

# Machine Learning

Date 23/08/23

Learning : Improving with experience (E) at some task.  
(+) (+)

Eg:-

Learn to play checkers.

T → Play checkers

P → % of game won

E → Opportunity to play game.

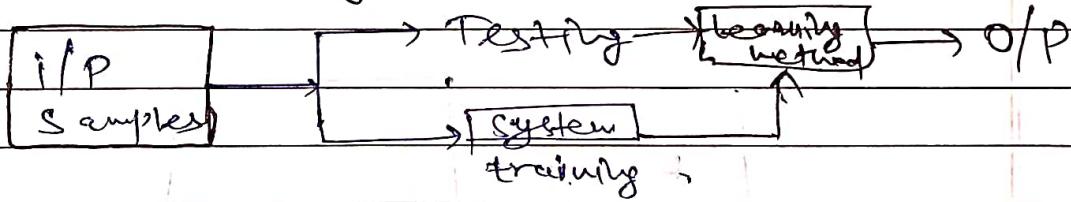
H → Past record / History.

M.L :- Optimizes a performance criteria by learning and using example data or past experience.

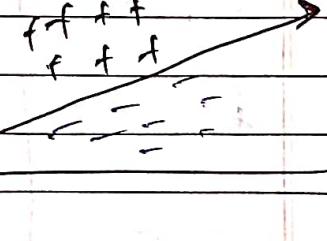
Role of Statistics (Inference from a sample)

Pole of G.S (Computer Science) efficient algorithms to solve optimization problem representing and evaluating the model from inference.

M.L framework



→ M.L has three components (i) representation, (ii) evolution, (iii) optimization.



Class

## Classification Matrix

Date: \_\_\_\_\_ Page: \_\_\_\_\_

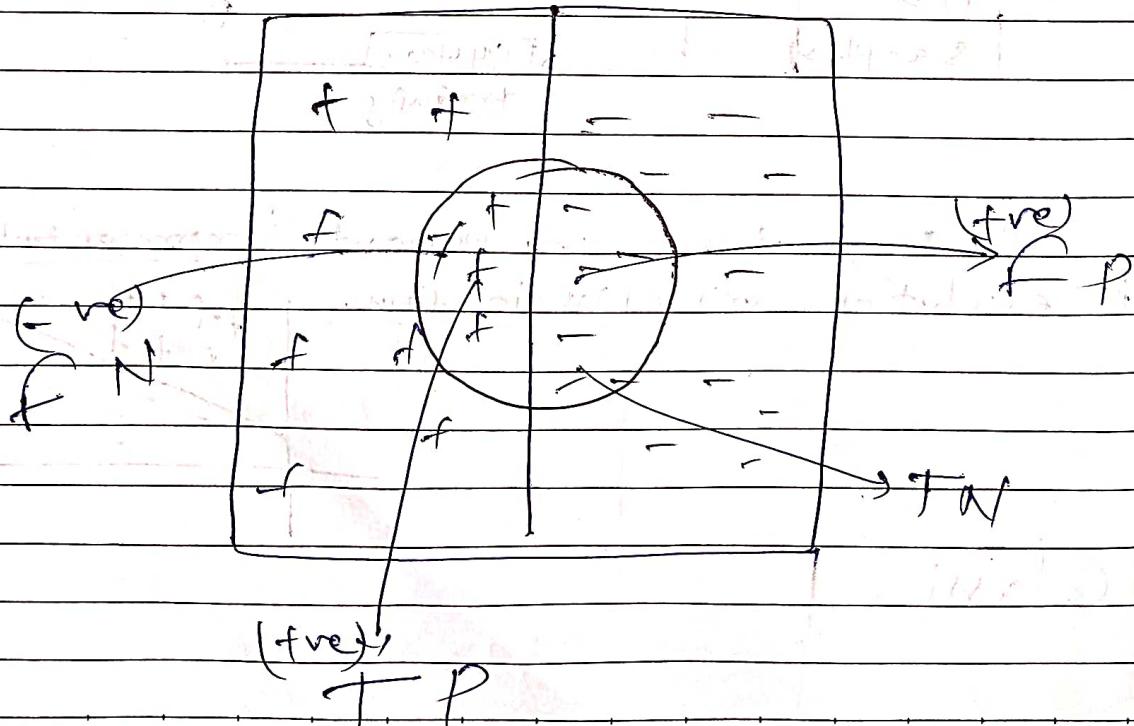
→ Performance measure for ML Classification Algorithm.

True +ve → When you predict an observation belongs to class and it actually <sup>does</sup> belongs to that class.

True -ve → When you predict an observation does not belongs to a class and it actually does not belongs to that class.

False +ve → Occurs when you predict observation belongs to one class (positive) when it really it does not.

False -ve → Occurs when you predict observation belongs to one class (negative) when it does not.



Classification

measure

Date: \_\_\_\_\_ Page: \_\_\_\_\_

Precision =  $\frac{TP}{TP + FP}$

$$\frac{TP}{TP + FP}$$

Recall =  $\frac{TP}{TP + FN}$

$$\frac{TP}{TP + FN}$$

Precision

How many selected items are relevant.

It is defined as precision of relevant ex. ( $TP$ ) among all the examples which are predicted to belong in a certain class.

No. of correctly classified the examples divided by total no. of examples that are classified as true.

Precision =  $\frac{TP}{TP + FP}$

Recall

How many relevant items are selected

T.P rate / Sensitivity is defined as

fraction of examples that were predicted to belong to a class w.r.t all of the examples are truly belongs to in the class.

Recall is the no. of correctly classified true examples divided by total no. of actual true example

M the test said set.

Ex: Calculating accuracy, precision, recall  
for the following model that classifies tumors as malignant (the +ve class)  
or benign the (-ve class)

True +ve (TP)	False Neg
---------------	-----------

Really +ve	Really -ve
------------	------------

ML model predicted +ve	ML model predicted -ve
------------------------	------------------------

No. of TP result = 1	No. of FP result
----------------------	------------------

True -ve	FP
----------	----

Really +ve	Really -ve
------------	------------

ML model predicted: malignant	ML model predicted: not malignant
-------------------------------	-----------------------------------

No. of TP result = 1	No. of FP result
----------------------	------------------

FN

Really +ve	TN
------------	----

ML model predicted: Benign	Really -ve
----------------------------	------------

No. of FN result = 8	No. of TN = 90
----------------------	----------------

Accuracy =  $\frac{TP + TN}{TP + TN + FN + FP} = \frac{1 + 90}{1 + 90 + 8 + 1} = 0.8$

precision =  $\frac{TP}{TP + FP} = \frac{1}{1 + 1} = 0.5$

Recall =  $\frac{TP}{TP + FN}$

Confusion matrix

Predicted

		Predicted	
		T	F
Actual	T	TP	FP
	F	FN	TN

(Positive)

		Actual	
		+	-
Predicted	+	TP	FP
	-	FN	TN

Accuracy: It is defined as the % of correct prediction for test data. It can be calculated by dividing the no. of correct prediction by the no. of total prediction.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$F_1$  value / score / measure

P - precision  
R - recall

$$F_1 = \frac{2Pr}{P+R}$$

$$P = 0.5$$

$$R = 0.11$$

$$F_1 = \frac{2 \times 0.5 \times 0.11}{0.5 + 0.11}$$

$$\frac{0.11}{0.61}$$

$$F_1 = \frac{11}{61}$$

$\rightarrow F_1$  is harmonic mean of precision & recall.

Predicted

	True	False
Actual True	$T_p = 1$	$F_N = 99$
" False	$F_P = 0$	$F_N = 0$

$$P = \frac{T_p}{T_p + F_p} = \frac{1}{1 + 0} = 1 = 100\%$$

$$R = \frac{T_p}{T_p + F_N} = \frac{1}{1 + 99} = \frac{1}{100} = 0.01 = 1\%$$

$$F_1 = \frac{2 \times 1 \times \frac{1}{100}}{2} = \frac{0.02}{2} = \frac{1}{100} = 1\%$$

(Recall - TPR)

Date: \_\_\_\_\_ Page: \_\_\_\_\_

Specificity  $\rightarrow \frac{TN}{TN+FP}$

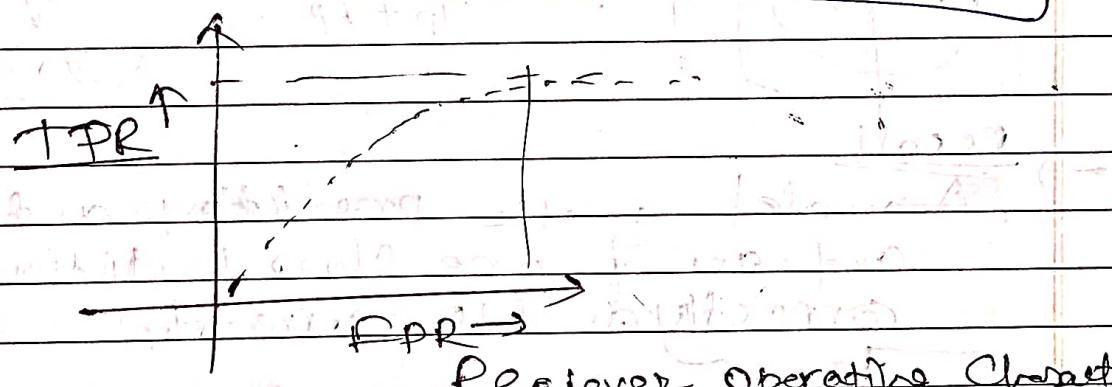
False positive rate  $\rightarrow 1 - \text{Specificity}$

False negative rate  $\rightarrow \frac{FN}{TN+FN}$

~~$\frac{TP + FN + TN}{TN + FN}$~~

Accuracy  $\rightarrow \frac{TP + TN}{TP + TN + FP + FN}$

$$\text{FPR} = \frac{FP}{TN+FN}$$



Receiver operating characteristic curve

(ROC curve)

It is a graph showing the performance of classification model at all classification thresholds. This curve plots two parameter TPR & FPR.

Area under ROC curve (AUROC)

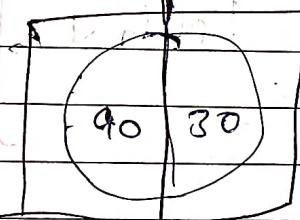
from sklearn.metrics import precision\_score

Regression Performance matrix

Mean Square error :

Precision

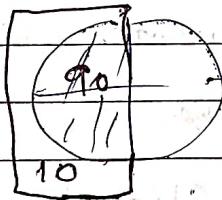
→ A model make prediction and predicts 120 examples as belonging to minority Class 90 of which are correct, 30 of which are incorrect



$$P = \frac{T_p}{T_p + F_p} = \frac{90}{90 + 30} = 75\%$$

recall

→ A model make prediction and predicts and 90 of the Class prediction are correctly and 10 incorrectly



$$\text{Recall} = \frac{90}{90 + 10} = \frac{90}{100} = 90\%$$

A model make prediction and predicts

For example for first minority class where 50 are correct and 20 are incorrect if predict 10 for 2nd class with 99 correct and 51 incorrect final precision for H.

$$P_H = \frac{50}{70} / \frac{10}{124} = \frac{99}{150}$$

$$\text{Precision} = \frac{T_p + 99}{T_p + 100}$$

Date: \_\_\_\_\_ Page: \_\_\_\_\_

$$\begin{bmatrix} x_i \\ A \\ I^2 \end{bmatrix}$$

Strain tensor matrix  $\underline{\underline{\epsilon}}$   
Multiplication  $\underline{\underline{\epsilon}} \underline{\underline{A}} \underline{\underline{x}}$

Numpy, Scipy, matplotlib, Pandas, Sklearn

Sklearn

All the libraries  
matrix

If I + IC - NLP

Monday, Tuesday

$$P = \frac{T_p}{T_p + F_p}$$

True, Wed, Thu  
Fri, Sat

Suppose

- Q) 1000 patient get tested for flu. Of them 900 are actually healthy and 100 are actually sick. For Sick people or a test were done they got +ve for 620 true and 380 -ve. For healthy people same test are done 180 -ve and 8820 got -ve. You have to formulate confusion matrix and find out precision & recall.

Predicted  $\rightarrow$

	+	-
Sick	$T_p = 620$	$F_p = 380$
Healthy	$F_N = 180$	$8820$

(multiple h(x))

Hypothesis Space & Inductive Bias

Inductive learning: It involves using evidence to determine the outcome

(Learn from examples/past experience/etc.)

→ Induction: deriving the function from given data

→ Specific data to (General model)

→ Given a collection of example ( $x^{(i)}, f(x^{(i)})$ ),  $i \in \mathbb{N}$

If a  $f(x)$ , return  $h(x)$  that approximates  $f(x)$   
where  $f(x)$  is called true function which  
correctly maps the input space 'x' to  
the output space  $\gamma$  & approximating  
function  $h(x)$  is called hypothesis function

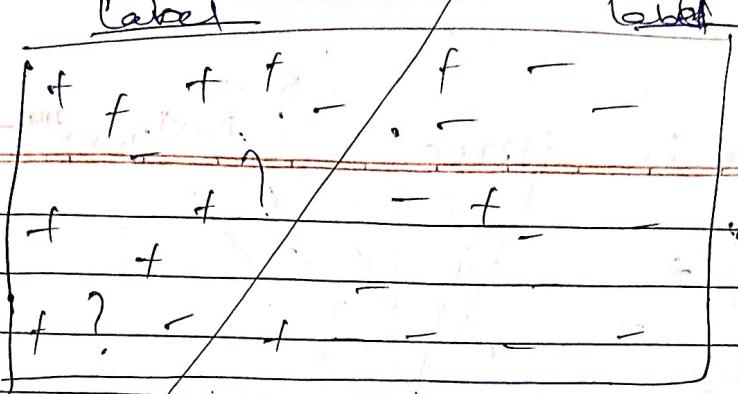
→ Given examples of a function ( $x, f(x)$ )

→ Predict function  $f(x')$  for new example ( $x'$ )

Hypothesis Space ( $H$ ) is the set of all  
legal hypothesis that describe using feature  
you have chosen  $\boxed{h \in H}$

→ Supervised learning works by a device  
that explores hypothesis space.

→ Each setting of a parameter in machine  
is a different hypothesis about the  
function it needs if we pass to  
S/P vectors.



Date: \_\_\_\_\_ Page: \_\_\_\_\_

$f(x)$

(Example of the hypothesis function)

Given data or examples, find the function

$$\begin{array}{c|c} x_1 & \text{Unknown} \\ x_2 & \text{function} \\ x_3 & \\ x_4 & \end{array} \rightarrow Y = f(x_1, x_2, x_3, x_4)$$

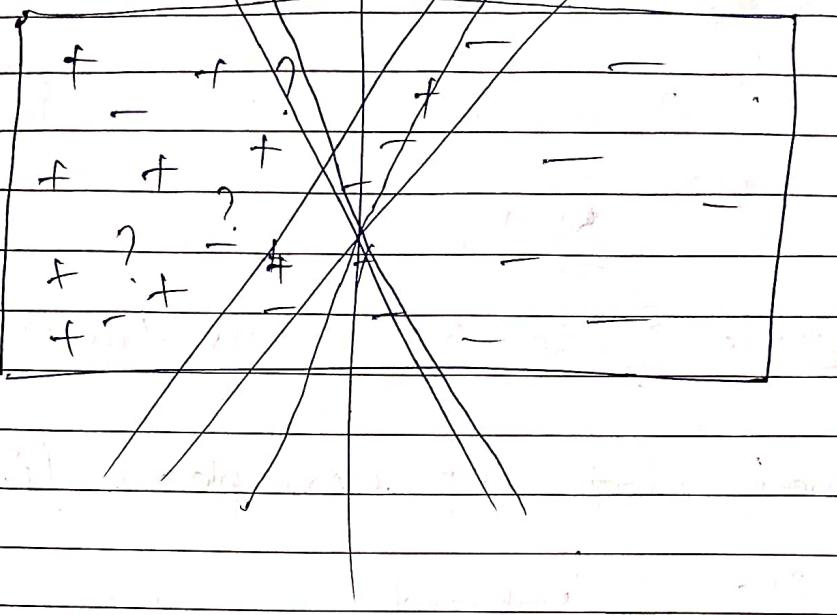
example)

	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	0	0	1	0	0
2	0	1	0	0	0
3	0	0	1	1	1
4	1	0	0	1	1
5	0	1	1	0	0
6	1	1	0	0	0
7	0	1	0	1	0

## Hypothesis Space

$h_1(x)$   
 $h_2(x)$   
 $h_3(x)$

Date: \_\_\_\_\_ Page: \_\_\_\_\_



→ How many distinct Boolean functions of  $n$  boolean attributes  $\rightarrow 2^n$

→ How many distinct

→ Inductive Bias →

→ Need to make assumption & example experience alone doesn't allow us to make conclusion about unseen data instances

→ Inductive Bias : Explicit or implicit assumption (s) about what kind of model is wanted.  
→ The constraint(s) on the hypothesis space is called inductive bias.

BiasIV PreferencesV restrictions

(i) Restrictions: Limit the hypothesis space (ex: book at rules - only small 'h' belongs  $H$  ( $h \in H$ ) are allowed  
(Language bias))

(hypothesis Hypothesis)  
function Space)

(ii)  $\rightarrow$  Preference: propose ordering on hypothesis space / includes all possible hypothesis,

Search bias

Restriction Bias: Hypothesis space restriction bias we restrain the language of (greater hypothesis) hypothesis space ex: K-DNF OR of AND  
sum of product

K-CNF AND of OR  
product of sum

The algorithm

Search through only the subset of possible hypothesis (& complete hypothesis space)

Yet searches these space completely.

This type of bias called restriction & language bias, bias. Because no. of possible hypothesis is restricted.

Preference Bias: It is an order or unit measure matter that serve as best to relation or preference in hypothesis space this algorithm search incomplete through the set of possible hypothesis preferentially select that lead to small decision tree. This type of bias called preference or search Bias.

~~Ex. (Occam's razor)~~ - We prefer a simple formula for ~~a lot~~.

~~Principle of minimal description length  
(An extension of Occam's razor)~~

The best hypothesis is one that minimizes the total length of hypothesis and description of exception of this hypothesis.

## Occam's Razor

### i) Razor

Given two model with same generalization error the simpler is preferred because simplicity desirable itself.

### (ii) Razor

Given two model with same training set error the simpler one should be preferred because it is likely to have better generalization error.

## Cross Validation

### Leave One Out method

Hold Out set - The available data set is divided into two disjoint sets.

(i) Training set -  $d_{tr}$  {for learning a model}

(ii) Test data set -  $d_{te}$  {for testing a model}

→ Training set should not be used in testing and the test set should not be used in learning

(Test set - hold set)

### Leave One Out Cross Validation

The "Leave one out" method is used when the data set is very small; this is a special case of cross validation. Each fold of the cross validation uses only a single test example and all the rest data set is training.

$d_1$	$d_2$	$d_3$	$d_4$	$d_5$
-------	-------	-------	-------	-------

Training -  $d_2 d_3 d_4$   
Testing -  $d_1$

Training -  $d_1 d_3 d_4 d_5$   
Testing -  $d_2$

Training -  $d_1 d_2 d_3 d_5$   
Testing -  $d_4$

Training -  $d_1 d_2 d_3 d_4$   
Testing -  $d_5$

Training -  $d_1 d_3 d_4 d_5$   
Testing -  $d_2$

Training -  $d_1 d_2 d_4 d_5$   
Testing -  $d_3$

→ If original dataset has 'm' examples  
this 'm' field crosses ~~fold~~ validation.

→ 'n' fold cross validation - (Lernmont)

n-fold

Validation set → The selectable

data is divided into 3 subset

(i) Training set

(ii) Validation set

(iii) Test set

It is used frequently for standard parameter in learning algorithm.

In such cases the values give the best accuracy of the validation which are used by parameters.

Cross validation for estimating as well.

02/Sept/2023

Underfitting

Date:

Page:

Overfitting

## Overfitting & Underfitting

Signal & Noise: The underlying pattern that your machine learning aims to learn from data.

Noise: The irrelevant and random data in data set. If data set is large the ML model will learn the signal. It is better when data set is small.

Overfitting: The ML model is very complex. In such a case model learns noise in training data perform very well on it. However when you use to model to test other data set, the model does not perform well gives high error.

→ This model has huge variance.

Underfitting: The ML model is very simple. The model has very few features and regularized way more than needed. Huge bias less variance in its prediction which lead to large error.

Bias: Error caused because the model can not represent the concept.

Variance: Error caused because the learning algorithm overfits to small change (noise) in the fix training data.

Linear Regression - (i) Regression Analysis

Regression Analysis → It is statistical technique for linear relationship.

→ Dependent Variable & Independent Variable

Regression

Regression It is used in more and more in Analytical research that affects daily life from how much we pay automobile insurance to what ads appear in social media.

Prediction

- Linear regression is standard mathematical technique for predicting numerical outcome
- The Linear model is one of an important model of parametric model.

$$Y = C + mx$$

Dependent Variable

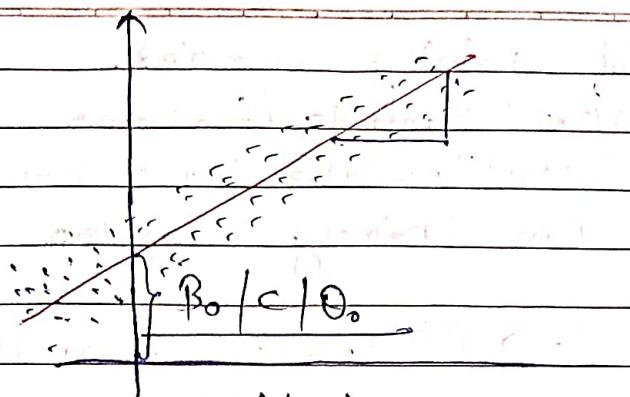
Independent Variable

$$\text{P.M.L} \rightarrow Y = \beta_0 + \beta_1 x$$

{  $\beta_0, \beta_1$  - parameters }

$$h_0(x) = \theta_0 + \theta_1 x$$

Estimate the values of  $\theta_0$  &  $\theta_1$



We begin by considering linear regression

$$f(x, \theta) = \theta_0 + \theta_1 x, \text{ where } \theta = [\theta_0, \theta_1]^T$$

are parameters

We need to set.

We can measure the prediction loss in term of square error.

Cost function / ~~or~~  $\text{Loss}(y, \hat{y})$

$$= (y - \hat{y})^2$$

↓  
 actual value      ↓  
 predicted value

$x$	$y$	$\theta_0 = 0$	$\theta_0 = 10$
100	100	$\theta_1 = 1$	$\theta_1 = 0$
200	200	$\theta_1 = 1$	$\theta_1 = 0.5$
300	300	$\theta_1 = 1$	$\theta_1 = 0.33$

$$\text{Or } 1/3 \times 100 \approx 100 \quad \frac{1}{3} \times 280$$

Performance Measure  $\frac{(100-100)^2}{n}$  cost function

Cost function / loss =  $(y - \hat{y})^2$

Mean Square Error

Square root mean error.

$n = \text{no. of training examples}$

$x^{(i)} = \text{"input" variable / features}$

$y^{(i)} = \text{"Output" variable / target example}$

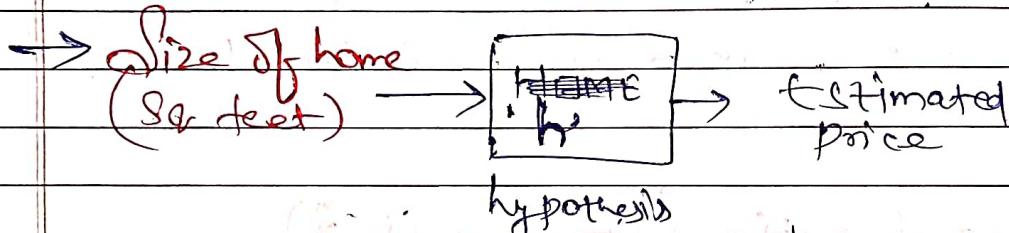
$(x, y) \rightarrow \text{One training example}$

$(x^{(i)}, y^{(i)}) \rightarrow \text{i}^{\text{th}} \text{ training example}$

### Example

$$x^{(1)} = 450, y^{(1)} = 350$$

$$x^{(1)} = 280, y^{(1)} = ?.$$



→ 'h' maps from  $x^{(i)}$  to  $y^{(i)}$

→ Representation of  $h$

$$\boxed{h_0(x) = \theta_0 + \theta_1 x}$$

values of parameters

$$\theta_0 = 0, \theta_1 = 0.5$$

$$\theta_0 = 1.5, \theta_1 = 0$$

$$\theta_0 = 1, \theta_1 = 1.5$$

(i) Choose  $\theta_0, \theta_1$  such that  $h_0(x)$  is close to 'y' for training example  $(x, y)$

(ii) Minimize  $(\theta_0, \theta_1) (h_0(x) - y)^2$

(iii) Minimize

$$\frac{1}{2} \sum_{i=1}^n (h_0(x^{(i)}) - y^{(i)})^2$$

↓  
part of  $x^{(i)}$

Optimize  $\xrightarrow{\text{Minimize}} \xrightarrow{\text{Maximize}}$  either

### Gradient Descent method

$$\left| \text{Minimize } (\theta_0, \theta_1) \frac{1}{2n} \sum_{i=1}^n (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \right|$$

Idea: find approx parameters  $(\theta_0, \theta_1)$  which minimize the cost function.

- Start with some  $\theta_0, \theta_1$
- update the value of  $\theta_0, \theta_1$  to minimize the cost function
- Stop till get good result or no change in  $\theta_0, \theta_1$ .

### Gradient Descent method.

- Initialize with random  $\theta_0, \theta_1$
- Calculate the gradients of cost function w.r.t parameters of cost function ( $\theta$ )

$$\frac{\partial}{\partial \theta_0} = \frac{1}{n} \sum_{i=1}^n (h_\theta(x) - y)^2$$

$$\frac{\partial}{\partial \theta_1} = \frac{1}{n} \sum_{i=1}^n \{ h_\theta(x) - y \}^2$$

- update the parameters

$$(\theta_0)_{\text{new}} = (\theta_0) - \alpha \cdot \frac{\partial}{\partial \theta_0} \left\{ \frac{1}{n} \sum_{i=1}^n (h_\theta(x) - y)^2 \right\}$$

$$(\theta_1)_{\text{new}} = (\theta_1) - \alpha \cdot \frac{\partial}{\partial \theta_1} \left\{ \frac{1}{n} \sum_{i=1}^n (h_\theta(x) - y)^2 \right\}$$

- Stop till get good result or no change in  $\theta_0, \theta_1$ .
- Learning Parameters

$$\left| -0.05 \leq \alpha \leq 1.0 \right|$$

→ If value of  $\alpha$  is too small then gradient descent can be converge too slow.

→ If value of  $\alpha$  is very large then gradient descent can be passed the minima and may fail to converge.

$$\frac{d\mathcal{L}(\theta)}{d\theta_0} \geq (h_\theta(x) - y)$$

$$\frac{d\mathcal{L}(\theta)}{d\theta_1} = \frac{d}{d\theta_1} \left( \frac{1}{2} (h_\theta(x) - y)^2 \right)$$

$$= (h_\theta(x) - y) \frac{d}{d\theta_1} \left( \frac{1}{2} (\theta_0 + \theta_1 x) - y \right)$$

$$\frac{d(\text{Cost function})}{d\theta_i} = (h_\theta(x) - y) x_i$$

$$x = [1, 2, 3, 4, 5]$$

$$y = [3, 6, 7, 11, 15]$$

$$\theta_1 = ?, \quad \theta_0 = ?$$

Learning rate ( $\alpha$ ) = 0.5

$$\hat{y} = \theta_0 + \theta_1 x \quad (x=1, y=3)$$

for  $x=1$

$$\hat{y}=12 = h_\theta(x), \quad y=3$$

$$\frac{d\mathcal{L}(\theta)}{d\theta_0} = (h_\theta(x) - y)$$

$$\theta_0 = \theta_0 - \alpha \frac{d}{d\theta_0} (h_\theta(x) - y)^2 = \frac{(12 - 3)}{5}$$

$$y = \theta_0 + \theta_1 x$$

Date: \_\_\_\_\_ Page: \_\_\_\_\_

(1)  $x = [1, 2, 3, 4, 5]$

$$y = [3, 6, 2, 11, 15] \text{ learning rate}$$

$$\theta_0 = 2, \theta_1 = 5$$

$$x < 0.5$$

$$\frac{\partial C(\theta)}{\partial \theta_0} = (h_\theta(x) - y)$$

$$= \frac{\partial}{\partial \theta_0} \left( \frac{1}{2} h_\theta(x) - y \right)^2$$

$$= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_0} (h_\theta(x) - y)$$

$$= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x - y)$$

$$\frac{\partial C(\theta)}{\partial \theta_0} = (h_\theta(x) - y) (\theta_0 + \theta_1 x - y)$$

$$\frac{\partial C(\theta)}{\partial \theta_1} = (h_\theta(x) - y) x_i$$

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (h_\theta(x) - y)^2$$

$$\theta_1 := \theta_1 - \alpha (h_\theta(x) - y) x_i$$

$X \geq 1, Y \geq 3$  Date: \_\_\_\_\_ Page: \_\_\_\_\_

$$Y \geq 0, f_0(x) \quad 0, 5$$

$$Y \geq h_0(x) \geq 0, 5$$

$$h_0(x) = f_0(x) - 1$$

$$h_0(x) = 12, Y \geq 3$$

$$\begin{cases} X \geq 2 \\ Y \geq 6 \end{cases}$$

$$h_0(x) = 5f_0(x) - 2$$

$$h_0(x) = 19, Y \geq 6$$

$$\begin{cases} X \geq 7 \\ Y \geq 2 \end{cases}$$

$$h_0(x) = 26, Y \geq 7$$

$$\begin{cases} X \geq 4 \\ Y \geq 11 \end{cases}$$

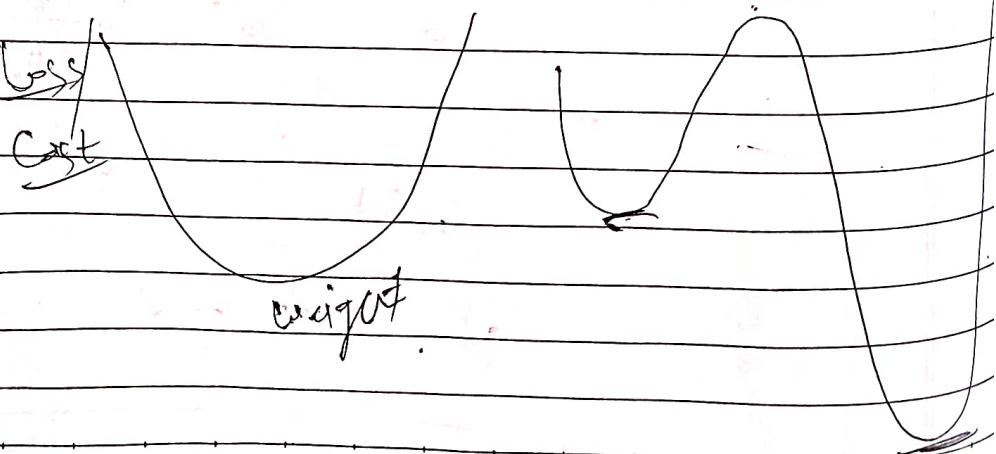
$$h_0(x) = 33, Y \geq 11$$

$$\begin{cases} X \geq 11 \\ Y \geq 15 \end{cases}$$

$$h_0(x) = 40, Y \geq 15$$

Learning loss  
Cost  
Rate

weight



① 1946) = 1946) 1946) 1946) 1946)

Epoch = 0

09/09/23

Cost Square error $(y - \hat{y})^2$	Cost $\frac{\partial}{\partial \theta_0}$	Cost $\frac{\partial}{\partial \theta_1}$
81	9	$9x1$
169	13	$13x1$
261	19	$19x1$
481	22	$22x1$
625	25	$25x1$
		<u>Sum = 76</u>
		<u>Sum = 5</u>

$$\theta_{(0)}^{\text{new}} = \theta_{(0)}^{\text{(old)}} - \frac{\alpha d}{m}$$

$$= 3 - \frac{1}{5} \times 12.5$$

$$\theta_1^{\text{(new)}} = \theta_1^{\text{(old)}} - \frac{\alpha d}{m}$$

$$= 7 - \frac{1}{5} \times 6$$

Epoch = 1

$$\theta_0^{\text{(new)}}$$

$$\theta_1^{\text{(new)}}$$

113 Master

1000 Dry weight Method 1000g  
Calcd. 6.9g

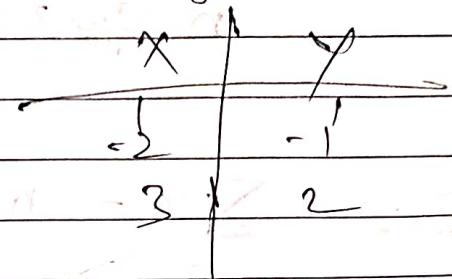
1000 - 212.1 = 787.9

1000 - 212.1 = 787.9

1000 - 212.1 = 787.9

R	I	P	Q
1	13	6.9	9
1	14	7.0	10
1	15	7.1	11
1	16	7.2	12
1	17	7.3	13
1	18	7.4	14
1	19	7.5	15
1	20	7.6	16
1	21	7.7	17
1	22	7.8	18
1	23	7.9	19
1	24	8.0	20
1	25	8.1	21
1	26	8.2	22
1	27	8.3	23
1	28	8.4	24
1	29	8.5	25
1	30	8.6	26
1	31	8.7	27
1	32	8.8	28
1	33	8.9	29
1	34	9.0	30
1	35	9.1	31
1	36	9.2	32
1	37	9.3	33
1	38	9.4	34
1	39	9.5	35
1	40	9.6	36
1	41	9.7	37
1	42	9.8	38
1	43	9.9	39
1	44	10.0	40
1	45	10.1	41
1	46	10.2	42
1	47	10.3	43
1	48	10.4	44
1	49	10.5	45
1	50	10.6	46
1	51	10.7	47
1	52	10.8	48
1	53	10.9	49
1	54	11.0	50
1	55	11.1	51
1	56	11.2	52
1	57	11.3	53
1	58	11.4	54
1	59	11.5	55
1	60	11.6	56
1	61	11.7	57
1	62	11.8	58
1	63	11.9	59
1	64	12.0	60
1	65	12.1	61
1	66	12.2	62
1	67	12.3	63
1	68	12.4	64
1	69	12.5	65
1	70	12.6	66
1	71	12.7	67
1	72	12.8	68
1	73	12.9	69
1	74	13.0	70
1	75	13.1	71
1	76	13.2	72
1	77	13.3	73
1	78	13.4	74
1	79	13.5	75
1	80	13.6	76
1	81	13.7	77
1	82	13.8	78
1	83	13.9	79
1	84	14.0	80
1	85	14.1	81
1	86	14.2	82
1	87	14.3	83
1	88	14.4	84
1	89	14.5	85
1	90	14.6	86
1	91	14.7	87
1	92	14.8	88
1	93	14.9	89
1	94	15.0	90
1	95	15.1	91
1	96	15.2	92
1	97	15.3	93
1	98	15.4	94
1	99	15.5	95
1	100	15.6	96
1	101	15.7	97
1	102	15.8	98
1	103	15.9	99
1	104	16.0	100
1	105	16.1	101
1	106	16.2	102
1	107	16.3	103
1	108	16.4	104
1	109	16.5	105
1	110	16.6	106
1	111	16.7	107
1	112	16.8	108
1	113	16.9	109
1	114	17.0	110
1	115	17.1	111
1	116	17.2	112
1	117	17.3	113
1	118	17.4	114
1	119	17.5	115
1	120	17.6	116
1	121	17.7	117
1	122	17.8	118
1	123	17.9	119
1	124	18.0	120
1	125	18.1	121
1	126	18.2	122
1	127	18.3	123
1	128	18.4	124
1	129	18.5	125
1	130	18.6	126
1	131	18.7	127
1	132	18.8	128
1	133	18.9	129
1	134	19.0	130
1	135	19.1	131
1	136	19.2	132
1	137	19.3	133
1	138	19.4	134
1	139	19.5	135
1	140	19.6	136
1	141	19.7	137
1	142	19.8	138
1	143	19.9	139
1	144	20.0	140
1	145	20.1	141
1	146	20.2	142
1	147	20.3	143
1	148	20.4	144
1	149	20.5	145
1	150	20.6	146
1	151	20.7	147
1	152	20.8	148
1	153	20.9	149
1	154	21.0	150
1	155	21.1	151
1	156	21.2	152
1	157	21.3	153
1	158	21.4	154
1	159	21.5	155
1	160	21.6	156
1	161	21.7	157
1	162	21.8	158
1	163	21.9	159
1	164	22.0	160
1	165	22.1	161
1	166	22.2	162
1	167	22.3	163
1	168	22.4	164
1	169	22.5	165
1	170	22.6	166
1	171	22.7	167
1	172	22.8	168
1	173	22.9	169
1	174	23.0	170
1	175	23.1	171
1	176	23.2	172
1	177	23.3	173
1	178	23.4	174
1	179	23.5	175
1	180	23.6	176
1	181	23.7	177
1	182	23.8	178
1	183	23.9	179
1	184	24.0	180
1	185	24.1	181
1	186	24.2	182
1	187	24.3	183
1	188	24.4	184
1	189	24.5	185
1	190	24.6	186
1	191	24.7	187
1	192	24.8	188
1	193	24.9	189
1	194	25.0	190
1	195	25.1	191
1	196	25.2	192
1	197	25.3	193
1	198	25.4	194
1	199	25.5	195
1	200	25.6	196
1	201	25.7	197
1	202	25.8	198
1	203	25.9	199
1	204	26.0	200
1	205	26.1	201
1	206	26.2	202
1	207	26.3	203
1	208	26.4	204
1	209	26.5	205
1	210	26.6	206
1	211	26.7	207
1	212	26.8	208
1	213	26.9	209
1	214	27.0	210
1	215	27.1	211
1	216	27.2	212
1	217	27.3	213
1	218	27.4	214
1	219	27.5	215
1	220	27.6	216
1	221	27.7	217
1	222	27.8	218
1	223	27.9	219
1	224	28.0	220
1	225	28.1	221
1	226	28.2	222
1	227	28.3	223
1	228	28.4	224
1	229	28.5	225
1	230	28.6	226
1	231	28.7	227
1	232	28.8	228
1	233	28.9	229
1	234	29.0	230
1	235	29.1	231
1	236	29.2	232
1	237	29.3	233
1	238	29.4	234
1	239	29.5	235
1	240	29.6	236
1	241	29.7	237
1	242	29.8	238
1	243	29.9	239
1	244	30.0	240
1	245	30.1	241
1	246	30.2	242
1	247	30.3	243
1	248	30.4	244
1	249	30.5	245
1	250	30.6	246
1	251	30.7	247
1	252	30.8	248
1	253	30.9	249
1	254	31.0	250
1	255	31.1	251
1	256	31.2	252
1	257	31.3	253
1	258	31.4	254
1	259	31.5	255
1	260	31.6	256
1	261	31.7	257
1	262	31.8	258
1	263	31.9	259
1	264	32.0	260
1	265	32.1	261
1	266	32.2	262
1	267	32.3	263
1	268	32.4	264
1	269	32.5	265
1	270	32.6	266
1	271	32.7	267
1	272	32.8	268
1	273	32.9	269
1	274	33.0	270
1	275	33.1	271
1	276	33.2	272
1	277	33.3	273
1	278	33.4	274
1	279	33.5	275
1	280	33.6	276
1	281	33.7	277
1	282	33.8	278
1	283	33.9	279
1	284	34.0	280
1	285	34.1	281
1	286	34.2	282
1	287	34.3	283
1	288	34.4	284
1	289	34.5	285
1	290	34.6	286
1	291	34.7	287
1	292	34.8	288
1	293	34.9	289
1	294	35.0	290
1	295	35.1	291
1	296	35.2	292
1	297	35.3	293
1	298	35.4	294
1	299	35.5	295
1	300	35.6	296
1	301	35.7	297
1	302	35.8	298
1	303	35.9	299
1	304	36.0	300
1	305	36.1	301
1	306	36.2	302
1	307	36.3	303
1	308	36.4	304
1	309	36.5	305
1	310	36.6	306
1	311	36.7	307
1	312	36.8	308
1	313	36.9	309
1	314	37.0	310
1	315	37.1	311
1	316	37.2	312
1	317	37.3	313
1	318	37.4	314
1	319	37.5	315
1	320	37.6	316
1	321	37.7	317
1	322	37.8	318
1	323	37.9	319
1	324	38.0	320
1	325	38.1	321
1	326	38.2	322
1	327	38.3	323
1	328	38.4	324
1	329	38.5	325

→ Consider the set of points  $(1, 1)$ ,  $(-2, -1)$ , &  $(3, 2)$ . Plot these points & the least square regression line on the same graph.



### Matrix Notation

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Cost function:  $\frac{1}{n} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2$

$$\equiv \frac{1}{n} \left\| \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \right\|^2$$

$$= \frac{1}{n} \|Y - X\theta\|^2$$

$$\frac{d}{d\theta} \frac{1}{n} \|Y - X\theta\|^2 = \frac{d}{d\theta} \frac{1}{n} (Y - X\theta)^T (Y - X\theta)$$

$$\frac{d}{d\theta} \frac{1}{n} \|Y - X\theta\|^2 = \frac{d}{d\theta} \frac{1}{n} (Y^T Y - 2Y^T X\theta + \theta^T X^T X\theta)$$

### ③ Inverse Method

Date \_\_\_\_\_ Page \_\_\_\_\_

$$\Rightarrow -2x^T y + 2x^T x \theta = 0$$

Feature  $\theta$

$$x^T y = x^T x \theta$$

Normal Equation:  $\theta = (x^T x)^{-1} x^T y$

→ horse's saddle

12/09/2023

Toggle (data set) (CSF)

reinforcement (Support Vector Machine)

Numpy, Pandas, Matplotlib  
(create library file) (plotting) (graph)

McKearn, Scipy

Assignment - (i) Load the dataset for graphical representation  
 (ii) Shape  
 (iii) Description

Date: 19/09/2023

## 1) Linear Regression with multiple variables

- $X_1, X_2, X_3, X_4$  and more

- Examples of multiple variables (Housing data set)

	price	area	bedrooms	bathrooms	stories
d <sub>1</sub>	133000	7420	4	1	2
d <sub>2</sub>	1225000	8960	4	4	4
d <sub>3</sub>	1825000	8960	5	3	5
	avg.				

→ Here y = price

$$X [x_1, x_2, x_3, x_4]$$

$x_1 = \text{area}$

Rigid  
lessened

$x_2 = \text{bed rooms}$

Elastic

(Lessor - least absolute shrinkage)

$x_3 = \text{bathrooms}$

Elastic

& selection parameters

$x_4 = \text{stories}$

Stiffer  
matrix

→  $X^i$ : input features of  $i$ th training example

→  $X_j^i$ :  $j$ th features of  $i$ th training example.

$$\rightarrow h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

→ Generalize hypothesis for multiple variables

$$\rightarrow h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$X = \begin{bmatrix} x \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\rightarrow h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\rightarrow [\theta_0 \quad \theta_1 \quad \theta_2 \quad \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$h_0(x) = \theta^T x$$

### Gradient descent Method

From this status

### Regularization

Solve

Date \_\_\_\_\_

Page \_\_\_\_\_

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\left[ \begin{matrix} \theta^0 & \theta^1 & \theta^2 & \dots & \theta^n \end{matrix} \right] \left[ \begin{matrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \right]$$

$$\left. \right\} h_{\theta}(x) = \theta^T(x)$$

Gradient Descent Method

Generalizing with

Regularization -

Defining the loss function

# Regularization

Penalty

Date: Shrink  
reduce  
Bell-Hobby

## Regularization (Shrinkage in stats)

- Solution to overfitting
- Reduce number of features
- Manually select features to keep
- Model selection algorithm.

### Regularization

- Keep ~~all~~ features but reduced magnitude or value of parameters  $\theta_j$ .
- Works well when we have a lot of features.

### Linear Regression with multiple variables

- Small value for parameter
- Less prone to overfitting
- Simpler hypothesis.

### Add regularization parameter to cost function

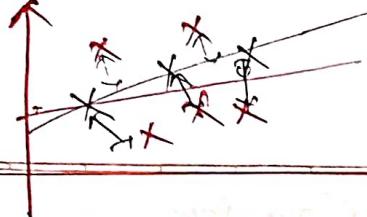
( $\lambda$ ) to shrink parameters.

→ If  $\lambda$  (lambda) is set to an extremely large value this would result in under fitting



↳ High bias

→ Only penalize thetas from 'i', not for '0'



X - training data  
X - testing data  
Date: \_\_\_\_\_ Page: \_\_\_\_\_

reduce stiffness of curve

→ for reducing stiffness or generalized model we have to add penalty to cost function,

### Objective

- Shrink the coefficients or weights of features in the model.
- Getting rid of high degree of polynomial features from the model.
- The idea of regularization revolves around modifying the loss function / Cost function  $C_\theta(\theta)$  in particular, we add a regularization term that enforces some specified properties of model parameter.

$$h_\theta(x) = h_0(x) + \lambda \text{Reg}(\theta)$$

- There are ~~three~~ two types of regularization
- Lasso ( $\ell_1$  norm)
  - Ridge ( $\ell_2$  norm)
- elastic net  
Combination of Lasso & Ridge).

## L1 Lasso Regularization

→ We need to calculate a regularization term that penalize parameter magnitudes for our cost/loss function, we will again use MSE (mean square error)

→ Regularized Cost function / loss function

$$\text{Cost function (L1)} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 + \frac{\lambda}{m} \sum_{j=1}^n |\theta_j|$$

(Loss function)  
↓  
↓  
(Normal)

$$\text{Regularization} = \frac{\lambda}{2} \sum_{j=1}^m |\theta_j|$$

$$= \frac{\lambda}{2} \sum_{j=1}^m |\theta_j|^2$$

$$= \frac{\lambda}{2} \|\theta\|_2^2$$

Diagram

MLE

Probability

Distribution

## Maximum Likelihood Estimation (MLE)

MLE is method that determines values for the parameters of a model. The parameters value are found such that the process described model produced the data that is actually observed.

### MLE Principle :

→ Choose parameters that maximize the likelihood function.

### I Minimise

→ A coin is flipped 100 times. Given that there is 55 heads, find the maximum likelihood estimate for the probability 'p' of heads on a single toss.

$$nCr(p)^r \cdot (1-p)^{n-r}$$

Probability of occurrence of head =

$$P(55 \text{ heads}) = \frac{100}{55} C_{55} p^{55} (1-p)^{45} = \frac{55}{200}$$

$$55,100 C_{55} p^{54} (1-p)^{45} - 45,100 C_{55} p^{55} (1-p)^{44}$$

$$(1-p)^{44} 100 C_{55} p^{54} ((1-p) \cdot 55 - 45 \cdot p)$$

$$(-55p + 55 - 45p)$$

$$= 100 C_{55} (1-p)^{54} p^{54} (-100p + 55) \Rightarrow$$

$$p = 55/100 = 0.55$$

Date: \_\_\_\_\_ Page: \_\_\_\_\_  
Experiment! flip the coin 100 times &  
Count the no. of heads.

Note! The data is the result of experiment  
in this case it is 55 heads.

Parameters of interest: We are interested  
in the value of unknown parameter  $p$ .

Likelihood or Likelihood function! This is  
 $(\text{data}/p)$

Note! It is a function of both the  
data & parameter  $p$ . In this case  
the likelihood is  $p^{55}(1-p)^{45}$ .

$$Y = \log p^{55} (1-p)^{45}$$

$$\log Y$$

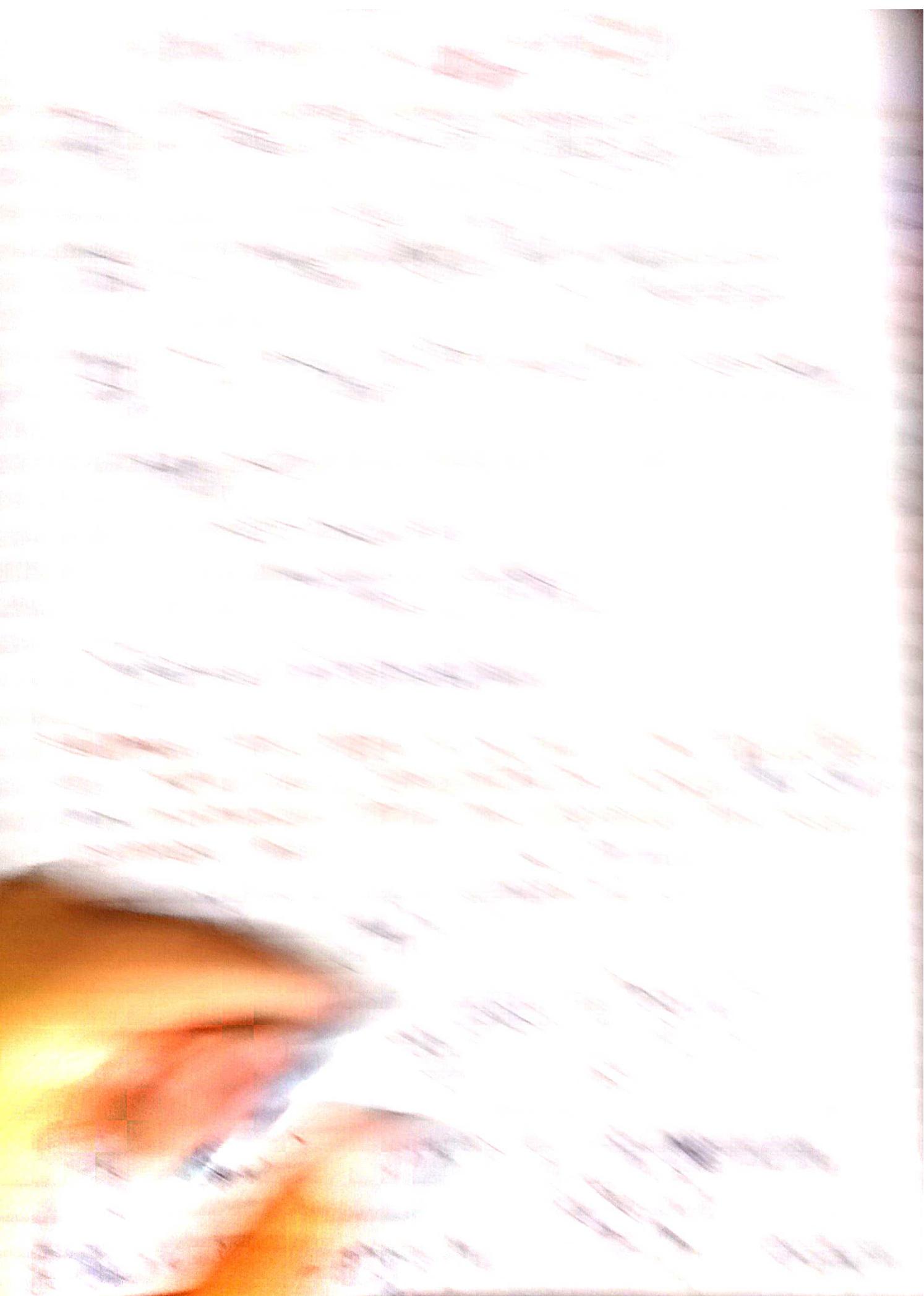
$$\log Y = \log 55 + \log p + \log (1-p)^{45}$$

$$\log Y = \log 55 + 55 \log p + 45 \log (1-p)$$

(differentiate wrt.  $p$ )

$$\frac{dy}{dp} = 0 + \frac{55}{p} + \frac{45}{1-p}$$

$$\frac{dy}{dp} = \log p^{55} (1-p)^{45} \left\{ \frac{55}{p} - \frac{45}{1-p} \right\}$$



$$\frac{dy}{d(\mu)} = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{(x_i-\mu)^2}{2\sigma^2}} \sum \frac{(x_i-\mu)}{\sigma^2} = 0$$

$$\frac{dy}{d(\sigma)} = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{(x_i-\mu)^2}{2\sigma^2}} \left( \frac{(x_i-\mu)^2}{\sigma^2} + \frac{1}{2\pi\sigma^2} \right) = 0$$

$$\boxed{\mu = \bar{x}}$$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\sigma^4 + \frac{(x_i - \bar{x})^2}{\sigma^3}$$

$$X_i = 9, 9.5, 11.5, \dots \quad \mu = \bar{x} = \frac{9+9.5+11}{3} = 10$$

$$f(x_i/\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\bar{x})^2}{2\sigma^2}} \quad \bar{x} = \frac{29.5}{3} = 9.833$$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i^2 - (\bar{x})^2 - 2\bar{x}(x_i))}$$

$$\begin{aligned} x &= \{x_i\} \\ n &= \{x_i\}^n \end{aligned}$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i^2 - n(\bar{x})^2)$$

$$\frac{1}{n} (81 + 90.25 + 101) \rightarrow 3x \frac{(9.83)}{24.5} = 29.0$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Marginal Probability :- The probability of an event irrespective of the outcomes of the other random variables eg.  $P(A)$ .

Q) In the other words, it is known as simple probability which refers to the probability of occurrence of a single event  $\text{eg. } P(A), P(B)$

Joint Probability :- The probability of two (or more) simultaneous events eg.  $P(A \text{ and } B)$  or  $P(A \text{ or } B)$

→ It is often described in terms of event A & B from two dependent random variables eg.  $X \& Y$ , the joint probability is often summarized as just outcomes eg. A & B

Conditional probability :- The probability of two (more) event given the occurrence of another

eg:  $P(A \text{ given } B)$  or  $P(A|B)$ .

Q)

Gender	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	Total
Male	30	80	90	200
Female	20	40	60	120
Total	50	120	150	300

(i) find out probability of total no. of female in students

$$P(F) = \frac{120}{300} = \frac{1}{3}$$

Date: \_\_\_\_\_ Page: \_\_\_\_\_  
 Q) probability that next male student with grade  
 find out

$$P(M) = \frac{30}{200} = 0.15 \text{ about } 15\%$$

$$P(M) = \frac{30}{200} = \frac{1}{10} = 0.1$$

Q) Considered as data is given that student  
 is male then what the probability

$$P(M_A) = \frac{20}{200} = \frac{1}{10} = 0.15$$

$$P(A|B)$$

$$P(A|M) = \frac{P(A \cap M)}{P(M)} = \frac{30}{200} = \frac{3}{20}$$

~~$$P(A|B) = \frac{P(A \cap B)}{P(B)} \Rightarrow P(A|B) \cdot P(B) = P(A \cap B)$$~~

~~$$P(B|A) \cdot P(A) \Rightarrow P(B|A) \cdot P(A) = P(A \cap B)$$~~

~~$$P(A \cap B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$~~

Baye's theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$A_1 \cap B, A_2 \cap B, A_3 \cap B$$

$$P(B) = A_1 \cap B + A_2 \cap B + A_3 \cap B$$

$$P(B) = P(B|A_1) \cdot P(A_1) + P(B|A_2) \cdot P(A_2) + P(B|A_3) \cdot P(A_3)$$

$$P(B) = \underbrace{P(B|A_i) P(A_i)}_{\sum P(B|A_i) \cdot P(A_i)}$$

Baye's Theorem  $\rightarrow$  Principled way of calculating a conditional probability without the joint probability term result

$P(A|B)$  is referred to as posterior probability and  $P(A)$

is referred to as prior probability.

$P(B|A)$   $\rightarrow$  is referred to as likelihood.

$P(B)$   $\rightarrow$  is referred to as evidence.

Likelihood  $\times$  prior

Posterior:  $\frac{\text{Likelihood} \times \text{prior}}{\text{Evidence}}$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Posterior ( $P$ )

$\rightarrow$  What is  $P$ , that there is fire given that there is

a. Smoker ( $S$ )  $P(F|S) = P$

$A_1 \cap B$ ,  ~~$A_1 \cup B$~~ ,  $A_2 \cap B$ ,  $A_3 \cap B$

$$P(B) = A_1 \cap B + A_2 \cap B + A_3 \cap B$$

$$P(B) = P(B|A_1) \cdot P(A_1) + P(B|A_2) \cdot P(A_2) + P(B|A_3) \cdot P(A_3)$$

$$P(B) = P(B|A_i) P(A_i)$$

$$\sum_{i=1}^3 P(B|A_i) \cdot P(A_i)$$

Baye's Theorem  $\rightarrow$  Principle way of calculating a conditioned-probability without the joint probability the result

$P(A|B)$  is referred to as posterior probability and  $P(A)$

is referred to prior probability.

$P(B|A)$  is referred likelihood.

$P(B)$  is referred to as evidence.

Likelihood  $\times$  prior  
posterior = evidence

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Posterior ( $F$ )

$\Rightarrow$  what is  $P_F$  that there is the given that there is

a. smoker ( $S$ )  $P(F|S) = P_F$

$P(A)$  = Prior

$P(\text{smoker})$  = evidence

$P(B|A)$  = Likelihood

$P(A|B)$  = Posterior

- (Q) There are 3 boxes labeled 'A', 'B', 'C'.  
 Box A contains 2 Red & 3 Black balls  
 Box B contains 3 Red & 1 Black  
 & Box C contains 1 Red & 4 Black

	Red	Black
Box - A	2	3
Box - B	3	1
Box - C	1	4

There are 3 boxes are identical and have equal probability of getting picked considered as that a red ball is chosen then what is probability that this red ball was picked out of box A.

$P(R/A)$  =  $P(A)$

$P(R|A)$  =  $P(RP)$

$P(R)$

(2)  
5

Date: \_\_\_\_\_ Page: \_\_\_\_\_

$$P(A|R) = P(R|A) \cdot P(A)$$

$$\{ P(R|A) \cdot P(A)$$

$$= \frac{2}{5} \cdot \frac{1}{3}$$

$$\frac{2}{5} \cdot \frac{1}{3} + \frac{3}{5} \cdot \frac{1}{3} + \frac{1}{5} \cdot \frac{1}{3}$$

$$= \frac{2}{15}$$

$$\frac{2}{10} + \frac{3}{12} + \frac{1}{15}$$

$$\rightarrow \frac{2}{18}$$

$$\frac{24 + 45 + 12}{180}$$

$$= \frac{21}{180}$$

$$= \frac{2x + x_{12}}{27} = \frac{28}{9}$$

$$\frac{28}{15} / \frac{81}{180} = \frac{2x \times 180 + 2}{27}$$

$$\frac{81 \times 15}{27}$$

$$= \frac{8}{27}$$

Name based Classification

↳ Classification

P(A), Name