

CREDIT CARD FRAUD DETECTION USING MACHINE LEARNING

TEAM MEMBERS:

- **Anubhav Gupta (102303789)**
- **Chirag Bansal (102303700)**
- **Shrey Rawat (102303669)**
- **Dushyant Singh (102303807)**

ABSTRACT

- DETECTING FRAUDULENT TRANSACTIONS IS CRUCIAL DUE TO RISING GLOBAL CREDIT CARD FRAUD.
- MACHINE LEARNING MODELS PROVIDE FASTER, MORE ADAPTIVE FRAUD DETECTION COMPARED TO TRADITIONAL METHODS.
- THE PROJECT BUILDS A REAL-TIME FRAUD DETECTION SYSTEM USING HISTORICAL TRANSACTION DATA.

INTRODUCTION

- GROWTH OF ONLINE TRANSACTIONS INCREASED THE RISK OF FRAUD.
- TRADITIONAL METHODS (RULES/HEURISTICS) ARE INSUFFICIENT.
- MACHINE LEARNING USES DATA PATTERNS TO DETECT FRAUD DYNAMICALLY.

PROBLEM STATEMENT

FRAUD CAUSES
BILLIONS OF
DOLLARS IN
LOSSES
ANNUALLY.

HARD TO DETECT DUE TO:

- HIGH VOLUME OF TRANSACTIONS
- CONSTANTLY EVOLVING FRAUD PATTERNS

OBJECTIVES

01

- DEVELOP A ROBUST MACHINE LEARNING MODEL FOR FRAUD DETECTION.

02

- MINIMIZE FALSE NEGATIVES (FRAUD TRANSACTIONS MARKED AS NORMAL).

03

- IMPROVE FINANCIAL SECURITY FOR USERS AND INSTITUTIONS.

METHODOLOGY

- DATA COLLECTION
- DATA PREPROCESSING
- MODEL SELECTION
- MODEL EVALUATION

DATA COLLECTION

- We used a **HISTORICAL DATASET CONTAINING CREDIT CARD TRANSACTION RECORDS.**
- **EACH RECORD INCLUDES FEATURES LIKE TOTAL INCOME, AMOUNT CREDIT, AND A TARGET INDICATING FRAUD OR LEGITIMATE.**

DATASET SOURCE - [HTTPS://WWW.KAGGLE.COM/DATASETS/MISHRA5001/CREDIT-CARD/DATA](https://www.kaggle.com/datasets/mishra5001/credit-card/data)

DATA PREPROCESSING

- HANDLING MISSING VALUES
- ENCODING CATEGORICAL FEATURES
- SPLITTING DATA (TRAIN AND TEST SET)

HANDLING MISSING VALUES

- **NULL VALUES ARE CHECKED USING `df.isnull().sum()` COMMAND**
- **NULL RECORDS ARE REMOVED USING `dropna()` COMMAND**

```
#Checking for null values
df.isnull().sum()

✓ 0.0s
```

TARGET	0
NAME_CONTRACT_TYPE	0
CODE_GENDER	0
FLAG_OWN_CAR	0
AMT_INCOME_TOTAL	0
AMT_CREDIT	0
AMT_ANNUITY	12
AMT_GOODS_PRICE	278
DAYS_BIRTH	0
DAYS_EMPLOYED	0
NAME_EDUCATION_TYPE	0
OCCUPATION_TYPE	96391
OBS_30_CNT_SOCIAL_CIRCLE	1021
DEF_30_CNT_SOCIAL_CIRCLE	1021
AMT_REQ_CREDIT_BUREAU_DAY	41519
AMT_REQ_CREDIT_BUREAU_WEEK	41519
AMT_REQ_CREDIT_BUREAU_MON	41519
CNT_FAM_MEMBERS	2
REGION_RATING_CLIENT	0
REG_REGION_NOT_WORK_REGION	0
DAYS_LAST_PHONE_CHANGE	1
dtype: int64	

ENCODING CATEGORICAL FEATURES

- Categorical features were identified by checking the data types of columns. Columns with data type object

```
ctgr_features= df.select_dtypes(include=["object"]).columns  
print(ctgr_features)  
Index(['NAME_CONTRACT_TYPE', 'CODE_GENDER', 'FLAG_OWN_CAR',  
       'NAME_EDUCATION_TYPE', 'OCCUPATION_TYPE'],  
      dtype='object')
```

- `pd.get_dummies()` is used to encode the categorical features into separate numerical/boolean columns

SPLITTING DATA (TRAIN AND TEST SET)

- Dataset is divided into features and target set
- Imbalanced data is handled using SMOTE
- The dataset was split into training(80%) and testing(20%) sets using the `train_test_split()` function .

