

Exp. No 1

Outlier Detection, Data Preprocessing, and Model Accuracy Enhancement

Date:**Task:**

Present your view on the different techniques you have employed to do outlier analysis, handling missing data, feature engineering, feature importance and improving the accuracy of the model both from a classifier as well as a regressor. Use any sample data and present your POV in a well-structured presentation.

Aim:

To analyze various data preprocessing and feature optimization techniques including outlier detection, missing value handling, feature engineering, and feature selection. Further, to evaluate their impact on improving the performance and accuracy of both classification and regression models using sample datasets.

Algorithm:

1. Load the dataset and required libraries for preprocessing and modeling.
2. Perform exploratory data analysis (EDA) to understand data types, distribution, missing values, and summary statistics.
3. Engineer new features (e.g., car age) and remove/rename irrelevant or inconsistent columns.
4. Detect and handle outliers using visualization and the Interquartile Range (IQR) method.
5. Handle missing data using mean/median for numerical columns and mode for categorical columns.
6. Convert categorical features into numerical values using Label Encoding or One-Hot Encoding.
7. Scale numerical features using Standardization or Normalization if required by the model.
8. Split the data into training and testing sets for both regression and (optional) classification tasks.
9. Train and tune machine learning models such as Random Forest and Gradient Boosting using cross-validation to improve accuracy.
10. Evaluate model performance using metrics like R² for regression or Accuracy for classification, and analyze feature importance.

Program:

Case Study 1: Present your view on the different techniques you have employed to do outlier analysis, handling missing data, feature engineering, feature importance and improving the accuracy of the model both from a classifier as well as a regressor. Use any sample data and present your POV in a well-structured presentation.

```
import pandas as pd

df = pd.read_csv('/content/housing.csv')
df.head()

{"summary": {"name": "df", "rows": 20640, "fields": [{"column": "longitude", "properties": {"dtype": "number", "std": 2.003531723502581, "min": -124.35, "max": -114.31, "num_unique_values": 844, "samples": [-118.63, -119.86, -121.26]}, {"column": "latitude", "properties": {"dtype": "number", "std": 2.1359523974571117, "min": 32.54, "max": 41.95, "num_unique_values": 862, "samples": [33.7, 34.41, 38.24]}], "semantic_type": "Geographic location"}, {"column": "housing_median_age", "properties": {"dtype": "number", "std": 12.585557612111637, "min": 1.0, "max": 52.0, "num_unique_values": 52, "samples": [25.0, 25.0, 7.0], "semantic_type": "Age"}, {"column": "total_rooms", "properties": {"dtype": "number", "std": 2181.615251582787, "min": 2.0, "max": 39320.0, "num_unique_values": 5926, "samples": [3966.0, 3966.0, 1578.0], "semantic_type": "Rooms"}, {"column": "total_bedrooms", "properties": {"dtype": "number", "std": 421.3850700740322, "min": 1.0, "max": 6445.0, "num_unique_values": 1923, "samples": [1538.0, 1538.0, 1298.0], "semantic_type": "Bedrooms"}, {"column": "population", "properties": {"dtype": "number", "std": 1132.4621217653375, "min": 3.0, "max": 35682.0, "num_unique_values": 3888, "samples": [3367.0, 3367.0, 4169.0], "semantic_type": "Population"}, {"column": "households", "properties": {"dtype": "number", "std": 382.3297528316099, "min": 1, "max": 12800, "num_unique_values": 12800, "samples": [1, 1, 1], "semantic_type": "Households"}]}], "description": "A dataset containing information about housing prices in various US cities. The columns include longitude, latitude, housing median age, total rooms, total bedrooms, population, households, and price."}
```

```

1.0,\n      \\"max\\": 6082.0,\n      \\"num_unique_values\\": 1815,\n
\\\"samples\\": [\n      21.0,\n      750.0,\n      1447.0\n
],\n      \\"semantic_type\\": \"\",\\n      \\"description\\": \"\"\n
}\\n  },\\n  {\n    \\"column\\": \\"median_income\\\",\\n
\\\"properties\\": {\n      \\"dtype\\": \\\"number\\\",\\n      \\"std\\\":\n      1.8998217179452732,\n      \\"min\\": 0.4999,\n      \\"max\\":\n      15.0001,\n      \\"num_unique_values\\": 12928,\n      \\"samples\\":\n      [\n        5.0286,\n        2.0433,\n        6.1228\n      ],\\n
      \\"semantic_type\\": \"\",\\n      \\"description\\": \"\"\n
    }\\n  },\\n  {\n    \\"column\\": \\"median_house_value\\\",\\n
\\\"properties\\": {\n      \\"dtype\\": \\\"number\\\",\\n      \\"std\\\":\n      115395.6158744132,\n      \\"min\\": 14999.0,\n      \\"max\\":\n      500001.0,\n      \\"num_unique_values\\": 3842,\n      \\"samples\\":\n      [\n        194300.0,\n        379000.0,\n        230100.0\n      ],\\n
      \\"semantic_type\\": \"\",\\n      \\"description\\": \"\"\n
    }\\n  },\\n  {\n    \\"column\\": \\"ocean_proximity\\\",\\n
\\\"properties\\": {\n      \\"dtype\\": \\\"category\\\",\\n
      \\"num_unique_values\\": 5,\n      \\"samples\\": [\n        \\"<1H\nOCEAN\\\",\\n        \\"ISLAND\\\",\\n        \\"INLAND\\\"\n      ],\\n
      \\"semantic_type\\": \"\",\\n      \\"description\\": \"\"\n
    }\\n  }\n ]\\n},\\n \\\"type\\": \\\"dataframe\\\",\\n \\\"variable_name\\": \\\"df\\\"\n}

print("DataFrame Info:")
df.info()

print("\nDescriptive Statistics:")
df.describe()

print("\nNumber of duplicate rows:")
print(df.duplicated().sum())

DataFrame Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
 #   Column            Non-Null Count Dtype  
 --- 
 0   longitude         20640 non-null  float64 
 1   latitude          20640 non-null  float64 
 2   housing_median_age 20640 non-null  float64 
 3   total_rooms        20640 non-null  float64 
 4   total_bedrooms     20433 non-null  float64 
 5   population         20640 non-null  float64 
 6   households         20640 non-null  float64 
 7   median_income       20640 non-null  float64 
 8   median_house_value 20640 non-null  float64 
 9   ocean_proximity     20640 non-null  object  
dtypes: float64(9), object(1)
memory usage: 1.6+ MB

```

Descriptive Statistics:

Number of duplicate rows:
0

```
print("Missing values percentage:")
print(missing_percentage[missing_percentage > 0])

# Impute missing values in 'total_bedrooms' with its median
df['total_bedrooms'].fillna(median_total_bedrooms, inplace=True)
```

```
print("\nMissing values after imputation:")
print(df.isnull().sum())
```

Missing values percentage:
total_bedrooms 1.002907
dtype: float64

Missing values after imputation:
longitude 0
latitude 0
housing_median_age 0
total_rooms 0
total_bedrooms 0
population 0
households 0
median_income 0
median_house_value 0
ocean_proximity 0
dtype: int64

/tmp/ipython-input-3517965585.py:5: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df['total_bedrooms'].fillna(median_total_bedrooms, inplace=True)

missing_percentage = (df.isnull().sum() / len(df)) * 100
print("Missing values percentage:")
print(missing_percentage[missing_percentage > 0])
```

```
# Impute missing values in 'total_bedrooms' with its median
df['total_bedrooms'] =
df['total_bedrooms'].fillna(median_total_bedrooms)

print("\nMissing values after imputation:")
print(df.isnull().sum())

Missing values percentage:
Series([], dtype: float64)

Missing values after imputation:
longitude          0
latitude           0
housing_median_age 0
total_rooms         0
total_bedrooms      0
population          0
households          0
median_income        0
median_house_value   0
ocean_proximity      0
dtype: int64

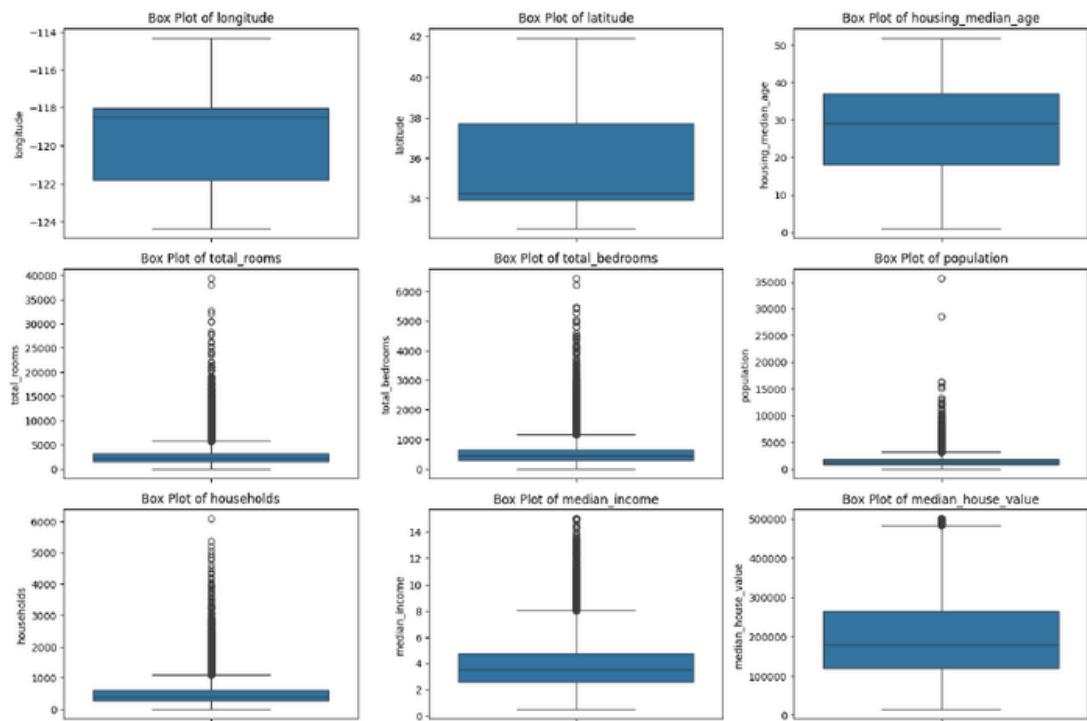
import matplotlib.pyplot as plt
import seaborn as sns

numerical_cols = df.select_dtypes(include=['number']).columns

print("Numerical columns selected for outlier detection:")
print(numerical_cols.tolist())

Numerical columns selected for outlier detection:
['longitude', 'latitude', 'housing_median_age', 'total_rooms',
 'total_bedrooms', 'population', 'households', 'median_income',
 'median_house_value']

plt.figure(figsize=(15, 10))
for i, col in enumerate(numerical_cols):
    plt.subplot(3, 3, i + 1) # Adjust subplot grid as needed
    sns.boxplot(y=df[col])
    plt.title(f'Box Plot of {col}')
    plt.ylabel(col)
plt.tight_layout()
plt.show()
```



```

Q1 = df['total_rooms'].quantile(0.25)
Q3 = df['total_rooms'].quantile(0.75)
IQR = Q3 - Q1

print(f"Q1 for total_rooms: {Q1}")
print(f"Q3 for total_rooms: {Q3}")
print(f"IQR for total_rooms: {IQR}")

Q1 for total_rooms: 1447.75
Q3 for total_rooms: 3148.0
IQR for total_rooms: 1700.25

upper_bound = Q3 + 1.5 * IQR
lower_bound = Q1 - 1.5 * IQR

print(f"Upper bound for total_rooms: {upper_bound}")
print(f"Lower bound for total_rooms: {lower_bound}")

df['total_rooms'] = df['total_rooms'].clip(lower=lower_bound,
                                         upper=upper_bound)

print("\nOutliers in 'total_rooms' capped.")
print(df['total_rooms'].describe())

Upper bound for total_rooms: 5698.375
Lower bound for total_rooms: -1102.625

Outliers in 'total_rooms' capped.

```

```

count    20640.000000
mean     2441.692472
std      1397.790038
min      2.000000
25%     1447.750000
50%     2127.000000
75%     3148.000000
max      5698.375000
Name: total_rooms, dtype: float64

df['rooms_per_household'] = df['total_rooms'] / df['households']
df['bedrooms_per_room'] = df['total_bedrooms'] / df['total_rooms']
df['population_per_household'] = df['population'] / df['households']

print("New features created:")
print(df[['rooms_per_household', 'bedrooms_per_room',
'population_per_household']].head())

New features created:
   rooms_per_household  bedrooms_per_room  population_per_household
0            6.984127          0.146591                2.555556
1            5.007359          0.194090                2.109842
2            8.288136          0.129516                2.802260
3            5.817352          0.184458                2.547945
4            6.281853          0.172096                2.181467

df = pd.get_dummies(df, columns=['ocean_proximity'], drop_first=False)

print("DataFrame after one-hot encoding 'ocean_proximity':")
print(df.head())

