

Name - Harshita Singh  
Sec - K

Univ Rollno - 2315000950  
Class Rollno - (31)

Course / Branch - B.Tech (CSE)

### ASSIGNMENT - 1 (Machine Learning)

Ans 1 - Used-car dataset → Regression or classification?

- Task given : Predict whether a car will sell above its listing price → this is a binary classification problem (Yes = Above price, No = Not above price).
- Justification :
  - Classification : The output is discrete (0/1). Model predicts probabilities (e.g., 0.8 chance the car sells above price).
  - Regression alternative : If instead we tried to predict the actual selling price as a continuous variable, then we'd use regression.
- Outputs & Metrics :
  - Classification : output = probability or class label.  
Metrics : Accuracy, Precision, Recall, F1-score, ROC-AUC.
  - Regression : Output = numerical price prediction.  
Metrics : RMSE, MAE,  $R^2$ .

Ans 2 - Handling Missing BMI & Glucose Values

- Options :
  1. Mean imputation : Replace missing values with the mean. Works if data is normally distributed without outliers.
  2. Median imputation : Replace with the median. Better when distributions are skewed or when extreme outliers exist (common for BMI and glucose).
  3. Dropping records : Appropriate only if :
    - Missing values are rare,
    - The records are missing completely at random (MCAR).
- Clinical relevance :
  - BMI and glucose often have skewed distributions, with important clinical thresholds. Using median imputation preserves robustness.

- Mean imputation could distort results if a few patients have extreme values.
- Dropping should be avoided if missingness is systematic.

Ans 3 :- High Error / Underfitting (high bias) - Model is performing poorly even on training.

Steps to improve :-

- \* Feature Engineering.
- \* Use a Stronger algo.
- \* Increase model capacity / complexity, tune, hyper parameters.
- \* Check data quality / labels, add more information feature.

Ans 4 :- One-hot Encoding Example

Input:

| <u>Name</u> | <u>Department</u> |
|-------------|-------------------|
| Alice       | HR                |
| Bob         | Engineering       |
| Charlie     | HR                |
| Dana        | Sales             |

After one-hot encoding Department :

| <u>Name</u> | <u>Dept_Engineering</u> | <u>Dept_HR</u> | <u>Dept_Sales</u> |
|-------------|-------------------------|----------------|-------------------|
| Alice       | 0                       | 1              | 0                 |
| Bob         | 1                       | 0              | 0                 |
| Charlie     | 0                       | 1              | 0                 |
| Dana        | 0                       | 0              | 1                 |



Ans 5 (distances) :

$$\begin{aligned} d(E, A) &= \sqrt{(6-8)^2 + (7-6)^2} = \sqrt{5} = 2.23 && \text{pass} \\ d(E, B) &= \sqrt{(6-5)^2 + (7-4)^2} = \sqrt{10} = 3.16 && \text{fail} \\ d(E, C) &= \sqrt{(6-7)^2 + (7-5)^2} = \sqrt{5} = 2.23 && \text{pass} \\ d(E, D) &= \sqrt{(6-3)^2 + (7-2)^2} = \sqrt{34} = 5.8 && \text{fail} \end{aligned}$$

$k = 3$  nearest : A (pass), C (pass), B (fail) .

prediction : pass (2 vs 1)

Ans 6 PCA on points  $(1,0), (0,1), (2,2), (4,4)$

| $x$ | $y$ | $(x-\bar{x})$ | $(y-\bar{y})$ | $(x-\bar{x})(y-\bar{y})$ | $(x-\bar{x})^2$ | $(y-\bar{y})^2$ |
|-----|-----|---------------|---------------|--------------------------|-----------------|-----------------|
| 1   | 0   | -0.75         | -1.75         | 1.31                     | 0.56            | 3.06            |
| 0   | 1   | -1.75         | -0.75         | 1.31                     | 3.06            | 0.56            |
| 2   | 2   | 0.25          | 0.25          | 0.06                     | 0.06            | 0.06            |
| 4   | 4   | 2.25          | 2.25          | 5.06                     | 5.06            | 5.06            |
|     |     |               |               | 7.74                     | 8.74            | 8.74            |

# mean  $\bar{x} = 1.75$   
 $\bar{y} = 1.75$

# Covariance Matrix :-

$$\begin{bmatrix} \text{Cov}(x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Cov}(y) \end{bmatrix}$$

$$\text{Cov}(x) = \frac{1}{3} \times 8.74 = 2.91$$

$$\text{Cov}(x, y) = \frac{1}{3} \times 7.74 = 2.58 = \text{Cov}(y, x)$$

$$\text{Cov}(y) = \frac{1}{3} \times 8.74 = 2.91$$

$$\text{Covariance Matrix} = \begin{bmatrix} 2.91 & 2.58 \\ 2.58 & 2.91 \end{bmatrix}$$

### # Eigen Values

For a symmetric matrix of the form  $\begin{bmatrix} a & c \\ c & a \end{bmatrix}$ , eigen values are  $a+c$  and  $a-c$ .

here,  $a = 2.91$ ,  $c = 2.58$

$$\bullet \lambda_1 = a + c = 2.91 + 2.58 = 5.5$$

$$\bullet \lambda_2 = a - c = 2.91 - 2.58 = 0.33 \ (\approx 1/3)$$

Interpretation / dimensionality reduction :

Total variance  $= \lambda_1 + \lambda_2 = 5.83$ . The first principal component ( $\lambda_1 = 5.5$ ) captures  $\approx 5.5 / 5.83 = 94.3\%$  of the variance

PCA reduces dimensionality by projecting data onto the top  $k$  principal components (here  $k=1$ ) which retain the largest possible fraction of variance - projecting onto the first eigenvector keeps  $\approx 94\%$  of the variability while reducing  $2D \rightarrow 1D$ .

