

TECHNO MAIN SALT LAKE

(FORMERLY TECHNO INDIA, SALT LAKE)

Name.....Debasmita Lakshmi.....

Roll No.13030822008.....Stream.....CSE-AIMI.....

Subject ..Machine Learning Applications (PCA/ML-601).....Semester.....6.....

Invigilator's Signature[Signature].....Date.....

Part A

1. Two most common supervised learning tasks are - Classification & Regression.
2. The purpose of Validation Set is to reduce the problem of overfitting and underfitting and by improving generalization it reduce generalization error.
3. There are only one model parameter (θ) and bias term (θ_0) with a single feature variable.
4. AUC value of perfect classifier is 1.
5. Recall is more important for Spam detection System.

Part B

9. Confusion Matrix :- Confusion Matrix is a comprehensive model evaluation or performance measure matrix which is visualized as a tabular form and counts the number of actual outputs versus predicted output.

Actual	P	TP	FN
	N	FP	TN
		P	N
		Predicted	

TP - True Positive is model is predicted positive classes correctly.

FP - False Positive is model predicted positive but actual is negative (Type I error).

FN - False negative is model predicted negative but actual is positive. (Type II error).

TN - True negative is model predicted negative classes correctly.

$TP + FP = \text{Total +ve predictions}$

$FN + TN = \text{Total -ve predictions}$

$TP + FN = \text{Total actual positives}$

$FP + TN = \text{Total actual negatives}$

$TP + FP + FN + TN = \text{Total predictions}$

□ Importance :-

① Confusion matrix provides how the model is performing in unseen data, helps to evaluate its generalization capability.

② It helps to calculate other performance measure metrics

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

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

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

ROC Curve = Plot of TPR and FPR at different threshold.

PR Curve = Curve of Precision & Recall.

In given classification problem, $TN = 82$

$FP = 3$,

$FN = 5$,

$TP = 10$.

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 3} = \frac{10}{13}$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 5} = \frac{10}{15} = \frac{2}{3}$$

$$\text{False negative rate} = \frac{FN}{FN + TP} = \frac{5}{5 + 10} = \frac{5}{15} = \frac{1}{3}$$

$$\text{False Positive rate} = \frac{FP}{FP + TN} = \frac{3}{3 + 82} = \frac{3}{85}$$

6.) Train-Test Split :- Train-Test Split is defined by dividing the whole dataset into two portions.

① Train Set :- For training the data to the model. (higher in ratio)

② Test Set :- For testing the model predictions and performance in unseen data. (smaller in ratio). (70:30)

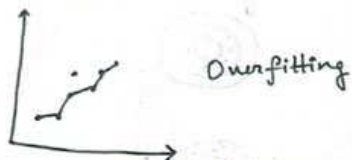
It is generally used in Supervised learning algorithms where we train the model by labelled data, and test the model by providing the unlabelled data, and compare the actual and predicted output.

• from sklearn, model selection import train-test-split.

Overfitting :- Overfitting occurs when model performs good in training Set and poor in testing Set.

• In overfitting model tries to memorize the data rather than generalizing.

- In overfitting, model has high variance ~~and low bias~~.
- Model tries to be too complex that it tries to predict all points in training but due to high variance it fails to predict unseen data.



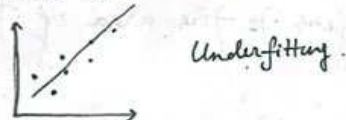
Underfitting :- Underfitting occurs when model performs poor in training and testing data.

- In underfitting model does not able to extract the true patterns of data and have its own path (far from actual output).

- Here, model has high bias and ~~low variance~~.

- ~~Class imbalance leads something high bias and~~

- Model is too simple or linear and can't capture non-linear relations.



□ Ways to Prevent.

- Use ensemble learning techniques (Bagging & Boosting)
- Use Regularization methods (Ridge, Lasso, Elastic Net).
- Standardization and Normalization.
- For Overfitting, do feature selection and for underfitting use kernel trick (SVM) and model with non-linear channels.

7.) Bias :- Bias refers to the difference between actual and predicted output. It defines how close the generalizing hypothesis is of true hypothesis. (Low-Bias :- Near to actual value). (High Bias :- Far from actual, happens due to class imbalance, sparse dataset).

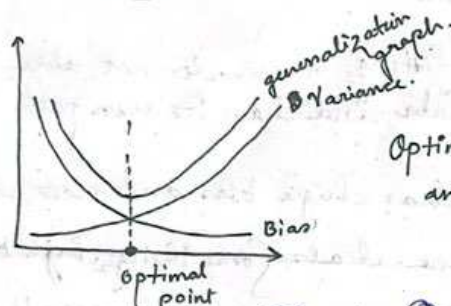
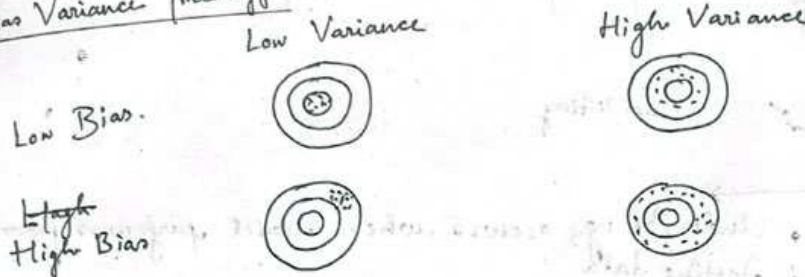
Variance :- Variance refers how model predictions change by changing the features/feature vectors. It refers how scattered the predicted values are (Low Variance → data (predicted) are not scattered (close to each other)). High Variance → predicted points are scattered).

□ Reducing Way to Reduce Them. :-

- To reduce bias, we need to balance the dataset by doing oversampling of minor class and under sampling of major class.
- do ensemble methods like ~~boosting~~ random forest, bagging.

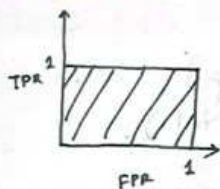
To reduce Variance we need to do Boosting techniques. We do feature selection (Regularization techniques).

Bias Variance Tradeoff

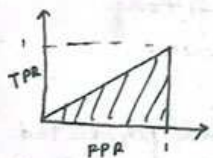


Optimal point is where bias and Variance is minimum. (Low bias and low Variance)

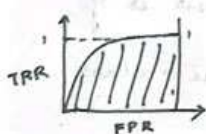
10. ROC Curve: ROC (Receiver Operating Characteristics) refers to the graph of TPR vs FPR at different thresholds. AUC: AUC (Area under Curve) refers to the area of ROC Curve.



Perfect Classifier ($AUC = 1$).



Random Classifier ($AUC = 0.5$)

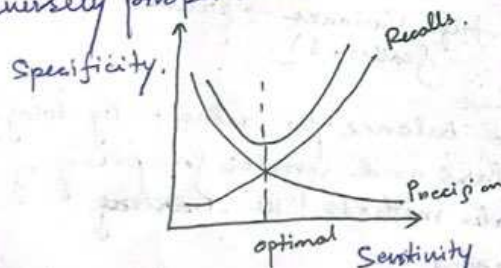


Random Classifier ($0.5 \leq AUC \leq 1$).

Precision Recall Tradeoff It defines by the Curve of Precision and Recall. It plots precision and recalls. Precision (Specificity) & Recall (Sensitivity) are inversely proportions.

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

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



It defines the change of FP (Precision) and FN (Recall).

Optimal points defines where precision and Recall is minimum