

Subject:

Date:

10 9

Machine Learning

Phase 0: Prerequisites

Mathematics for understanding Models

Debug effectively, Train as a neural
scientist

Scikit Learn

Mean: The average of a set of numbers

Median: The middle value in a sorted
set of numbers.

Mode: The most frequent value in a
set of numbers.

Percentages, Ratios, Basic algebra,

Plotting x vs y, Probability, vectors,

Standard deviation:

Subject : _____

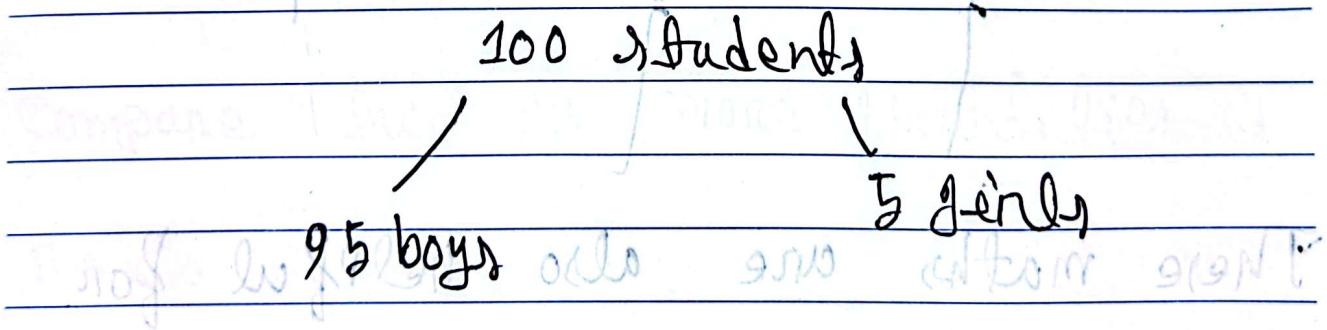
Date : _____

$$\text{Mean} = (\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_n) / n$$

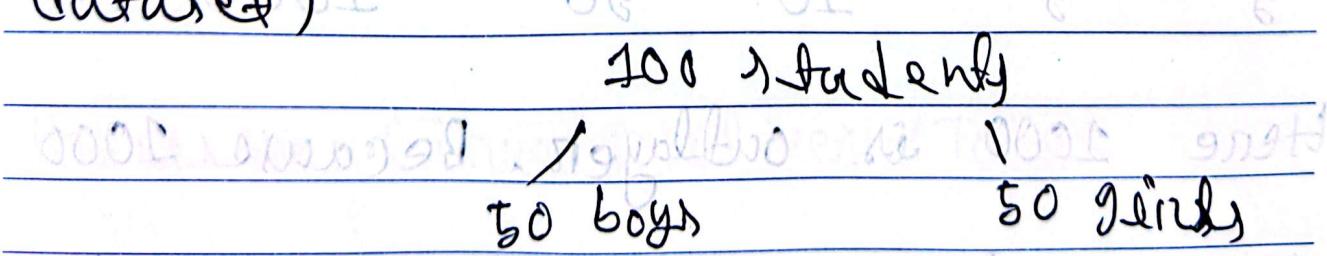
median = middle value (if odd count),
if even (middle 2 average)

mode = most repeated value.

If our dataset is unbalanced, our model can give us wrong prediction.



(This is an example of unbalanced dataset)



(A balanced dataset).

Compare actual value vs Predicted

Value (3000) bbo & 1000 - 1000 = 2000

Students	Actual Mark	Predicted Mark	Error
A	60	80	20
B	80	85	- 5

These marks are also helpful for

analyzing error (Ansible 0.0)

5 5 10 50 100.0 (Ansible 0.0)

Here 1000 is outlier. Because 1000 is way too much higher compared to other values.

(Ansible 0.0)

Subject: _____

Date: _____

Too much low values can be also outliers. In this example it can be -100.

Cleaning outliers helps the model to work efficiently.

Percentages and Ratio: Expressing numbers as parts of 100 as Percentage.

Compare two or more quantities in ratio.