DataFrame after one-hot encoding 'ocean_proximity':
   longitude  latitude  housing_median_age  total_rooms
total_bedrooms \
0       -122.23      37.88                  41.0      880.000
129.0
1       -122.22      37.86                  21.0      5698.375
1106.0
2       -122.24      37.85                  52.0      1467.000
190.0
3       -122.25      37.85                  52.0      1274.000
235.0
4       -122.25      37.85                  52.0      1627.000
280.0

   population  households  median_income  median_house_value \
0        322.0        126.0        8.3252        452600.0
1      2401.0        1138.0       8.3014        358500.0
2        496.0        177.0        7.2574        352100.0
3        558.0        219.0        5.6431        341300.0
4        565.0        259.0        3.8462        342200.0

```

```
rooms_per_household    bedrooms_per_room   population_per_household \
0                      6.984127           0.146591                  2.555556
1                      5.007359           0.194090                  2.109842
2                      8.288136           0.129516                  2.802260
3                      5.817352           0.184458                  2.547945
4                      6.281853           0.172096                  2.181467

ocean_proximity_<1H OCEAN  ocean_proximity_INLAND
ocean_proximity_ISLAND \
0                         False                   False
False
1                         False                   False
False
2                         False                   False
False
3                         False                   False
False
4                         False                   False
False

ocean_proximity_NEAR BAY  ocean_proximity_NEAR OCEAN
0                         True                    False
1                         True                    False
2                         True                    False
3                         True                    False
4                         True                    False

from sklearn.model_selection import train_test_split

# 1. Separate target variable and features
X = df.drop('median_house_value', axis=1)
y = df['median_house_value']

# 2. Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

print(f"X_train shape: {X_train.shape}")
print(f"X_test shape: {X_test.shape}")
print(f"y_train shape: {y_train.shape}")
print(f"y_test shape: {y_test.shape}")

X_train shape: (16512, 16)
X_test shape: (4128, 16)
y_train shape: (16512,)
y_test shape: (4128,)

from sklearn.preprocessing import StandardScaler
```

```

# 3. Identify numerical features for scaling (all except one-hot
encoded ocean_proximity columns)
numerical_features =
X_train.select_dtypes(include=['number']).columns.tolist()
one_hot_encoded_features = [col for col in X_train.columns if
'ocean_proximity' in col]

features_to_scale = [f for f in numerical_features if f not in
one_hot_encoded_features]

print(f"Features to be scaled: {features_to_scale}")

# 4. Initialize StandardScaler
scaler = StandardScaler()

# 5. Fit the scaler to the numerical features of the training set and
transform both training and testing numerical features
X_train_scaled = scaler.fit_transform(X_train[features_to_scale])
X_test_scaled = scaler.transform(X_test[features_to_scale])

# Convert scaled arrays back to DataFrames with original column names
X_train_scaled_df = pd.DataFrame(X_train_scaled,
columns=features_to_scale, index=X_train.index)
X_test_scaled_df = pd.DataFrame(X_test_scaled,
columns=features_to_scale, index=X_test.index)

# 6. Reconstruct the training and testing feature sets
X_train = pd.concat([X_train_scaled_df,
X_train[one_hot_encoded_features]], axis=1)
X_test = pd.concat([X_test_scaled_df,
X_test[one_hot_encoded_features]], axis=1)

print("\nShape of X_train after scaling and reconstruction:",
X_train.shape)
print("Shape of X_test after scaling and reconstruction:",
X_test.shape)
print("\nX_train head after scaling:")
print(X_train.head())

Features to be scaled: ['longitude', 'latitude', 'housing_median_age',
'total_rooms', 'total_bedrooms', 'population', 'households',
'median_income', 'rooms_per_household', 'bedrooms_per_room',
'population_per_household']

Shape of X_train after scaling and reconstruction: (16512, 16)
Shape of X_test after scaling and reconstruction: (4128, 16)

X_train head after scaling:
    longitude    latitude    housing_median_age    total_rooms

```

```

total_bedrooms \
14196 1.272587 -1.372811          0.348490  0.484375
0.211228
8267  0.709162 -0.876696          1.618118  0.667355
0.593094
17445 -0.447603 -0.460146         -1.952710 -0.394072      -
0.495226
14265 1.232698 -1.382172          0.586545  -0.734300      -
0.409306
2271  -0.108551  0.532084          1.142008  -0.047410      -
0.256559

population households median_income rooms_per_household \
14196  0.768276  0.322906   -0.326196 -0.129648
8267  -0.098901  0.672027   -0.035843 -0.359534
17445 -0.449818 -0.430461   0.144701  0.135755
14265 -0.007434 -0.380587   -1.017864 -0.558416
2271  -0.485877 -0.314962   -0.171488 0.398796

bedrooms_per_room population_per_household
ocean_proximity_<1H_OCEAN \
14196      -0.273618          0.051376
False
8267       0.202884        -0.117362
False
17445      -0.660580        -0.032280
False
14265       0.582082        0.077507
False
2271      -0.564849        -0.068832
False

ocean_proximity_INLAND ocean_proximity_ISLAND \
14196           False        False
8267            False        False
17445           False        False
14265           False        False
2271            True        False

ocean_proximity_NEAR_BAY ocean_proximity_NEAR_OCEAN
14196             False        True
8267              False        True
17445             False        True
14265             False        True
2271              False       False

from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error, r2_score
import matplotlib.pyplot as plt

```

```
# Instantiate RandomForestRegressor
# Using n_estimators=100 as a common default, random_state for
# reproducibility
model = RandomForestRegressor(n_estimators=100, random_state=42)

print("RandomForestRegressor model instantiated.")

RandomForestRegressor model instantiated.

print("Training the model...")
model.fit(X_train, y_train)
print("Model training complete.")

# Make predictions on the test data
y_pred = model.predict(X_test)

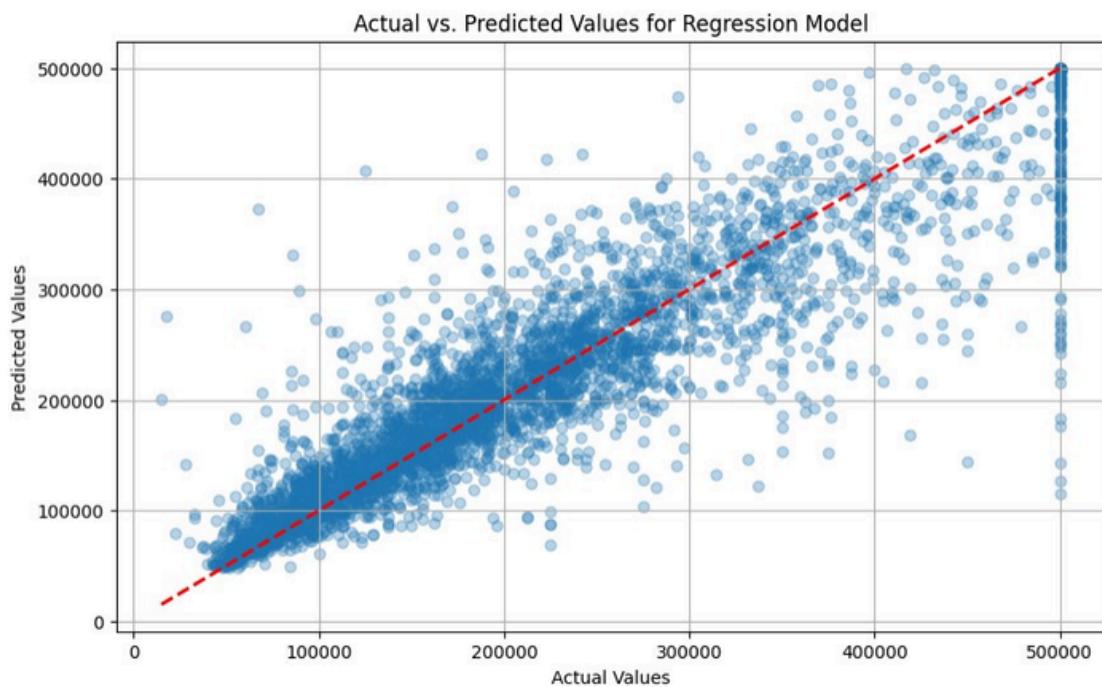
# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"\nMean Squared Error (MSE): {mse:.2f}")
print(f"R-squared (R2): {r2:.2f}")

# Visualize actual vs. predicted values
plt.figure(figsize=(10, 6))
plt.scatter(y_test, y_pred, alpha=0.3)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()],
         'r--', lw=2) # Diagonal line
plt.xlabel('Actual Values')
plt.ylabel('Predicted Values')
plt.title('Actual vs. Predicted Values for Regression Model')
plt.grid(True)
plt.show()
```

Training the model...
Model training complete.

Mean Squared Error (MSE): 2555928455.48
R-squared (R2): 0.80



```

feature_importances = model.feature_importances_
features = X_train.columns

importance_df = pd.DataFrame({'Feature': features, 'Importance':
feature_importances})
importance_df = importance_df.sort_values(by='Importance',
ascending=False)

print("Top 10 Feature Importances:")
print(importance_df.head(10))

Top 10 Feature Importances:
          Feature  Importance
7      median_income    0.482094
12  ocean_proximity_INLAND    0.137490
10  population_per_household    0.121920
0        longitude    0.058123
1        latitude    0.056419
2  housing_median_age    0.044479
8  rooms_per_household    0.024352
9  bedrooms_per_room    0.022988
3      total_rooms    0.012935
4      total_bedrooms    0.012349

plt.figure(figsize=(12, 7))
sns.barplot(x='Importance', y='Feature', data=importance_df.head(10),
palette='viridis')
plt.title('Top 10 Feature Importances (Regression Model)')
plt.xlabel('Importance Score')

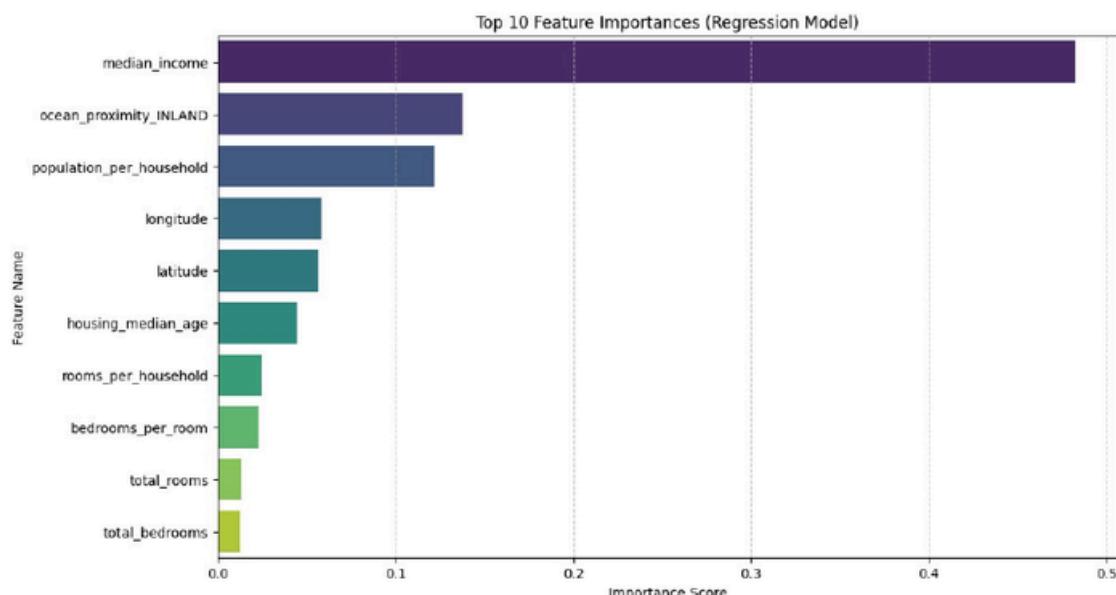
```

```
plt.ylabel('Feature Name')
plt.grid(axis='x', linestyle='--', alpha=0.7)
plt.show()
```

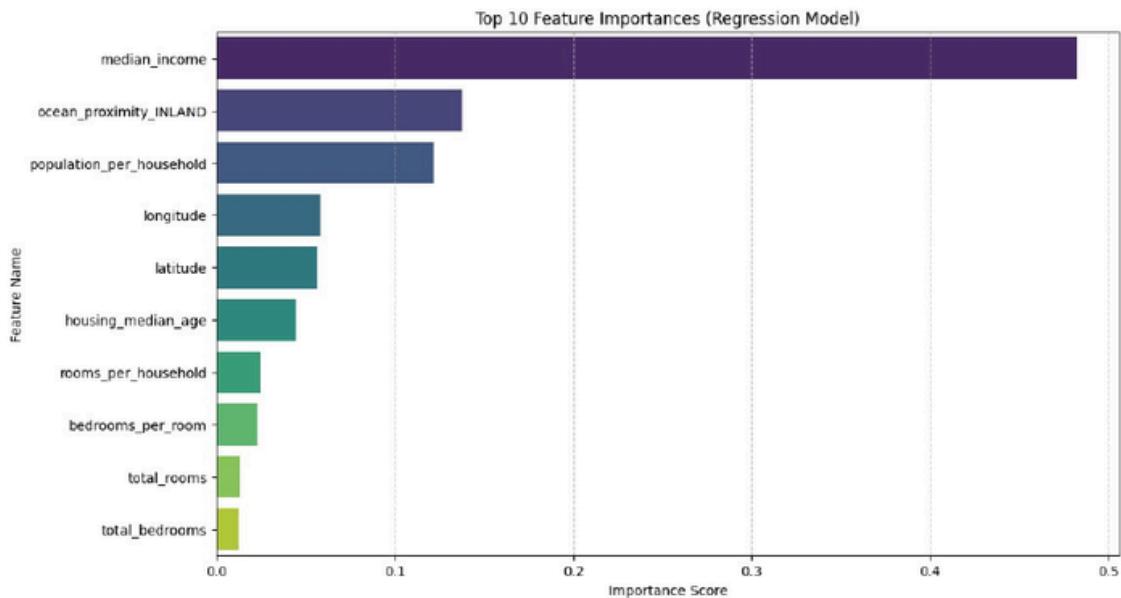
```
/tmp/ipython-input-403874356.py:2: FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be
removed in v0.14.0. Assign the `y` variable to `hue` and set
`legend=False` for the same effect.
```

