**U-1** Describe ML process flow with oppropriate diagram.
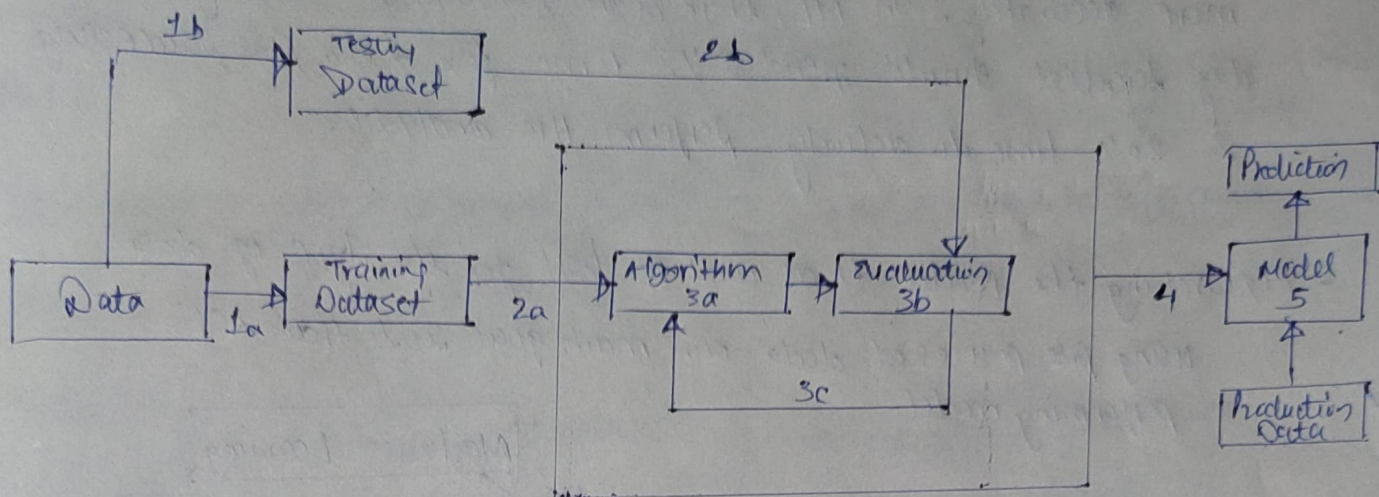


fig(1):- overview of ML process flow.

ML workflow in 3 stages:-

1. Gathering data

2. Data Pre-processing

3. Researching the model that will be best for the type of data.

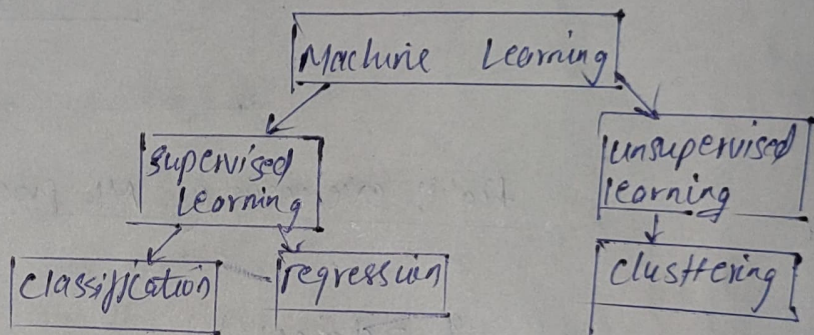4. Training & testing the model.

5. Evaluation.

- **Gathering of data :-**
  - It depends on the type of project we desire to make, if we want to make an ML project that uses real-time data, then we can build an IOT system that using different Sensor data. The datu set can be collected from various sources such as file, database, Sensor and many other.

19630Y02g

- Data pre - processing :-
            It helps in building machine learning models
  more accurately. In ML, there is an 80/20 rule. Every
  day scientist should spend 80% time for data pure processing
  & 20% time to actually perform the analysis.

- Researching the model that will be best for the type of data.
      using pre-processed data our main goal is to train best
      performing model.



- Training & testing the model on data :-
            we split the model in 3 section which are training data,
  validation data and testing data.

- Evaluation :- Its an integral part of model development process, helps
            to find the best model that represent our data &
      how well the chosen model will work in future.

Page-301

PATEZ
301029

Assignment - 01

Distributed system.

Page - 03

es-02: Differtiate → A F E S T D D D—

(a) Analysis hand caye Storg Tool Databa Data lite
method cuse Sea node nottru

| | Structured Data | Unstructured Data |
|---|---|---|
| formatts | Several formats | A huge variety of formats |
| Data model | Pre-defined/not flexible | Not pre-defined / flexible |
| Storages | Data warehouses | Data lakes |
| Databases | SQL Relational Databases | NOSQL Non-relational databases |
| Ease of Search | Easy to search | Difficult to Search |
| Data nature | Quantitive | Qualitative |
| Analysis method | • Classification <br> • Regression <br> • Data clustering | • Data stacking <br> • Data mining |
| Tools and technologies | • RDBMS <br> • CRM <br> • OLAP <br> • OLTP | • NOSQL DBMS <br> • AI-driven tools <br> • Data storage architecture <br> • Data visulization tools. |

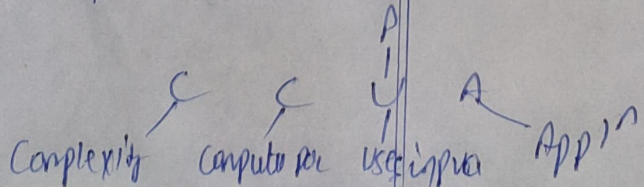| (b) | Online Machine learning | offline Machine learning |
|---|---|---|
| Complexity | More complex | Tess complex |
| Computatioinal power | More Computaional power is required | fewer computaional power is required |

C C U A
Complexity compute per useinpva App^n

| | | |
|---|---|---|
| - Use in Production | Harder to implement & control because the production model changes in real-time according to its data feed. | Easier to implement because offline learning provides engineers with more time to perfect the model before deployment. |
| - Application | used where new data patterns are constantly required (eg. weather prediction tools.) | used where data application data patterns remains constant and don't have sudden concept drifts (eg. image classification) |

## Ques-03  Define bias-variance trade-off?

The goal of any supervised machine learning algorithm is to achieve low bias and low variance. In return the algorithm should achieve good prediction performance.