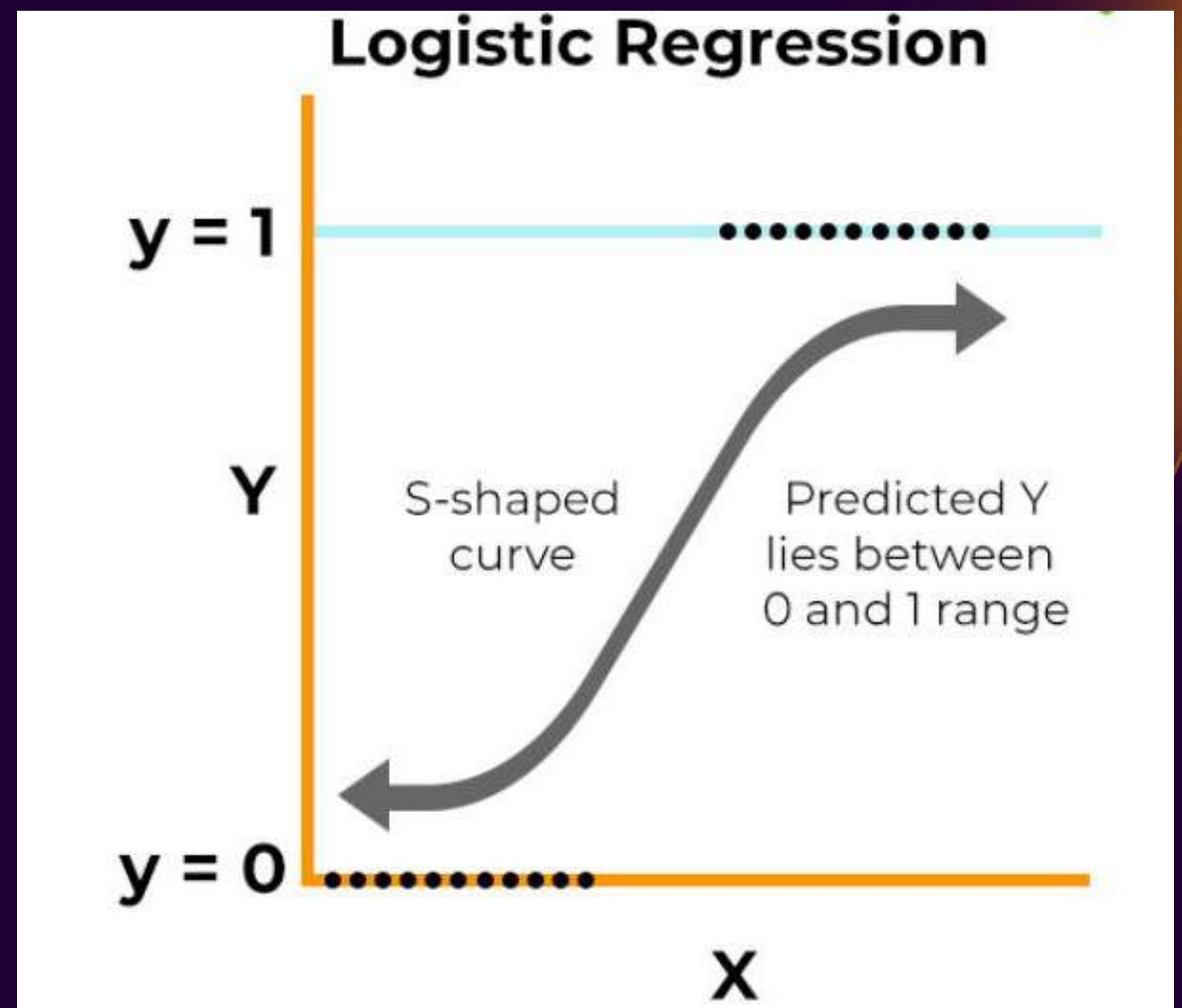
METHODOLOGY

MODEL SELECTION AND TRAINING

- LOGISTIC REGRESSION
- DECISION TREE
- RANDOM FOREST
- K-NEAREST NEIGHBORS
(KNN)

LOGISTIC REGRESSION

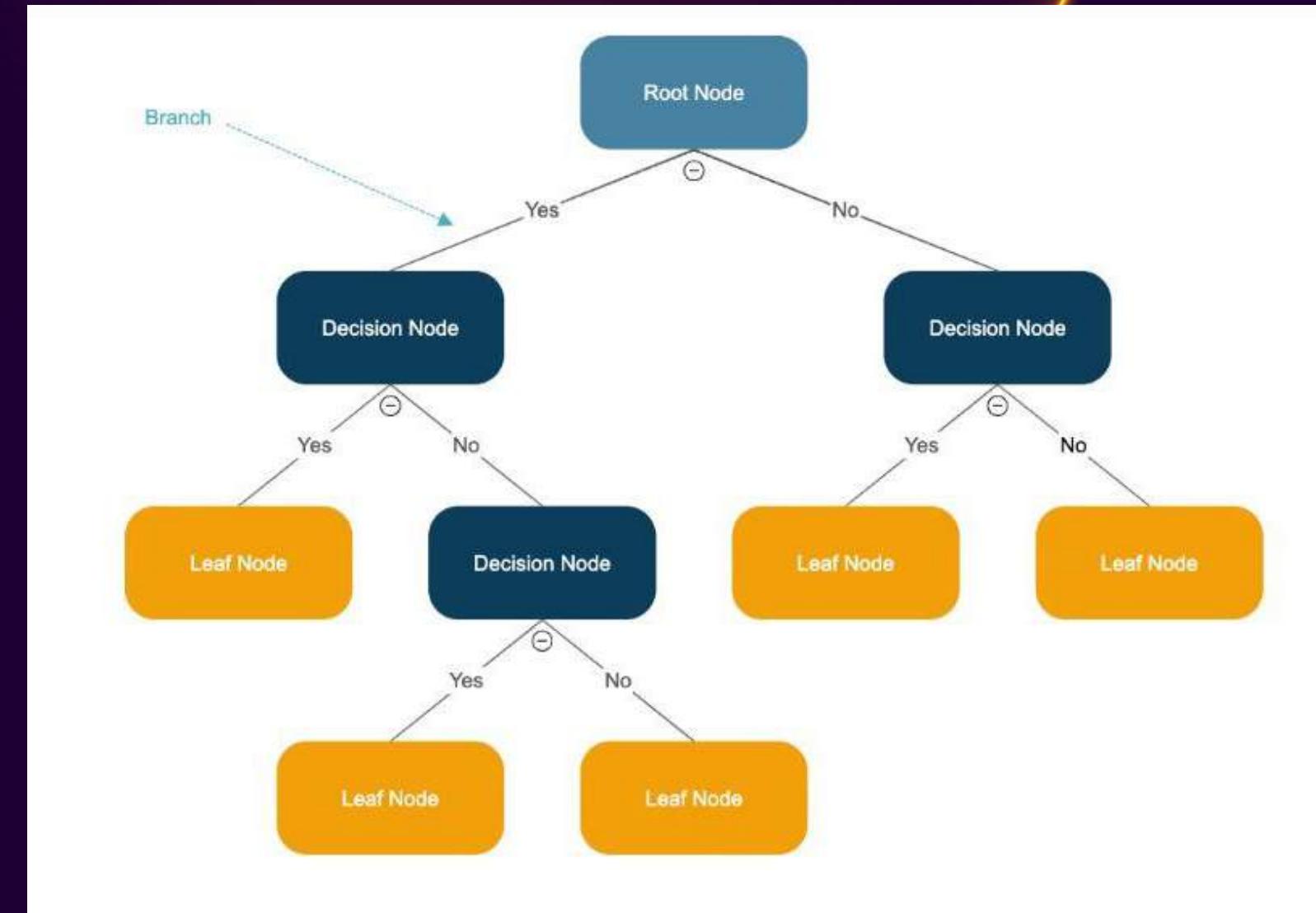
- LOGISTIC REGRESSION WAS USED AS A BASELINE MODEL FOR BINARY CLASSIFICATION.
- IT PREDICTS THE PROBABILITY OF FRAUD USING A LINEAR DECISION BOUNDARY BASED ON LOGISTIC(SIGMOID) FUNCTION.