Accuracy = (Correct prediction / Total pred) $\times 100$

Precision = (True positives / Total predictive positives)

Recall = True positives / Total actual positives

Subject: _____

Date: 12 9

Basic algebra

$$y = mx + b$$

Ex: 500 daily income

1000 joining bonus (Fixed)

0 - 1000 1 - 1500 2 - 2000 3 - 2500

y = Final salary predicted, (output / result)

x = days worked, (input)

m = 500 (day, rate of change), \rightarrow slope

b = 1000 bonus, (y-intercept point)

$$y = 500x + 1000$$

u days

$$y = 500 \times 4 + 1000 = 3000$$

Answer (3000) written out → 3000

subject: _____

Date: ___-___-___

Drafting x vs y

Hours of studies

Marks

1

40

2

50

3

60

4

75

 $\frac{5}{5}$

90

x = independent variable

y = dependent variable

Pattern here = More hours = More Marks

We have to find a pattern

Vectors

P. 30 of book

\mathbf{g}_3	R_1	R_2	R_3	possible to represent
\mathbf{g}_2	0.1			E
\mathbf{g}_1	0.2			S
[2, 4]	2F			P
	0.5			T

Vector means a list of numbers

that describes details about a thing

H. w

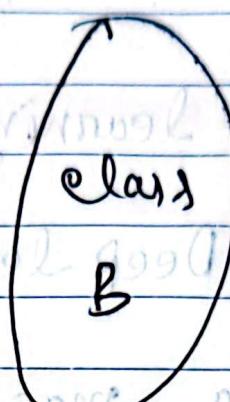
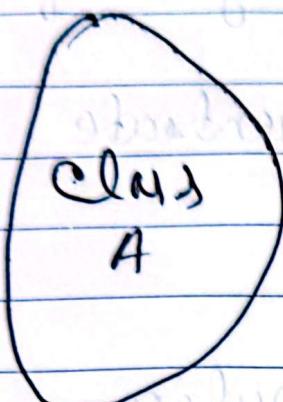
(160 / 60) = height and weight

A single data can be stored as

Vector

A whole dataset = Bunch of vectors

Standard Deviation:



70, 72, 68

30, 70, 100

Close data

Spread data

Standard deviation measures whether

data is closed to each other or

spread.

To analyze dataset we use standard deviation. consistency

data prep, data normalization and

Cleaning outliers, We use SD.

Subject:

Date:

Phase ① H-Foundation of ML

What is ML? Different types of ML, Scikit-Learning, Differentiate AI, ML and Deep Learning.

ML = When a machine/computer

learns from past data, pattern to make prediction or decision.

Machine learns from historical/past data.

AI is layer 1 or outer layer.

AI is a kind of technique that helps a machine to act smart.

Deep learning: Deep learning try to mimic a human brain using neuron.

Self driving is an example of DL.

AI → Entire school - with rules

ML → A class → without rules

DL → A robot with a brain.

Scikit learning helps to predict.

Types of ML

i) Supervised Learning: Ideal for tasks with

Labeled data where the desired

output is known.

ii) Unsupervised Learning: Suitable for

identifying patterns in unlabelled

Data without explicit guidance.

Subject: _____

Date:

(iii) Reinforcement Learning): Best for

Training agents to make decisions
through interaction with an environment.

Supervised learning: Price prediction

Number, Age prediction

Tools: Scikit learning: Linear and Logistic regression, KNN, Decision Tree

I/O + O/P

Unsupervised learning: Clustering based on pattern. Only input. Machine will

try to find pattern/group of objects.

K-means, PCA, DBSCAN.

We care; Marketing, Fraud detection

Subject: _____

Date:

Reinforcement Learning: Machine is learning

learning over the time. Try lennon.

Ex: Games, Tesla, Robotics.

Action + Reward.

Algo: Dynamic Programming, Monte Carlo,

Temporal Difference.

Core ML concepts

Input (x), Output (y) x is always

capital, y is small.

Model: A formula that helps to

Predict.

Training: Understand the pattern.

Testing: Giving new data for expected

Prediction.

Prediction: Final answer.

HEARTS

Subject: _____

Date:

Scikit Learn gives us ready made ML models. It is a python library.

i) Mat - wodleg. lib library

Phase 2 : Data Pre Processing

(Preparing Data)

Raw Data



Handle missing data



Encoding Categorical values



Feature Scaling



Select Data



Prepared Data

Process of task will proceed in below

Naive Bayes

Decision Tree

Subject: .

Date:

Handle missing data; using Pandas.

Encoding Categorical Values & Categories to numbers. Text - Numerical data.

Feature scaling; Standard scalar or

min, max scalers.

Split data for achieving fairness.

