**NAME – PRIYANSHI ASATI**

**SAP ID – 500108525**

**ROLL NO – R2142220982**

**BATCH- 4 (AIML-NON HONS.)**

## PREDICTIVE ANALYTICS LAB- 5

**Customer Churn Prediction (regression)**

- **Dataset: Telecom customer data (e.g., from Kaggle):**
  - **Data preprocessing (handling missing values, outliers, feature scaling).**

  **CODE:**

```
import pandas as pd

from sklearn.impute import SimpleImputer

from sklearn.preprocessing import StandardScaler

import matplotlib.pyplot as plt

import seaborn as sns

import numpy as np

data = pd.read_csv('/content/Telecom_customer
churn.csv.zip')

num_cols = data.select_dtypes(include=['float64',
'int64']).columns

imputer_num = SimpleImputer(strategy='mean')

data[num_cols] =
imputer_num.fit_transform(data[num_cols])
```

```
cat_cols =
data.select_dtypes(include=['object']).columns

imputer_cat =
SimpleImputer(strategy='most_frequent')

data[cat_cols] =
imputer_cat.fit_transform(data[cat_cols])

z_scores = np.abs((data[num_cols] -
data[num_cols].mean()) / data[num_cols].std())

threshold = 3

data[num_cols] = np.where(z_scores > threshold,
data[num_cols].mean(), data[num_cols])

scaler = StandardScaler()

data[num_cols] =
scaler.fit_transform(data[num_cols])

data.head()
```

**OUTPUT:**

| | rev_Mean | mou_Mean | totmrc_Mean | da_Mean | ovrmou_Mean | ovrrev_Mean | vceovr_Mean | dat |
|---|---|---|---|---|---|---|---|---|
| 0 | -0.960957 | -0.600041 | -1.100190 | -0.365203 | -0.565426 | -0.576005 | -0.570269 | |
| 1 | 0.072225 | 0.019651 | -0.372904 | -0.365203 | -0.162033 | -0.082038 | -0.069650 | |
| 2 | -1.177109 | -1.091560 | -1.368688 | -0.570643 | -0.565426 | -0.576005 | -0.570269 | |
| 3 | -0.529038 | -1.098028 | -0.344884 | -0.570643 | -0.565426 | -0.576005 | -0.570269 | |
| 4 | 0.002437 | 0.226018 | 1.310940 | -0.570643 | -0.565426 | -0.576005 | -0.570269 | |

5 rows × 100 columns

| datovr_Mean | roam_Mean | change_mou | ... | forgntvl | ethnic | kid0_2 | kid3_5 | kid6_10 |
|---|---|---|---|---|---|---|---|---|
| -0.194986 | -0.273458 | -0.744487 | ... | -0.283288 | N | U | U | U |
| -0.194986 | -0.273458 | 2.785990 | ... | -0.283288 | Z | U | U | U |
| -0.194986 | -0.273458 | 0.038925 | ... | -0.283288 | N | U | Y | U |
| -0.194986 | -0.273458 | 0.053006 | ... | -0.283288 | U | Y | U | U |
| -0.194986 | -0.273458 | 0.257820 | ... | -0.283288 | I | U | U | U |

| kid11_15 | kid16_17 | creditcd | eqpdays | Customer_ID |
|---|---|---|---|---|
| U | U | Y | -0.089967 | -1.732033 |
| U | U | Y | -0.598883 | -1.731999 |
| U | U | Y | 0.040132 | -1.731964 |
| U | U | Y | 0.040132 | -1.731930 |
| U | U | Y | 0.217065 | -1.731895 |