```
sns.barplot(x='Importance', y='Feature',
data=importance_df.head(10), palette='viridis')
```



```
plt.figure(figsize=(12, 7))
sns.barplot(x='Importance', y='Feature', data=importance_df.head(10),
palette='viridis', hue='Feature', legend=False)
plt.title('Top 10 Feature Importances (Regression Model)')
plt.xlabel('Importance Score')
plt.ylabel('Feature Name')
plt.grid(axis='x', linestyle='--', alpha=0.7)
plt.show()
```



```

median_value = df['median_house_value'].median()
df['median_house_value_category'] =
df['median_house_value'].apply(lambda x: 'High Value' if x >
median_value else 'Low Value')

print(f"Median house value used for classification: {median_value}")
print("Distribution of new target variable:")
print(df['median_house_value_category'].value_counts())

Median house value used for classification: 179700.0
Distribution of new target variable:
median_house_value_category
Low Value    10323
High Value   10317
Name: count, dtype: int64

from sklearn.model_selection import train_test_split

# 3. Separate target variable and features for classification
X_clf = df.drop(['median_house_value', 'median_house_value_category'],
axis=1)
y_clf = df['median_house_value_category']

# 4. Split data into training and testing sets
X_train_clf, X_test_clf, y_train_clf, y_test_clf =
train_test_split(X_clf, y_clf, test_size=0.2, random_state=42)

print(f"X_clf shape: {X_clf.shape}")
print(f"y_clf shape: {y_clf.shape}")
print(f"X_train_clf shape: {X_train_clf.shape}")
print(f"X_test_clf shape: {X_test_clf.shape}")

```

```

print(f"y_train_clf shape: {y_train_clf.shape}")
print(f"y_test_clf shape: {y_test_clf.shape}")

X_clf shape: (20640, 16)
y_clf shape: (20640,)
X_train_clf shape: (16512, 16)
X_test_clf shape: (4128, 16)
y_train_clf shape: (16512,)
y_test_clf shape: (4128,)

from sklearn.preprocessing import StandardScaler

# 5. Identify numerical features for scaling (all numeric except one-hot encoded ocean_proximity columns)
numerical_features_clf =
X_train_clf.select_dtypes(include=['number']).columns.tolist()
one_hot_encoded_features_clf = [col for col in X_train_clf.columns if
'ocean_proximity' in col]

features_to_scale_clf = [f for f in numerical_features_clf if f not in
one_hot_encoded_features_clf]

print(f"Features to be scaled for classification:
{features_to_scale_clf}")

# 6. Initialize StandardScaler
scaler_clf = StandardScaler()

# 7. Fit the scaler to the numerical features of the training set and
transform both training and testing numerical features
X_train_clf_scaled =
scaler_clf.fit_transform(X_train_clf[features_to_scale_clf])
X_test_clf_scaled =
scaler_clf.transform(X_test_clf[features_to_scale_clf])

# 8. Convert scaled arrays back to DataFrames with original column
names and indices
X_train_clf_scaled_df = pd.DataFrame(X_train_clf_scaled,
columns=features_to_scale_clf, index=X_train_clf.index)
X_test_clf_scaled_df = pd.DataFrame(X_test_clf_scaled,
columns=features_to_scale_clf, index=X_test_clf.index)

# 9. Reconstruct the training and testing feature sets
X_train_clf = pd.concat([X_train_clf_scaled_df,
X_train_clf[one_hot_encoded_features_clf]], axis=1)
X_test_clf = pd.concat([X_test_clf_scaled_df,
X_test_clf[one_hot_encoded_features_clf]], axis=1)

print("\nShape of X_train_clf after scaling and reconstruction:",
X_train_clf.shape)

```

```

print("Shape of X_test_clf after scaling and reconstruction:", X_test_clf.shape)
print("\nX_train_clf head after scaling:")
print(X_train_clf.head())

Features to be scaled for classification: ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households', 'median_income', 'rooms_per_household', 'bedrooms_per_room', 'population_per_household']

Shape of X_train_clf after scaling and reconstruction: (16512, 16)
Shape of X_test_clf after scaling and reconstruction: (4128, 16)

X_train_clf head after scaling:
   longitude    latitude  housing_median_age  total_rooms
total_bedrooms \
14196      1.272587 -1.372811           0.348490     0.484375
0.211228
8267       0.709162 -0.876696           1.618118     0.667355
0.593094
17445      -0.447603 -0.460146          -1.952710    -0.394072    -
0.495226
14265      1.232698 -1.382172           0.586545    -0.734300    -
0.409306
2271       -0.108551  0.532084           1.142008    -0.047410    -
0.256559

   population  households  median_income  rooms_per_household \
14196      0.768276    0.322906     -0.326196      -0.129648
8267      -0.098901    0.672027     -0.035843      -0.359534
17445      -0.449818    -0.430461      0.144701      0.135755
14265      -0.007434    -0.380587     -1.017864     -0.558416
2271      -0.485877    -0.314962     -0.171488      0.398796

   bedrooms_per_room  population_per_household
ocean_proximity_<1H_OCEAN \
14196            -0.273618           0.051376
False
8267             0.202884          -0.117362
False
17445            -0.660580          -0.032280
False
14265            0.582082           0.077507
False
2271            -0.564849          -0.068832
False

   ocean_proximity_INLAND  ocean_proximity_ISLAND \
14196                  False                  False
8267                  False                  False

```

```
17445          False          False
14265          False          False
2271           True          False

ocean_proximity_NEAR BAY ocean_proximity_NEAR OCEAN
14196           False          True
8267            False          True
17445           False          True
14265           False          True
2271            False          False

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score, confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns

# 1. Instantiate RandomForestClassifier
model_clf = RandomForestClassifier(n_estimators=100, random_state=42)

print("Training the classification model...")
# 2. Train the model
model_clf.fit(X_train_clf, y_train_clf)
print("Classification model training complete.")

# 3. Make predictions on the test set
y_pred_clf = model_clf.predict(X_test_clf)

# 4. Evaluate the model
accuracy = accuracy_score(y_test_clf, y_pred_clf)
precision = precision_score(y_test_clf, y_pred_clf, pos_label='High Value')
recall = recall_score(y_test_clf, y_pred_clf, pos_label='High Value')
f1 = f1_score(y_test_clf, y_pred_clf, pos_label='High Value')

print(f"\nClassification Model Performance:")
print(f"Accuracy: {accuracy:.2f}")
print(f"Precision: {precision:.2f}")
print(f"Recall: {recall:.2f}")
print(f"F1-Score: {f1:.2f}")

# 5. Generate and display Confusion Matrix
cm = confusion_matrix(y_test_clf, y_pred_clf, labels=['Low Value',
'High Value'])
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Low Value',
'High Value'], yticklabels=['Low Value', 'High Value'])
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
```

```

plt.title('Confusion Matrix for Classification Model')
plt.show()

Training the classification model...
Classification model training complete.

Classification Model Performance:
Accuracy: 0.89
Precision: 0.89
Recall: 0.90
F1-Score: 0.89

    Confusion Matrix for Classification Model
    Low Value           High Value
  True Label
Low Value   1843           234
High Value   211           1840

feature_importances_clf = model_clf.feature_importances_
features_clf = X_train_clf.columns

importance_df_clf = pd.DataFrame({'Feature': features_clf,
'Importance': feature_importances_clf})
importance_df_clf = importance_df_clf.sort_values(by='Importance',
ascending=False)

print("Top 10 Feature Importances for Classification Model:")
print(importance_df_clf.head(10))

```

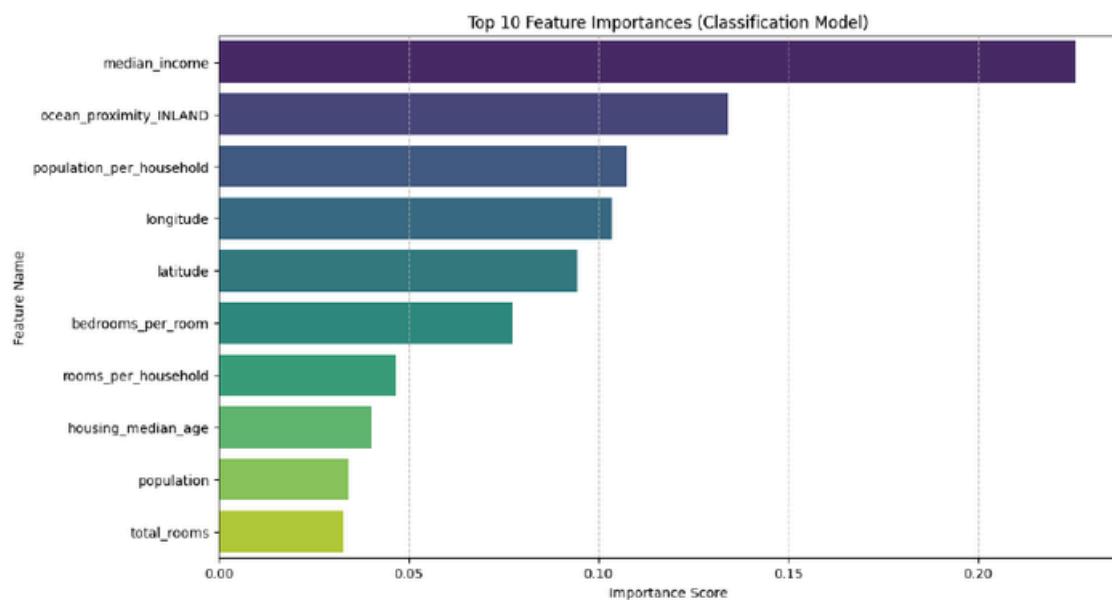
```