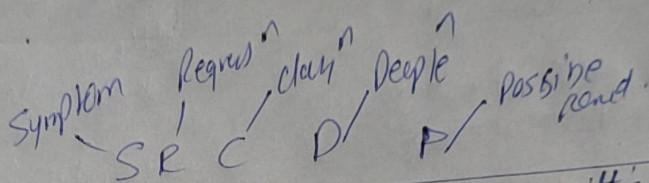The parameterization of ML is often a battle to balance out bias & variance

example of configuring the bias-variance trade-off for specific algorithm.

- The K-nearest algorithm has low bias & high variance. but the trade-off can be changed by increasing the value of K which increases the no. of neighbours that contribute the prediction and in return, increases the bias of the model.

There is no escaping the relationship between bias and variance in ML.

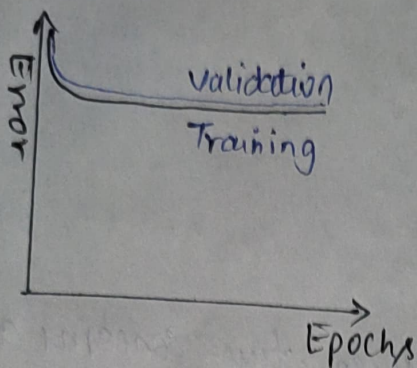$$\boxed{Bias \propto \frac{1}{variance}}$$

There is a trade-off at play between these two concerns and the algorithms you choose and the way you choose to configure them are finding different balances this trade-off for your problem.
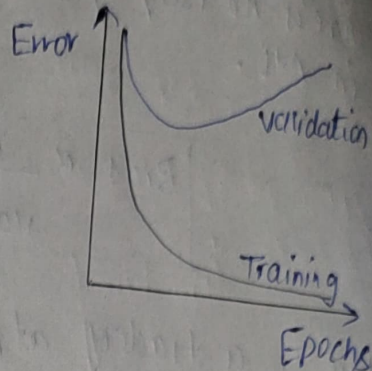
Symptom  Regres⌃  clan⌃  Deeple      possible
         S R      C      D    P       Rand.

**Ques-04 Differentiate:**

|  | Underfitting | Overfitting |
|---|---|---|
| Symptoms | • high training error<br>• training error close to test error<br>• high bias | • very low training error<br>• training error much lower than test error<br>• high variance |
| Regression illustration |  |  |
| classification illustration |  |  |

| Deep. Learning illustration |   Error / Epochs with validation and Training curves |   Error / Epochs with validation and Training curves |
|---|---|---|
| Possible remedies | • Complexity model<br>• Add more feature<br>• Train longer | • Perform regulonization<br>• Get more data |

_End._

(10 marks) extra topics.
→ cross validation

↳ Cross-Validation :

- It's a technique to in which we train our model using the subset of the data-set and then evaluate using the complemantary subset of data.

- There are three steps involved in cross validation

1. Reserve some portion of sample data-set.

2. Using the rest data-rest set to train the model.

3. test the model using the reserve portion of data set.

- Why we use Cross Validation
    - to test stability of model
    - we can't just fit our model on the training dataset.
    - we need a particular sample of dataset, which is not
      part of training dataset

Method used for Cross Validation
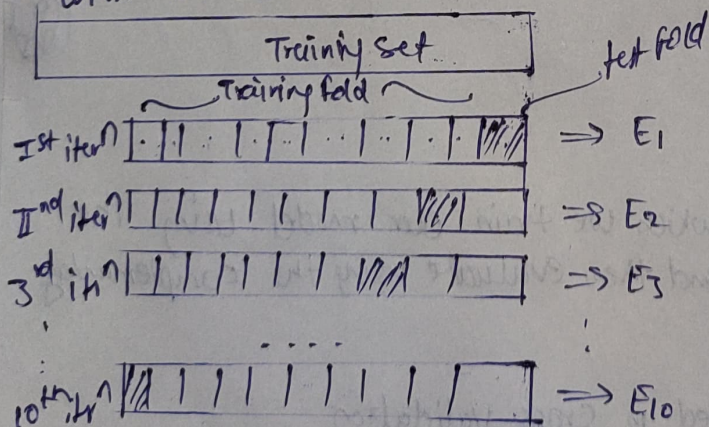
① K- fold cross- validation
    - divides the input dataset in K-groups of sample of equal
      sizes. , Those sample are called folds.

Steps for K-Fold
    - split the input dataset into K groups.
    - for each group-
        -: use one group as reserve or test data set
        -: use remaining group as training dataset.
        -: fit the model on the training set &
           evaluate the performance of model usip
           the test set.

K-fold cross valid^n
with K = 10



$$\Rightarrow E_1$$
$$\Rightarrow E_2$$
$$\Rightarrow E_3$$
$$\Rightarrow E_{10}$$

$$E = \frac{1}{10} \sum_{i=1}^{10} E_i$$

② Stratified K-fold Cross validation

→ Best approach to deal w bias & variance

→ works on stratification concept
→ similar to K-fold with some little changes.

→ rearranging the data so that each fold or group is a good representative of complete dataset.