

# **MCSE Orange Level Problem**

**Name: Adarsh D SRN: PES2UG24AM010**

**Section: A (AIML)**

## Python Code Used in Analysis:

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import StandardScaler

from sklearn.linear_model import LogisticRegression

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import classification_report,
confusion_matrix, roc_auc_score, roc_curve

# Load the dataset

df = pd.read_csv('loan_data.csv')

# Data Cleaning

# Handle missing values

df['employment_length'].fillna(df['employment_length'].median(),
), inplace=True)
```

```
df['annual_income'].fillna(df['annual_income'].median(),  
inplace=True)

df.dropna(subset=['loan_status'], inplace=True)

# Outlier Detection and Treatment

def remove_outliers(df, column):  
  
    Q1 = df[column].quantile(0.25)  
  
    Q3 = df[column].quantile(0.75)  
  
    IQR = Q3 - Q1  
  
    lower_bound = Q1 - 1.5 * IQR  
  
    upper_bound = Q3 + 1.5 * IQR  
  
    return df[(df[column] >= lower_bound) & (df[column] <=  
upper_bound)]  
  
df = remove_outliers(df, 'annual_income')  
  
df = remove_outliers(df, 'loan_amount')
```

```
# Data Preprocessing

# Convert categorical variables to numerical

df['loan_status'] = df['loan_status'].map({'Approved': 1,
'Denied': 0})

df = pd.get_dummies(df, columns=['employment_type',
'loan_purpose'], drop_first=True)

# Feature Engineering

df['debt_to_income_ratio'] = df['monthly_debt'] /
(df['annual_income'] / 12)

df['loan_to_income_ratio'] = df['loan_amount'] /
df['annual_income']

# Exploratory Data Analysis

# Correlation matrix

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

sns.heatmap(df.corr(), annot=True, cmap='coolwarm', center=0)

plt.title('Correlation Matrix')

plt.savefig('correlation_matrix.png')

plt.close()
```

```
# Distribution plots

fig, axes = plt.subplots(2, 2, figsize=(15, 10))

sns.histplot(df['annual_income'], kde=True, ax=axes[0, 0])

axes[0, 0].set_title('Annual Income Distribution')

sns.histplot(df['loan_amount'], kde=True, ax=axes[0, 1])

axes[0, 1].set_title('Loan Amount Distribution')

sns.histplot(df['credit_score'], kde=True, ax=axes[1, 0])

axes[1, 0].set_title('Credit Score Distribution')

sns.countplot(x='loan_status', data=df, ax=axes[1, 1])

axes[1, 1].set_title('Loan Status Distribution')

plt.tight_layout()

plt.savefig('distributions.png')

plt.close()

# Statistical Analysis

# Group statistics by loan status
```



```
# Scale features

scaler = StandardScaler()

X_train_scaled = scaler.fit_transform(X_train)

X_test_scaled = scaler.transform(X_test)

# Logistic Regression Model

log_reg = LogisticRegression(random_state=42, max_iter=1000)

log_reg.fit(X_train_scaled, y_train)

y_pred_log = log_reg.predict(X_test_scaled)

y_pred_proba_log = log_reg.predict_proba(X_test_scaled)[:, 1]

print("Logistic Regression Results:")

print(classification_report(y_test, y_pred_log))

print(f"ROC AUC Score: {roc_auc_score(y_test, y_pred_proba_log):.4f}")

# Random Forest Model

rf_model = RandomForestClassifier(n_estimators=100,
random_state=42, max_depth=10)
```

```
rf_model.fit(X_train_scaled, y_train)

y_pred_rf = rf_model.predict(X_test_scaled)

y_pred_proba_rf = rf_model.predict_proba(X_test_scaled)[:, 1]

print("\nRandom Forest Results:")

print(classification_report(y_test, y_pred_rf))

print(f"ROC AUC Score: {roc_auc_score(y_test,
y_pred_proba_rf):.4f}")

# Feature Importance

feature_importance = pd.DataFrame({

    'feature': X.columns,

    'importance': rf_model.feature_importances_


}).sort_values('importance', ascending=False)

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

sns.barplot(x='importance', y='feature',
data=feature_importance.head(10))

plt.title('Top 10 Feature Importances')
```

```
plt.savefig('feature_importance.png')

plt.close()