$$P(y = 1 | \mathbf{x}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

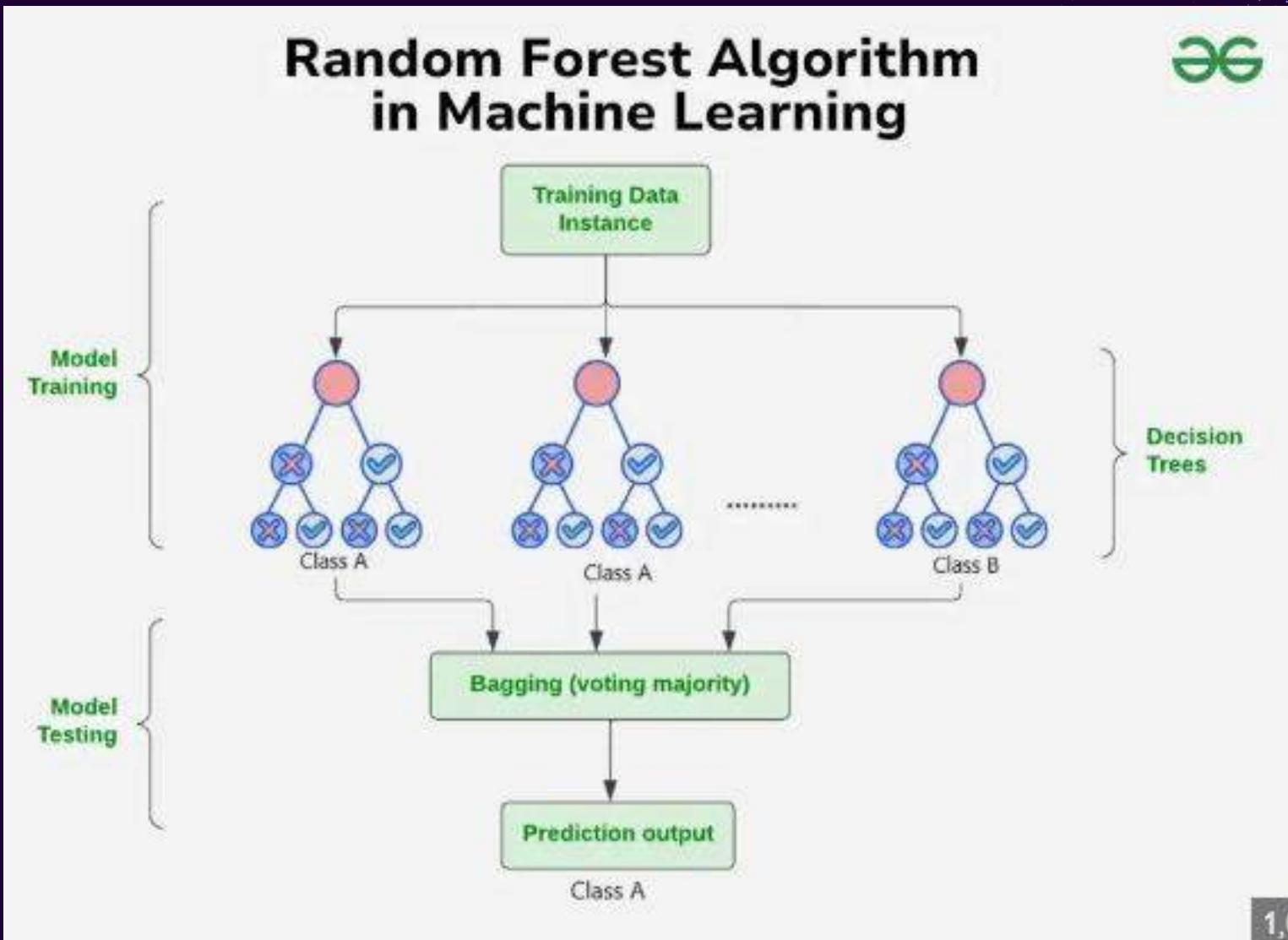
DECISION TREE

- A DECISION TREE SPLITS DATA BASED ON FEATURE VALUES TO CLASSIFY TRANSACTIONS.
- It builds a tree-like structure where each node represents a decision based on a feature.
- EACH LEAF NODE CORRESPONDS TO A PREDICTED CLASS (FRAUD OR NON-FRAUD).



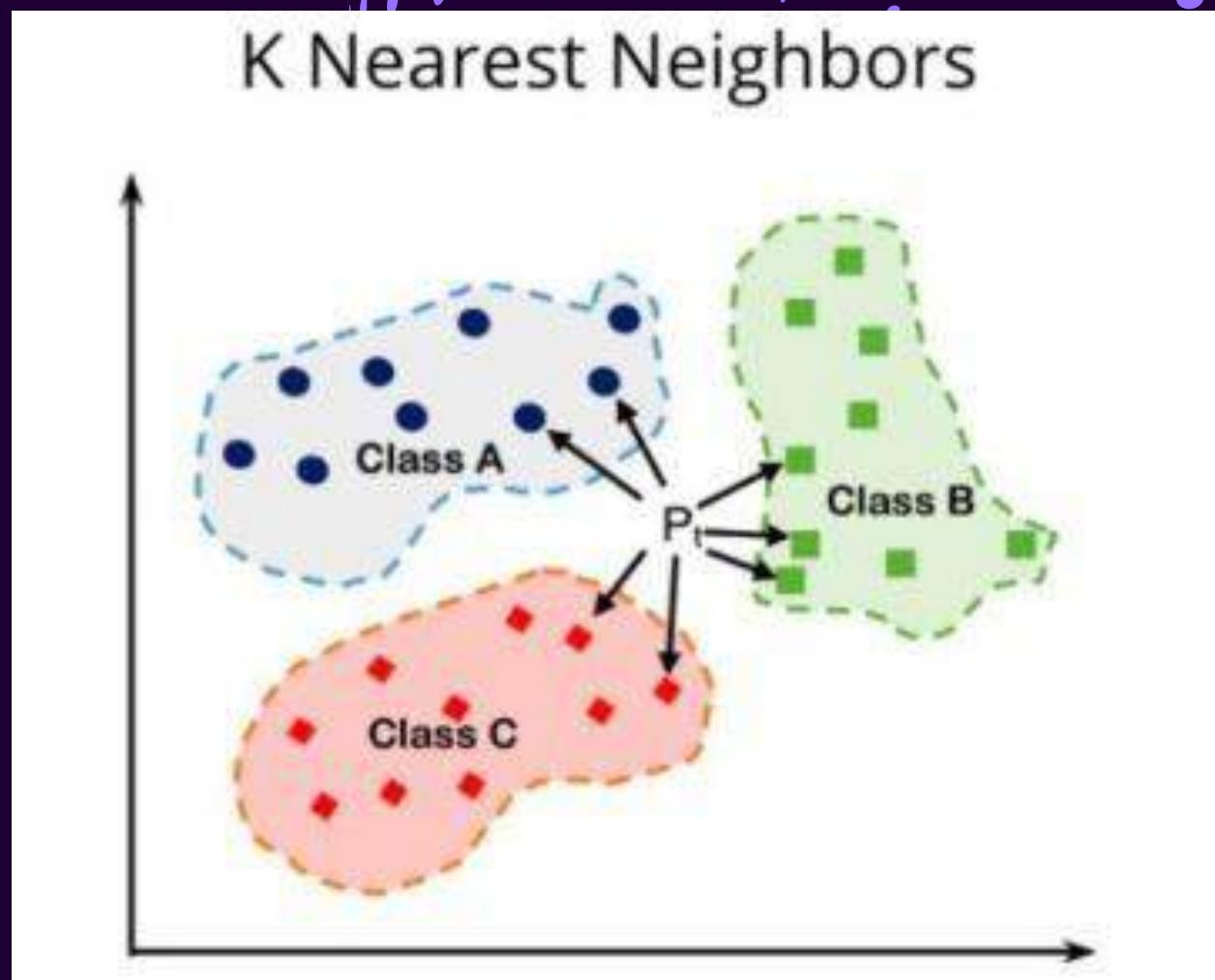
RANDOM FOREST

- RANDOM FOREST IS AN ENSEMBLE MODEL MADE OF MULTIPLE DECISION TREES.
- EACH TREE IS BUILT using a random subset of data and features, and the final prediction is based on the majority vote



K-NEAREST NEIGHBORS

- K-NEAREST NEIGHBORS IS A SUPERVISED MACHINE LEARNING ALGORITHM USED for classification and regression.
- IT CLASSIFIES A NEW DATA POINT BASED ON THE MAJORITY CLASS OF ITS K NEAREST NEIGHBORS IN THE FEATURE SPACE.