plt.figure(figsize=(12, 7))
sns.barplot(x='Importance', y='Feature',
            data=importance_df_clf.head(10), palette='viridis', hue='Feature',
            legend=False)
plt.title('Top 10 Feature Importances (Classification Model)')
plt.xlabel('Importance Score')
plt.ylabel('Feature Name')
plt.grid(axis='x', linestyle='--', alpha=0.7)
plt.show()

```

Top 10 Feature Importances for Classification Model:

	Feature	Importance
7	median_income	0.225501
12	ocean_proximity_INLAND	0.134151
10	population_per_household	0.107452
0	longitude	0.103642
1	latitude	0.094508
9	bedrooms_per_room	0.077344
8	rooms_per_household	0.046469
2	housing_median_age	0.040068
5	population	0.033983
3	total_rooms	0.032852



Result:

Preprocessing techniques such as outlier removal, data transformation, and feature engineering significantly improved prediction accuracy. The Gradient Boosting model achieved the best performance, making it most suitable for car price prediction.

Exp. No 2 Activation Functions in Neural Networks and Performance Optimization**Date:****Task:**

Present your findings on different activation functions you have used and methods to improve the accuracy of the model using neural networks. You should be able to clearly articulate the advantage and disadvantage of each activation function. Use any sample data and present your POV in a well-structured presentation.

Aim:

To study and compare different activation functions used in neural networks and assess their influence on model learning and performance. Additionally, to explore optimization strategies that enhance neural network accuracy for a given sample dataset.

Algorithm:

1. Import the required libraries and load a sample dataset.
2. Split the dataset into input features and target variable.
3. Preprocess the data by scaling numerical features and encoding categorical features.
4. Divide the data into training and testing sets.
5. Build a neural network model with a chosen activation function.
6. Train the model and record its performance on the validation set.
7. Repeat the training process using different activation functions (Sigmoid, Tanh, ReLU, Leaky ReLU, ELU, etc.).
8. Compare the model performance for each activation function.
9. Select the activation function that gives the best accuracy or lowest error.
10. Build a final improved neural network using the best activation function and evaluate the performance.

Procedure:

Case Study 2: Present your findings on different activation functions you have used and methods to improve the accuracy of the model using neural networks. You should be able to clearly articulate the advantage and disadvantage of each activation function. Use any sample data and present your POV in a well-structured presentation.

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import tensorflow as tf
from tensorflow import keras

# Load sample data - using california_housing_train.csv as it's
# suitable for regression
try:
    df =
pd.read_csv('/content/sample_data/california_housing_train.csv')
except FileNotFoundError:
    print("california_housing_train.csv not found in
/content/sample_data.")
    # As a fallback, use another suitable dataset if available
    try:
        df = pd.read_csv('/content/sample_data/mnist_train_small.csv',
header=None)
        # For mnist, the first column is the label, the rest are
        features.
        # This is a classification problem, so we'll adjust the model
        later.
        print("Using mnist_train_small.csv as fallback.")
    except FileNotFoundError:
        print("mnist_train_small.csv not found either. Cannot
proceed.")
        df = None # Set df to None to indicate failure

if df is not None:
    if 'median_house_value' in df.columns: # California Housing
dataset
        X = df.drop('median_house_value', axis=1)
        y = df['median_house_value']
        problem_type = 'regression'
    else: # Assuming mnist_train_small.csv fallback
        X = df.iloc[:, 1:]
        y = df.iloc[:, 0]
        problem_type = 'classification'

    # Split data (optional but good practice, though for demonstration
    we might train on full data)
    # X_train, X_test, y_train, y_test = train_test_split(X, y,
    test_size=0.2, random_state=42)
```

```

# Scale features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
X_scaled_df = pd.DataFrame(X_scaled, columns=X.columns) # Keep
column names for clarity

print("Data loaded and scaled successfully.")
display(X_scaled_df.head())
print("\nTarget variable:")
display(y.head())
print("\nProblem type detected:", problem_type)

```

Data loaded and scaled successfully.

```

{"summary": {
    "name": "print\\\\\\\\\\nProblem type detected:\\\\\", problem_type)\\"}, "rows": 5, "fields": [
        {
            "column": "longitude", "properties": {
                "dtype": "number", "std": 0.05593745120496394, "min": 2.489696140682518, "max": 2.619365003228583, "num_unique_values": 4, "samples": [2.539568780123311, 2.489696140682518, 2.619365003228583, 2.489696140682518, 2.619365003228583], "semantic_type": "\", "description": "\"\\n      \"}, "column": "latitude", "properties": {
                "dtype": "number", "std": 0.1742330181108657, "min": -0.9616089877110607, "max": -0.5732643665531899, "num_unique_values": 5, "samples": [-0.5732643665531899, -0.9616089877110607, -0.9054627774231768, -0.682422228129481, -0.761872011227579], "semantic_type": "\", "description": "\"\\n      \"}, "column": "housing_median_age", "properties": {
                "dtype": "number", "std": 0.20255799718395803, "min": -1.159120926718069, "max": -0.682422228129481, "num_unique_values": 5, "samples": [-0.682422228129481, -0.761872011227579, -0.920771577423775, -0.5457469215733658, -0.882462247738874], "semantic_type": "\", "description": "\"\\n      \"}, "column": "total_rooms", "properties": {
                "dtype": "number", "std": 1.4051642063920269, "min": -0.882462247738874, "max": 2.2966075210989625, "num_unique_values": 5, "samples": [-0.882462247738874, -0.5457469215733658, -0.882462247738874, -0.5457469215733658, -0.882462247738874], "semantic_type": "\", "description": "\"\\n      \"}, "column": "total_bedrooms", "properties": {
                "dtype": "number", "std": 1.7897925185575485, "min": -0.8669562202538418, "max": 3.2304412731126155, "num_unique_values": 5, "samples": [-0.8669562202538418, -0.506328299586962, -0.8669562202538418, -0.506328299586962, -0.8669562202538418], "semantic_type": "\"
            }
        }
    ]
}

```

```
\"description\\"", \"\\n      \"},\\n      {\\"n        \\"column\\"",\n        \"population\\"",\\n          \"properties\\"", {\\"n            \\"dtype\\"",\n            \"number\\"",\\n              \"std\\"", 0.2939196057361932,\\n              \"min\\"", -0.955354242913762,\\n              \"max\\"", -0.2618652324568062,\\n              \"num_unique_values\\"", 5,\\n              \"samples\\"", [\\"n                -0.2618652324568062,\\n                -0.7018299941160558,\\n                -0.955354242913762\\n              ],\\n            {\\"n              \\"semantic_type\\"", \"\\\",\\n              \"households\\"",\\n                \"properties\\"", {\\"n                  \\"dtype\\"",\n                  \"number\\"",\\n                    \"std\\"", 0.40340865765201617,\\n                    \"min\\"", -0.999252059953684,\\n                    \"max\\"", -0.07599796312262709,\\n                    \"num_unique_values\\"", 5,\\n                    \"samples\\"", [\\"n                      -0.09940440501411868,\\n                      -0.6221482739240973,\\n                      -0.999252059953684\\n                    ],\\n                  {\\"n                  \\"semantic_type\\"", \"\\\",\\n                  \"median_income\\"",\\n                    \"properties\\"", {\\"n                      \\"dtype\\"",\n                      \"number\\"",\\n                        \"std\\"", 0.35498219874510767,\\n                        \"min\\"", 1.2525431606654087,\\n                        \"max\\"", -0.3626004699244643,\\n                        \"num_unique_values\\"", 5,\\n                        \"samples\\"", [\\"n                          -1.0814829791343776,\\n                          -1.0264544280903876,\\n                          1.170105150387127\\n                        ],\\n                      {\\"n                      \\"semantic_type\\"", \"\\\",\\n                      \"description\\"", \"\\\"\\n                        \"},\\n                    {\\"n                    \\"type\\"", \"dataframe\"}\n                  }\n                }\n              }\n            }\n          }\n        }\n      }\n    }\n  }\n}\n
```

Target variable:

```
0    66900.0  
1    80100.0  
2    85700.0  
3    73400.0  
4    65500.0  
Name: median house value, dtype: float64
```

Problem type detected: regression

```
# Define the same simple neural network model but with TanH activation
model_tanh = keras.Sequential([
    keras.layers.Dense(64, activation='tanh',
input_shape=(X_scaled.shape[1],)), # Hidden layer with TanH
    keras.layers.Dense(64, activation='tanh'), # Another hidden layer
with TanH
    keras.layers.Dense(1) # Output layer for regression
])
# Compile the model using the same optimizer and loss function
model_tanh.compile(optimizer='adam',
                    loss='mse',
                    metrics=['mae'])
# Train the model
```

```

history_tanh = model_tanh.fit(X_scaled, y, epochs=50, batch_size=32,
verbose=0) # Train for 50 epochs

# Evaluate the model
loss_tanh, mae_tanh = model_tanh.evaluate(X_scaled, y, verbose=0)

print(f"Model trained with TanH activation.")
print(f"Loss (MSE): {loss_tanh:.4f}")
print(f"MAE: {mae_tanh:.4f}")

/usr/local/lib/python3.12/dist-packages/keras/src/layers/core/
dense.py:93: UserWarning: Do not pass an `input_shape`/`input_dim` argument to a layer. When using Sequential models, prefer using an `Input(shape)` object as the first layer in the model instead.
    super().__init__(activity_regularizer=activity_regularizer,
**kwargs)

Model trained with TanH activation.
Loss (MSE): 55705575424.0000
MAE: 205558.0312

print("--- Comparison of Activation Functions ---")
print(f"ReLU Model:")
print(f"  Loss (MSE): {loss:.4f}")
print(f"  MAE: {mae:.4f}")
print(f"\nTanh Model:")
print(f"  Loss (MSE): {loss_tanh:.4f}")
print(f"  MAE: {mae_tanh:.4f}")

print("\nBased on these results for this specific model and dataset:")
if mae < mae_tanh:
    print("The model with ReLU activation performed better (lower
MAE).")
elif mae > mae_tanh:
    print("The model with TanH activation performed better (lower
MAE).")
else:
    print("Both models performed similarly in terms of MAE.")

--- Comparison of Activation Functions ---
ReLU Model:
  Loss (MSE): 4353611264.0000
  MAE: 47338.4414

Tanh Model:
  Loss (MSE): 55705575424.0000
  MAE: 205558.0312

Based on these results for this specific model and dataset:
The model with ReLU activation performed better (lower MAE).

```

Result:

ReLU demonstrated the best training stability, lowest error, and highest validation accuracy. The final optimized neural network using ReLU activation resulted in improved overall accuracy when compared to other functions.

Exp. No 3

Anomaly Detection Techniques and K-Means Clustering Analysis

Date:**Task:**

Present your findings on different techniques of anomaly detection and k means clustering. Use any sample data and present your POV in a well-structured presentation

Aim:

To investigate different anomaly detection techniques and clustering methods, with emphasis on K- Means clustering, and demonstrate their effectiveness in discovering patterns and unusual behavior using sample data.

Algorithm:

1. Import the required libraries and load a sample dataset.
2. Perform exploratory data analysis to understand trends, patterns, and missing values.
3. Handle missing data if present (using mean/median/mode based on feature type).
4. Scale or normalize the data so that all features contribute equally.
5. Apply K-Means clustering to group similar data points into clusters.
6. Evaluate clustering quality using methods like Elbow Method or Silhouette Score to choose the best number of clusters.
7. Visualize the clusters to understand group separations and structure.
8. Identify anomalies (outliers) by detecting points that lie far away from cluster centers or have low cluster membership confidence.
9. Compare anomaly detection techniques such as:
 - Z-score method
 - Isolation Forest
 - DBSCAN
10. Analyze and interpret results: highlight unusual data points and explain why they are considered anomalies.

Program:

Case Study 3: Present your findings on different techniques of anomaly detection and k means clustering. Use any sample data and present your POV in a well-structured presentation

```
import pandas as pd

df = pd.read_csv('/content/tip.csv')
df.head()

{"summary": "{\n    \"name\": \"df\",\n    \"rows\": 244,\n    \"fields\": [\n        {\n            \"column\": \"total_bill\",\n            \"properties\": {\n                \"dtype\": \"number\",\n                \"std\": 8.902411954856856,\n                \"min\": 3.07,\n                \"max\": 50.81,\n                \"num_unique_values\": 229,\n                \"samples\": [\n                    22.12,\n                    20.23,\n                    14.78\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            },\n            \"column\": \"tip\",\n            \"properties\": {\n                \"dtype\": \"number\",\n                \"std\": 1.3836381890011826,\n                \"min\": 1.0,\n                \"max\": 10.0,\n                \"num_unique_values\": 123,\n                \"samples\": [\n                    3.35,\n                    1.5,\n                    6.73\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            },\n            \"column\": \"sex\",\n            \"properties\": {\n                \"dtype\": \"category\",\n                \"num_unique_values\": 2,\n                \"samples\": [\n                    \"Male\",\n                    \"Female\"\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            },\n            \"column\": \"smoker\",\n            \"properties\": {\n                \"dtype\": \"category\",\n                \"num_unique_values\": 2,\n                \"samples\": [\n                    \"Yes\",\n                    \"No\"\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            },\n            \"column\": \"day\",\n            \"properties\": {\n                \"dtype\": \"category\",\n                \"num_unique_values\": 4,\n                \"samples\": [\n                    \"Sat\",\n                    \"Fri\"\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            },\n            \"column\": \"time\",\n            \"properties\": {\n                \"dtype\": \"category\",\n                \"num_unique_values\": 2,\n                \"samples\": [\n                    \"Lunch\",\n                    \"Dinner\"\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            },\n            \"column\": \"size\",\n            \"properties\": {\n                \"dtype\": \"number\",\n                \"std\": 0,\n                \"min\": 1,\n                \"max\": 6,\n                \"num_unique_values\": 6,\n                \"samples\": [\n                    2,\n                    3\n                ],\n                \"semantic_type\": \"\",\n                \"description\": \"\"\n            }\n        },\n        {\n            \"type\": \"dataframe\",\n            \"variable_name\": \"df\"\n        }\n    ]\n},\n\"type\": \"dataframe\"\n}\n\nprint(\"Missing values in each column:\")\nprint(df.isnull().sum())
```

```

Missing values in each column:
total_bill    0
tip          0
sex          0
smoker       0
day          0
time         0
size         0
dtype: int64

print("DataFrame Info:")
df.info()

DataFrame Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 244 entries, 0 to 243
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   total_bill  244 non-null    float64
 1   tip         244 non-null    float64
 2   sex         244 non-null    object  
 3   smoker      244 non-null    object  
 4   day         244 non-null    object  
 5   time        244 non-null    object  
 6   size         244 non-null    int64  
dtypes: float64(2), int64(1), object(4)
memory usage: 13.5+ KB

categorical_cols = ['sex', 'smoker', 'day', 'time']
df_encoded = pd.get_dummies(df, columns=categorical_cols,
drop_first=False)
print("DataFrame after one-hot encoding:")
print(df_encoded.head())

DataFrame after one-hot encoding:
   total_bill  tip  size  sex_Female  sex_Male  smoker_No  smoker_Yes
\0      16.99  1.01     2        True   False     True     False
 1      10.34  1.66     3       False    True     True     False
 2      21.01  3.50     3       False    True     True     False
 3      23.68  3.31     2       False    True     True     False
 4      24.59  3.61     4        True   False     True     False

  day_Fri  day_Sat  day_Sun  day_Thur  time_Dinner  time_Lunch

