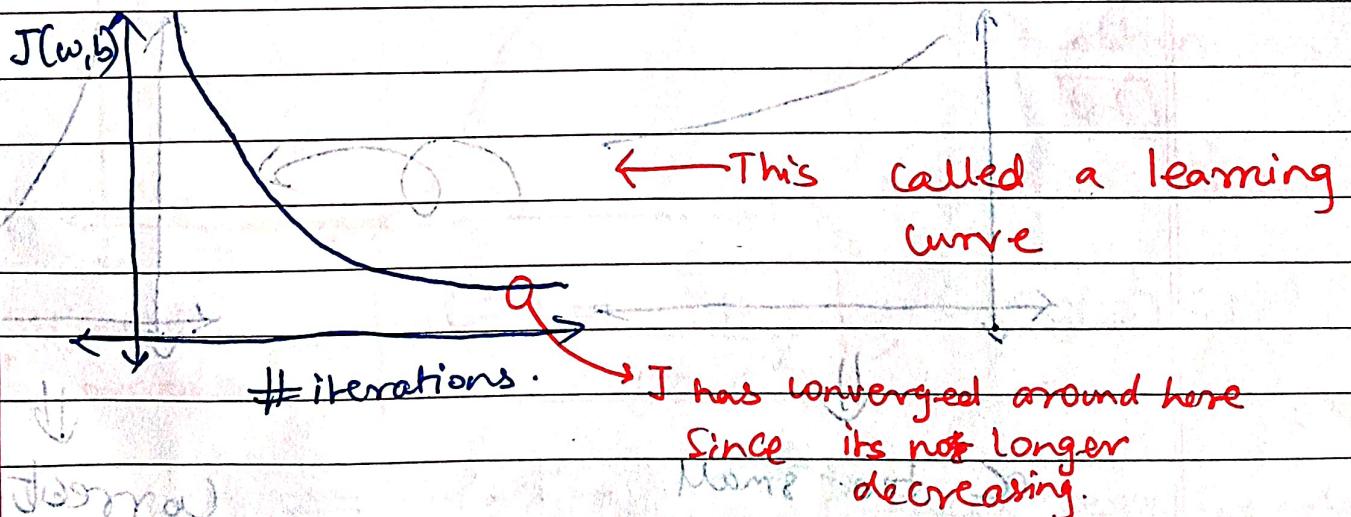
$$\begin{aligned} -1 \leq x_j &\leq 1 \\ -3 \leq x_j &\leq 3 \\ -0.3 \leq x_j &\leq 0.3 \end{aligned} \quad \left. \begin{array}{l} \text{acceptable} \\ \text{ranges.} \end{array} \right\}$$

$$\begin{aligned} 0 \leq x_j &\leq 3 \\ -2 \leq x_j &\leq 0.5 \end{aligned} \quad \left. \begin{array}{l} \text{Okay, no rescaling} \\ \text{needed.} \end{array} \right\}$$

$-100 \leq x_2 \leq 100 \rightarrow \text{too large, rescale}$

$-0.0001 \leq x_4 \leq 0.0001 \rightarrow \text{too small, rescale.}$

→ To make sure GD is working correctly  
plot a graph of cost function  $J(w, b)$   
after each iteration.

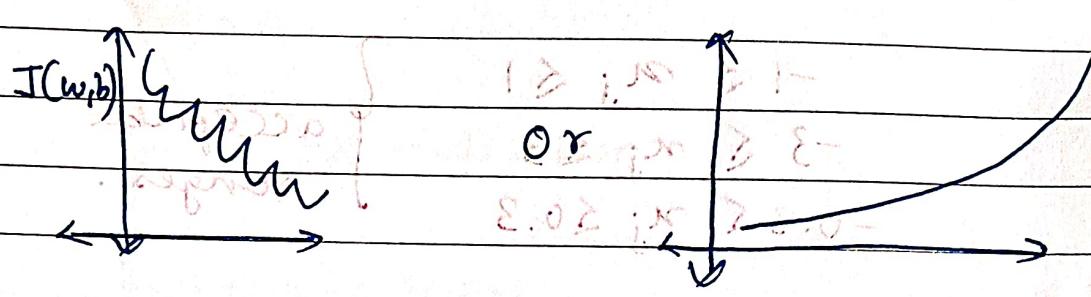


→ Automatic convergence test

$$\text{let } \epsilon = 0.001$$

if  $J(w, b)$  decreased by  $\leq \epsilon$   
then  $J$  has converged.

## # selecting $\alpha$



→ either  
big or

$$w_i = w_i + \alpha d_i$$

→  $\alpha$  is too large ← use minus sign

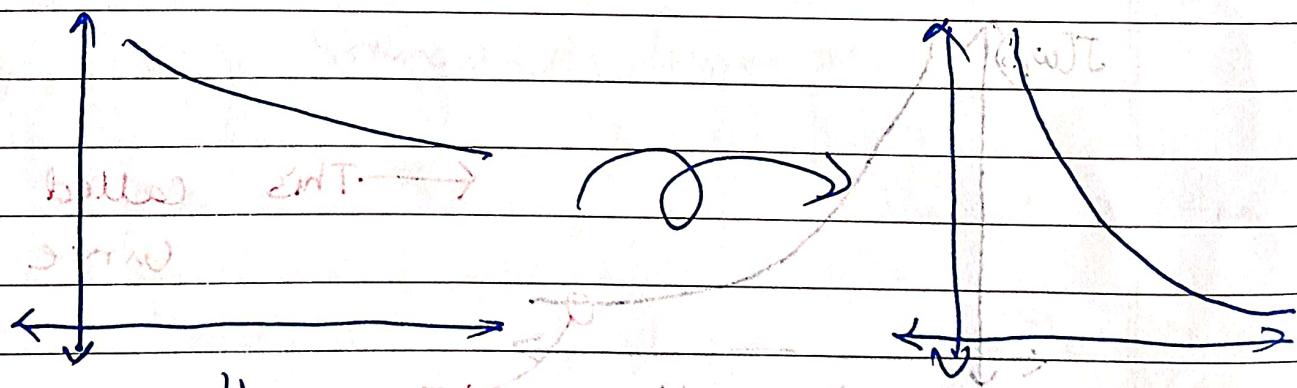
too large  $w_i = w_i - \alpha d_i$ ,

→ or  $\alpha$  is too large

Ideal values of  $\alpha$  for the start of S

(d<sub>i+1</sub>)<sub>T</sub> mentioned → sharp & flat

... 0.001, 0.01, 0.1, 1, 10, ...



so from graph 1.5 & T is written off

incorrect alpha

$\alpha$  too small

Correct  $\alpha$

$$\begin{aligned} 7.29 & \text{ corresponds to } \alpha = 0.001 \\ 100.0 & = 3 \cdot 49 \end{aligned}$$

$\alpha$  is chosen so that  $(d_{i+1})_T$  is

as small as possible and  $T$  is next

## WEEK 3

\* Classification

→ limited output options for "y"

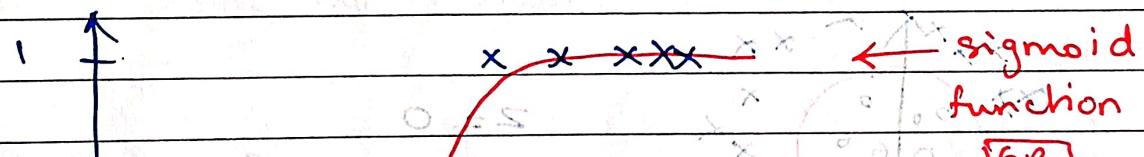
binary classification: y can only be one of 2 values

→ Binary classification problems are solved using  
Logistic regression.

→ although it is named  
in logistic regres<sup>sion</sup>, it is still  
output is either 0 or 1

→ linear main S line hit karne hai

logistic mein S-shaped curve



[GR]

logistic function

Sigmoid  
f<sup>n</sup> formula

$$g(z) = \frac{1}{1 + e^{-z}} ; 0 < g(z) < 1$$

Let  $z = \vec{w} \cdot \vec{x} + b$  → linear regr. formula

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

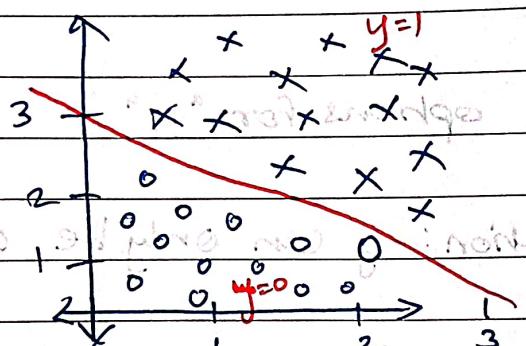
(error)  
loss

using logistic formula

→ Its output is the probability that  
class is 1.

## # Decision Boundary

android 128GB



→ It is the boundary

$$z = \vec{w} \cdot \vec{x} + b = 0$$

→ which means the result is ~~zero~~.

prior belief  $Z=0$  at a decision with a neutral value is boundary. (neither 0 nor 1).

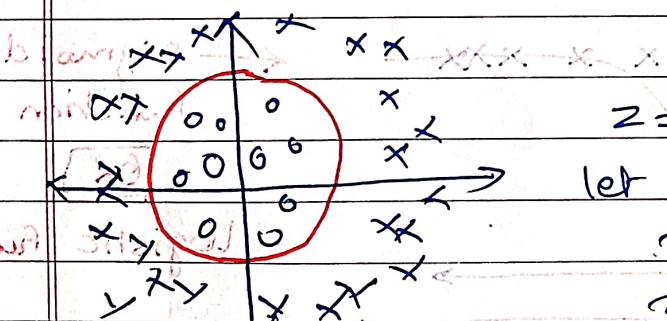
Bernard, Dr. H. J. Green, Mrs. G.

$$4.162 \cdot \vec{x} + b = 0$$

• Odele:  $\text{and the SVD is } \vec{w} \cdot \vec{x} + b = 0$  (so 0 weight in the sum)

$$x_1 + x_2 - 3 = 0$$

$$5\pi/3 \quad \text{Locality} + d_2 = 3\pi/3 \quad \text{Dilution}$$



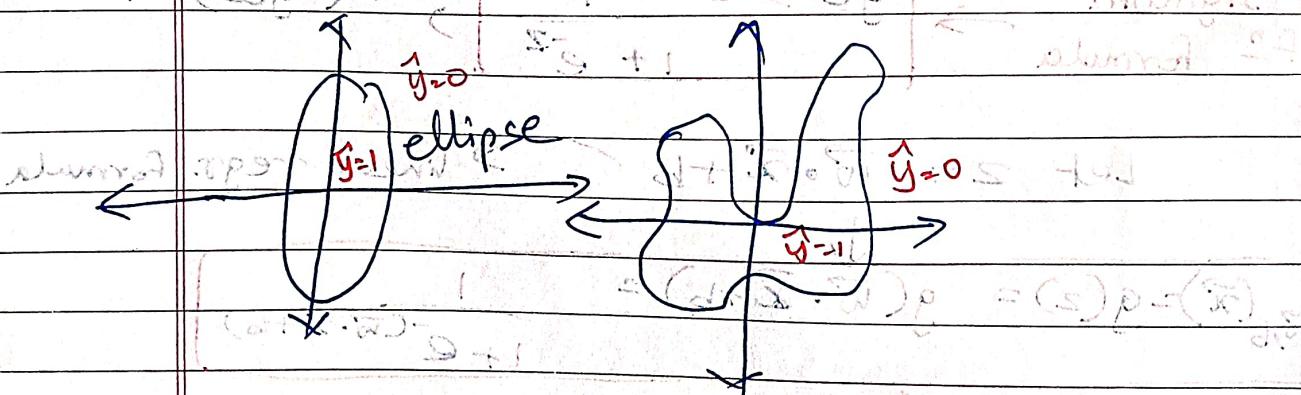
卷之三

$$\text{let } z = w_1 x_1 + w_2 x_2 + b$$

$$x_1^2 + x_2^2 - 1 = 0$$

$$x_1^2 + x_2^2 = 1 \leftarrow \text{circle of}$$

## Radius 1



more complex boundaries can be formed

## Cost function for logistic regression

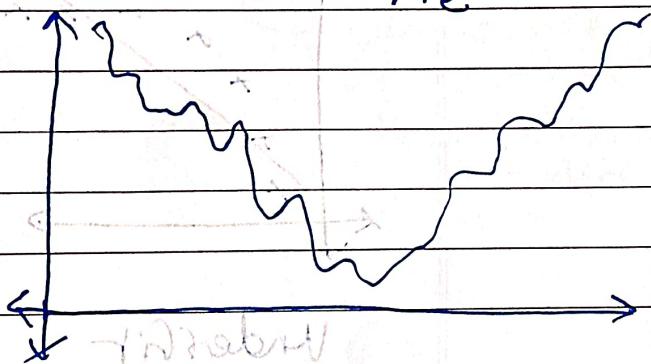
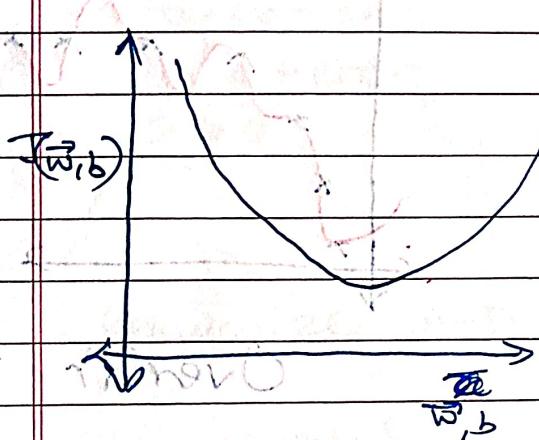
→ The square error cost function does not work → loss of linear

$$\text{linear loss: } J(w, b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^{(i)}) - y^{(i)})^2$$

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

logistic regression

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



Only one minima always multiple local minima.

using the same cost  $f$  is not ideal.