RESULTS

- MODEL PERFORMANCE
- CONFUSION MATRIX
- ROC CURVE
- PRECISION RECALL
CURVE

MODEL PERFORMANCE

ACCURACY

- THE PROPORTION OF CORRECT PREDICTIONS (BOTH TRUE POSITIVES AND TRUE NEGATIVES) OUT OF ALL PREDICTIONS.

$$\text{Accuracy} = \frac{\text{Correct predictions}}{\text{All predictions}}$$

Precision

- THE PROPORTION OF CORRECT POSITIVE PREDICTIONS (FRAUDULENT TRANSACTIONS) OUT OF ALL PREDICTED POSITIVES.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

MODEL PERFORMANCE

RECALL

- THE PROPORTION OF CORRECT POSITIVE PREDICTIONS (FRAUDULENT TRANSACTIONS) OUT OF ALL ACTUAL POSITIVES (FRAUDULENT TRANSACTIONS).

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

F1-SCORE

- THE HARMONIC MEAN OF PRECISION AND RECALL, USED WHEN YOU NEED A BALANCE BETWEEN THE TWO.

$$\begin{aligned}\text{F1 Score} &= \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \\ &= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}\end{aligned}$$

MODEL PERFORMANCE

LOGISTIC REGRESSION

	precision	recall	f1-score	support
0	0.72	0.69	0.70	33343
1	0.70	0.73	0.72	33343
accuracy			0.71	66686
macro avg	0.71	0.71	0.71	66686
weighted avg	0.71	0.71	0.71	66686

DECISION TREE

	precision	recall	f1-score	support
0	0.84	0.98	0.90	33343
1	0.97	0.81	0.89	33343
accuracy			0.90	66686
macro avg	0.91	0.90	0.89	66686
weighted avg	0.91	0.90	0.89	66686

RANDOM FOREST

	precision	recall	f1-score	support
0	0.89	0.97	0.93	33343
1	0.97	0.88	0.92	33343
accuracy			0.92	66686
macro avg	0.93	0.92	0.92	66686
weighted avg	0.93	0.92	0.92	66686

K-NEAREST NEIGHBORS

	precision	recall	f1-score	support
0	0.89	0.92	0.91	33343
1	0.92	0.88	0.90	33343
accuracy			0.90	66686
macro avg	0.90	0.90	0.90	66686
weighted avg	0.90	0.90	0.90	66686

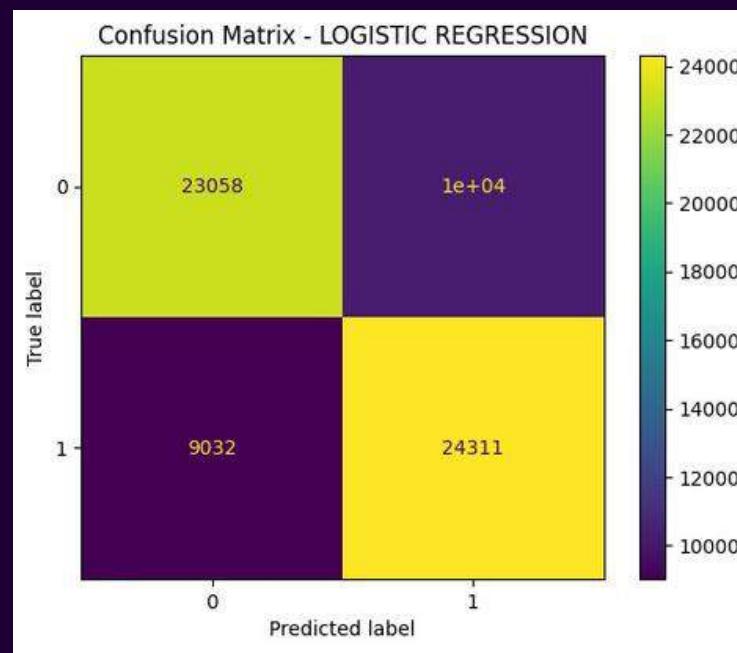
CONFUSION MATRIX

- A CONFUSION MATRIX IS A TABLE USED TO EVALUATE THE PERFORMANCE OF A CLASSIFICATION MODEL BY COMPARING THE PREDICTED CLASS TO THE ACTUAL CLASS.
- IT SUMMARIZES THE RESULTS OF CLASSIFICATION AND HELPS TO VISUALIZE THE PERFORMANCE WITH RESPECT TO TRUE POSITIVES, FALSE POSITIVES, TRUE NEGATIVES, AND FALSE NEGATIVES.

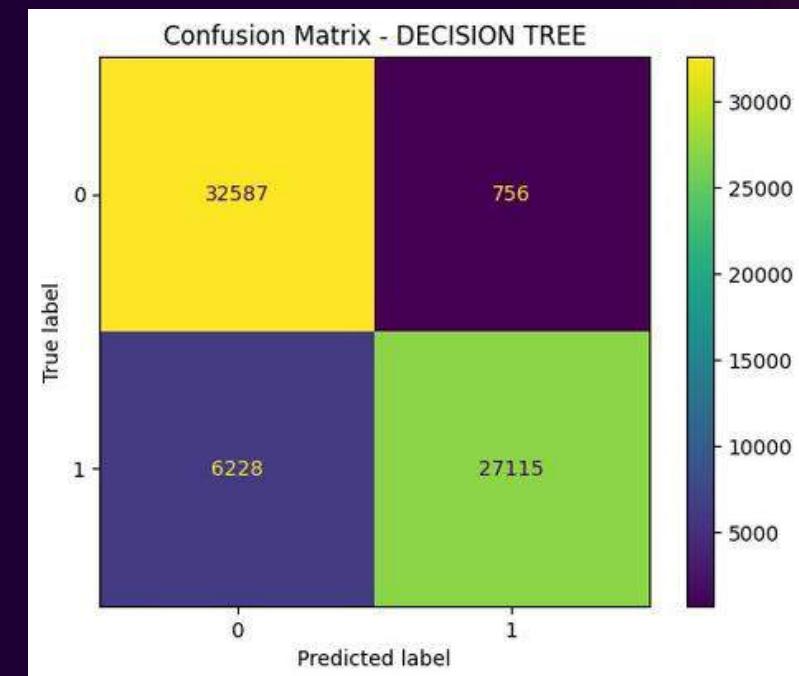
		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