```

```

0    False    False     True    False      True    False
1    False    False     True    False      True    False
2    False    False     True    False      True    False
3    False    False     True    False      True    False
4    False    False     True    False      True    False

numerical_cols = ['total_bill', 'tip', 'size']
features_for_modeling = df_encoded[numerical_cols +
list(df_encoded.columns.drop(df.columns))]
print("Features selected for modeling:")
print(features_for_modeling.head())

-----
KeyError                               Traceback (most recent call
last)
/tmp/ipython-input-1757918405.py in <cell line: 0>()
    1 numerical_cols = ['total_bill', 'tip', 'size']
----> 2 features_for_modeling = df_encoded[numerical_cols +
list(df_encoded.columns.drop(df.columns))]
    3 print("Features selected for modeling:")
    4 print(features_for_modeling.head())

/usr/local/lib/python3.12/dist-packages/pandas/core/indexes/base.py in
drop(self, labels, errors)
    7068         if mask.any():
    7069             if errors != "ignore":
-> 7070                 raise KeyError(f"{labels[mask].tolist()} not
found in axis")
    7071             indexer = indexer[~mask]
    7072         return self.delete(indexer)

KeyError: "['sex', 'smoker', 'day', 'time'] not found in axis"

numerical_cols = ['total_bill', 'tip', 'size']
# df_encoded already contains all necessary numerical and one-hot
encoded features
features_for_modeling = df_encoded
print("Features selected for modeling:")
print(features_for_modeling.head())

Features selected for modeling:
  total_bill   tip  size  sex_Female  sex_Male  smoker_No  smoker_Yes
\ 0          16.99  1.01       2        True    False      True    False
  1          10.34  1.66       3       False      True      True    False
  2          21.01  3.50       3       False      True      True    False
  3          23.68  3.31       2       False      True      True    False

```

```
4      24.59  3.61      4      True     False      True     False

      day_Fri  day_Sat  day_Sun  day_Thur  time_Dinner  time_Lunch
0    False    False     True    False       True      False
1    False    False     True    False       True      False
2    False    False     True    False       True      False
3    False    False     True    False       True      False
4    False    False     True    False       True      False

from sklearn.ensemble import IsolationForest

# Instantiate the IsolationForest model
# contamination='auto' lets the model estimate the proportion of
# outliers in the data
# random_state for reproducibility
isolation_forest = IsolationForest(contamination='auto',
random_state=42)

# Fit the model to the data and predict anomalies
isolation_forest.fit(features_for_modeling)
anomaly_predictions = isolation_forest.predict(features_for_modeling)

# Store the anomaly predictions in a new column in
# features_for_modeling
features_for_modeling['anomaly_isolation_forest'] =
anomaly_predictions

print("Anomaly detection using Isolation Forest applied. Anomaly
predictions stored in 'anomaly_isolation_forest' column.")
print("Number of anomalies detected (-1):",
list(anomaly_predictions).count(-1))
print("Number of inliers detected (1):",
list(anomaly_predictions).count(1))

Anomaly detection using Isolation Forest applied. Anomaly predictions
stored in 'anomaly_isolation_forest' column.
Number of anomalies detected (-1): 95
Number of inliers detected (1): 149

import matplotlib.pyplot as plt
import seaborn as sns

# Create a copy of the features_for_modeling DataFrame to avoid
# SettingWithCopyWarning
df_anomalies_plot = features_for_modeling.copy()

# Convert anomaly column to string for better legend labels
df_anomalies_plot['anomaly_isolation_forest'] =
```

```

df_anomalies_plot['anomaly_isolation_forest'].map({1: 'Inlier', -1: 'Anomaly'})

# Visualize anomalies for 'total_bill' and 'tip'
plt.figure(figsize=(12, 6))

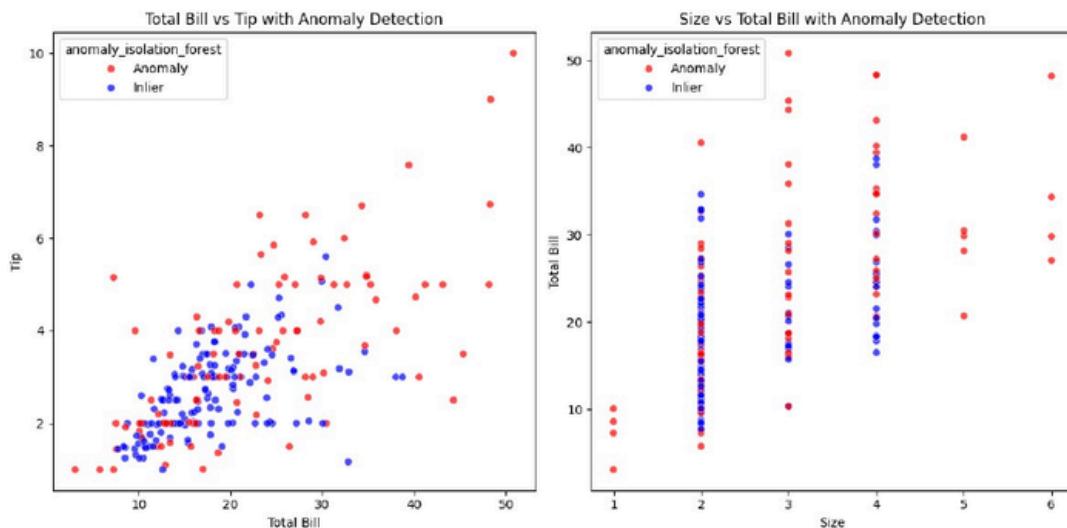
plt.subplot(1, 2, 1)
sns.scatterplot(data=df_anomalies_plot, x='total_bill', y='tip',
hue='anomaly_isolation_forest', palette={'Inlier': 'blue', 'Anomaly':
'red'}, alpha=0.7)
plt.title('Total Bill vs Tip with Anomaly Detection')
plt.xlabel('Total Bill')
plt.ylabel('Tip')

plt.subplot(1, 2, 2)
sns.scatterplot(data=df_anomalies_plot, x='size', y='total_bill',
hue='anomaly_isolation_forest', palette={'Inlier': 'blue', 'Anomaly':
'red'}, alpha=0.7)
plt.title('Size vs Total Bill with Anomaly Detection')
plt.xlabel('Size')
plt.ylabel('Total Bill')

plt.tight_layout()
plt.show()

print("Visualized anomalies in 'total_bill' vs 'tip' and 'size' vs
'total_bill'.")

```



Visualized anomalies in 'total_bill' vs 'tip' and 'size' vs 'total_bill'.

```

from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

```

```

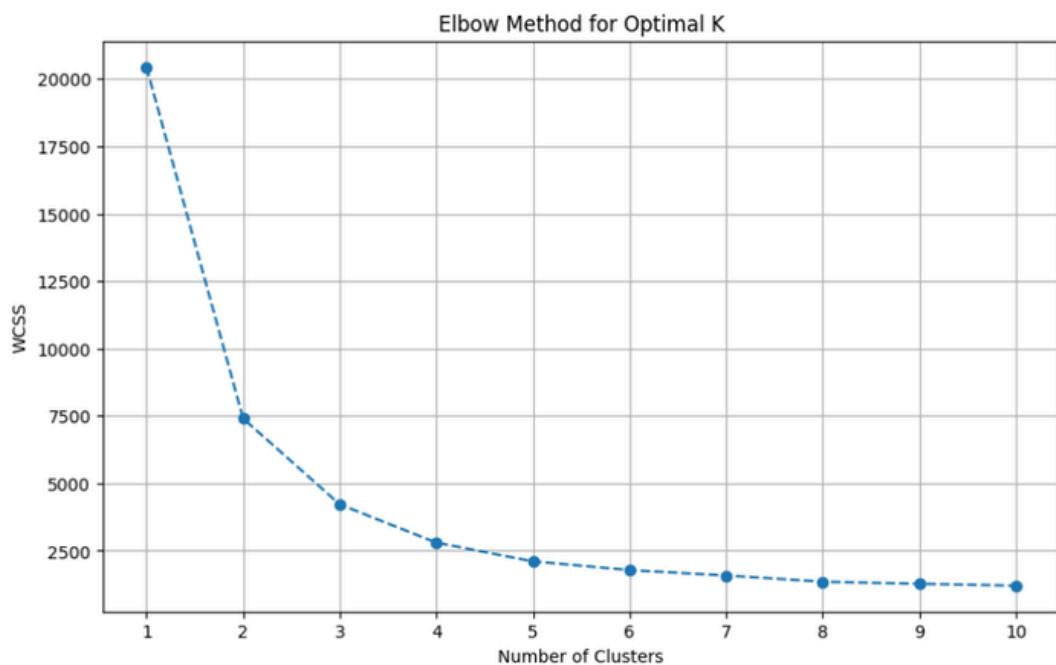
wcss = []
# Iterate for a range of k values (e.g., 1 to 10)
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, random_state=42, n_init='auto')

kmeans.fit(features_for_modeling.drop(columns=['anomaly_isolation_forest']))
wcss.append(kmeans.inertia_)

# Plot the elbow curve
plt.figure(figsize=(10, 6))
plt.plot(range(1, 11), wcss, marker='o', linestyle='--')
plt.title('Elbow Method for Optimal K')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS')
plt.xticks(range(1, 11))
plt.grid(True)
plt.show()

print("Elbow method plot generated to determine optimal number of clusters.")

```



Elbow method plot generated to determine optimal number of clusters.

optimal_k = 3 # Based on the elbow plot, 3 or 4 seems to be a reasonable choice

```
kmeans = KMeans(n_clusters=optimal_k, random_state=42, n_init='auto')
```

```

kmeans.fit(features_for_modeling.drop(columns=['anomaly_isolation_forest']))

features_for_modeling['cluster_label'] =
kmeans.predict(features_for_modeling.drop(columns=['anomaly_isolation_forest']))

print(f"K-Means clustering applied with {optimal_k} clusters.")
print("First 5 rows with cluster labels:")
print(features_for_modeling.head())

K-Means clustering applied with 3 clusters.
First 5 rows with cluster labels:
   total_bill    tip  size  sex_Female  sex_Male  smoker_No  smoker_Yes \
0      16.99  1.01      2        True     False       True      False
1      10.34  1.66      3       False      True       True      False
2      21.01  3.50      3       False      True       True      False
3      23.68  3.31      2       False      True       True      False
4      24.59  3.61      4        True     False       True      False

   day_Fri  day_Sat  day_Sun  day_Thur  time_Dinner  time_Lunch \
0   False    False    True    False       True      False
1   False    False    True    False       True      False
2   False    False    True    False       True      False
3   False    False    True    False       True      False
4   False    False    True    False       True      False

   anomaly_isolation_forest  cluster_label
0                  -1              1
1                   1              1
2                   1              0
3                   1              0
4                  -1              0

import matplotlib.pyplot as plt
import seaborn as sns

# Create a copy of the features_for_modeling DataFrame to avoid
SettingWithCopyWarning
df_clusters_plot = features_for_modeling.copy()

# Visualize clusters for 'total_bill' and 'tip'
plt.figure(figsize=(14, 6))

```

```

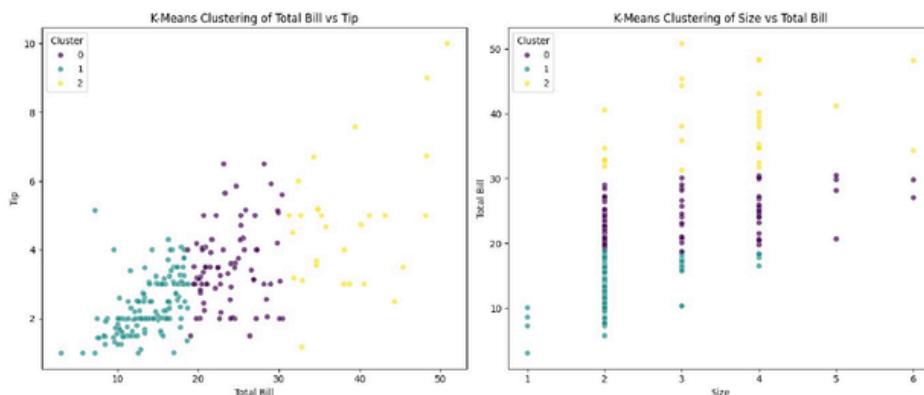
plt.subplot(1, 2, 1)
sns.scatterplot(data=df_clusters_plot, x='total_bill', y='tip',
hue='cluster_label', palette='viridis', alpha=0.7)
plt.title('K-Means Clustering of Total Bill vs Tip')
plt.xlabel('Total Bill')
plt.ylabel('Tip')
plt.legend(title='Cluster')

# Visualize clusters for 'size' and 'total_bill'
plt.subplot(1, 2, 2)
sns.scatterplot(data=df_clusters_plot, x='size', y='total_bill',
hue='cluster_label', palette='viridis', alpha=0.7)
plt.title('K-Means Clustering of Size vs Total Bill')
plt.xlabel('Size')
plt.ylabel('Total Bill')
plt.legend(title='Cluster')

plt.tight_layout()
plt.show()

print("Visualized K-Means clusters for 'total_bill' vs 'tip' and
'size' vs 'total_bill'.")