Methods for Cleaning data (Missing Value)

i) Drop Na

ii) Fill Na (Value)

iii) isnull. sum

(aabb - bb) freq

count like this new avoid aabb

Subject: _____

Date: _____

A simple program by Pandas

```
import Pandas as pd
```

```
data = pd.read_csv('titanic.csv')  
# some data
```

```
df = pd.DataFrame(data)
```

```
print("Original Dataframe is", df)
```

```
print(df.isnull().sum())
```

It will provide us how many data

is missing.

```
df_drop = df.dropna()
```

```
print(df_drop)
```

drop rows with any null values.

Subject : _____

Date : - -

For falling;

```
df['Age'].fillna(df['Age'].mean(), inplace = True)
```

How will fill missing value

Always try to fell data using mean

on median when working with Numerical

data: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GM12878> (1M) 2019

For categorized data, we made

Do not drop data blindly, try to fill.

Encoding Categorical Values

i) Label encoding.

ii) One-hot encoding.

Encoding values

i) Non numerical values (strings, enums)

into numeric features the model can use

Most ML algorithms expect numeric

continuous inputs. Raw string doesn't

work.

Label encoding

gender

encoded

Male

0

Female

1

Replace each category with a number.

Subject: _____

Date: - -

apply: Gender Column

Let's apply Label Encoding:

```
from sklearn.preprocessing import LabelEncoder
```

```
import pandas as pd
```

```
df = pd.read_csv('sample.csv')
```

```
df_Label = df.copy() # making a copy
```

```
le = LabelEncoder()
```

```
df_Label = LabelEncoder()
```

```
df_Label['Gender_Encoded'] = le.fit_transform
```

```
(df_Label['Gender'])
```

```
Print('Label Encoded')
```

```
Print(df_Label[['Name', 'Gender', 'Gender  
-Encoded']])
```

Subject: _____

Date: _____

One-hot encoding turns a categorical feature into a set of binary indicator features.

- i) Original feature City with values
[Dhaka, Noakhali, Comilla, Tangail]
- ii) One-hot features: City = Dhaka, City = Noakhali, City = Tangail
- iii) A row with City = Noakhali becomes
[0, 1, 0]
 $\text{df_encoded} = \text{pd.get_dummies(df_label, columns=[\text{'City'}])}$
`print('In One-Hot Encoded Data (City)')`
`print(df_encoded)`

Feature Scaling: It is the process of

transforming numerical features so they share a similar range or distribution.

Min-Max scaling (Normalization): $x' = \frac{x - \text{min}}{\text{max} - \text{min}}$

/ ($\text{max} - \text{min}$), maps to $[0, 1]$.

Splitting data into two parts

i) One part will use to train model

ii) Another part for testing

import pandas as pd

from sklearn.preprocessing import

StandardScaler, MinMaxScaler

from sklearn.model_selection import

train_test_split

Subject: Machine LearningDate: 12-12-2022

$$\text{Standard Scaler}, z = \frac{x - \mu}{\sigma}$$

x = Actual Value, μ = Mean

σ = Standard deviation of a column.

Column.

$$\text{MinMaxScaler} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

x = Actual Value, x_{\max} = Max Value

of that column, x_{\min} = Min Value

of that column.

Train-Test-Splitting to divide the

dataset into training, testing and

testing set. It is used to evaluate

accuracy of machine learning model.

212A - 102 - 1001

Phase 3: Supervised ML

Supervised means we will train model with output.

- i) Predict numbers ii) Categories Predict

- Collect label data → Train Model → Make Predictions → Evaluate performance → Refine Model.

ML Algo and Training Process

Regression

Linear Regression

Classification

Logistic

KNN

Decision Tree

Model Training

• fit()

• predict()

Subject: _____

Date:

We use regression when we have to

Predict numbers. Ex: Marks, Price

We use classification when we have

To predict a category. Ex: Mail

Spam.

Linear Regression:

1- Finds a pattern

42000 person hrs off IM

2- Straight line

Regression

3- Line.

Linear Regression is a ML model which

is used for predict numbers using

a straight line relationship between

input x and output y .

Q1 Q2 Q3 Q4

Classification is about predicting a discrete label (Category), not a continuous number like Linear Regression.

Ex: Pass vs Fail, Spam vs not Spam.

Logistic Regression is a ML model used for predicting a binary (0, 1 / Yes, No) output.