CONFUSION MATRIX

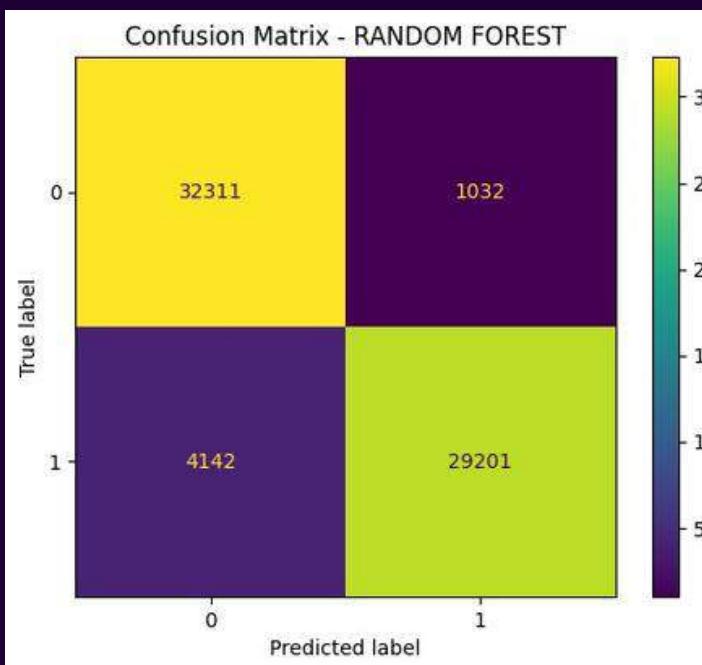
LOGISTIC REGRESSION



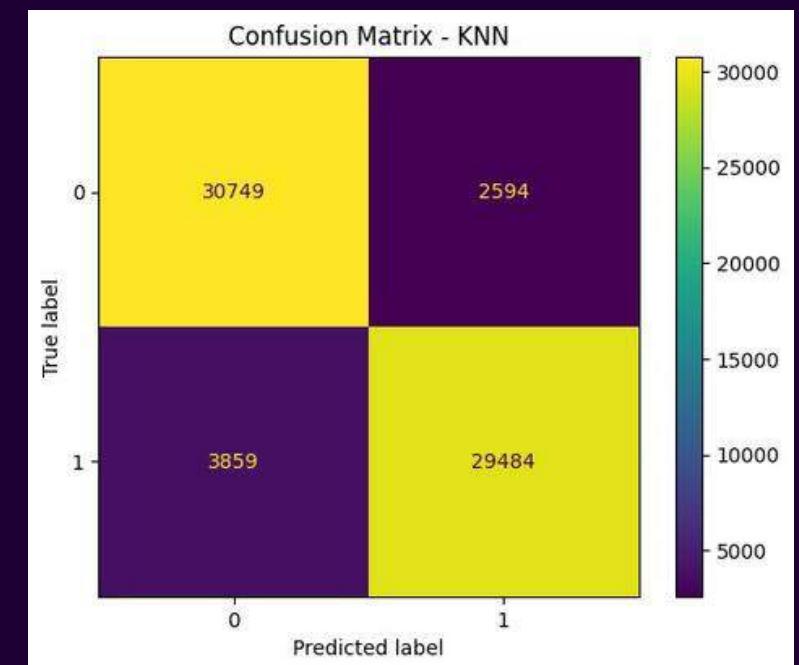
DECISION TREE



RANDOM FOREST

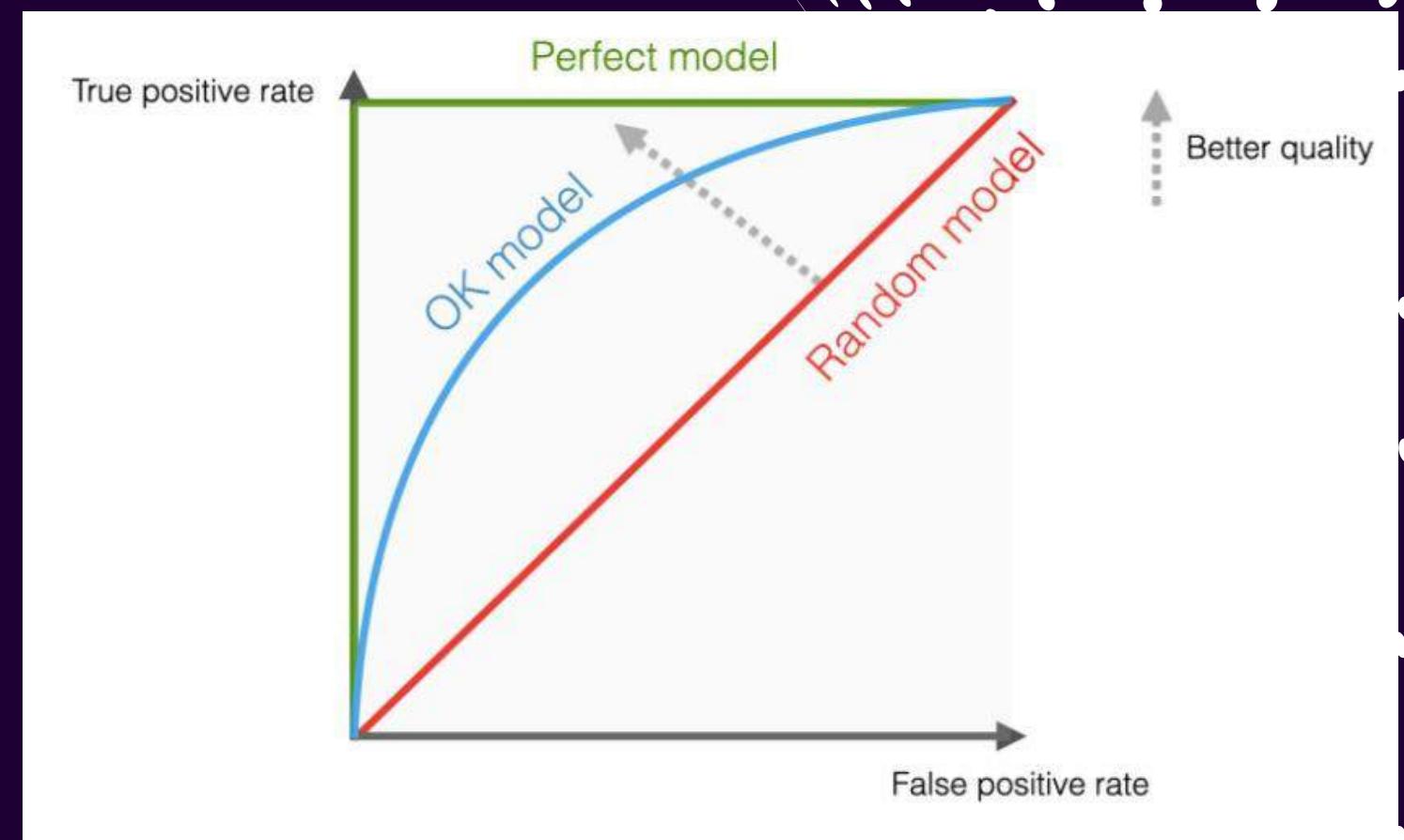


K-NEAREST NEIGHBORS



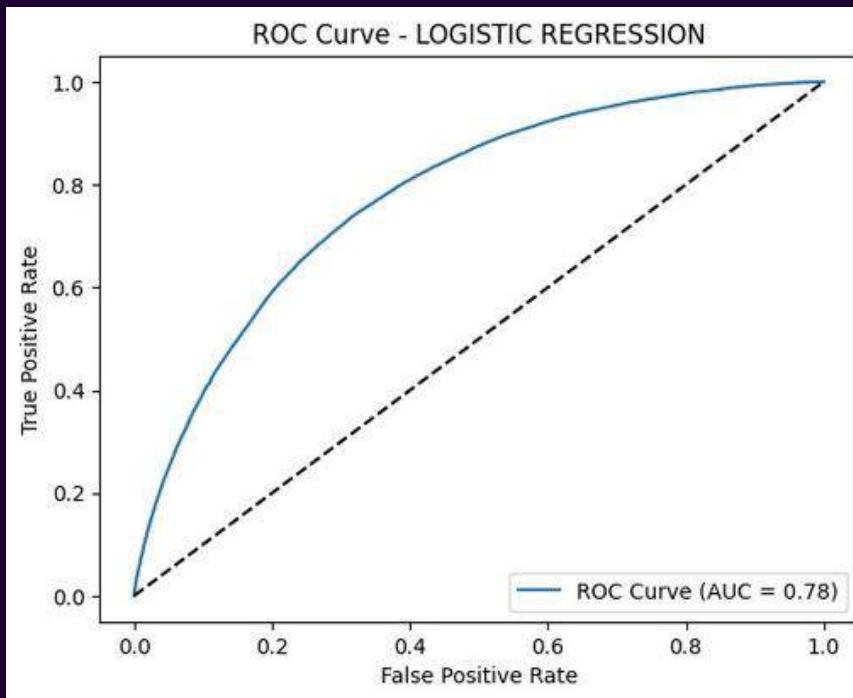
ROC CURVE

- THE RECEIVER OPERATING CHARACTERISTIC (ROC) CURVE IS A GRAPHICAL REPRESENTATION OF A CLASSIFIER'S PERFORMANCE ACROSS ALL CLASSIFICATION THRESHOLDS.
- IT PLOTS THE TRUE POSITIVE RATE (RECALL) AGAINST THE FALSE POSITIVE RATE FOR VARIOUS THRESHOLD VALUES.