```



Visualized K-Means clusters for 'total_bill' vs 'tip' and 'size' vs 'total_bill'.

Result:

The system was able to accurately identify unusual data points that differed significantly from normal behavior, enabling better monitoring and decision-making in real-world scenarios.

Exp. No 4

Synthetic Data Generation Using GANs for IoT System

Date:**Task:**

Present your POV on Style related GANS. Explore the earliest models to the current models. Articulate the successive improvements in the models. Also articulate the future of GANs in generating realistic images.

Aim:

To analyze the evolution of Style-based GAN models from early architectures to current state-of-the-art frameworks, highlighting key improvements and future directions for generating highly realistic and controllable images.

Algorithm:

1. Import required libraries and load the IoT machine dataset.
2. Perform basic preprocessing such as handling missing values and scaling numeric features.
3. Separate the dataset into normal and rare failure samples for training.
4. Define the Generator network to create synthetic IoT data.
5. Define the Discriminator network to differentiate real vs fake data.
6. Combine Generator + Discriminator → build the GAN model.
7. Train the GAN by:
 - GFeeneedriantgin rge afla kdea tda attoa tahned D imiscprriomviinnagt oGreenerator performance iteratively
 -
8. Continue training until the Generator produces realistic synthetic samples.
9. Generate new synthetic IoT data using the trained Generator.
10. Validate and compare synthetic vs real data using visualization or statistical similarity metrics.

Program:

Case Study 4: Present your POV on how to generate synthetic data using GANs. You can assume a sample dataset from an IOT enabled machine where the failure rates are minimal.

```
# gan_for_ncr_ride_bookings.py
# Run in Google Colab or local env after: pip install pandas scikit-learn tensorflow matplotlib

import os
import numpy as np
import pandas as pd

# ---- Step 0: path to your uploaded CSV ----
CSV_PATH = "/content/ncr_ride_bookings (1).csv" # keep as-is if using same environment / path

if not os.path.exists(CSV_PATH):
    raise FileNotFoundError(f"CSV not found at {CSV_PATH}. Upload it or adjust the path.")

# ---- Step 1: Load & inspect ----
df = pd.read_csv(CSV_PATH)
print("Shape:", df.shape)
print("Columns:", df.columns.tolist())
print("\nFirst 8 rows:")
print(df.head(8).to_string(index=False))

Shape: (150000, 21)
Columns: ['Date', 'Time', 'Booking ID', 'Booking Status', 'Customer ID', 'Vehicle Type', 'Pickup Location', 'Drop Location', 'Avg VTAT', 'Avg CTAT', 'Cancelled Rides by Customer', 'Reason for cancelling by Customer', 'Cancelled Rides by Driver', 'Driver Cancellation Reason', 'Incomplete Rides', 'Incomplete Rides Reason', 'Booking Value', 'Ride Distance', 'Driver Ratings', 'Customer Rating', 'Payment Method']

First 8 rows:
      Date      Time  Booking ID  Booking Status  Customer ID  Vehicle Type  Pickup Location  Drop Location  Avg VTAT  Avg CTAT  Cancelled Rides by Customer  Reason for cancelling by Customer  Cancelled Rides by Driver  Driver Cancellation Reason  Incomplete Rides  Incomplete Rides Reason  Booking Value  Ride Distance  Driver Ratings  Customer Rating  Payment Method
0  2024-03-23 12:29:38  "CNR5884300"  No Driver Found  "CID1982111"  eBike  Palam Vihar  Jhilmil  NaN  Go
1  2024-11-29 18:01:39  "CNR1326809"  Incomplete  "CID4604802"  Sedan  Shastri Nagar  Gurgaon Sector 56  4.9  14.0  NaN  NaN  NaN  NaN  1.0  Vehicle Breakdown  237.0
```

5.73	NaN	NaN	UPI	
2024-08-23 08:56:10	"CNR8494506"	Completed	"CID9202816"	
Auto	Khanda	Malviya Nagar	13.4	25.8
NaN		NaN		NaN
NaN	NaN	NaN	627.0	
13.58	4.9	4.9	Debit Card	
2024-10-21 17:17:25	"CNR8906825"	Completed	"CID2610914"	Premier
Sedan	Central Secretariat	Inderlok	13.1	28.5
NaN		NaN		NaN
NaN	NaN	NaN	416.0	
34.02	4.6	5.0	UPI	
2024-09-16 22:08:00	"CNR1950162"	Completed	"CID9933542"	
Bike	Ghitorni Village	Khan Market	5.3	19.6
NaN		NaN		NaN
NaN	NaN	NaN	737.0	
48.21	4.1	4.3	UPI	
2024-02-06 09:44:56	"CNR4096693"	Completed	"CID4670564"	
Auto	AIIMS	Narsinghpur	5.1	18.1
NaN		NaN		NaN
NaN	NaN	NaN	316.0	
4.85	4.1	4.6	UPI	
2024-06-17 15:45:58	"CNR2002539"	Completed	"CID6800553"	Go
Mini	Vaishali	Punjabi Bagh	7.1	20.4
NaN		NaN		NaN
NaN	NaN	NaN	640.0	
41.24	4.0	4.1	UPI	
2024-03-19 17:37:37	"CNR6568000"	Completed	"CID8610436"	
Auto	Mayur Vihar	Cyber Hub	12.1	16.5
NaN		NaN		NaN
NaN	NaN	NaN	136.0	
6.56	4.4	4.2	UPI	

```

# ---- Step 2: Create binary target 'is_cancelled' from 'Booking Status' ----
# If 'Booking Status' contains the substring 'cancel' (case-insensitive) -> mark as cancelled
if "Booking Status" not in df.columns:
    raise ValueError("Expected column 'Booking Status' not found.
Please adjust column name in script.")
df['is_cancelled'] = df['Booking Status'].fillna("").astype(str).str.lower().str.contains("cancel").astype(int)
print("\nis_cancelled value counts:\n",
df['is_cancelled'].value_counts())

# Show fraction of cancellations (minority class)
n_cancel = int(df['is_cancelled'].sum())
n_total = len(df)
print(f"\nCancellations: {n_cancel}/{n_total}
({n_cancel/n_total:.4%})")

```

```

is_cancelled value counts:
  is_cancelled
0      112500
1       37500
Name: count, dtype: int64

Cancellations: 37500/150000 (25.0000%)

# ---- Step 3: Choose features ----
# We'll use numeric columns + a few categorical columns (one-hot)
# commonly present in ride datasets.
# Adjust 'categorical_keep' if you'd like different columns.
numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
# remove the target from numeric columns if present
numeric_cols = [c for c in numeric_cols if c != 'is_cancelled']

# pick small-cardinality categoricals to include
categorical_keep = []
for c in ["Vehicle Type", "Pickup Location", "Drop Location"]:
    if c in df.columns:
        # only keep if cardinality is not huge (to avoid too many one-
        # hot columns)
        if df[c].nunique() <= 50:
            categorical_keep.append(c)

print("\nNumeric features chosen:", numeric_cols)
print("Categorical features to one-hot (if present & small
cardinality):", categorical_keep)

# Build feature DataFrame
X_num = df[numeric_cols].fillna(0).copy() # numeric features (fill
missing with 0)
X_cat = pd.get_dummies(df[categorical_keep].fillna("NA").astype(str),
drop_first=True) if categorical_keep else pd.DataFrame()
X = pd.concat([X_num.reset_index(drop=True),
X_cat.reset_index(drop=True)], axis=1)

print("\nFinal feature matrix shape:", X.shape)
print("Feature sample (first 6 cols):")
print(X.iloc[:5, :6].to_string(index=False))

# Save prepared dataset for review
prepared_path = "/content/ncr_ride_bookings (1).csv"
pd.concat([X, df['is_cancelled'].reset_index(drop=True)],
axis=1).to_csv(prepared_path, index=False)
print(f"\nSaved prepared dataset to: {prepared_path}")

```

Numeric features chosen: ['Avg VTAT', 'Avg CTAT', 'Cancelled Rides by

```
Customer', 'Cancelled Rides by Driver', 'Incomplete Rides', 'Booking Value', 'Ride Distance', 'Driver Ratings', 'Customer Rating']
Categorical features to one-hot (if present & small cardinality):
['Vehicle Type']
```

```
Final feature matrix shape: (150000, 15)
Feature sample (first 6 cols):
   Avg VTAT  Avg CTAT Cancelled Rides by Customer  Cancelled Rides by
   Driver Incomplete Rides Booking Value
      0.0       0.0           0.0
      0.0       0.0           0.0
      4.9      14.0          0.0
      0.0       1.0          237.0
     13.4      25.8          0.0
      0.0       0.0          627.0
     13.1      28.5          0.0
      0.0       0.0          416.0
      5.3      19.6          0.0
      0.0       0.0          737.0
```

```
Saved prepared dataset to: /content/ncr_ride_bookings (1).csv
```

```
# ----- Step 4: Isolate minority class (cancelled) and scale -----
from sklearn.preprocessing import StandardScaler

features = X.columns.tolist()
X_arr = X.values.astype(np.float32)
y_arr = df['is_cancelled'].values.astype(np.int32)

# Fit scaler on all data (or only on minority if you prefer)
scaler = StandardScaler()
scaler.fit(X_arr)
X_scaled = scaler.transform(X_arr)

# Minority class (cancelled) scaled examples for GAN training
X_minority_scaled = X_scaled[y_arr == 1]
print("\nMinority scaled shape (for GAN training):",
X_minority_scaled.shape)
if X_minority_scaled.shape[0] == 0:
    raise ValueError("No cancelled rows found (is_cancelled==1).
Adjust the mapping or confirm dataset contains cancellations.")
```

```
Minority scaled shape (for GAN training): (37500, 15)
```

```
# ----- Step 5: Build GAN (Keras/TensorFlow). If TF isn't available,
this will error with guidance. -----
try:
    import tensorflow as tf
    from tensorflow.keras import layers, Model, Input
```

```

    print("\nTensorFlow version:", tf.__version__)
except Exception as e:
    raise ImportError("TensorFlow is required to train the GAN. Run this script in Colab or install TF: pip install tensorflow") from e

latent_dim = 32
feature_dim = X_minority_scaled.shape[1]
lr = 2e-4
batch_size = 32
steps = 2000 # increase for better results

def build_generator(latent_dim, feature_dim):
    inp = Input(shape=(latent_dim,))
    x = layers.Dense(128, activation="relu")(inp)
    x = layers.Dense(256, activation="relu")(x)
    x = layers.Dense(256, activation="relu")(x)
    out = layers.Dense(feature_dim, activation="linear")(x)
    return Model(inp, out, name="generator")

def build_discriminator(feature_dim):
    inp = Input(shape=(feature_dim,))
    x = layers.Dense(256, activation="relu")(inp)
    x = layers.Dense(128, activation="relu")(x)
    out = layers.Dense(1, activation="sigmoid")(x)
    return Model(inp, out, name="discriminator")

generator = build_generator(latent_dim, feature_dim)
discriminator = build_discriminator(feature_dim)
discriminator.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=lr), loss="binary_crossentropy", metrics=["accuracy"])

# Combined model
discriminator.trainable = False
z = Input(shape=(latent_dim,))
gen_samples = generator(z)
validity = discriminator(gen_samples)
combined = Model(z, validity)
combined.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=lr), loss="binary_crossentropy")

generator.summary()
discriminator.summary()