It is a classification model unlike linear regression below.

from sklearn.linear_model import LogisticRegression

$X = [[1], [2], [3], [4], [5]]$ # hours studied

$y = [0, 0, 1, 1, 1]$ # 0 is fail, 1 is pass

model = LogisticRegression()

model.fit(X, y)

hours = float(input("Enter Hours:"))

result = model.predict([1 hours]) [0]

if result == 1:

print("You are likely to pass")

else:

print("You may fail")

KNN - K nearest neighbours

KNN is a simple, instance based

supervised ML algorithm used for

Classification

i) Choose K: Select the number of

neighbours (k) e.g. k = 3 (an odd number)

ii) Calculate distance: For a given

query point, calculate the distance

to all points in the training

iii) Find Neighbors: Identify the k points

closest to the query point.

iv) Vote: The query point is assigned to the class most common among its k nearest neighbors.

Taking decision based on nearest data point.

Suppose, in your friend group, 7 friends like Football. So, there is a possibility that 8th friend can start liking Football.

KNN is not suitable for big dataset.

Subject: Fruit Example

Date:

Weight(gm)	Size(cm)	Fruit
180	7	Apple
200	7.5	Apple
250	8	Orange
300	8.5	Orange
330	9	Orange
360	9.5	Orange

from `sklearn.neighbors import KNN`

`KNeighborsClassifier`

`X = [] # some input`

`y = [] # output based on input - Model`

`model = KNeighborsClassifier(n_neighbors=3)`

`model.fit(X, y)`

`# Take Value of X from user`

`prediction = model.predict([[X]])[0]`

`# Print answer`

Decision Tree

A decision tree is a supervised ML

algorithm used for classification and

regression tasks. It works like by splitting

data into branches based on feature

values, creating a tree like model of

decisions.

[Income ?]

High

Yes

Not high

[Age < 40 ?]

Yes

No

No

NO

age Income buys a sports car

age	Income	buys a sports car
25	High	Yes
40	Low	NO
30	High	Yes
50	Medium	NO
35	Low	NO

Overfitting :- Overfitting is when a

ML model learns the training data

too closely and loses the ability
to generalize to new data.

(gives poor performance on unseen

examples)

Underfitting is when a model

cannot learn the relationship

between input and output well enough,

leading to low accuracy.

Perfect fit is perfect and accurate

for model.

Subject:

Date: 17 9 2015

Phase 4: Model Evaluation and METRICS

SKLearn Metrics: Master the use of

'sklearn.metrics' for comprehensive model evaluation.

Regression Metrics Learn to evaluate

regression models MAE, MSE and RMSE

Confusion Matrix:

Classification Metrics:

Accuracy: (%) of total right prediction

$$\text{Acc} = \frac{\text{Correct_pred}}{\text{Total_pred}} \times 100$$

Precision: It measures how many of the

items predicted as positives are actually

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

TP = True Positives, FP = False positives.

Ex: If a spam filter detects 10 emails as spams and 8 of them are actual spams, the precision is $8/10 = 0.8$ (80%).

Recall: Recall measures how many of the actual positive items were correctly identified by the model.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN}) \quad \text{FN} = \text{False Negative}$$

Ex: If there are 12 spam emails in total, and the filter correctly detects 8 of them, recall is $8/12 = 0.67$ (67%).

F1 score: It is a metric used in ML to evaluate the performance of a classification model by

Subject: _____

Date: _____

Combining both precision and recall

into a single number.

$$F1 \text{ score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Range: 0 - 1. It is useful when we need to balance False positives and False negatives and when dataset is imbalanced.

Confusion Matrix: A confusion matrix

is a table used in ML to

evaluate the performance of a

classification algorithm. It shows

how many predictions were correct

and incorrect, and breaks them down by

each class.

Actual/Predicted	Spam	Not spam
spam	3	1
Not spam	2	4

MAP, MSE, RMSE

Student	Real Mark	Model Mark	Mse
A	90	85	5
B	60	70	10
C	80	70	10
D	100	95	5

i) MAE (Mean Absolute Error):

The average of the absolute differences between actual and predicted values.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

y_i is actual value,
 \hat{y}_i is the predicted value.

1) Take the mistake difference

2) Remove the common signs

3) add, 4) divide

ii) MSE (Mean Squared Error): The average

of the squared differences between actual and predicted values.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Useful when large errors are particularly undesirable.

Subject : _____

Date :

1) Square mistakes

2) Add

3) Divide by N

RMSE (Root Mean Squared Error) : Just

The square root of MSE

useful - Penalizes large errors more than

MAE. These 3 are most important.

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

A

Subject: _____

Date:

3 Step Formula:

1 - X and Y Numeric?

Scatter Plot Relationships

2 - In one column a category

Bar Plot / Count Plot

3 - Want to see a distribution?

Histogram / KDE Plot / Box Plot

Problem Statement; Final Exam Score