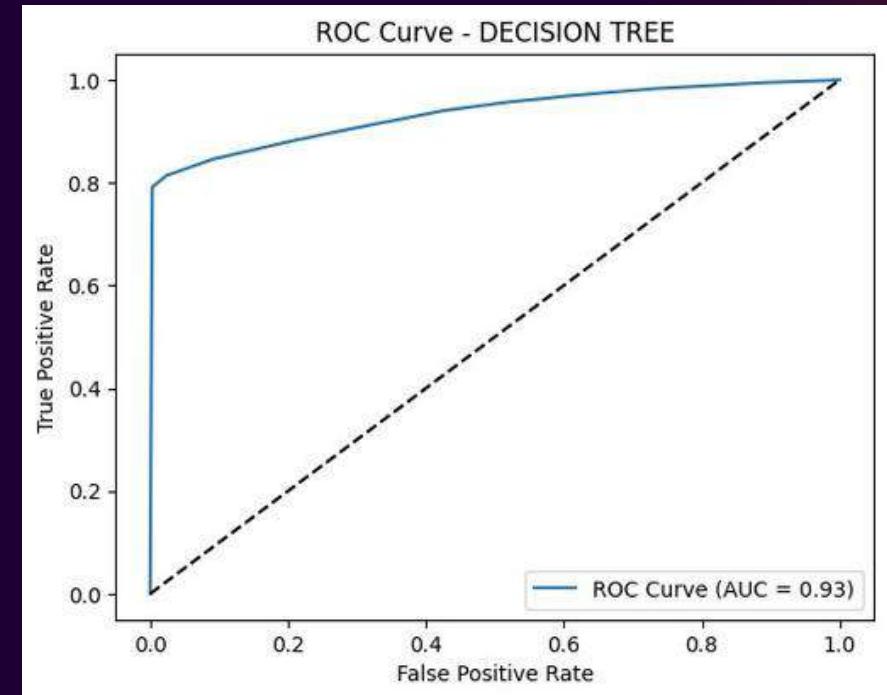


ROC CURVE

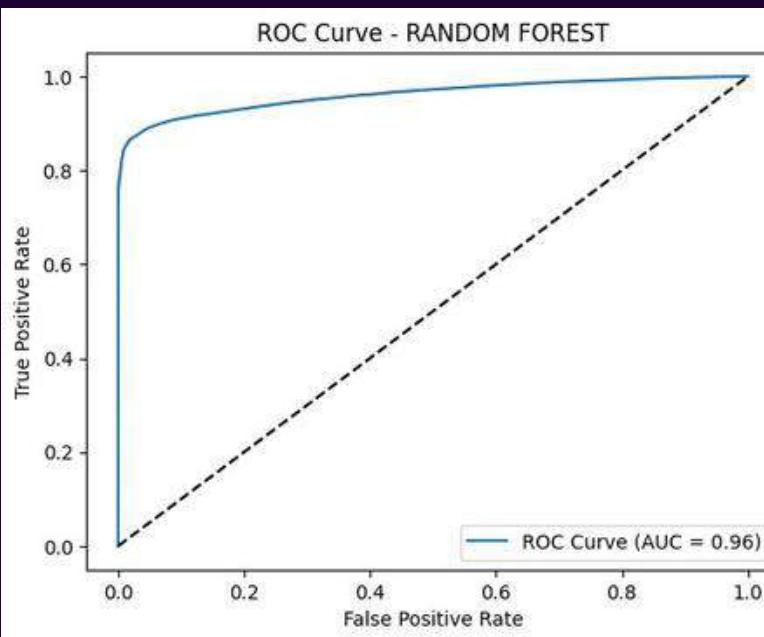
LOGISTIC REGRESSION



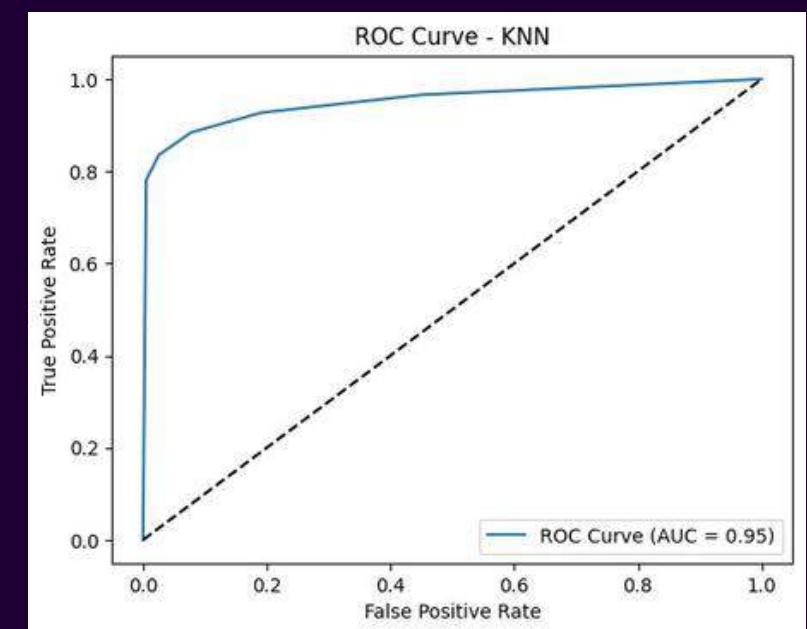
DECISION TREE



RANDOM FOREST

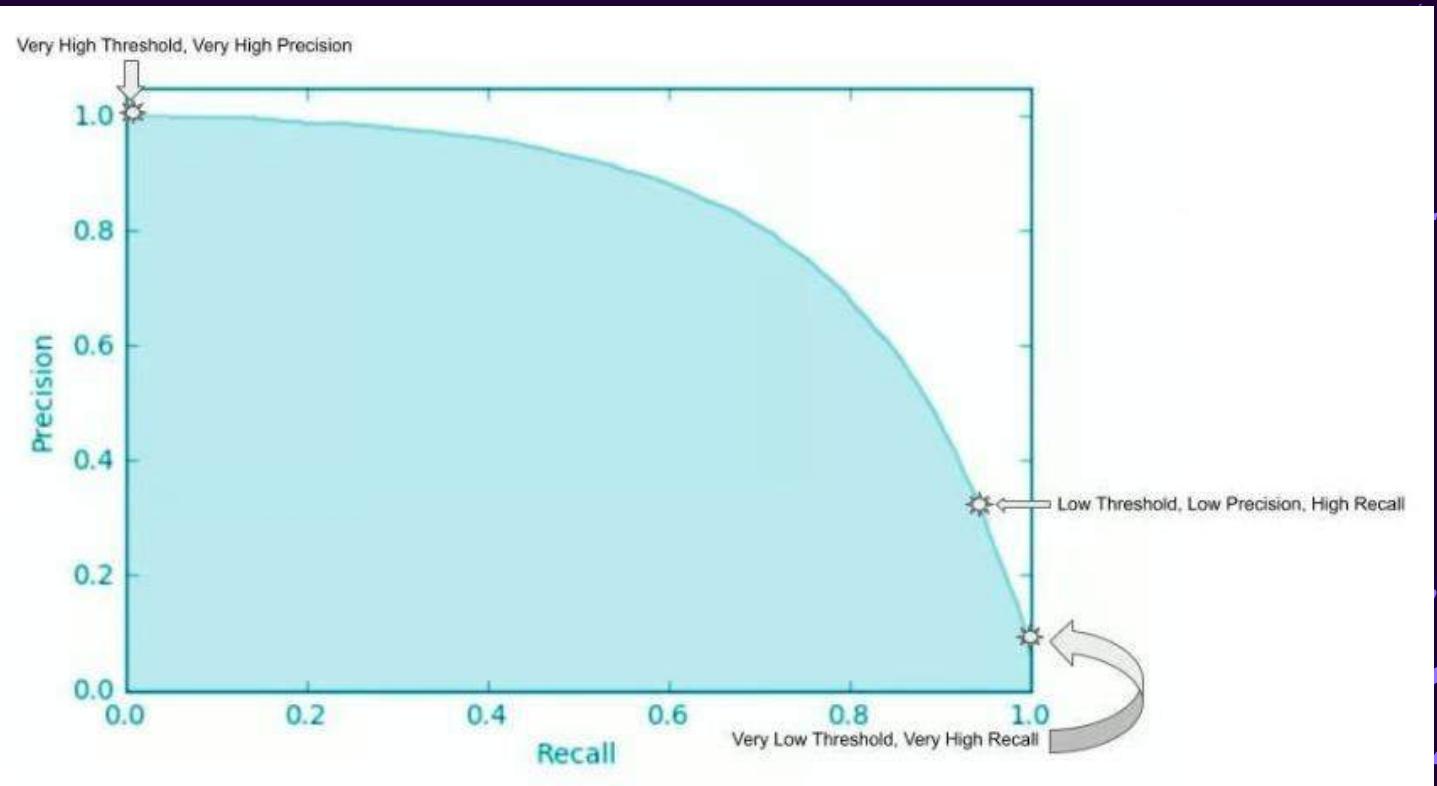


K-NEAREST NEIGHBORS



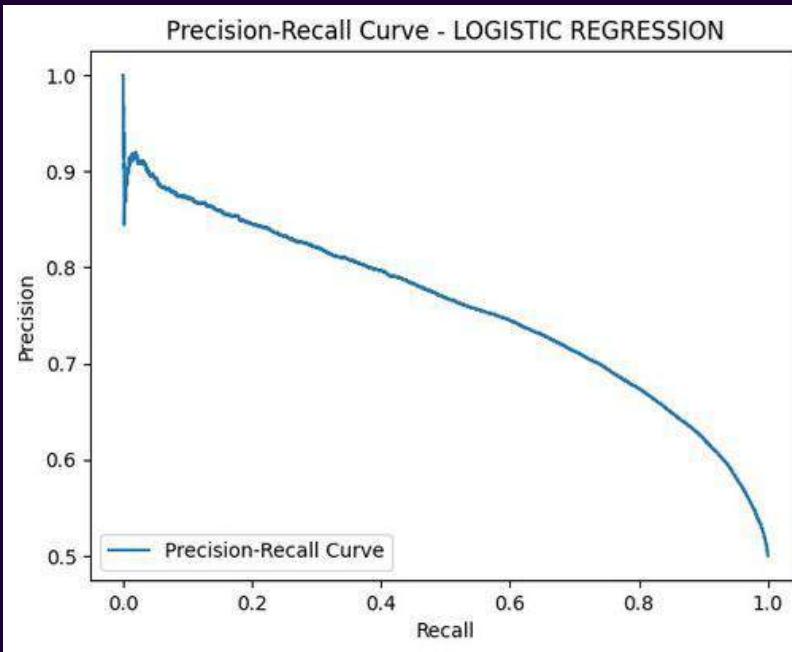
PRECISION RECALL CURVE

- THE PRECISION-RECALL CURVE (PR CURVE) IS A GRAPHICAL REPRESENTATION OF THE RELATIONSHIP BETWEEN PRECISION AND RECALL ACROSS DIFFERENT THRESHOLD VALUES.
- IT IS PARTICULARLY USEFUL FOR EVALUATING CLASSIFIERS ON IMBALANCED DATASETS, SUCH AS CREDIT CARD FRAUD DETECTION, WHERE FRAUDULENT TRANSACTIONS ARE MUCH RARER THAN LEGITIMATE ONES.

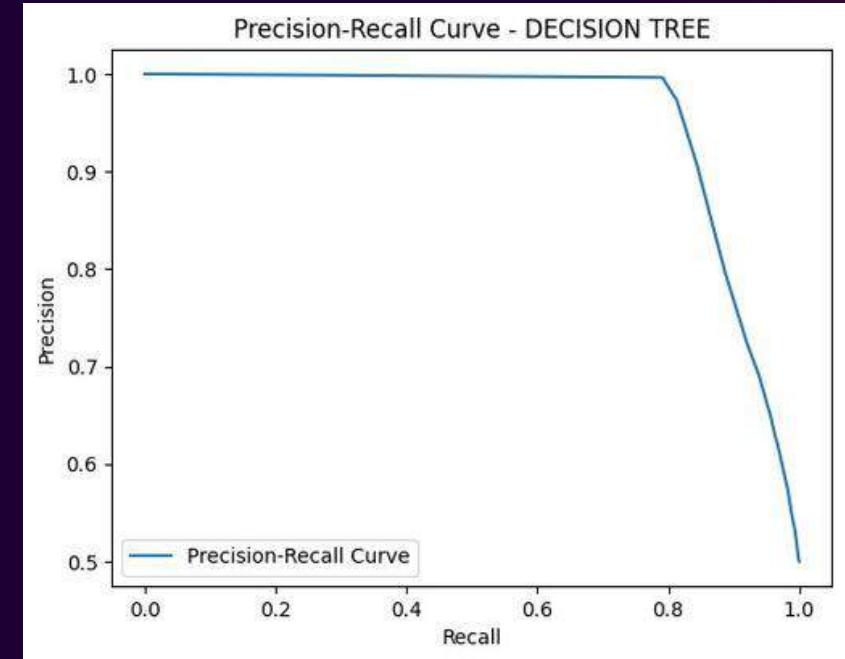


PRECISION RECALL CURVE

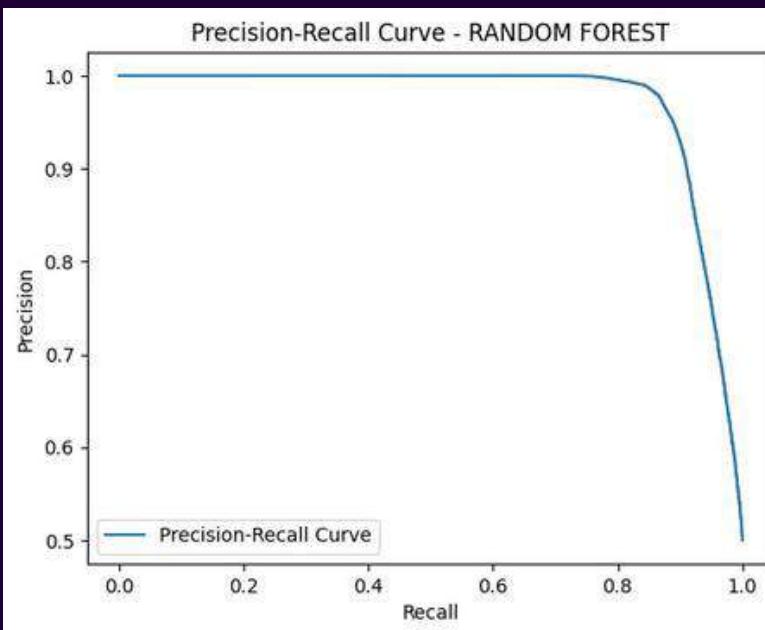
LOGISTIC REGRESSION



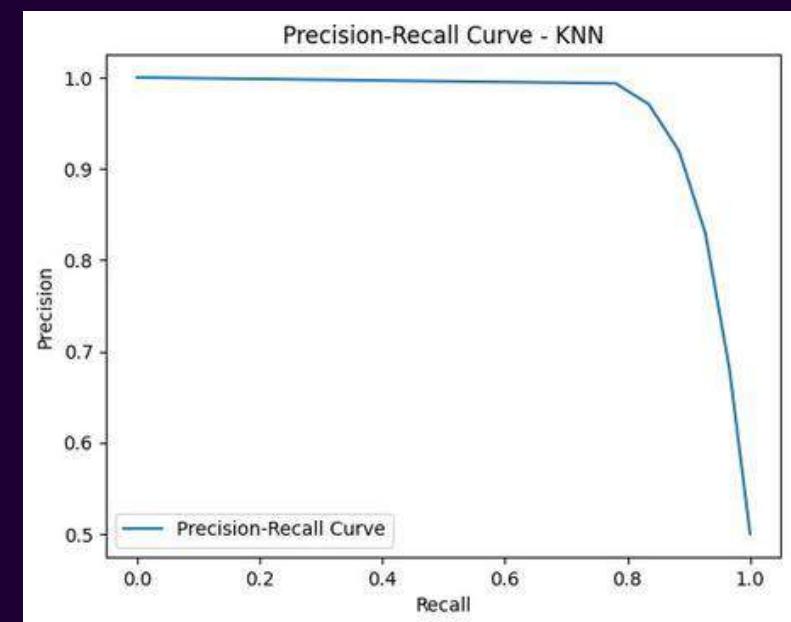
DECISION TREE



RANDOM FOREST



K-NEAREST NEIGHBORS



HYPERPARAMETER TUNING

- HYPERPARAMETER TUNING HELPS FIND THE BEST SET OF PARAMETERS FOR EACH MODEL, IMPROVING ITS ABILITY TO DETECT FRAUDULENT TRANSACTIONS.
- TECHNIQUES LIKE GRID SEARCH OR RANDOMIZED SEARCH ARE USED TO EXPERIMENT WITH VARIOUS HYPERPARAMETERS AND ENHANCE MODEL ACCURACY.

	model	best_score	best_params
0	decision_tree	0.915027	{'max_depth': 25}
1	random_forest	0.936740	{'max_depth': 20, 'n_estimators': 100}
2	logistic_regression	0.582882	{'max_iter': 1000, 'solver': 'liblinear'}
3	knn	0.922030	{'knn_n_neighbors': 9, 'knn_p': 1, 'knn_wei...

CONCLUSION

- **MACHINE LEARNING IMPROVES FRAUD DETECTION ACCURACY SIGNIFICANTLY.**
- Reduces financial risk for businesses and users.
- Dynamic, scalable solution adaptable to new fraud patterns.

THANK YOU