- **Exploratory data analysis (EDA) to identify potential predictors.**

  **CODE:**

```
plt.figure(figsize=(10, 6))

sns.heatmap(data.isnull(), cbar=False, cmap='viridis')

plt.title('Missing Data Heatmap')

plt.show()

plt.figure(figsize=(12, 8))

correlation_matrix = data[num_cols].corr()

sns.heatmap(correlation_matrix, annot=False, cmap='coolwarm', linewidths=0.5)

plt.title('Correlation Heatmap')

plt.show()

num_cols = data.select_dtypes(include=['float64', 'int64']).columns

data[num_cols].hist(bins=20, figsize=(16, 12), color='skyblue')

plt.suptitle('Distribution of Numeric Features')

plt.show()

plt.figure(figsize=(12, 8))
```

```python
sns.boxplot(data=data[num_cols], orient='h', palette='coolwarm')

plt.title('Boxplot of Numeric Features')

plt.show()

top_corr_features = correlation_matrix.abs().unstack().sort_values(ascending=False).drop_duplicates().head(5).index

top_corr_data = data[list(set([i[0] for i in top_corr_features]))]

sns.pairplot(top_corr_data)

plt.suptitle('Pairplot of Top Correlated Features', y=1.02)

plt.show()

cat_cols = data.select_dtypes(include=['object']).columns

for col in cat_cols:

    plt.figure(figsize=(8, 4))

    sns.countplot(x=data[col], palette='coolwarm')

    plt.title(f'Distribution of {col}')

    plt.xticks(rotation=90)

plt.show()
```
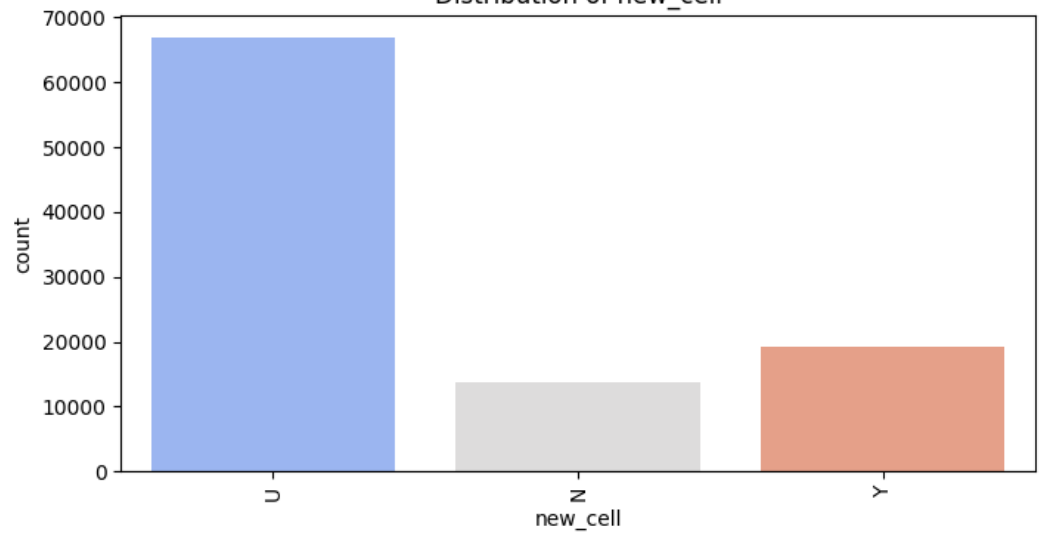
**OUTPUT:**