# ROC Curve Comparison

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

fpr_log, tpr_log, _ = roc_curve(y_test, y_pred_proba_log)

fpr_rf, tpr_rf, _ = roc_curve(y_test, y_pred_proba_rf)

plt.plot(fpr_log, tpr_log, label=f'Logistic Regression (AUC = {roc_auc_score(y_test, y_pred_proba_log):.4f})')

plt.plot(fpr_rf, tpr_rf, label=f'Random Forest (AUC = {roc_auc_score(y_test, y_pred_proba_rf):.4f})')

plt.plot([0, 1], [0, 1], 'k--', label='Random Classifier')

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive Rate')

plt.title('ROC Curve Comparison')

plt.legend()

plt.savefig('roc_curve.png')

plt.close()
```

```
# Confusion Matrix

fig, axes = plt.subplots(1, 2, figsize=(15, 5))

sns.heatmap(confusion_matrix(y_test, y_pred_log), annot=True,
fmt='d', cmap='Blues', ax=axes[0])

axes[0].set_title('Logistic Regression Confusion Matrix')

axes[0].set_xlabel('Predicted')

axes[0].set_ylabel('Actual')

sns.heatmap(confusion_matrix(y_test, y_pred_rf), annot=True,
fmt='d', cmap='Blues', ax=axes[1])

axes[1].set_title('Random Forest Confusion Matrix')

axes[1].set_xlabel('Predicted')

axes[1].set_ylabel('Actual')

plt.tight_layout()

plt.savefig('confusion_matrices.png')

plt.close()

print("\nAnalysis complete. All visualizations saved.")
```

# Orange Level Data Analysis Project

## 1. Project Workflow and Requirements

This project involves comprehensive data analysis following a structured workflow:

- **Data Collection:** Gathering loan application data from various sources
- **Data Cleaning:** Handling missing values, outliers, and inconsistencies
- **Data Preprocessing:** Transforming and preparing data for analysis
- **Exploratory Data Analysis:** Understanding patterns and relationships in the data
- **Statistical Analysis:** Applying statistical methods to derive insights
- **Visualization:** Creating meaningful visual representations of findings
- **Model Development:** Building predictive models for loan approval
- **Documentation:** Comprehensive reporting of methodology and results

## 2. Data Preprocessing Summary

The data preprocessing phase involved several critical steps:

- **Missing Value Treatment:** Identified and handled missing values using appropriate imputation techniques
- **Outlier Detection:** Used statistical methods to identify and treat outliers in numerical variables
- **Data Type Conversion:** Ensured all variables have appropriate data types for analysis
- **Feature Engineering:** Created new features to enhance model performance
- **Data Normalization:** Scaled numerical features for consistent analysis
- **Categorical Encoding:** Converted categorical variables into numerical format

## 3. Descriptive Statistics and Inferences

Key statistical findings from the dataset:

- **ApplicantIncome:** Mean income shows significant variation across applicants
- **CoapplicantIncome:** Distribution indicates varying levels of additional income support
- **LoanAmount:** Central tendency and spread reveal typical loan request patterns
- **Loan\_Amount\_Term:** Most loans follow standard term periods
- **Credit\_History:** Strong indicator of loan approval probability
- **Property\_Area:** Geographic distribution affects loan approval rates

## 4. Visualization Results

Multiple visualizations were created to understand data patterns:

- **Distribution Plots:** Histograms showing the distribution of numerical variables
- **Box Plots:** Identifying outliers and quartile ranges
- **Correlation Heatmaps:** Understanding relationships between variables
- **Bar Charts:** Categorical variable frequency distributions
- **Scatter Plots:** Exploring relationships between continuous variables

## 5. Q-Q Plot Analysis

The Q-Q (Quantile-Quantile) plot for ApplicantIncome was generated to assess normality:

- **Purpose:** Evaluate whether ApplicantIncome follows a normal distribution
- **Interpretation:** Deviations from the diagonal line indicate departures from normality
- **Findings:** The plot reveals the distribution characteristics and potential skewness
- **Implications:** Informs the choice of statistical tests and transformations needed

**Q-Q Plot for ApplicantIncome:**



## 6. Statistical Inferences and Model Findings

Statistical analysis revealed several key insights:

- **Hypothesis Testing:** Conducted tests to validate assumptions about the data
- **Significance Levels:** Identified statistically significant relationships between variables
- **Confidence Intervals:** Established ranges for parameter estimates
- **Model Performance:** Evaluated using appropriate metrics (accuracy, precision, recall)
- **Feature Importance:** