

TECHNO MAIN SALT LAKE

(FORMERLY TECHNO INDIA, SALT LAKE)

Name..... Kewal Rij

Roll No. 13030822018 Stream..... CSE (AIML)

Subject Machine Learning Semester..... 6th

(PCC-AIML601)

Invigilator's Signature  Date.....

(Part-A)

1. The two most Common tasks are:-
In age recognition and Customer-churn Analysis.
2. The purpose of Validation test Set are:-
After analysing the data through trainset and testing ~~set~~
through test set for multiple times to find the final
accuracy Validation Set is needed.
3. (i) Coefficient of the Variable,
(ii) Loss function i.e mean Square error (MSE)
(iii) Gradient decent.
- 4.) AUC value of perfect classifier $\tan \theta \geq 1$.
- 5.) Out of precision, recall, recall is more important for
a spam email detection System.

(Part-B)

6. Train-test-split:-

when we get a new dataset, we divide the dataset initially into train-test-split. i.e. 70-80% of data are kept initially for training and after training the data model, we use 30-20% of dataset for testing. On this basis, we train the model as well as test the model simultaneously, this division of dataset into 70-30 is splitting of data for the purpose of training and testing of data.

Overfitting of data: And Underfitting:

In Overfitting of data, we basically try to fit data more precisely and accurately, we try to cover each and every point possible, which is very strict in nature.

And, Underfitting of data says that, we are less concerned of considering data, we traverse in dataset, and whatever number of data we get is the possible data, But it does not give more accurate models, because sample size of data may decrease.

⑧ General algorithm that are available to minimize Cost function:-

- Naive Bayes algorithm
- Random Forest algorithm
- Decision trees

Cost functions depend on the regression model, we select, Cost function may reduce for some problem and also may increase for some problem,

Basically, Cost function is low for linear regression model and at the same time it is increased for logistic regression model.

In Linear regression, basically we simplify the problem with given datasets, which takes less time than Logistic regression problem, So, it will lower the Cost function and simultaneously it will increase the Cost function for Logistic regression problem.

⑨ Confusion matrix is the essential aspects of training model, we will get Confusion matrix for every model we train, It basically tells about the model in more details, without Confusion matrix we can't be able to tell about the precision, recall, false negative rate, false positive rate of model.

Given: $TN = 82$

$FP = 3$

$FN = 5$

$TP = 10$

$$\text{false negative rate} = \frac{FN}{TN + FN}$$

$$= \frac{5}{82 + 5}$$

$$= \frac{5}{87}$$

$$\boxed{FNR \approx 0.05}$$

$$\text{false positive rate} = \frac{FP}{TP + FP}$$

$$= \frac{3}{10 + 3}$$

$$= \frac{3}{13}$$

$$\boxed{FPR \approx 0.23}$$

$$\text{PRECISION} = \frac{10}{10 + 3 + 3 + 82} = \frac{10}{100} = 0.1$$

$$\text{Recall} = \frac{TP}{TN + TP} = \frac{10}{82 + 10} = \frac{10}{92} = 0.108$$

- ⑦ Bias and Variance of a Machine Learning model.
- Bias are dataset (how much deviated from the actual central tendency).
 - Variance is the (how dataset is or far from the central nature of a given model).

To reduce Bias and Variance, we should avoid:-

- overfitting and underfitting
- increase sample space of datasets
- more frequent train-test-split,

Bias-Variance trade off is basically, first we know how much deviation is there in the model, then we fix the model by knowing how far we are from the central tendency, & this process is called Bias-Variance trade off.