→ loss:  $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$  ↑ loss  $f$  for linear reg.

logistic loss function:

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} \text{min} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ \text{max} \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

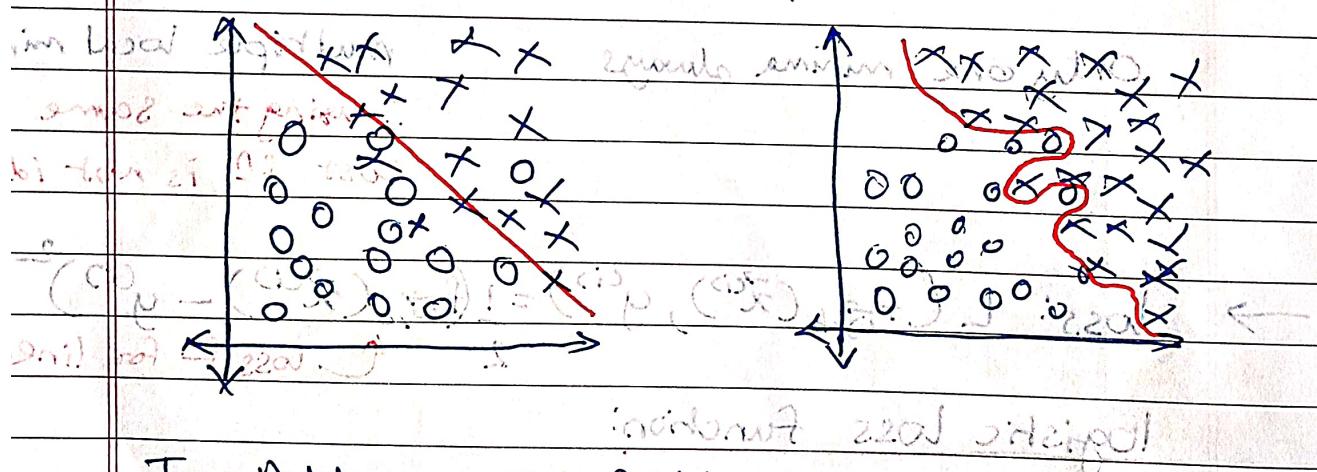
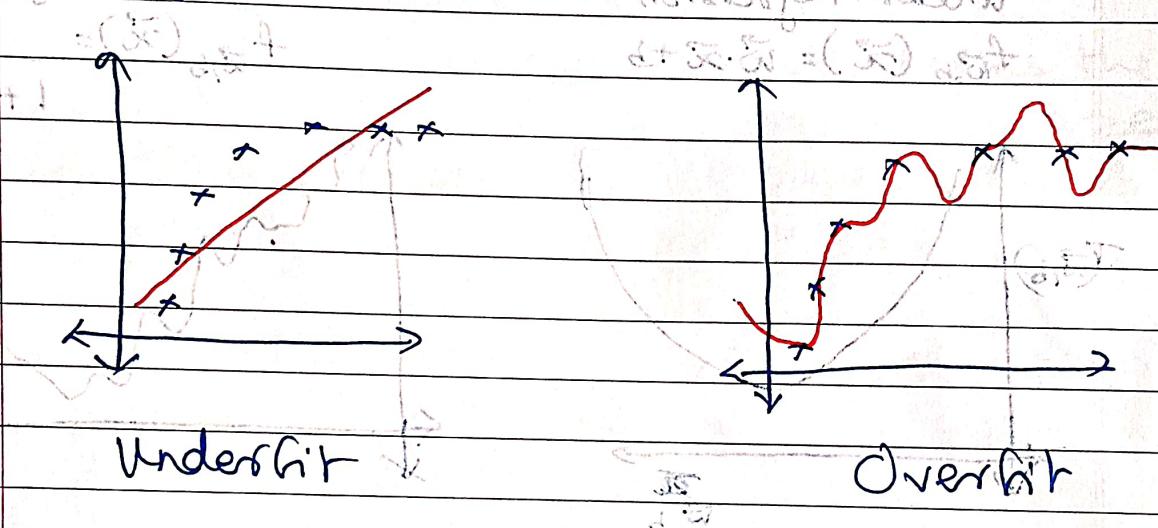
$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

Converts to convex. y can only be 1 or 0 in binary classification.

plug this loss  $f$  in  $J(w, b) = \frac{1}{m} \sum L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

→ Overfitting: When you fit the training set so well that the model gets worse at predicting values for new data

→ Underfit: When the model uses less features to fit data and results to high error



To Address Overfitting:

1) Collect more data

2) Select only a subset of features

3) Reduce size of parameters

- Regularization

$$(X^T X)^{-1} X^T Y = \hat{\beta}$$

## \* Regularization and cost function loss function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \text{Mean squared error}$$

$$\text{regularisation term} = \frac{\lambda}{2m} \sum_{j=1}^m w_j^2 + \frac{\lambda}{2m} b^2$$

$\lambda$ : regularization

a: parameter of regularization  
 $\lambda > 0$  or main reason not doing that  
 can include or exclude b bcs it makes very less diff.

## \* Regularized GD

repeat {

$$w_j = w_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m [(f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)}] + \frac{\lambda}{m} w_j \right]$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

} simultaneous update.

we don't add anything bcs we aren't changing how L is updated

\* Regularized cost  $J^h$  for logistic reg.

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{w,b}(\vec{x}^{(i)})) + (1-y^{(i)})$$

$$\log(1-f_{w,b}(\vec{x}^{(i)}))]$$

$$+ \frac{\lambda}{2m} \sum_{j=1}^n w_j^2]$$

→ Regularized GD for logistic reg. is  
the same for linear reg.

GD loss function

} loss func

$$+ \left[ \frac{\lambda}{m} (w_1^2 + w_2^2 + \dots + w_n^2) \right] \quad \text{if } \lambda > 0 \Rightarrow \text{GD loss} = J(w)$$

$$(m_y - (m_x)) \quad \text{if } \lambda = 0 \Rightarrow \text{GD loss} = J(w)$$

• analog loss function

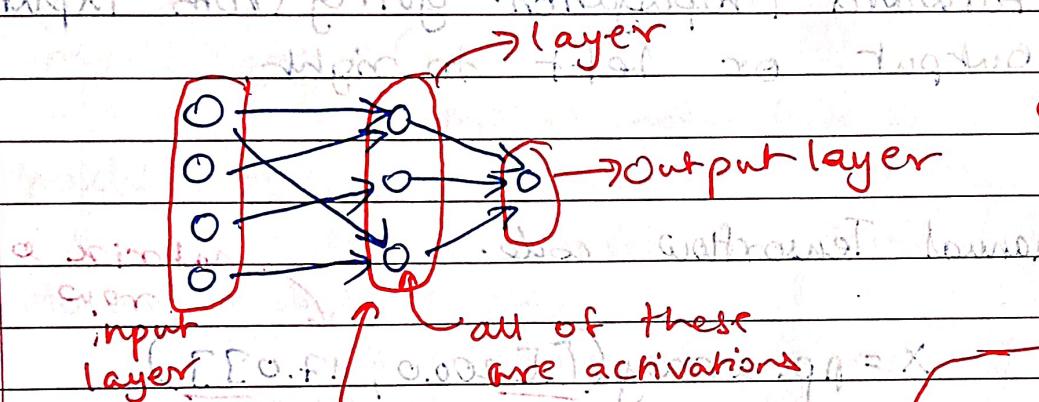
# Course 2 week 1

## Advanced Learning Algorithms

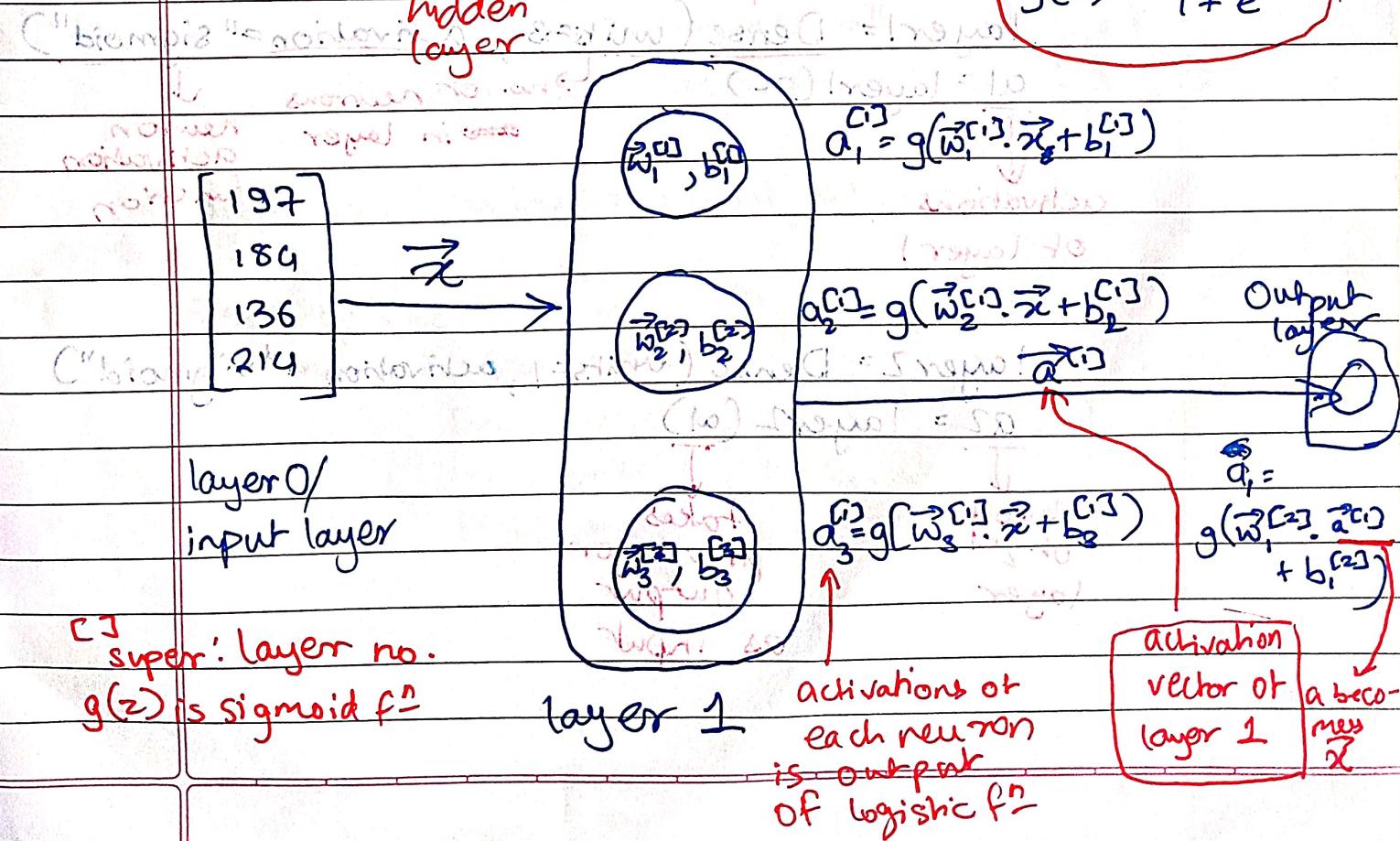
a neuron  $\rightarrow \circ$

input  $\rightarrow \circ \rightarrow$  probability

all neurons use a function for output  
so "neurons are just a bunch of logistic regression models bundled together"



Each layer is a vector



## # General Notations of one neuron's activation

$$a_j^{[l]} = g(w_j^{[l]} \cdot \vec{a}^{[l-1]} + b_j^{[l]})$$

$j^{\text{th}}$  neuron of  $l^{\text{th}}$  layer

$g(\cdot)$ : activation function

Forward pass

# Forward Propagation: going from input to output or left to right.

\* Manual Tensorflow code:

$$x = np.array([200.0, 17.0])$$

creates layer

layer1 = Dense(units=3, activation="sigmoid")

$$a1 = layer1(x)$$

no. of neurons

in layer

neuron activation

function

activations  
of layer1

$$a2 = layer2(a1)$$

Output  
of 2nd  
layer

takes  
prev layer  
output  
as input

Output  
layer

on input: input  
"A" shape 2x1(s)

## ~~Course~~ Course 2 week 2

### \* Auto tensorflow code

```
layer_1 = Dense(units=3, activation="sigmoid")
```

```
layer_2 = Dense(units=1, activation="sigmoid")
```

```
model = Sequential([layer_1, layer_2])
```

↳ joins layers automatically.

```
x = np.array([[200.0, 17.0], [212.0, 18.0]])
```

```
[.200.0, 17.0], [212.0, 18.0]])
```

```
y = np.array([1, 0, 0, 1])
```

```
model.compile(...)
```

↳ compiles model layers

```
model.fit(x, y)
```

↳ train model on given inputs & o/p

```
model.predict(x_new)
```

↳ predict o/p of new x input.

```
(x) model = Sequential([
```

```
Dense(units=3, activation="sigmoid"),
```

```
Dense(units=1, activation="sigmoid")])
```

### \* Steps in training a Neural Network model

#### Step 1: Make a model

```
model = Sequential([
```

```
Dense(units= , activation= )
```

```
Dense(units= , activation= )
```

```
Dense(units= , activation= )])
```

## Step 2: Compile / specify loss & cost \*

Binary cross entropy:

model.compile(loss="binary\_crossentropy") uses the same loss as in BinaryCrossentropy(); MeanSquaredError() is binary classification

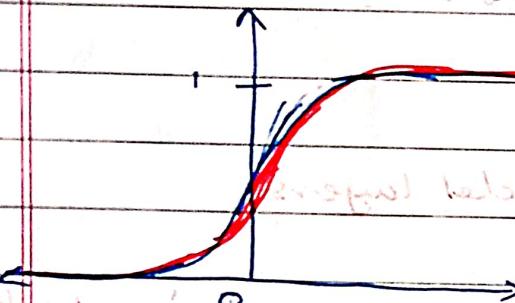
Step 3: fit model / minimize cost = fit()

fit(x, y, epochs=1000)

model.fit(x, y, epochs=1000)

## \* Types of activation functions

(ReLU, tanh, ...)