Missing Data Heatmap



Correlation Heatmap

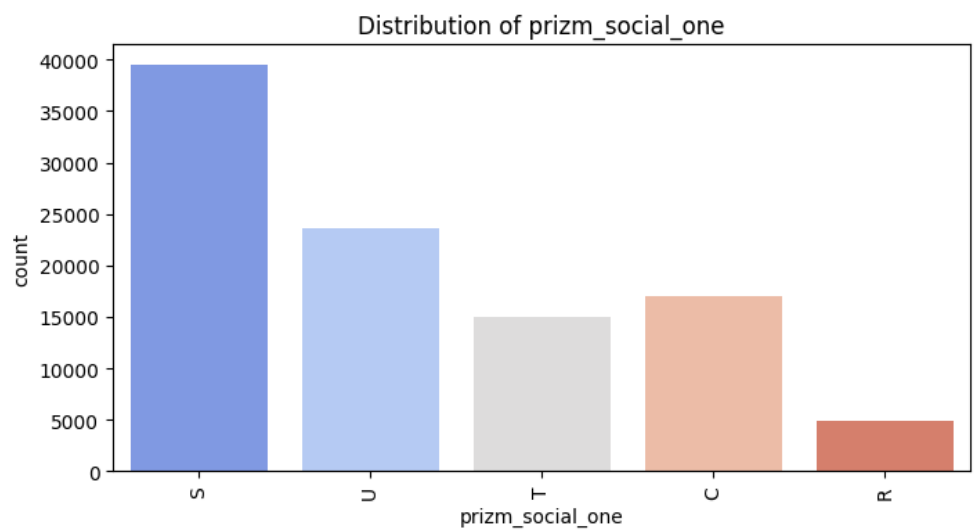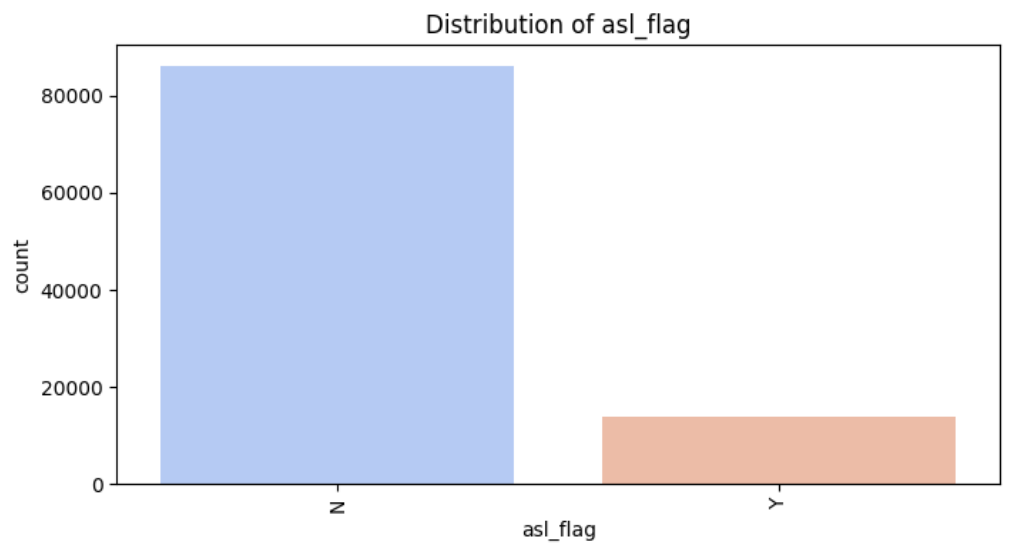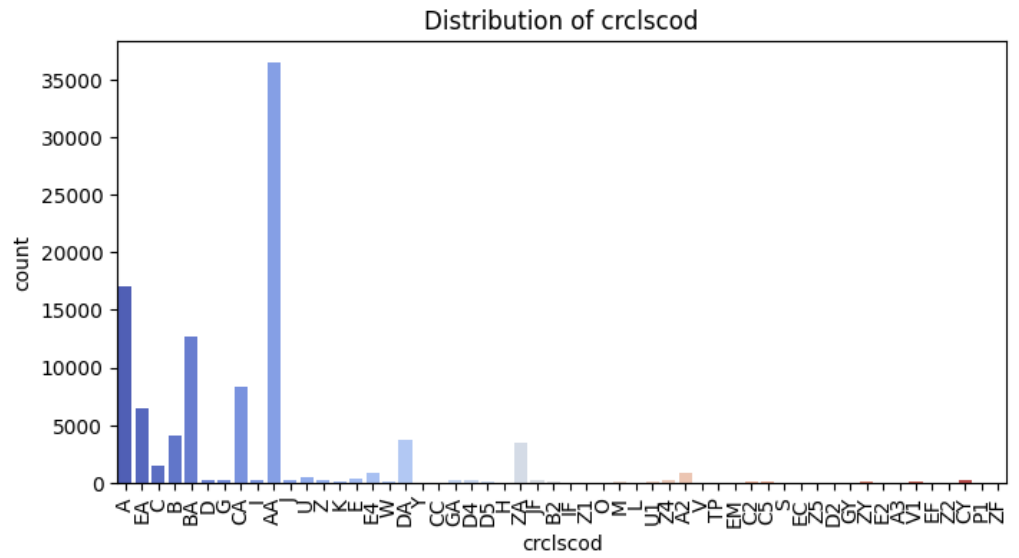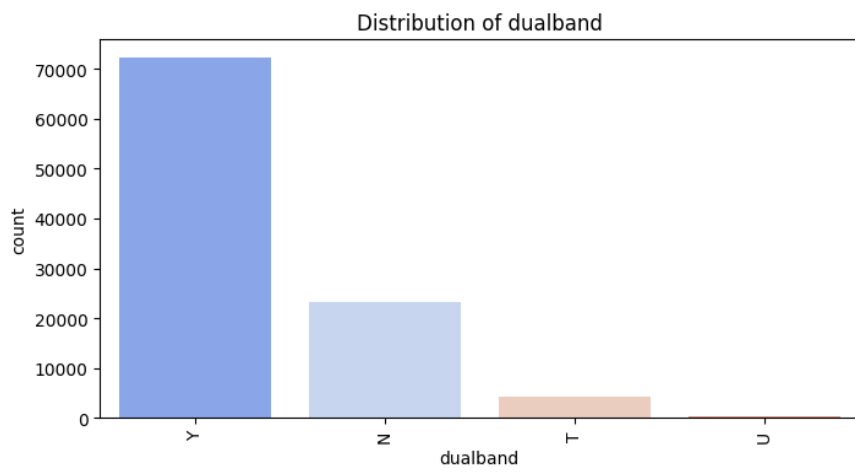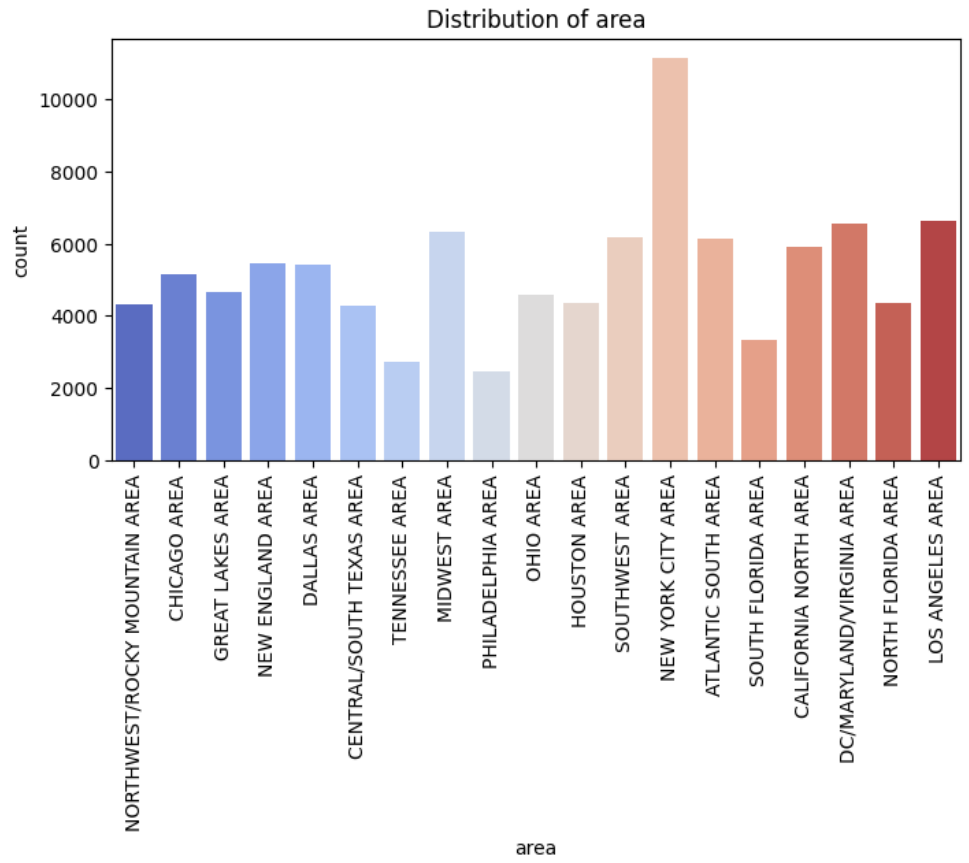Distribution of Numeric Features