### Ans 7. logistic regression (bank churn)

- logistic function maps linear score  $z$  into probability :

$$p = \frac{1}{(1 + e^{-z})}$$

- Positive coefficient for support calls  $\rightarrow$  more calls = higher churn probability (odds  $\uparrow$  by  $e^{\text{coef}}$ ).



Ans 8 Confusion matrix metrics (TP = 50, FN = 20, FP = 15, TN = 85)

- Accuracy = 79.4 %
- Precision = 76.9 %
- Recall = 71.4 %
- Specificity = 85 %
- F1 = 74.1 %

In critical health monitor → Recall more important. {

Ans 9 Smartphone resale price

- Task = predict exact resale price → Regression.
- output : numeric price ; Metrics : RMSE, MAE,  $R^2$ .
- If predicting price ranges → Classification.
- Output : class labels ; Metrics : Accuracy, Precision, Recall, F1.

Ans 10 . Missing creatinine & hemoglobin

- Mean imputation : Symmetric, no outliers.
- Median imputation : Skewed data / outliers.
- Drop records : very few missing, MCAR.

Ans 11 . Sentiment classifier (54% train, 56% val)

- Error type : underfitting.
- Fix : richer features (TF-IDF, n-grams, embeddings)
- stronger models (SVM, boosting, transformers), hyperparameter tuning, more / better data.

Ans 12.

| Item     | Stationary | Office Supplies | Art |
|----------|------------|-----------------|-----|
| Pen      | 1          | 0               | 0   |
| Notebook | 0          | 1               | 0   |
| Eraser   | 1          | 6               | 0   |
| Marker   | 0          | 0               | 1   |

Ans 13. KNN ( $k=3$ )  
 distance  $\rightarrow P_1 = \sqrt{10}$ ,  $P_2 \rightarrow \sqrt{10}$ ,  $P_3 \rightarrow \sqrt{10}$ ,  $P_4 = 5$   
 Nearest 3:  $P_1, P_2, P_3 \rightarrow$  pass  $\rightarrow 2$ , fail  $= 1$   
 $\rightarrow$  predict pass.

Ans 14 Covar. Matrix  $\rightarrow$  Eigen values show variance dir.  
 PCA keeps axis with max eigenvalue.

Ans 15 Sigmoid maps linear sum  $\rightarrow$  Prob.  
 Negative BMI coefficient = higher BMI + complications Prob.

Ans 16 Confusion Matrix (40, 10, 20, 130)  
 Accuracy = 85%.  
 Precision =  $40 / 60 \rightarrow 66.7\%$ .  
 Recall =  $40 / 50 \rightarrow 80\%$ .  
 Specificity =  $130 / 150 \rightarrow 86.71\%$ .  
 F1 Score = 72.7%.

for fraud alerts  $\rightarrow$  Recall more imp.

Ans 17. Classification (above / below median).  
 Regression only if predicting exact rent.  
 Metrics: Reg  $\rightarrow$  RMSE, MAE.  
 Classification  $\rightarrow$  Accuracy / F1.



Ans 18. Mean — Normal distm.  
 Median — Skewed / outliers.  
 Drop — small missing set & not crucial clinically.

Ans 19. Error — Underfitting  
 Fix — Richer features, Stronger, models, tuning.

Ans 20

| <u>Title</u> | <u>Sci-fi</u> | <u>Romance</u> | <u>Drama</u> |
|--------------|---------------|----------------|--------------|
| Inception    | 1             | 0              | 0            |
| Titanic      | 0             | 1              | 0            |
| Joker        | 0             | 0              | 1            |
| Up           | 0             | 0              | 0            |

Ans 21 KNN ( $k=3$ )  
 Distance from  $X_5(6,6)$ :  
 $X_1 = \sqrt{2}$ ,  $X_2 = \sqrt{13}$ ,  $X_3 = \sqrt{10}$ ,  $X_4 = \sqrt{12}$   
 Nearest 3:  $X_1$  (pass),  $X_3$ , (pass)  
 $X_2$  (tail)  $\rightarrow$  predict pass.

Ans 22 Covariance Eigenvalues  $\rightarrow$  One large, One small.

PCA keeps axis with larger eigenvalues  
 $\rightarrow$  max variance retained.

Ans 24. Accuracy  $\rightarrow (45+125)/200 \rightarrow 85\%$   
 Precision  $\rightarrow 45/70 \rightarrow 64.3\%$   
 Recall  $\rightarrow 45/50 \rightarrow 90\%$   
 Specificity  $\rightarrow 125/150 \rightarrow 83.3\%$   
 $F1 \rightarrow 75\%$   
 for fraud System  $\rightarrow$  Recall more