(ELU, ReLU) sigmoid:  $g(z) = \frac{1}{1 + e^{-z}}$

(...) sigmoid function

$$\text{sigmoid function: } g(z) = \frac{1}{1 + e^{-z}}$$

not sigmoid having no bias, range: 0 to 1  
(ReLU) binary function

ReLU: (rectified linear unit)

(rectified linear unit)

(\* binary function, e.g. if)  $g(z) = \max(0, z)$

if  $z < 0$  if  $z > 0$

$g(z) = 0$   $g(z) = z$  (bias straight forward is prioritized in backprop)

linear activation f:

1) identity:  $g(z) = z$

2) non-linear:  $g(z) = \vec{w} \cdot \vec{x} + b$

= identity:  $g(z) = z$

= identity:  $g(z) = z$

Binary step fn:

range: 0 to 1

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

if  $z < 0$ : 0 if  $z > 0$

$$g(z) = 0 \quad g(z) = 1$$

tanh (hyperbolic tangent)

range: -1 to 1

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

Softmax:

range: 0 to 1

$$g(z) = \frac{e^z}{\sum_{i=1}^n e^{z_i}}$$

Choosing the right fn based on

(i) Regression: Linear activation

(ii) Binary classification: Sigmoid

(iii) Multiclass classification: Softmax

(iv) Multilabel classification: sigmoid

(v) CNN: ReLU

(vi) RNN: Tanh or sigmoid

## \* Cost of logistic & Softmax

logistic regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x})$$

$$a_0 = g(z) = 1 - a_1 = P(y=0|\vec{x})$$

$$\text{loss} = -y \log a_1 - (1-y) \log (1-a_1)$$

$$L = -y \log a_1 - (1-y) \log a_2$$

$J(\vec{w}, b)$  = average loss

softmax regression

$$a_i = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$0 \leq a_i \leq \frac{e^{z_i}}{e^{z_1} + e^{z_2} + e^{z_3} + \dots + e^{z_N}}$$

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

$$a_N = \frac{e^{z_N}}{\sum e^{z_i}}$$

$$e^{z_1} + e^{z_2} + e^{z_3} + \dots + e^{z_N}$$

$$P(y=N|\vec{x})$$

loss ( $a_1, \dots, a_N, y$ )

$$= \begin{cases} -\log a_1 & \text{if } y=1 \\ -\log a_2 & \text{if } y=2 \\ \vdots \\ -\log a_N & \text{if } y=N \end{cases}$$

## \* MNIST with softmax regression

import tensorflow as tf

from tensorflow.keras import Sequential

from tensorflow.keras.layers import Dense

mnist = keras.datasets.mnist

model = Sequential()

model.add(Dense(units=256, activation='relu'))

model.add(Dense(units=10, activation='relu'))

Dense (units=10, activation='softmax'))

binary crossentropy

Step 1  
pecifying  
model

Step 2

from tensorflow.keras.losses import

~~losses~~ SparseCategoricalCrossentropy()

~~import tensorflow as tf~~ # Import TensorFlow and its submodules under tf

model.compile(loss=SparseCategoricalCrossentropy())

Step 3

model.fit(x, y, epochs=100)

Model training completed in 20s

# Improved Implementation of softmax

Since last layer is linear instead of softmax

Dense (units=10, activation='linear')

→ use from\_logits

model.compile(sparse\_categorical\_crossentropy(

from\_logits=True)

→ predict not using predict()

model.predict(x).argmax(1) # Prediction

z = model(x)

softmax = tf.nn.softmax(z)

(softmax - original softmax)

## \* Adam algorithm

François Fleuret (2011). Some optimization methods

→ Adam is an adaptive learning rate algorithm designed to improve training speeds in deep neural networks and reach convergence quickly.

→ It's an alternative to gradient descent and works faster than gradient descent.

→ If  $w_j$  (or  $b_j$ ) keeps moving in the same direction, increase  $\alpha_j$ .

→ If  $w_j$  (or  $b_j$ ) keeps oscillating, reduce  $\alpha_j$ .

→ MNIST Adam optimizer

model.compile(optimizer='adam')

optimizer = tf.keras.optimizers.Adam(

( $\times$ ) learning\_rate = 1e-3),

( $\times$ ) loss = tf.keras.losses.SparseCategoricalCrossentropy(

from\_logits=True))

## \* Additional layer types

### (i) Convolutional Layer

Creates a smaller gram (2x2) from original

→ Each neuron only looks at part of the previous layer's outputs

Why?

→ faster computation

→ Need less training data  
(less prone to overfitting)

→ A neural network made using these layers is called Convolutional Neural Network (CNN)

## Course 2 week 3

### Deciding what to try! Nescitabot

#### Options

- high variance → Get more training examples
- high variance → Try smaller set of features
- high bias → Try getting additional features
- high bias → Try adding polynomial features
- high bias → Try decreasing  $\lambda$
- high variance → Try increasing  $\lambda$ 
  - anything else?
  - (writing more or writing less)

Selecting what to do to debug ur model

It is important about underfitting & overfitting

(overfitting) (underfitting) (middle)

(i) If  $J_{train}(\vec{w}, b)$  is low but  
 $J_{test}(\vec{w}, b)$  is high,  
 Then the model is overfit.

→ Train, test, valid are 3 sets to divide our data into.

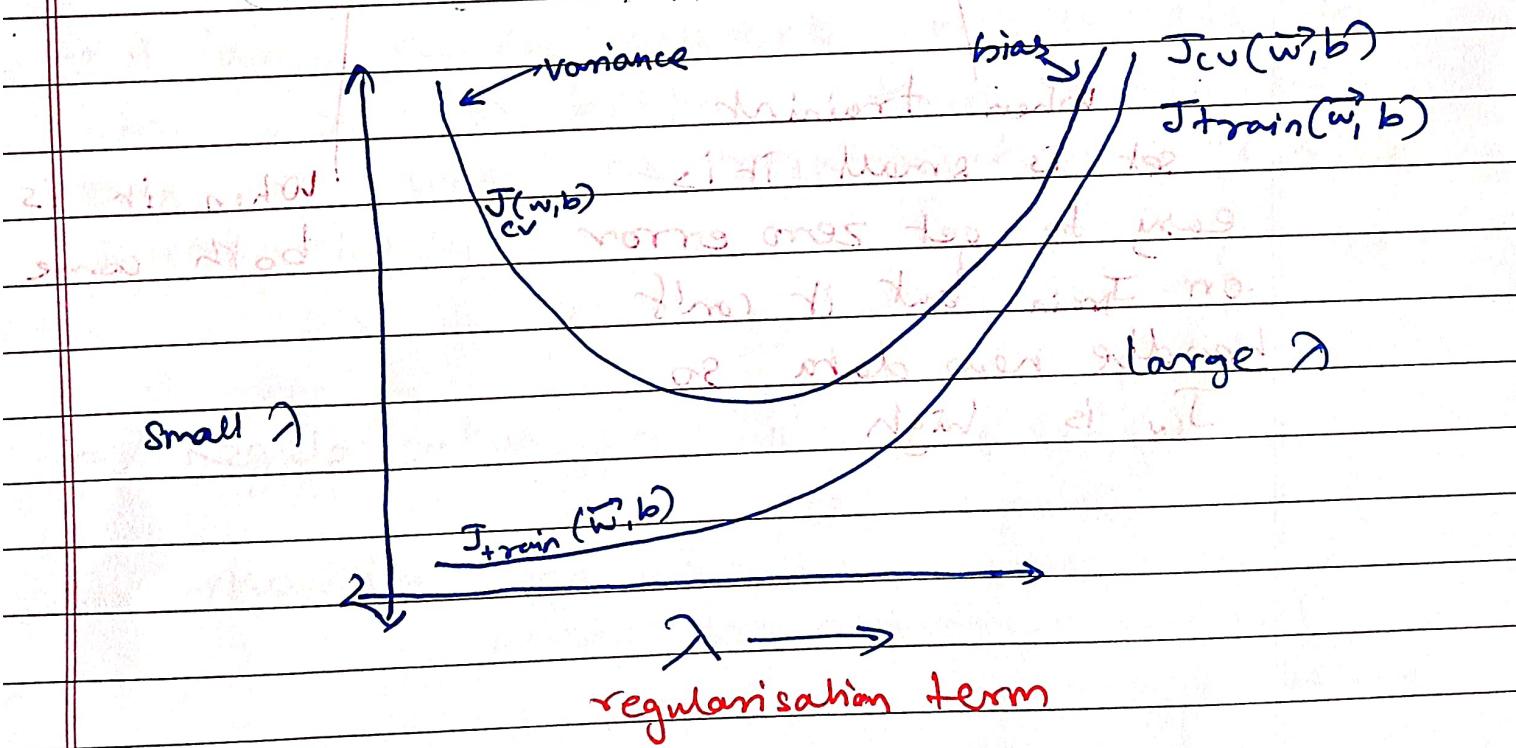
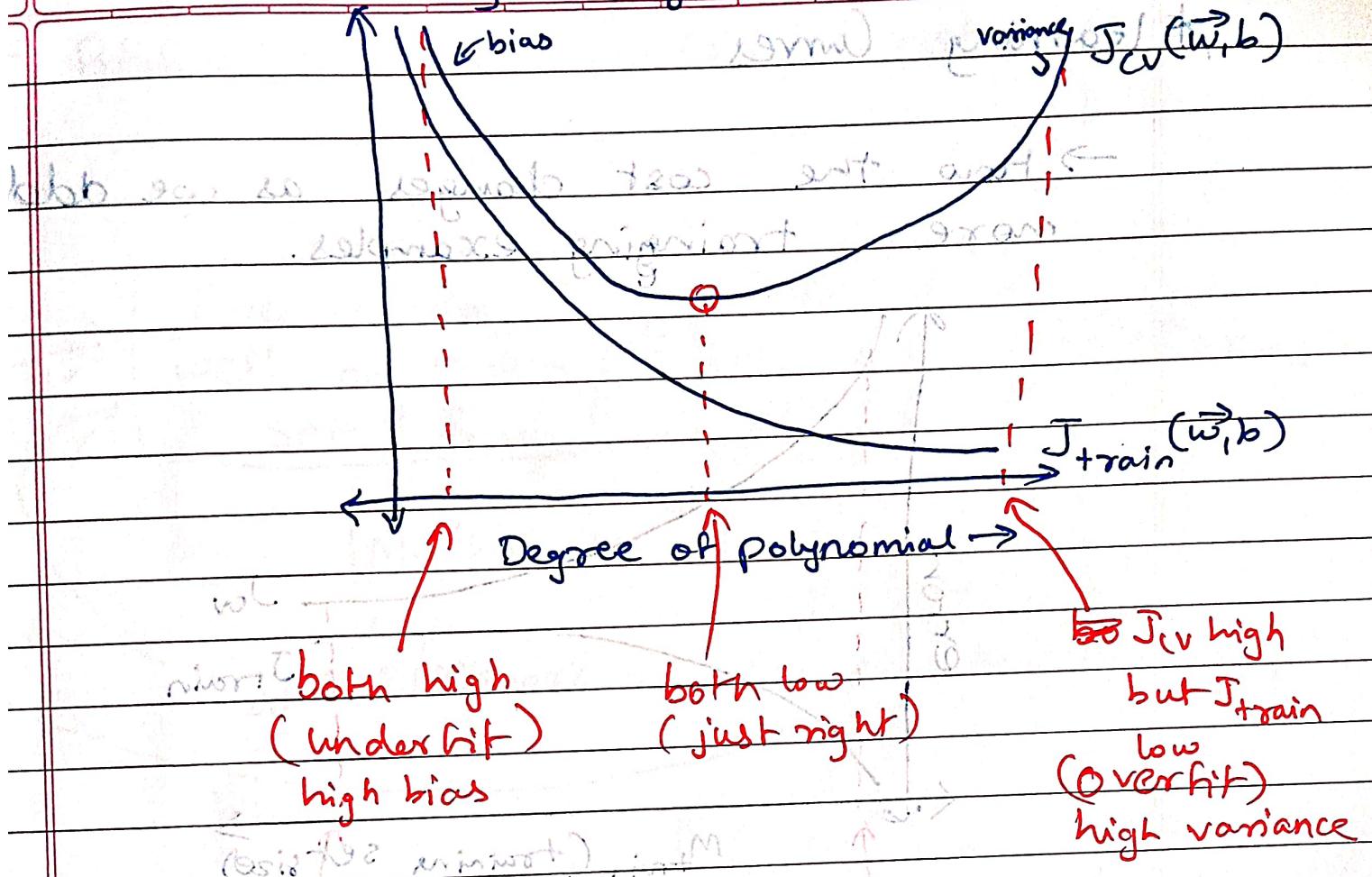
first train ~~model~~ on different models  
 and compare with valid.

then pick model with least valid  
 wr error.

then use test to check performance

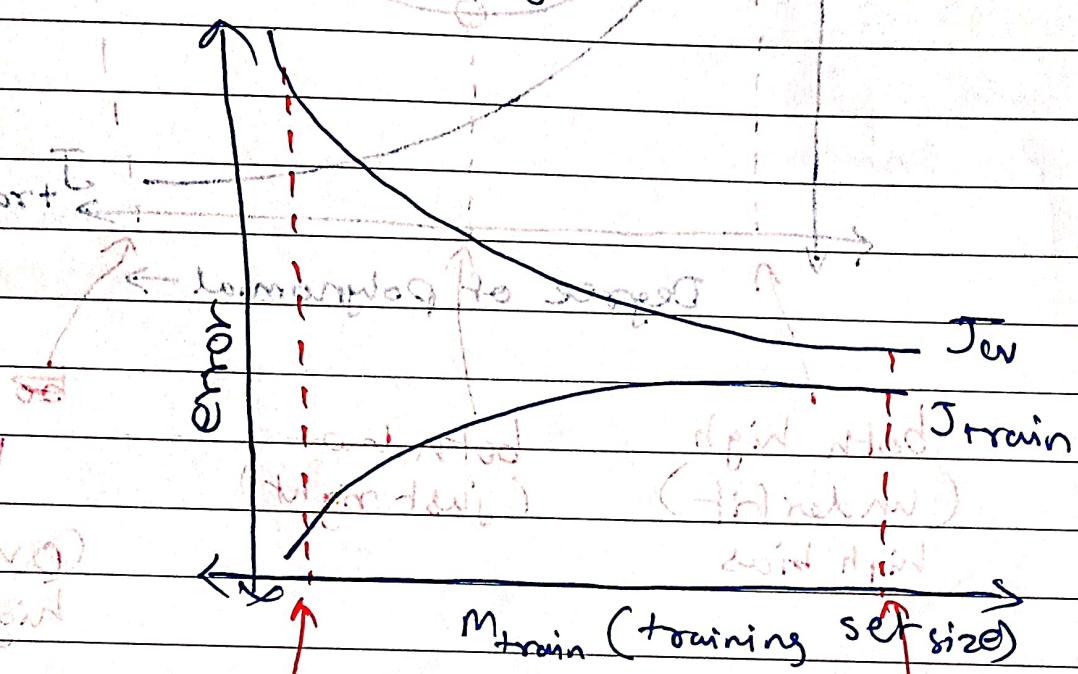
If the  $J_{train}$  &  $J_{cv}$  is low, then the model is right.

# Selecting degree of polynomial



↳ Learning curves

→ How the cost changes as we add more training examples.



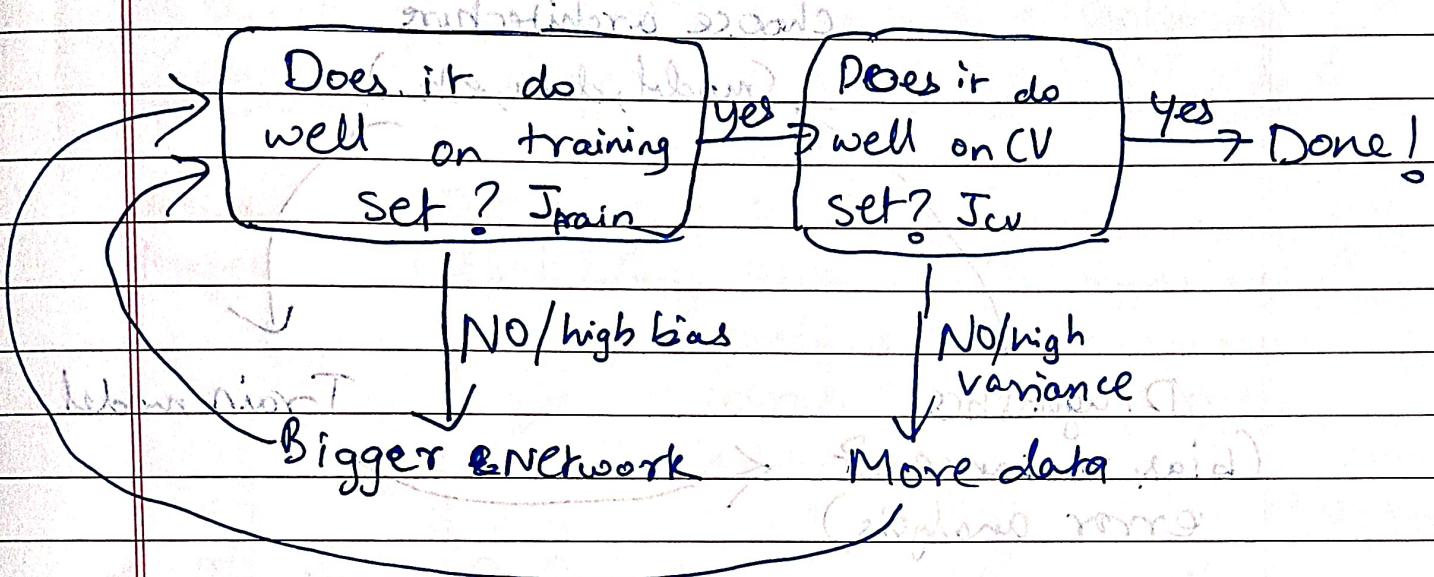
when training set is small it is easy to get zero error on  $J_{train}$  but it can't handle new data so

$J_{test}$  is high

When it is more both come closer

most often happens

## \* Bias & variance in Neural networks



→ A large neural network will usually do as well or better than a smaller one so long as regularisation is chosen appropriately.

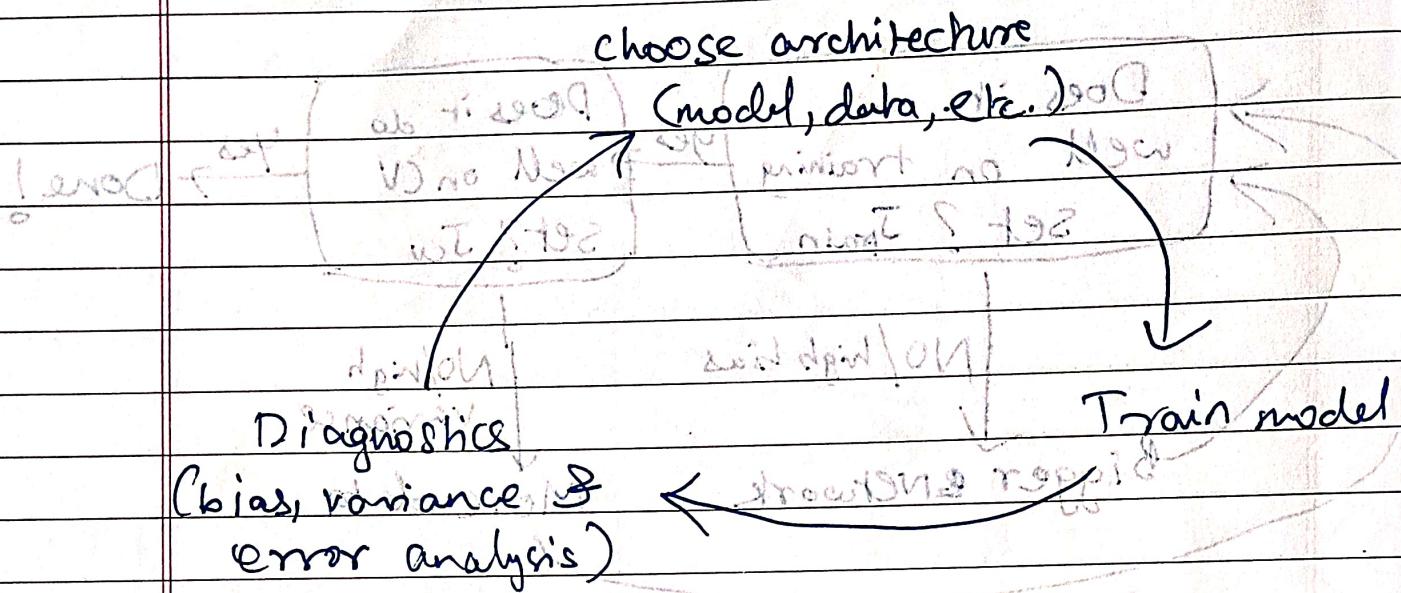
### ⇒ Regularisation in NNs

model = Sequential()  
 $\quad \quad \quad \text{Dense(units=, activation=, kernel\_regularisation=)}$

$$\text{regularisation=} L2(0.01)$$

$\text{Dense}(V=, a=, \text{kernel\_regularisation=} L2(0.01))$ ,  
 $\text{Dense}(V=, a=, u) \quad D$

## \* Iterative loop of machine learning



# Data Augmentation: modifying an existing training

example ~~to~~ <sup>to</sup> create a new training example

new ~~new~~ is next added to now ~~now~~

\* Transfer Learning

→ Downloading neural network parameters pretrained on a large dataset with same input type and own application.

→ further ~~train~~ <sup>fine-tune</sup> the network on ~~new~~ <sup>new</sup> data.

new ~~new~~

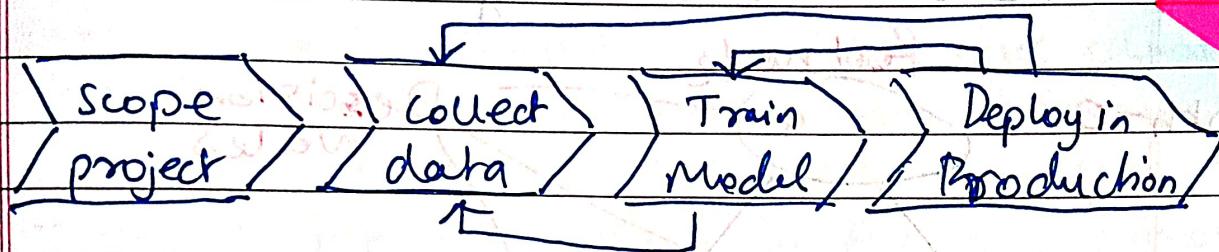
(10, 0) ~~1~~

(1, 0) ~~1~~ = ~~new~~ <sup>new</sup> for ~~new~~ <sup>new</sup>  $\rightarrow$   $\{1, 0\}$  ~~new~~ <sup>new</sup>

$\{5 \times 10, 0\}$  ~~new~~ <sup>new</sup>

## Machine Learning

\* All steps of a machine learning project.



Define Project      Define & Collect data      Training, error analysis & iterative improvement      Deploy, monitor & maintain system

\* Precision & Recall.

$y=1$  in presence of rare class we want to detect.

		Actual Class		Precision
		True Positive	False Pos.	$\frac{\text{True Pos.}}{\text{True Pos.} + \text{False Pos.}}$
Predicted Class	1	True Positive	False Pos.	$\frac{\text{True Pos.}}{\text{True Pos.} + \text{False Pos.}}$
	0	False Neg.	True Neg.	

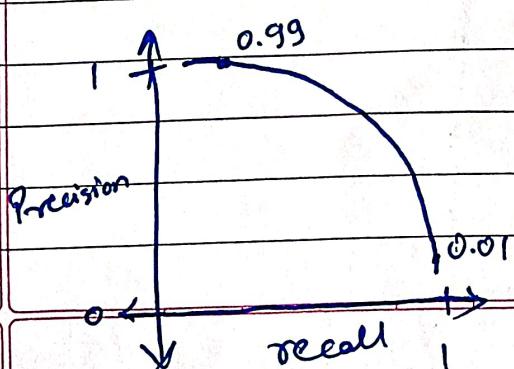
$$\text{Precision} = \frac{\text{True Pos.}}{\text{True Pos.} + \text{False Pos.}}$$

Recall

$$\text{Recall} = \frac{\text{True Pos.}}{\text{True Pos.} + \text{False Neg.}}$$

F<sub>1</sub> Score:

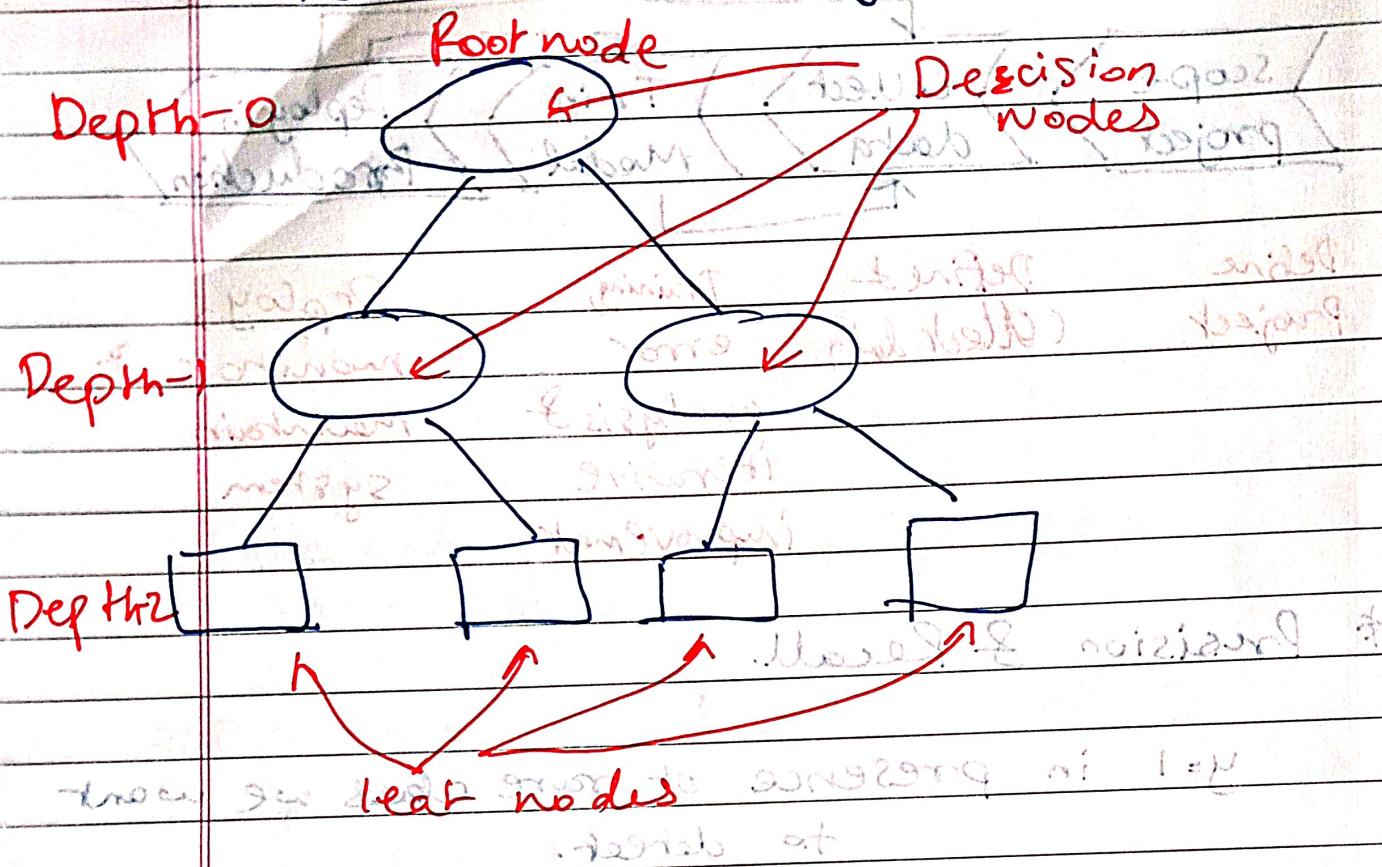
$$F_1 \text{ score} = \frac{1}{\frac{1}{2}(\frac{1}{P} + \frac{1}{R})} = \frac{2PR}{P+R}$$



Harmonic mean

Course 2 week 6

\* features take categorical or discrete values

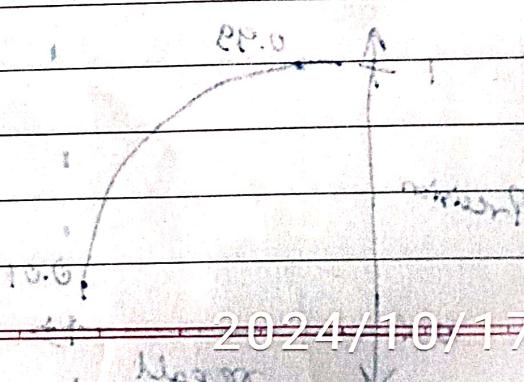


\* A Decision tree is a non-parametric supervised machine learning algorithm.

Output =  $\text{out}_{\text{out}}$   
 output of the Decision tree: How to choose what feature to split on at each node?

- Maximize Purity (minimize impurity)

$$\text{Gini} = \frac{1}{2} + \frac{1}{2} \cdot \left( \frac{1}{3} + \frac{2}{3} \right) = \frac{1}{2}$$



## # Decision 2: When does stop splitting.

- When a node is 100% one class.
- When splitting a node will result in the tree exceeding maximum depth
- When improvements in purity score are below a threshold.
- When number of examples in a node is below a threshold.

## # Measuring purity

→ Entropy as a measure of purity

$P_0, P_1 \rightarrow 2 \text{ classes (in binary) H classification}$

$$P_0 = 1 - P_1 \quad (\text{minist. soft margin})$$

$$H(P_1) = -P_1 \log_2(P_1) - P_0 \log_2(P_0)$$

$$H(P_1) = -P_1 \log_2(P_1) - (1-P_1) \log_2(1-P_1)$$

Entropy function

High purity means low impurity

Entropies are 0.0 if there is no impurity

Entropies are 1.0 if there is no purity

Max entropy = 1.0

Min entropy = 0.0

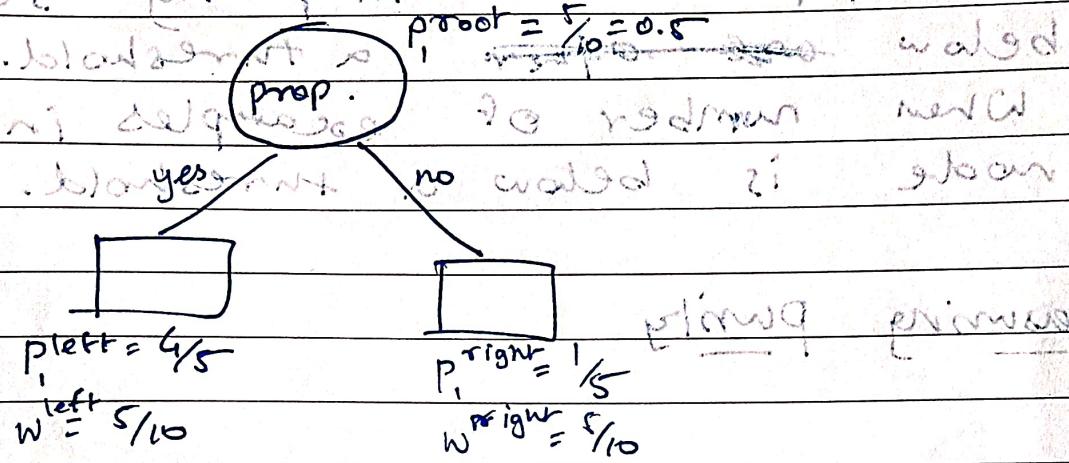
Blurred boundary at 0.5

2024/10/17 12:33

## \* Choosing good split node : Standard H

→ A split is chosen based on information gain. Now there is a profit & loss.

Highest information gain is chosen.



Information gain

$$\text{information gain} = H(p_{\text{root}}) - (w_{\text{left}} H(p_{\text{left}}) + w_{\text{right}} H(p_{\text{right}}))$$

## \* Decision Tree Training

(Start with all examples at the root node)

- Calculate information gain for all possible features, & pick the one with the highest information gain
- Split dataset according to selected feature, & create left & right branches of the tree.
- keep repeating splitting until stopping criteria is met:
  - o When a node is 100% one class.
  - o When splitting a node will result in the tree exceeding maximum depth
  - o Information gain from additional splits is less than threshold

→ When the no. of examples in a node is below a threshold.

# Tree Ensemble: collection of decision trees.

→ Random forest algorithm.

Given a training set of size  $m$ :

for  $b=1$  to  $B$ :

    Use sampling with replacement to create  
    a new training set of size  $m$ .  
    Train a new decision tree on this dataset.

Algorithm: At each node, when choosing a feature

to split, if  $n$  features are available,  
choose a random subset of  $k < n$  features  
and allow the algorithm to only choose  
from that subset of features.

$$k = \sqrt{n}$$

practical:  $k$  is often  $\sqrt{n}$

Want to prevent overfitting by using  $k < n$ .

overfitting occurs when the algorithm

can always predict the target value

of a given observation + variation

with a probability of 1.0.

You need more samples. But it's not always

possible to get more data.

XGBoost: following to cut out method  
Balanced tree model 21 Nov 2023

- most popular way of building trees.
- used in ML competitions

most popular from material ←  
Boosted trees intuition:

in ISR to add priority to noisy

Given training set of size  $m$ :

for  $b=1$  to  $B$ : : 8 at 1st not

at Use sampling with replacement to

create a new training set of size

no most notables even a short

But instead of picking all examples

instead of picking with equal ( $1/m$ ) probability, make it

more likely to pick misclassified

→ to the examples from a previously trained

set of trees. If we do this

then train a new decision tree

$$n = 1$$

## XGBoost → extreme Gradient Boosting

- open source implementation of boosted trees.
- fast efficient implementation,
- good choice of default splitting criteria & criteria when no stop splitting
- built in regularization for overfitting.
- highly competitive algorithm for ML competitions

## When to use Decision Trees?

DT:

- Works well on tabular (structured) data.  
(not in images, audio, text)
- Fast to train & make predictions.

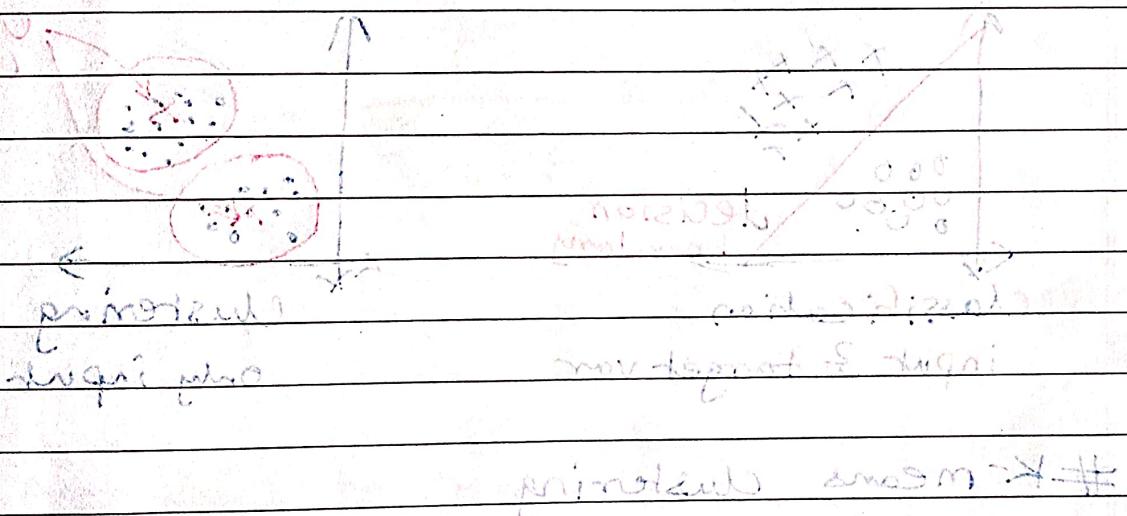
NN: • Good for non-linear data.

- Works well on all types of data.

- Slower than a DT.

- Works with transfer learning.

• Gabier: to learn together multiple NN.



Decision tree → good for tabular data.

# Course 3 Week 1

: TO

## Unsupervised Learning

(first, without supervision; task)

1) Clustering Algorithms

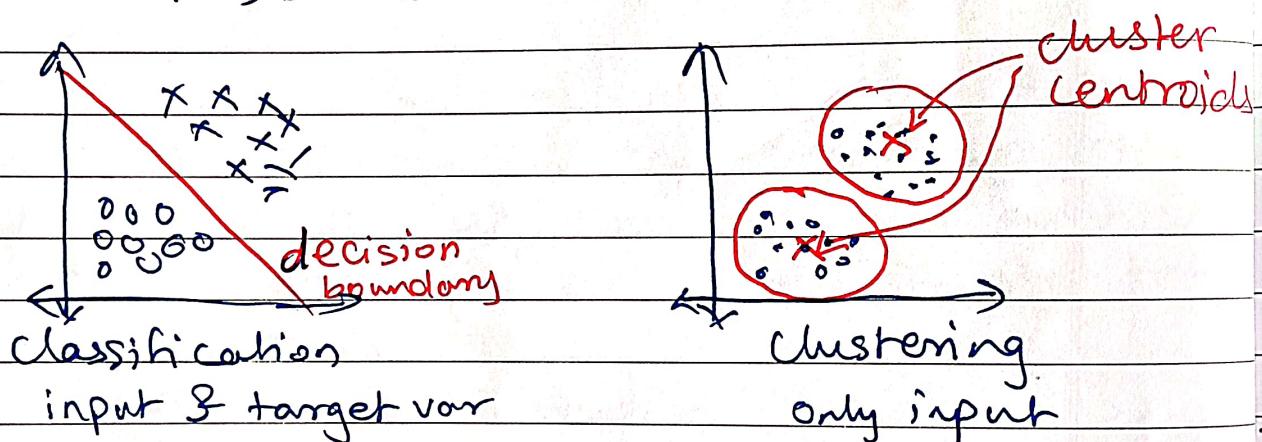
2) Anomaly detection

- Recommender systems (week 2)

3) Reinforcement learning (week 3)

### \* Clustering

→ Grouping similar data points together.



### # K-means clustering

- Randomly initialize K cluster centroids

$$\mu_1, \mu_2, \mu_3, \dots, \mu_K$$

with each being a vector of  $n$  features

repeat {

for  $i=1$  to  $m$

$c^{(i)} := \text{index}(\text{from } 1 \text{ to } K)$  of cluster centroid

closest to  $x^{(i)}$

$$\min_k \|x^{(i)} - \mu_k\|^2$$

it is squared.

for  $k=1$  to  $K$

$\mu_k := \text{average (mean) of points assigned}$   
to cluster  $k$

## # K-means optimization objective

$c^{(i)}$  = index of cluster ( $1, 2, \dots, k$ ) to which example  $x^{(i)}$  is currently assigned.

$\mu_k$  = cluster centroid

$\mu_{c(i)}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned.

### Cost function:

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$

$$\text{Algorithm} = \min_{\mu_1, \dots, \mu_k} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

It is called distortion function

## # Anomaly Detection

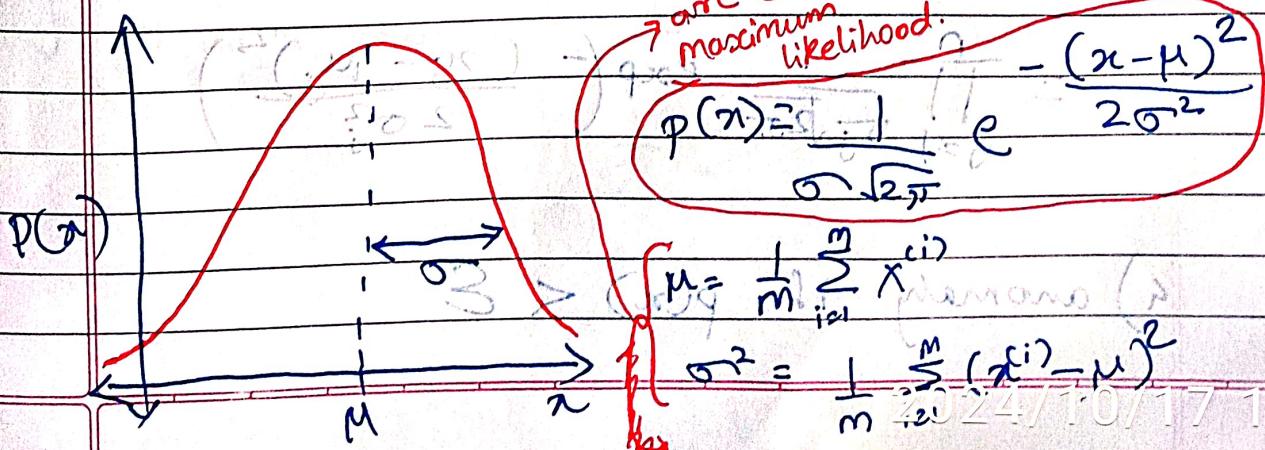
### # Gaussian (Normal distribution)

say  $x$  is a number

prob. of  $x$  is determined by a gaussian with a mean  $\mu$ , variance  $\sigma^2$

are called as maximum likelihood.

$$p(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



# Density Estimation notes in progress

Training set =  $\{\vec{x}_1^{(1)}, \vec{x}_2^{(2)}, \dots, \vec{x}_n^{(n)}\}$   
 Each example  $\vec{x}^{(i)}$  has  $n$  features

$$P(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

at most  $n$  features =  $(N, M)$

$$P(\vec{x}) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

## # Anomaly detection algorithm

1) choose parameters  $\mu_i$  that you think might be indicative of anomalous examples.

2) fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

3) Given example  $x$ , compute  $p(x)$

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$\frac{1}{\sigma_j} = \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

a) anomaly if  $p(x) < \epsilon$

## Course 3 week 2

### \* Recommender Systems

If you predicting for a single user,  
It is similar to linear regression.

fit a line  $wx+b$  based on user's choices and then predict what to recommend.

To learn parameters  $w^{(j)}, b^{(j)}$  for the user  $j$

cost  $f$  of lin. reg.

$$J(w^{(j)}, b^{(j)}) = \frac{1}{2} \sum_{i: r(i,j) \neq 1} (w^{(j)} \cdot x_i + b^{(j)} - y_{i,j})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k^{(j)})^2$$

Only if user has rated something

regularisation parameter term  
 $k = \text{no. of features}$

To learn parameters  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$  for all users;

$$J(w^{(1)}, w^{(2)}, \dots, w^{(n_u)}, b^{(1)}, b^{(2)}, \dots, b^{(n_u)}) = \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j) \neq 1} (w^{(j)} \cdot x_i + b^{(j)} - y_{i,j})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Cost function to learn  $\alpha^{(1)}, \dots, \alpha^{(m)}$ :  
for all features  $x^{(1)}, \dots, x^{(m)}$

$$\text{min}_{w, b} \sum_{i=1}^m \text{error}(w^{(i)}, b^{(i)}) = \sum_{i=1}^m \frac{1}{2} (y^{(i)} - w^{(i)} \cdot x^{(i)} + b)^2$$

inseu no log loss  $\rightarrow$  error  $\rightarrow$   $\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (x_k^{(i)})^2$   
feature  $x^{(i)}$  not true logistic

Remember of

Putting them together we get cost fn

$$\text{min}_{w, b} \sum_{i=1}^m \text{error}(w^{(i)}, b^{(i)})$$

$$\sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)})^2$$

$$\sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)})^2 + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)} \cdot x_k^{(i)} + b^{(i)} - y^{(i)})^2$$

$$\sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)})^2 + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)} \cdot x_k^{(i)} + b^{(i)} - y^{(i)})^2$$

$$\sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)})^2 + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)} \cdot x_k^{(i)} + b^{(i)} - y^{(i)})^2$$

$$\sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)})^2 + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^{n_k} r(i,j) (w_k^{(i)} \cdot x_k^{(i)} + b^{(i)} - y^{(i)})^2$$

#

Gradient descent (minimising loss function)

Collaborative filtering (ALS)

Using linear regression algo,

repeat {

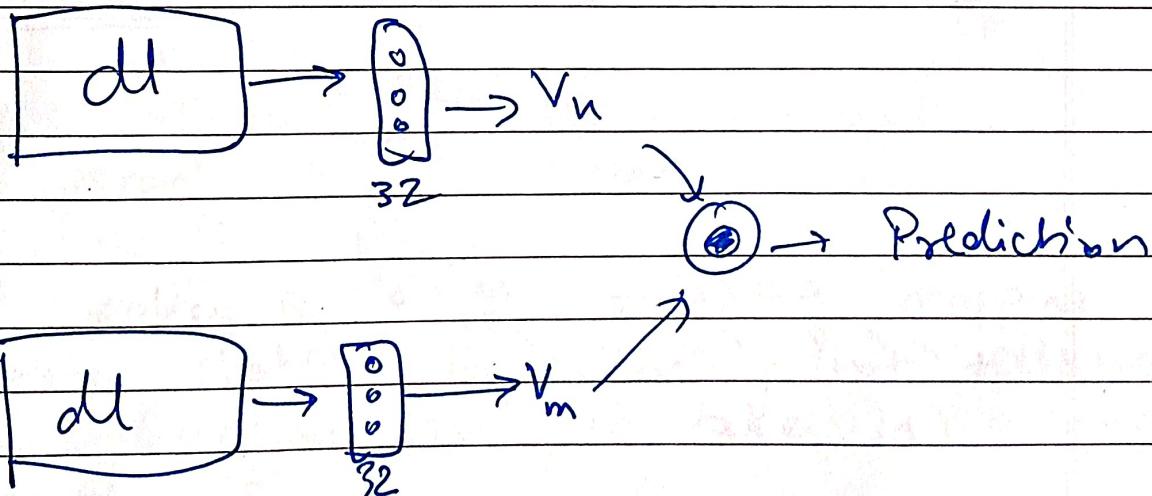
$$w_i^{(j)} = w_i^{(j)} - \alpha \frac{\partial}{\partial w_i^{(j)}} J(w, b, x)$$

$$b^{(j)} = b^{(j)} - \alpha \frac{\partial}{\partial b^{(j)}} J(w, b, x)$$

$$x_k^{(i)} = x_k^{(i)} - \alpha \frac{\partial}{\partial x_k^{(i)}} J(w, b, x)$$

Simultaneous update.

# Content-based filtering



Cost fn:

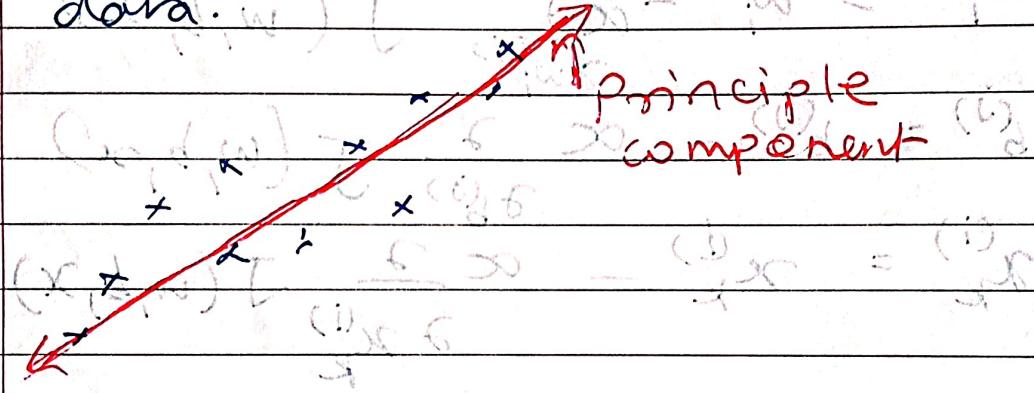
$$J = \sum_{(i,j) : r(i,j) = 1} (v_n^{(j)} \cdot v_m^{(i)} - y^{(i,j)})^2$$

+ NN regularisation term.

\* Principal component analysis  
 → It is an unsupervised learning algo

PCA algorithm: Take high dimension data and reduce it to 2D data for visualization.

→ find axis with maximum variance of data.



→ PCA is not linear regression.

principal second - third

$$NV \leftarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} U_1 \end{bmatrix}$$

numbered

$$NV \leftarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} U_2 \end{bmatrix}$$

: 120 120

$$S_{(i,j)} = m \cdot v \cdot v^T$$

$$\sum_{i=1}^n S_{(i,i)} = L$$

minimize  $L$

2024/10/17

## Course 3 week 3

### \* Reinforcement learning.

but is strongest & easiest / easiest of learning  $\leftarrow$  RL

when making move to start the game

$$x \longrightarrow y$$

no wrong or bad reward function  $\leftarrow$  RL

into state function  $\leftarrow$  RL

### Applications:

$\rightarrow$  controlling robots  $\leftarrow$  RL

$\rightarrow$  factory optimization

$\rightarrow$  financial (stock) trading

$\rightarrow$  Playing games (including video games)

discount factor: makes the algo impatient

$$\gamma = 0.9 \text{ or } 0.99 \text{ or } 0.999$$

### \* Policy ( $\pi$ )

Policy

state  $s$  is  $\xrightarrow{\pi}$  action  $a$

$\Rightarrow$  A policy is a  $\pi$   $\pi(s) = a$  mapping

from states to actions, that tells you

(what) action to take in a given state  $s$

$\rightarrow$  find a policy  $\pi$ , that tells you what action ( $a = \pi(s)$ ) to take in every state ( $s$ ) so as to maximize return.

## Lesson 8 MDP

## # Markov Decision Process (MDP)

→ Actions to be taken depends on the current state of the system only.

→ Which means that results prev or after this current state don't matter.

## # State-action value function

$Q(s, a)$  (return value you get)

from certain actions starting from state  $s$ .

- take action  $a$  once
- then behave optimally after

thus optimal opt with respect to next state

R.P.O → P.P.O → P.O =

The best possible return from state  $s$  is max <sub>$a$</sub>   $Q(s, a)$

The best possible action  $a$  in state  $s$  is the action  $a$  that gives max <sub>$a$</sub>   $Q(s, a)$

Suppose  $\pi = (\pi_a)$ . If  $\pi$  is a policy,  $\pi$  is a policy that says "Use first, avoid all others more often".

$$\Rightarrow Q_\pi(s, a) = R(s) + \gamma \max_{a'} Q_\pi(s, a')$$

which is called Bellman Equation  $(\pi = \pi)$

## \* Random (stochastic) environment

$$Q(s, a) = R(s) + \gamma \mathbb{E} [\max_{a'} Q(s', a')]$$

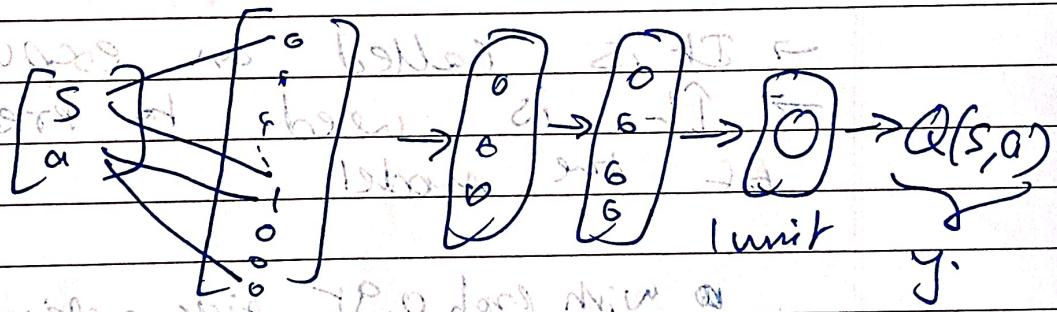
## \* Continuous state space:

current state is denoted by a vector

$$\begin{bmatrix} 0 \\ p \\ g \\ f \end{bmatrix}$$

vector of form  $s_i =$

current position of the system  
moving the Deep reinforcement learning.



Deep Q-network algorithm (DQN)

repeat {

→ take actions in the linear code. get  $(s, a, R(s), s')$

→ store 10,000 most recent  $(s, a, R(s), s')$  tuples

→ Train NN:

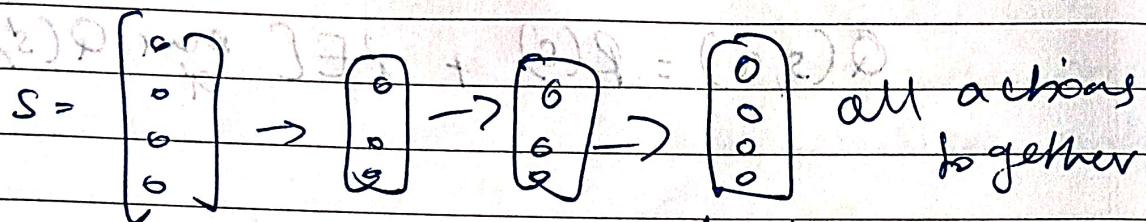
create training set of 10,000 eg using

$$x = (s, a) \quad \text{and} \quad y = R(s) + \gamma \max_{a'} Q(s', a')$$

Train  $Q_{\text{new}}$  such that  $Q_{\text{new}}(s, a) \approx y$ .

→ set  $Q = Q_{\text{new}}$ .

## \* Refined DQN (Sideloaded material) \*



: choose action by mixing

## \* E-greedy policy.

→ it's another refinement.

→ Epsilon greedy ~~etc~~ policy is used to take actions while learning.

→ make the algo randomly try new actions to make open to new options

→ It is called an exploration step.

→ It helps need to break preconceptions of the model.

• with prob 0.95 pick action  $a$  that maximizes  $Q(s, a)$

(w/ prob 0.05 pick random action)

gradually decrease  $\epsilon$  from

$$(0, \epsilon) \rightarrow (2, \epsilon) \rightarrow (1, \epsilon) \rightarrow (0, \epsilon) = x$$

prob  $(\epsilon)$  used with new weight

new  $\rightarrow 0.75 \epsilon$

\* Minibatch & soft updates  
→ more refinements.

minibatch: instead of using all eg of training set at each iteration, use only of a subset of the training set.

Soft update: don't change new Q drastically

$$w = 0.01 w_{\text{new}} + 0.99 w$$

$$b = 0.01 b_{\text{new}} + 0.99 b$$

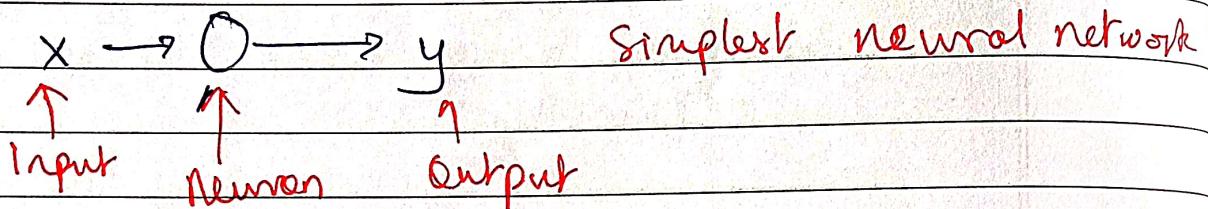
— X O X —

Jin

# Deep Learning Specialization

## Course 1: Neural Networks & Deep Learning

### Week 1: Introduction



### Week 2: Neural Networks Basic

~~shallow~~  
Week 3: Neural Networks ~~overview~~

ReLU is the most useful activation f $\sigma$

- Why do we need non-linear f $\sigma$ ?
- ① In order to learn interesting features
- ② If all layers are linear, then no matter layers or nodes we use, the overall model is just linear

### Week 4: Deep Neural Networks

\* All the weights should be initialized randomly instead of all weights to zeros.

## \* Parameters vs Hyperparameters.

→ Parameters:  $w^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots$

\* Hyperparameters: learning rate,  $\alpha$

# iterations, # no. of iterations.

# hidden layers

# hidden units

choice of activation functions.

Momentum term

Minibatch size

regularizations

\* Epochs: number of forward passes

\* Gradient descent: backpropagation with

\* Adam optimizer: backpropagation with

\* RMSprop: backpropagation with

\* Momentum: backpropagation with

\* Stochastic gradient descent: backpropagation with

\* Mini-batch gradient descent: backpropagation with

\* Conjugate gradient: backpropagation with

\* BFGS: backpropagation with

\* L-BFGS: backpropagation with

\* Newton's method: backpropagation with

\* Quasi-Newton methods: backpropagation with

\* Levenberg–Marquardt: backpropagation with

# Course 2: Improving Deep Neural Networks

## Week 1: Practical aspects of Deep learning

① → overfitting / high variance can be solved by regularization.

② → Dropout regularization is also very powerful.

Why does dropout work?

→ Intuition: Can't rely on any one feature, so have to spread our weights.  
→ This will shrink weights.

\* The more the nodes in a layer, the lesser should be the dropout probability. The lesser the nodes, the more the probability.

③ → Data augmentation can also help in overfitting.

④ → Early stopping

# Optimizing. ~~and a short Q&A session~~

① → Normalizing the inputs. ~~Normalizing~~  
 If decreases the steepness of the cost fn.

② → If you have a very deep neural network, the values of  $y$  can explode. This is called "exploding gradient".

③ → On the reverse if the values of weights is small then the value of  $y$  can vanish. This is called "vanishing gradient".

# weight initialization can help with this.

How's not we get rid of it for

almost larger  $n \rightarrow$  smaller  $w_i$

$$\text{Nor}(w_i) = \frac{2}{\sqrt{n}}$$

Sometimes standard deviation set equal

for ReLU  $\rightarrow w^{[l]} = \text{np.random.randn}(\text{shape})$   
 activation  $\rightarrow (\text{0.008} \geq |z|) \text{ then } \text{np.sqrt}(\frac{4}{n^{[l-1]}})$

For tanh  $\rightarrow$  we can't have one single input  
 activation

This is called a "Xavier initialization"

size standardization - small variance  $\Rightarrow$  large variance  $\Rightarrow$  small variance

personality US vs. UK vs. DE

## Week 2: Optimization algorithms

optimization algorithms help train models faster.

### ① Mini-Batch Gradient Descent.

Computing cost on a part of your data. This is efficient → this saves memory.

If mini-batch size = m : Batch gradient descent  
size = 1 : stochastic gradient descent

mini-batch size between that is  
not too big or too small  
gives the fastest learning

how to choose mini-batch size?

→ If target set is ( $\leq 2000$ ):

use batch grad desc.

typical minibatch sizes:

64, 128, 256, 512, 1024

\* Just make sure minibatch size fits in CPU / GPU memory

## # Exponentially weighted (moving) averages:

→ It is used to describe a time series

→ It decides how important the current observation is, as compared to previous observations.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$\beta$  is selected by user.

→ If you expand this for 100 eg,

$$V_{100} = 0.1 \theta_{100} + (0.1)(0.9) \theta_{99} + (0.1)(0.9)^2 \theta_{98} + (0.1)(0.9)^3 \theta_{97} + \dots$$

so this will decay exponentially.

It averages  $\propto \frac{1}{1-\beta}$  examples.

$$\Sigma = 1-\beta$$

## # Bias correction in Exponentially weighted averages.

For initial weights →

instead of  $V_t$ , calculate  $\frac{V_t}{1-\beta^t}$

## Gradient Descent with momentum

→ It uses the exponentially weighted averages:

momentum:

On iteration  $t$ :

compute  $d_w, d_b$  on current minibatch.

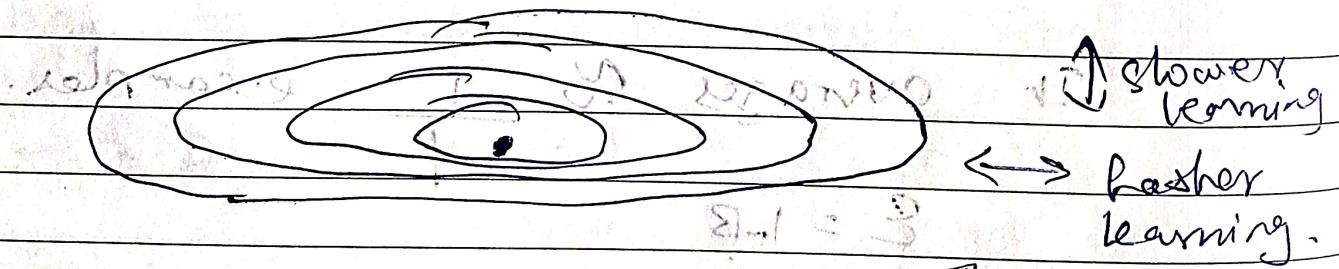
$$V_{dw} = \beta V_{dw} + (1-\beta) d_w$$

$$V_{db} = \beta V_{db} + (1-\beta) d_b$$

$$w = w - \alpha V_{dw}$$

$$b = b - \alpha V_{db}$$

→ It smooths out the grad descent, prevents overshooting.



this is what it does to prevent oscillations.



→ It helps to prevent oscillations.

↓ ↓ ↓

~~(a)~~

## RMS prop.

→ Root mean square prop.

On iteration t:

compute  $dW, db$  on current mini-batch.

$$S_{dw} = \beta S_{dw} + (1-\beta) dW^2$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$w = w - \alpha \frac{dW}{\sqrt{S_{dw} + \epsilon}}$$

$$\epsilon = 10^{-8}$$

$$b = b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

~~(b)~~

## Adam algorithm

→ combines momentum & rms prop.

→ Adaptive moment estimation

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0.$$

on iteration t:

compute  $dW, db$  on curr ~~iteration~~ mini-batch.

momentum }

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW$$

$$V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

momentum }

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2$$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

for adam,

$\alpha$ : needs to tune

$$\beta_1 = 0.9 (dw)$$

$$\beta_2 = 0.999 (dw)$$

$$\epsilon = 10^{-8}$$

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{(1-\beta_1^t)}, V_{db}^{\text{corrected}} = \frac{V_{db}}{(1-\beta_1^t)}$$

$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)}, S_{db}^{\text{corrected}} = \frac{S_{db}}{(1-\beta_2^t)}$$

$$w = w - \alpha$$

$$V_{dw} = \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}, b = b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

(5)

# learning rate decay.

→ slowly decrease learning over time.

1 epoch = 1 pass through data.

$$\alpha_k = \alpha_0 + \text{decay rate} \times k = \alpha_0 \times (1 + \text{decay rate} \times \text{Epoch No.})$$

$$\alpha_0 = 0.12 \quad \text{decay rate} = 0.1$$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05

Exponential decay:  $\alpha = \alpha_0 \times (1 - r)^k$

$$\alpha = \alpha_0 \times (1 - r)^{\text{Epoch Num}}$$

$$\alpha = \alpha_0 \times (1 - r)^k$$

$$\alpha_{(t+1)} = \alpha_t \times (1 - r)$$

( $t+1$ ) discrete staircase:

half the  $\alpha$  after

some time.

(6)

## Problem of local optima.

→ It's very unlikely that the grad desc gets stuck in a local optima in a higher dim space.

→ This problem is only possible in lower dim data.

$$\rightarrow X \in \mathbb{R}^n$$

## Week 3: Hyperparameter tuning

# appropriate scale for hyperparameters

$$n^{[L]} = 50 - 100 \text{ (number of samples)}$$

$$L = 2 - 4$$

$$\alpha = 0.0001 - 1 \text{ (use log scale.)}$$

$$\beta = 0.9 - 0.999 \text{ (using log scale of 1-\beta)}$$

## # Batch Normalisation

loop first step: calculate mean + std

$$\text{step 1: } \mu = \frac{1}{m} \sum x^{(i)}$$

loop second step: calculate variance

$$x = x - \mu$$

loop third step: calculate std

$$\sigma^2 = \frac{1}{m} \sum x^{(i)^2}$$

$$x = x / \sigma$$

problem: what about next layer? Solution:

→ we can normalize activations.

(given some intermediate values in NN

$$z^{(1)}, \dots, z^{(m)} : \text{mean} - \text{std}$$

$$\mu = \frac{1}{m} \sum z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum (z^{(i)} - \mu)^2$$

$$z_{\text{norm}} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$FF$$

$$\gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\beta - \mu$$

$$\text{given } z^{(i)} = z^{(i)}$$

this beta is not  
the momentum  
beta hyper parameter

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}} + \beta$$

learnable  
parameters of  
model

use  $\tilde{z}^{(i)}$  instead of  $z^{(i)}$

→ Sandwich batch norms b/w bet<sup>n</sup>. layers  
to normalize activations.

→ Batch norm also has a regularization side effect. b/c it ~~reduces~~<sup>increases</sup> the noise in values. This spreads weights across nodes.

→ Batch Norm increases learning speed.

### Course 3: Structuring Machine Learning Projects

#### Week 1: ML Strategy

Chain:

fit for training set well on cost  $f^n$

↓

fit dev set well on cost  $f^n$

↓

fit test set well on cost  $f^n$

↓

Performs well in real world.

## # Using a single number evaluation metric

Instead of trying to improve both precision or recall, just try improving the F1 score.

precision: of the examples recognised, how many are actual.

recall: out of actuals, how many recognised correctly.

F1 score = "average" of P & R

$$\left( \frac{2}{\frac{1}{P} + \frac{1}{R}} \right) \text{ Harmonic mean}$$

→ This applies everywhere

Instead of improving every metric, take a single metric as a benchmark for model performance.

Also compare your errors to human-level performance.

Learn from mistakes

Human learning: how do we model?

Human-level do

↓  
avoidable  
bias

Training error

↑  
Variance

Dev error

Train: bigger model

Train longer / better  
optimization algorithms.

NN architecture/hyperpara-

meters / search

RNN/CNN

more data

regularisation

(L1, L2)

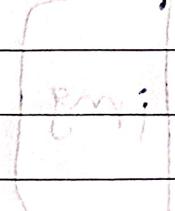
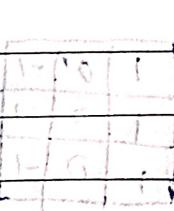
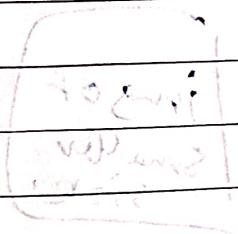
augmentation

hyperparameter

search

## 01 [ ] Week 2: ML strategy (if you don't have AI)

\* Build a simple first system quickly; then iterate



or foggy

data

high

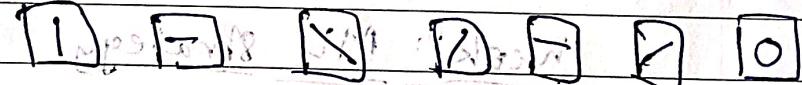
low

high

## Course 4: CNNs

### Week 1: Foundations of Convolutional Neural Networks

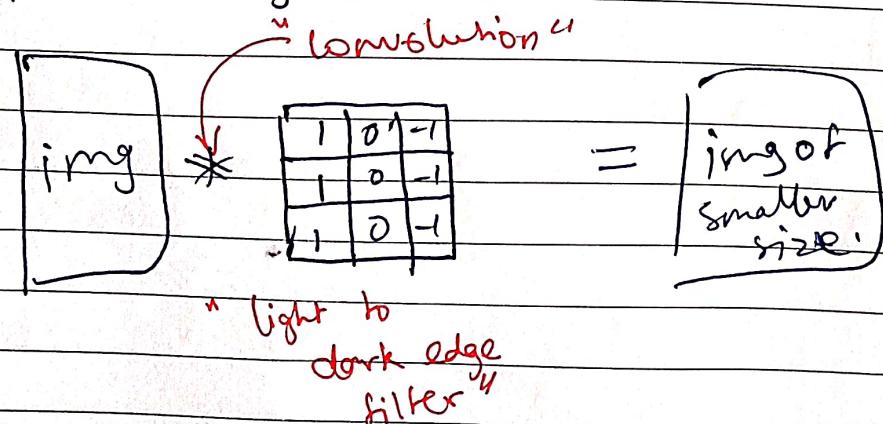
- Used to train neural networks with large inputs and large number of parameters.
- CNNs consist of kernels or "filters".
- Each kernel identifies a certain pattern in the image.

In first layer: 

In second layer: parts of objects - blurry

In third layer: objects themselves

- Vertical edge detection



→ Horizontal edge detection

$$\text{Image} \rightarrow \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix} = \text{Image}$$

→ Sobel Filter with 2x3 kernel

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ -1 & 0 & -1 \end{bmatrix} \leftarrow \text{gives more priority to center pixels.}$$

→ Scharr Filter

$$\begin{bmatrix} 3 & 0 & -3 \\ 16 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix} \leftarrow \text{another filter with same concept taken to extreme.}$$

→ Researchers used to handpick these kernel weights before but now we implement networks to "learn" the weights of the kernels.

## # Padding

without padding, introduce noise

- The pixels on the edge don't contribute a lot to the output.
- So we add a border to the image.

→ This preserves the size of the image & makes edge pixels contribute to the output.

1	0	1	0
0	1	0	1
1	0	1	0
0	1	0	1

## ★ Valid convolution: no padding

with stride 2

★ Same convolution: Pad so that output size is the same as the input.

→ Padding =  $\frac{f-1}{2}$  (size of kernel).

↳ always use padding 2 unless mentioned  
↳ always use stride 2 unless mentioned  
↳ filter size must be odd and smaller than input size

★ Always use odd size filters  
 $1 \times 1, 3 \times 3, 5 \times 5, 7 \times 7$ .

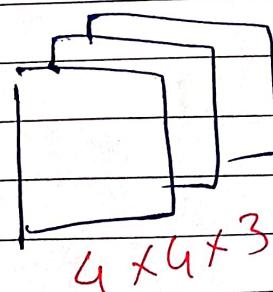
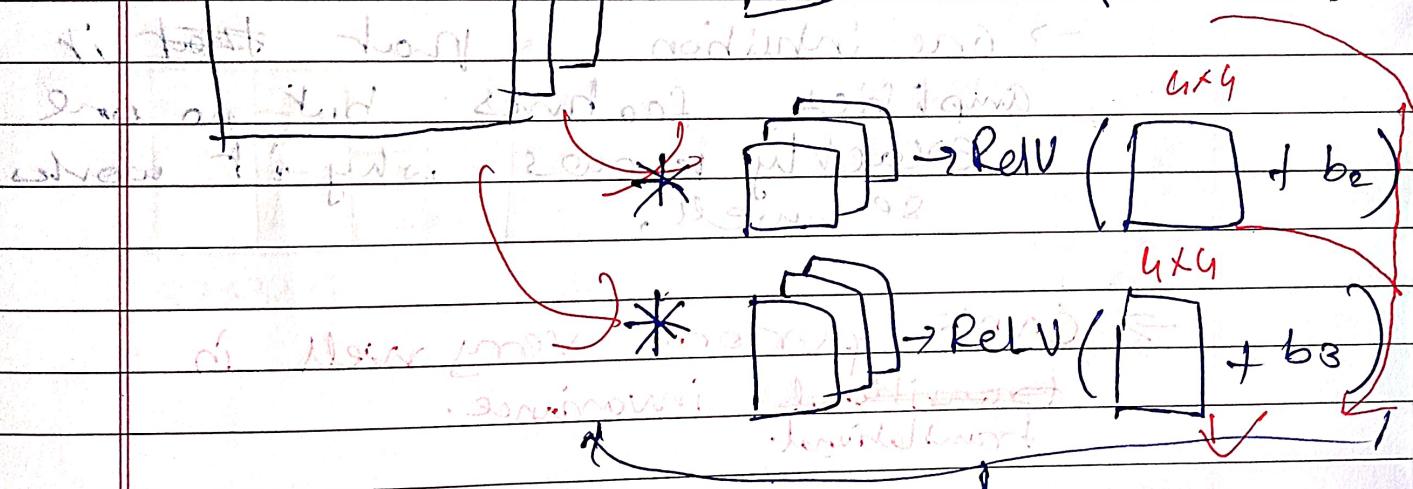
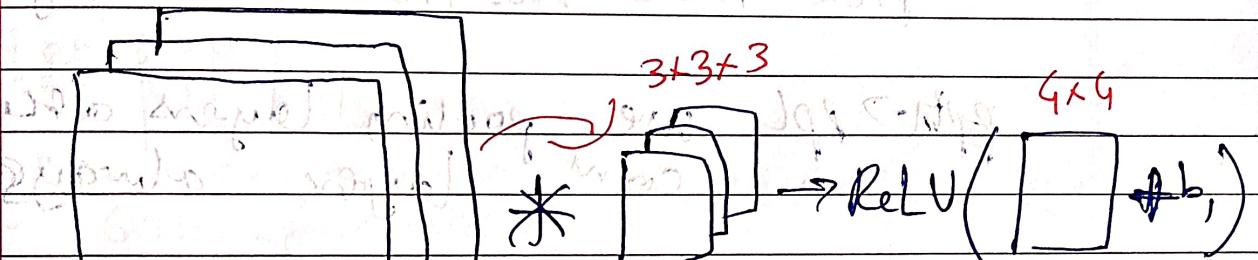
## Stride

Size of image after stride-

$$\text{nan } * \cdot f \cdot f = \left\lceil \frac{n + 2p - f}{s} + 1 \right\rceil$$

Padding P      Stride S :  $\times \left\lceil \frac{n + 2p - f}{s} + 1 \right\rceil$

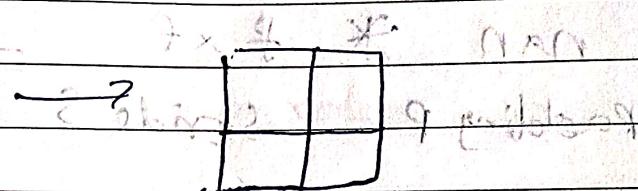
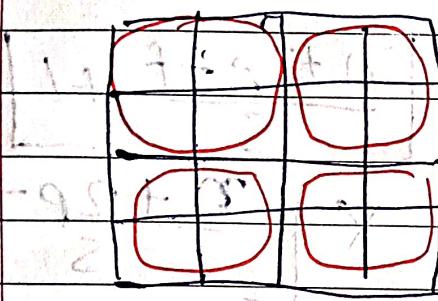
$6 \times 6 \times 3$



$4 \times 4 \times 3$

Pooling

Max Pool

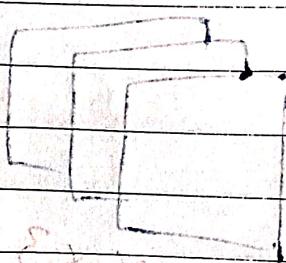


pick more no. from each region.

~~why~~  $\rightarrow$  ppl use pooling layers after convolution layer always.

$\rightarrow$  one intuition is that it amplifies features but no one especially knows why it works so well.

$\rightarrow$  CNNs perform very well in translational invariance.



Examp

# Course 4 week 2

## Deep convolutional models: Case studies.

~~Outline: LeNet-5~~

~~AlexNet~~

~~VGG Net~~

~~ResNets (152 layers)~~

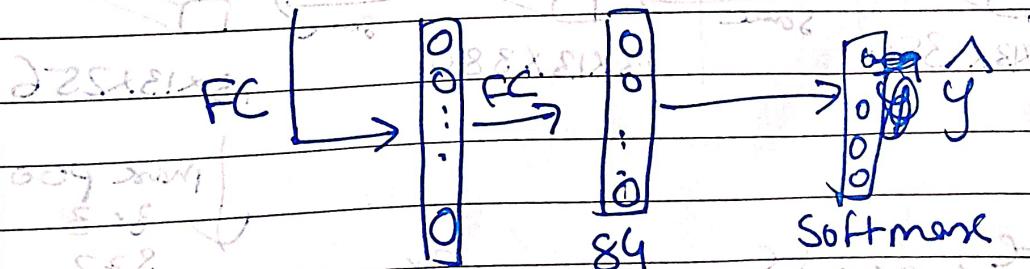
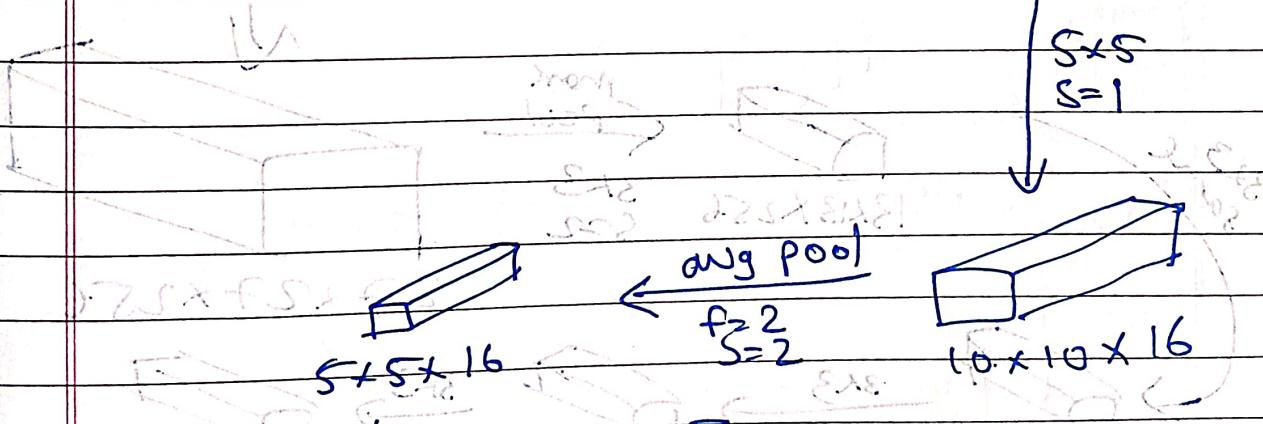
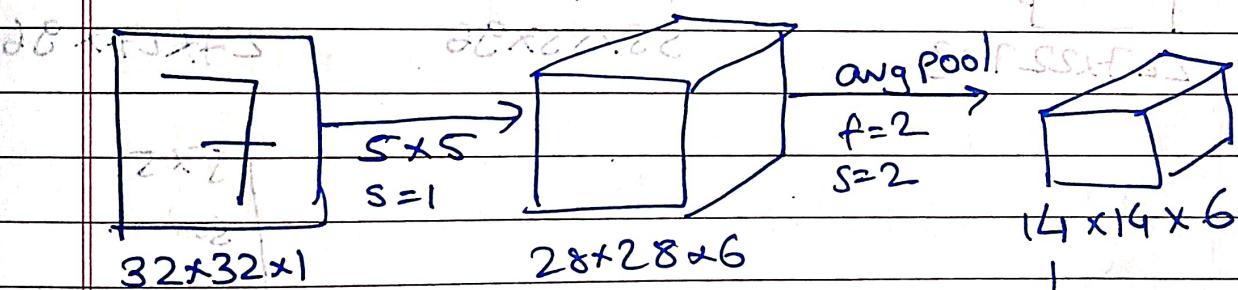
~~Inception Mod~~

### \* LeNet-5

(1998)

→ It was trained on grayscale images.

→ 60k parameters



$28 \times 28 \rightarrow 120 \rightarrow 10 \rightarrow 100$

softmax ←

Page

$n_H \downarrow n_W \downarrow$ , &  $n_C \uparrow$  as you go forward  
 when going forward goes  
 weight  $\downarrow$  bias  $\uparrow$

→ At that time ppl were using  
 Sigmoid / tanh activation more  
 than ReLU

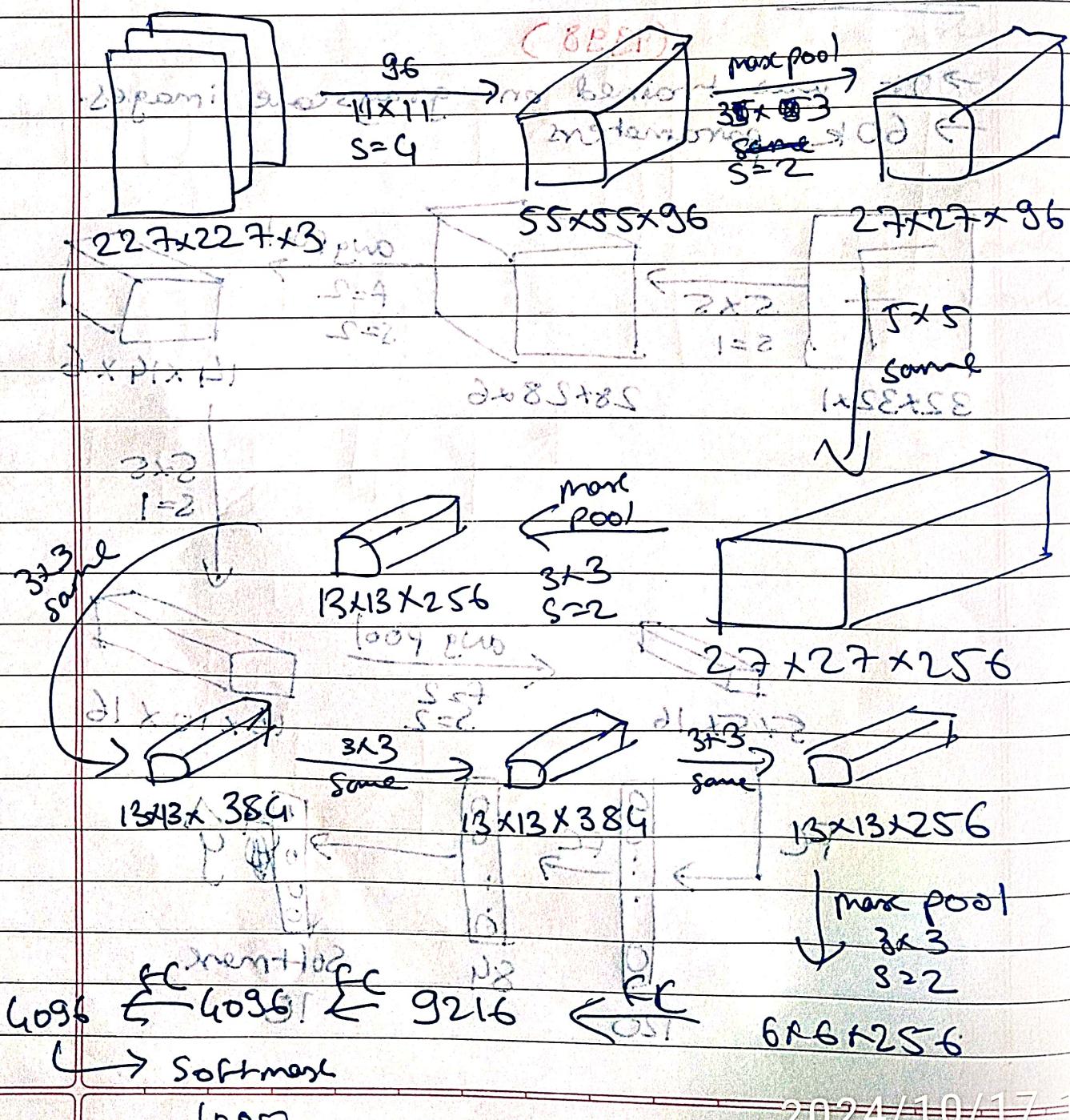
by ~~ReLU~~ V

\* AlexNet (2012)

~60M parameters

→ named after first author Alex Krizhevsky

2. AlexNet \*



- made use of ReLU
- Multiple GPUs
- New layer: Local response Normalisation (LRN)
- complicated & a lot of hyper parameters.

## \* VGG-16 (2015)

$\sim 138 \text{ M}$  parameters.

→ simpler

Conv =  $3 \times 3$  filter

$S = 1$ , same

Max pool =  $2 \times 2$ ,  $S = 2$

$224 \times 224 \times 3 \xrightarrow[\times 2]{\text{Conv 64}} 224 \times 224 \times 64$

$112 \times 112 \times 128 \xrightarrow[\times 2]{\text{Conv 128}}$   $\downarrow \text{pool}$   $112 \times 112 \times 64$

$56 \times 56 \times 128 \xrightarrow[\times 3]{\text{Conv 256}} 56 \times 56 \times 256 \xrightarrow{\text{Pool}} 28 \times 28 \times 256$   
 $\xrightarrow{\text{Conv 512}}$   
 $28 \times 28 \times 512$

$7 \times 7 \times 512 \xrightarrow{\text{Pool}}$   $14 \times 14 \times 512 \xleftarrow[\times 3]{\text{Conv 512}}$   $14 \times 14 \times 512 \xrightarrow{\text{Pool}}$

$4096 \xrightarrow{\text{FC}}$   $4096 \xrightarrow{\text{FC}}$   $4096 \rightarrow \text{softmax}$   
 $100.0$

→ No of filters

$64 \rightarrow 128 \rightarrow 256 \rightarrow 512$ .

→ VGG-19 is an even bigger model  
but performs almost as well as VGG-16