Boxplot of Numeric Features

Pairplot of Top Correlated Features


Distribution of new_cell

Distribution of crclscod

Distribution of asl_flag

Distribution of prizm_social_one

## Distribution of area



## Distribution of dualband

## Distribution of refurb_new



## Distribution of hnd_webcap



## Distribution of ownrent

## Distribution of dwlltype



## Distribution of marital



## Distribution of infobase

Distribution of HHstatin

Distribution of dwllsize

Distribution of ethnic

Distribution of kid3_5

Distribution of kid0_2

Distribution of kid6_10

Distribution of kid11_15



Distribution of kid16_17



Distribution of creditcd

- **Building a simple linear regression model.**

  **CODE:**

```python
X = data[['rev_Mean']]

y = data['totmrc_Mean']

X = X.dropna()

y = y[X.index]

X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

model = LinearRegression()

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

mse = mean_squared_error(y_test, y_pred)

r2 = r2_score(y_test, y_pred)

print(f'Mean Squared Error: {mse}')

print(f'R² Score: {r2}')

plt.figure(figsize=(10, 6))

plt.scatter(X_test, y_test, color='blue', label='Actual
Data')

plt.plot(X_test, y_pred, color='red', linewidth=2,
label='Regression Line')

plt.title('Simple Linear Regression: rev_Mean vs
totmrc_Mean')

plt.xlabel('rev_Mean')

plt.ylabel('totmrc_Mean')

plt.legend()

plt.show()
```
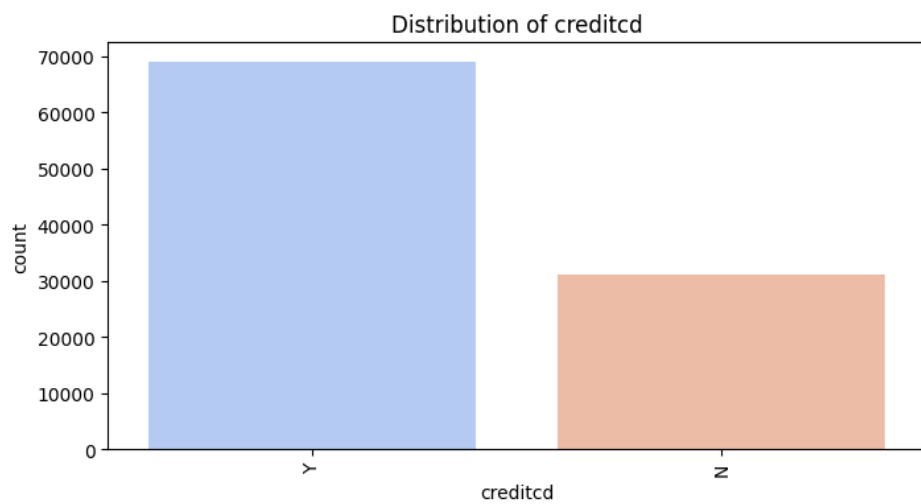
**OUTPUT:**

```
Mean Squared Error: 0.5998927894028158
R² Score: 0.3993436555753672
```


Simple Linear Regression: rev_Mean vs totmrc_Mean

- **Evaluating model performance using accuracy, precision, recall, F1-score.**

  **CODE:**

```
X = data[['rev_Mean'

y = data['churn']

X = X.dropna()

y = y[X.index]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = LogisticRegression()

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)

print(f'Accuracy: {accuracy}')

precision = precision_score(y_test, y_pred)
```

```python
print(f'Precision: {precision}')

recall = recall_score(y_test, y_pred)

print(f'Recall: {recall}')

f1 = f1_score(y_test, y_pred)

print(f'F1-Score: {f1}')

conf_matrix = confusion_matrix(y_test, y_pred)

print(f'Confusion Matrix:\n {conf_matrix}')

plt.figure(figsize=(6, 4))

sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', cbar=False)

plt.title('Confusion Matrix')

plt.ylabel('Actual Label')

plt.xlabel('Predicted Label')

plt.show()

print("\nClassification Report:")

print(classification_report(y_test, y_pred))
```

**OUTPUT:**

```
Accuracy: 0.5089066184956595
Precision: 0.6907894736842105
Recall: 0.010665312341289994
F1-Score: 0.021006301890056717
Confusion Matrix:
 [[10037    47]
 [ 9740   105]]
```

## Confusion Matrix



```
Classification Report:
              precision    recall  f1-score   support

           0       0.51      1.00      0.67     10084
           1       0.69      0.01      0.02      9845

    accuracy                           0.51     19929
   macro avg       0.60      0.50      0.35     19929
weighted avg       0.60      0.51      0.35     19929
```