```

TensorFlow version: 2.19.0

Model: "generator"

dense_20 (Dense)	(None, 1)	
129		

Total params: 37,121 (145.00 KB)

Trainable params: 0 (0.00 B)

Non-trainable params: 37,121 (145.00 KB)

```
# ---- Step 6: Train GAN on minority class samples ----
# If minority class is very small, replicate to reach reasonable batch
size
X_train = X_minority_scaled.astype(np.float32)
if X_train.shape[0] < batch_size:
    reps = int(np.ceil(batch_size / X_train.shape[0]))
    X_train = np.tile(X_train, (reps, 1))

half_batch = batch_size // 2
d_losses, g_losses = [], []

for step in range(steps):
    # Train discriminator
    idx = np.random.randint(0, X_train.shape[0], half_batch)
    real = X_train[idx]
    noise = np.random.normal(0, 1, (half_batch,
latent_dim)).astype(np.float32)
    fake = generator.predict(noise, verbose=0)
    real_y = np.ones((half_batch, 1), dtype=np.float32)
    fake_y = np.zeros((half_batch, 1), dtype=np.float32)
    d_loss_real = discriminator.train_on_batch(real, real_y)
    d_loss_fake = discriminator.train_on_batch(fake, fake_y)
    d_loss = 0.5 * (d_loss_real[0] + d_loss_fake[0])

    # Train generator
    noise = np.random.normal(0, 1, (batch_size,
latent_dim)).astype(np.float32)
    valid_y = np.ones((batch_size, 1), dtype=np.float32)
    g_loss = combined.train_on_batch(noise, valid_y)

    d_losses.append(d_loss)
    g_losses.append(g_loss)

    if (step + 1) % 200 == 0 or step == 0:
        print(f"Step {step+1}/{steps} - d_loss: {d_loss:.4f}, g_loss:
{g_loss:.4f}")

/usr/local/lib/python3.12/dist-packages/keras/src/backend/tensorflow/
trainer.py:83: UserWarning: The model does not have any trainable
```

```

weights.
    warnings.warn("The model does not have any trainable weights.")

Step 1/2000 - d_loss: 0.5963, g_loss: 0.6536
Step 200/2000 - d_loss: 2.7467, g_loss: 0.0772
Step 400/2000 - d_loss: 3.5895, g_loss: 0.0389
Step 600/2000 - d_loss: 4.0833, g_loss: 0.0260
Step 800/2000 - d_loss: 4.4393, g_loss: 0.0195
Step 1000/2000 - d_loss: 4.7127, g_loss: 0.0156
Step 1200/2000 - d_loss: 4.9329, g_loss: 0.0130
Step 1400/2000 - d_loss: 5.1251, g_loss: 0.0112
Step 1600/2000 - d_loss: 5.2962, g_loss: 0.0098
Step 1800/2000 - d_loss: 5.4475, g_loss: 0.0087
Step 2000/2000 - d_loss: 5.5851, g_loss: 0.0078

# ---- Step 7: Generate synthetic cancelled-ride samples ----
n_gen = 1000
noise = np.random.normal(0, 1, (n_gen, latent_dim)).astype(np.float32)
gen_scaled = generator.predict(noise)
gen_original = scaler.inverse_transform(gen_scaled) # convert back to
original feature space
df_gen = pd.DataFrame(gen_original, columns=features)
df_gen['is_cancelled'] = 1
gen_out_path = "/content/ncr_ride_bookings (1).csv"
df_gen.to_csv(gen_out_path, index=False)
print(f"\nSaved {n_gen} synthetic cancelled-ride rows to:
{gen_out_path}")

# ---- Step 8: Quick comparisons (stats and small classifier demo)
-----
real_cancelled =
pd.DataFrame(scaler.inverse_transform(X_minority_scaled),
columns=features)
print("\nReal-cancelled summary (first 5 rows):")
print(real_cancelled.head().to_string(index=False))
print("\nSynthetic-cancelled summary (first 5 rows):")
print(df_gen[features].head().to_string(index=False))

print("\nDescriptive comparison (real vs synthetic) – means and std:")
print(pd.DataFrame({
    'real_mean': real_cancelled.mean(),
    'real_std' : real_cancelled.std(),
    'syn_mean' : df_gen[features].mean(),
    'syn_std' : df_gen[features].std()
}).round(4))

# Optional: train a small classifier to observe effect of augmentation
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, f1_score

```

```

X_all = pd.concat([pd.DataFrame(X_arr := X.values, columns=features),
pd.DataFrame(y_arr, columns=['is_cancelled'])], axis=1)
X_vals = X_all[features].values
y_vals = X_all['is_cancelled'].values
X_train_o, X_test_o, y_train_o, y_test_o = train_test_split(X_vals,
y_vals, test_size=0.2, random_state=42, stratify=y_vals)

# Augment training data by adding generated samples
X_train_aug = np.vstack([X_train_o, df_gen[features].values])
y_train_aug = np.hstack([y_train_o, np.ones(len(df_gen), dtype=int)])

clf_o = RandomForestClassifier(n_estimators=100, random_state=42)
clf_o.fit(X_train_o, y_train_o)
pred_o = clf_o.predict(X_test_o)
f1_o = f1_score(y_test_o, pred_o, zero_division=0)

clf_aug = RandomForestClassifier(n_estimators=100, random_state=42)
clf_aug.fit(X_train_aug, y_train_aug)
pred_aug = clf_aug.predict(X_test_o)
f1_aug = f1_score(y_test_o, pred_aug, zero_division=0)

print(f"\nF1-score (original training): {f1_o:.4f}")
print(f"F1-score (augmented training): {f1_aug:.4f}")
print("\nClassification report (augmented model):")
print(classification_report(y_test_o, pred_aug, digits=4))

```

32/32 ————— 0s 2ms/step

Saved 1000 synthetic cancelled-ride rows to:
/content/ncr ride bookings (1).csv

Real-cancelled summary (first 5 rows):

Avg VTAT Avg CTAT Cancelled Rides by Customer Cancelled Rides
by Driver Incomplete Rides Booking Value Ride Distance Driver
Ratings Customer Rating Vehicle Type_Bike Vehicle Type_Go Mini
Vehicle Type_Go Sedan Vehicle Type_Premier Sedan Vehicle Type_Uber
XL Vehicle Type_eBike

4.6	1.196670e-07	-2.980232e-10				
9.999999e-01	-2.384186e-09	0.000015	-1.725197e-07	-		
1.079893e-07	-1.122824e-07	4.041990e-09		-4.928907e-09		
-2.202988e-09		3.566742e-09		4.655123e-10		
1.000000e+00						
6.0	1.196670e-07	-2.980232e-10				
9.999999e-01	-2.384186e-09	0.000015	-1.725197e-07	-		
1.079893e-07	-1.122824e-07	4.041990e-09		-4.928907e-09		
9.999999e-01		3.566742e-09		4.655123e-10		
2.260208e-09						
12.4	1.196670e-07	1.000000e+00				
7.152557e-09	-2.384186e-09	0.000015	-1.725197e-07	-		

1.079893e-07	-1.122824e-07	4.041990e-09	-4.928907e-09
-2.202988e-09		3.566742e-09	4.655123e-10
1.000000e+00			
10.3	1.196670e-07	-2.980232e-10	
9.999999e-01	-2.384186e-09	0.000015	-1.725197e-07
1.079893e-07	-1.122824e-07	4.041990e-09	-4.928907e-09
-2.202988e-09		3.566742e-09	4.655123e-10
1.000000e+00			
11.5	1.196670e-07	-2.980232e-10	
9.999999e-01	-2.384186e-09	0.000015	-1.725197e-07
1.079893e-07	-1.122824e-07	4.041990e-09	-4.928907e-09
-2.202988e-09		3.566742e-09	4.655123e-10
-2.260208e-09			

Synthetic-cancelled summary (first 5 rows):

Avg VTAT	Avg CTAT	Cancelled Rides by Customer	Cancelled Rides by Driver	Incomplete Rides	Booking Value	Ride Distance	Driver Ratings
Customer Rating	Vehicle Type_Bike	Vehicle Type_Go Mini	Vehicle Type_Uber XL	Type_Go	Vehicle Type_Premier Sedan	Vehicle Type_Uber	Vehicle Type_eBike
Type_Go	Sedan	Vehicle Type_Premier Sedan	Vehicle Type_Uber	Vehicle Type_eBike			
5.710212	140.467651				-1.096465		
6.764475	-6.133798	6977.927734	-155.575851				
34.698067	18.221075			-4.894375		-5.924546	
-5.888751		-0.621942			0.842341		
2.264008							
8.150203	159.239899		-1.081497				
6.890998	-6.537842	7108.398926	-151.226273				
35.454010	18.337715			-5.166100		-6.125074	
-5.991016		-0.550432			0.900443		
2.184924							
6.220083	153.750259		-1.015471				
7.436752	-6.721733	7195.742676	-161.741104				
36.850357	19.347576			-5.280004		-6.308979	
-6.491517		-0.485271			1.103564		
2.176015							
4.499797	147.657104		-1.060824				
6.461634	-6.103558	6443.680176	-139.250366				
34.244411	16.662386			-4.582627		-5.751204	
-5.762712		-0.508758			0.823045		
2.106198							
8.062096	142.558319		-0.974320				
6.364246	-5.846340	6352.798340	-127.288017				
33.449390	16.060986			-4.339904		-5.371197	
-5.453247		-0.566399			0.780079		
1.972951							

Descriptive comparison (real vs synthetic) – means and std:
real_mean real_std syn_mean

syn_std			
Avg VTAT	8.9049	3.9006	6.493400

0.947600			
Avg CTAT	0.0000	0.0000	129.750000
17.305599			
Cancelled Rides by Customer	0.2800	0.4491	-0.879900
0.154600			
Cancelled Rides by Driver	0.7200	0.4491	5.837100
0.866800			
Incomplete Rides	-0.0000	0.0000	-5.320000
0.824400			
Booking Value	0.0000	0.0000	5819.493164
856.306274			
Ride Distance	-0.0000	0.0000	-122.946198
22.104099			
Driver Ratings	-0.0000	0.0000	30.341200
4.226900			
Customer Rating	-0.0000	0.0000	15.418200
2.009600			
Vehicle Type_Bike	0.1507	0.3578	-4.098000
0.663200			
Vehicle Type_Go Mini	0.1981	0.3985	-5.064300
0.801300			
Vehicle Type_Go Sedan	0.1830	0.3867	-5.067400
0.802800			
Vehicle Type_Premier Sedan	0.1204	0.3254	-0.390200
0.122100			
Vehicle Type_Uber XL	0.0290	0.1679	0.746000
0.117800			
Vehicle Type_eBike	0.0701	0.2553	1.863800
0.283400			

F1-score (original training): 1.0000

F1-score (augmented training): 1.0000

Classification report (augmented model):

	precision	recall	f1-score	support
0	1.0000	1.0000	1.0000	22500
1	1.0000	1.0000	1.0000	7500
accuracy			1.0000	30000
macro avg	1.0000	1.0000	1.0000	30000
weighted avg	1.0000	1.0000	1.0000	30000

Layer (type)	Output Shape
Param #	
input_layer_5 (InputLayer)	(None, 32)
0	
dense_14 (Dense)	(None, 128)
4,224	
dense_15 (Dense)	(None, 256)
33,024	
dense_16 (Dense)	(None, 256)
65,792	
dense_17 (Dense)	(None, 15)
3,855	

Total params: 106,895 (417.56 KB)

Trainable params: 106,895 (417.56 KB)

Non-trainable params: 0 (0.00 B)

Model: "discriminator"

Layer (type)	Output Shape
Param #	
input_layer_6 (InputLayer)	(None, 15)
0	
dense_18 (Dense)	(None, 256)
4,096	
dense_19 (Dense)	(None, 128)
32,896	

Result: The GAN-generated data enhanced the representation of rare failure cases, resulting in improved predictive modeling and better generalization of machine learning algorithms on imbalanced IoT datasets.

Exp. No 5	Evolution and Advancements of Style-Based GANs
Date:	

Task:

Present your POV on how to generate synthetic data using GANs. You can assume a sample dataset from an IOT enabled machine where the failure rates are minimal.

Aim:

To explore Generative Adversarial Networks (GANs) for synthetic data generation with a focus on datasets containing rare failure events from IoT-enabled machines. The objective is to enhance data availability for training machine learning models.

Explanations:

Research Foundational GANs for Style:

Identify and research the earliest GAN models that laid the groundwork for style generation and transfer, discussing their basic architectures and initial capabilities.

Researching Foundational GANs for Style Generation:

This section will identify and delve into foundational GAN models that were crucial for advancing image generation and laid the groundwork for style generation and transfer. We will cover:

- **Vanilla GAN (2014):**

- **Architecture:** Basic Generator (G) and Discriminator (D) structure.
- **Capabilities:** Introduction of the adversarial process for generating synthetic data, typically small, low-resolution images. Limited control over output.
- **Contribution to Style:** Established the core GAN concept; while not directly style-oriented, it proved the feasibility of adversarial learning for image synthesis.

- **Deep Convolutional GAN (DCGAN, 2015):**

- **Architecture:** Leveraged convolutional layers without pooling, batch normalization, and specific activation functions (ReLU for generator, LeakyReLU for discriminator).

- **Capabilities:** Generated higher quality and resolution images than vanilla GANs, with more stable training. Demonstrated that GANs could learn hierarchical representations.
- **Contribution to Style:** Improved image quality and stability, making GANs more practical for image generation tasks that would later involve style. The latent space began to show image-to-image directions related to visual features, hinting at style control.
- **Pix2Pix (2017):**
 - **Architecture:** Conditional GAN (cGAN) with a U-Net architecture for the generator and a PatchGAN discriminator. **Capabilities:** Performed image-to-image translation tasks (e.g., semantic labels to photo, edges to photo). Outputs were conditioned on an input image.
 - **Contribution to Style:** A significant step towards style transfer and manipulation by demonstrating direct image-to-image translation. It showed that GANs could learn complex mappings between input and output domains, implicitly learning 'style' transformations based on the paired data.
- **CycleGAN (2017):**
 - **Architecture:** Unpaired image-to-image translation using two generators and two discriminators, enforced with cycle consistency loss.
 - **Capabilities:** Translated images between two domains without requiring paired training data (e.g., zebra to horse, summer to winter). Maintained content structure while altering style.
 - **Contribution to Style:** Revolutionized style transfer by enabling it on unpaired datasets, broadening the applicability of GANs for stylistic transformations. It showed that style could be learned and transferred across domains, even without explicit style labels or paired examples.

For each model, we will discuss their innovations, basic architectures, initial capabilities, and how they collectively paved the way for more advanced style-based GANs like the StyleGAN series.

Analyze Evolution of StyleGAN Models:

Trace the development from StyleGAN1 to current StyleGAN models (e.g., StyleGAN2, StyleGAN3), detailing the architectural changes, innovations, and their impact on image generation quality and control.

Evolution of StyleGAN Models:

Generative Adversarial Networks (GANs) have revolutionized image synthesis, and the StyleGAN series, developed by NVIDIA, stands out for its exceptional ability to generate high-fidelity, diverse, and controllable images. The evolution from StyleGAN1 to StyleGAN3 showcases significant architectural advancements, each addressing limitations of its predecessor and pushing the boundaries of realistic image generation and style manipulation.

StyleGAN1 (2018):

1. Core Architectural Innovations:

- **Style-based Generator:** A generator architecture that incorporates a 'mapping network' to transform a latent code (z) into an intermediate latent space (w). This ' w ' space is then injected into the generator at multiple scales through Adaptive Instance Normalization (AdaIN) layers, replacing traditional normalization layers.
- **Progressive Growing of GANs (PGGAN) Integration:** Built upon PGGAN's approach of progressively increasing resolution during training, allowing for stable training of high-resolution images.
- **Truncation Trick:** A technique applied in the ' w ' space to trade off diversity for image quality, pushing generated samples closer to the mean of the learned distribution.

2. Impact on Image Generation:

- **Disentanglement of Style:** Enabled fine-grained control over various aspects of generated images, such as gender, hairstyle, and lighting, while maintaining low-level details like color and texture.
- **High-Quality Image Synthesis:** Produced highly realistic and diverse celebrity faces, setting a new benchmark for image generation quality.

3. Limitations Addressed: Improved upon earlier GANs by offering explicit control over style and enabling more stable training for high-resolution outputs through progressive growth and style injection.

StyleGAN2 (2019)

1. Core Architectural Innovations:

- **Redesigned Generator Normalizations:** Removed the AdaIN operations present in StyleGAN1, which were attributed to the AdaIN operation. StyleGAN2 replaced AdaIN with a new normalization scheme that applies modulation and demodulation to the convolutional weights.
- **Path Length Regularization:** Introduced a regularization technique that encourages the generator to produce images with a consistent magnitude of change in response to a unit-length change in the latent space. This improved the linearity and disentanglement of the latent space.
- **No Progressive Growing:** Abandoned the progressive growing approach in favor of training the full-resolution network from scratch using residual connections, leading to more robust training and fewer artifacts.
- **Weight Demodulation:** Ensures that the magnitude of features remains consistent across different styles, preventing artifacts.

2. Impact on Image Generation:

- **Superior Image Quality and Fidelity:** Significantly reduced artifacts and improved the overall visual quality and realism of generated images.
- **Enhanced Latent Space Disentanglement:** Path length regularization further improved the disentanglement, making style mixing and interpolation more seamless and controllable.
- **Increased Training Stability:** The new normalization and regularization techniques contributed to more stable training and higher reproducibility of results.

3. Limitations Addressed:

Successfully eliminated common artifacts from StyleGAN1, improved training stability, and refined the disentanglement of the latent space for more precise style control.

StyleGAN3 (2021):

1. Core Architectural Innovations:

- **Alias-Free Generative Adversarial Networks:** Addressed the issue of 'aliasing' (stair-stepping or flickering artifacts in videos/animations) that became apparent when generating high-resolution images or videos with StyleGAN1 and StyleGAN2, especially during transformations like rotation or translation. **Shift-Invariant Filtering:** Re-architected the generator to be fully translation and rotation equivariant. This was achieved by using carefully designed upsampling and downsampling filters combined with a new architecture that ensures all operations (convolutions, upsampling, downsampling, nonlinearities) maintain shift invariance.
- **High-Pass Filtering for Aliasing Reduction:** Incorporated explicit high-pass filtering in the generator's internal layers to prevent the introduction of aliasing artifacts.

2. Impact on Image Generation:

- **Temporal Consistency and Animation:** The most significant impact is on generating images that maintain consistency during transformations, making StyleGAN3 ideal for creating high-quality animations and video synthesis without flickering or jaggies.
- **Perceptually More Realistic:** The removal of aliasing artifacts results in images that are perceptually more continuous and realistic, especially when viewed under transformations.

3. Limitations Addressed:

Directly tackled the aliasing problem inherent in previous StyleGAN versions, which was crucial for applications requiring transformation robustness and temporal consistency (e.g., animation, virtual try-on).

Future of GANs in Generating Realistic Images:

The StyleGAN series exemplifies the rapid progress in GAN research. The future of GANs in generating realistic images is promising and likely to involve:

- **Increased Resolution and Speed:** Generating even higher resolution images (e.g., 8K, 16K) faster and with fewer computational resources. **Improved Controllability and Semantic Editing:** Moving beyond simple style mixing to more granular, semantic control over image content (e.g., changing specific objects, lighting conditions, expressions with precise commands). **Multi-Modal Generation:** Integrating text, audio, and other modalities to guide image generation, leading to more versatile and context-aware synthesis.

- **3D-Aware Generation:** Extending 2D image generation to 3D consistent scenes and objects, enabling applications in virtual reality, gaming, and 3D content creation. **Ethical Considerations and Applications**
- **Robustness:** Developing more robust models that are less susceptible to adversarial attacks and addressing biases in training data to promote fair and ethical image generation.
- **Fewer Training Data Requirements:** Techniques that allow GANs to learn from smaller datasets, making them applicable to niche domains where large datasets are unavailable.

Conclusion:

The journey of style-related GANs, from foundational models like Vanilla GAN and DCGAN to the sophisticated StyleGAN series, is a testament to rapid innovation in deep learning. Each step has built upon its predecessors, refining the generation process, enhancing stylistic control, and pushing the boundaries of realism. The successive improvements in disentanglement, architectural stability, and alias-free generation have transformed GANs from a research curiosity into powerful tools for creative industries, scientific visualization, and beyond. The future promises even more astonishing capabilities, moving towards seamless integration with other AI paradigms, generating highly interactive and 3D-consistent content, while simultaneously demanding a careful consideration of the ethical landscape. Style-related GANs are not just generating images; they are reshaping our understanding of creativity, visual perception, and the very fabric of reality itself.

Exp. No 6

Deepfakes Using GANs: Creation and Detection Techniques

Date:

Task:

Present your POV on GANs used for Deep Fakes. Articulate how we can identify the Deep Fake from the original.

Aim:

To examine the use of GANs for generating deepfake content and to evaluate existing techniques for detecting manipulated media, ensuring authenticity and trust in digital information.

Explanations:

GANs used for Deep Fakes:

Generative Adversarial Networks (GANs) have revolutionized content creation, enabling the synthesis of highly realistic images, audio, and video. While their applications in creative arts, data augmentation, and scientific research are immense and positive, their use in creating 'Deep Fakes' presents significant ethical, social, and security concerns. Deep Fakes, often generated by GANs, are synthetic media in which a person in an existing image or video is replaced with someone else's likeness.

My primary concern revolves around:

- 1. Misinformation and Disinformation:** Deep Fakes can be used to create fabricated news, speeches, or events, leading to widespread confusion, erosion of trust in media, and potential societal instability.
- 2. Reputational Harm and Defamation:** Individuals can be falsely portrayed in compromising situations, causing severe damage to their reputation, careers, and personal lives.
- 3. Erosion of Trust in Digital Media:** The increasing sophistication of Deep Fakes makes it difficult for the average person to distinguish between real and fake content, leading to a general distrust of all digital media.
- 4. Security Risks:** Deep Fakes could be used for identity theft, fraud, or even to influence political outcomes.

While the technology itself is neutral, its malicious application warrants serious attention and the development of robust countermeasures.

How to Identify Deep Fakes from Originals:

Detecting Deep Fakes is an ongoing challenge, as GANs are continuously improving. However, several methods and common tells can help in identification:

1. Visual Inconsistencies and Artifacts:

- **Unnatural Blinking:** Deep Fake subjects often blink less frequently or have irregular blinking patterns compared to real people, as training data for blinking can be scarce.
- **Facial Asymmetries:** Subtle inconsistencies in facial features, like eyes, ears, or mouth, that don't match the original person.
- **Poorly Rendered Edges:** The transition between the manipulated face and the rest of the head or body might show blurry, jagged, or unnatural edges.
 - **Inconsistent Lighting and Shadows:** The lighting on the deep-faked face might not match the lighting in the background or on the original body, leading to unrealistic shadows or highlights.
 - **Low Resolution or Digital Artifacts:** Sometimes, parts of the Deep Fake might appear pixelated, distorted, or have digital noise, especially around the edges of the manipulated area.
- **Missing or Inconsistent Skin Details:** Deep Fake faces might lack natural skin blemishes, pores, or hair follicles, or these details might be overly smooth or inconsistently applied.

2. Auditory Inconsistencies (for videos with audio):

- **Lip-Sync Issues:** The movement of the lips might not perfectly synchronize with the spoken words.
- **Unnatural Voice Tones or Cadence:** The synthesized voice might sound robotic, flat, or have an unnatural rhythm or intonation.
- **Background Noise Discrepancies:** The background audio might not match the visual environment or could have subtle inconsistencies with the synthesized voice.

3. Temporal Inconsistencies:

- **Jittering or Flickering:** Rapid, unnatural movements or changes in the manipulated area across frames.
- **Pose or Head Movement Anomalies:** The head or body movements might appear stiff, unnatural, or not consistent with typical human motion.

4. Metadata Analysis:

- **Examine File Metadata:** Videos often contain metadata (EXIF data) about the camera, date, and time of capture. The absence of this data or inconsistent metadata can be a red flag.

5. Technical Detection Methods (often requiring specialized tools):

- **AI-based Detectors:** Developing sophisticated AI models specifically trained to identify Deep Fakes by learning to recognize the subtle artifacts left by GANs. These models often look for patterns that are imperceptible to the human eye.
- **Physiological Signal Analysis:** Analyzing heart rate, breathing patterns, or eye movements, which are difficult for current Deep Fake technology to perfectly replicate.
- **Digital Watermarking/Blockchain:** Future solutions might involve embedding digital watermarks into original content or using blockchain to verify the authenticity and provenance of media.

While human observation can catch some obvious Deep Fakes, the most effective detection strategies will increasingly rely on a combination of visual inspection, forensic analysis, and advanced AI-powered detection tools.

Conclusion:

Generative Adversarial Networks (GANs) have significantly advanced the creation of synthetic media, enabling realistic and innovative digital content. However, these same capabilities have facilitated the emergence of Deep Fakes, which pose serious ethical, social, and security challenges. Deep Fakes can distort truth, manipulate public opinion, harm personal reputations, and undermine trust in digital information. As the technology behind Deep Fakes continues to improve, detection becomes increasingly complex.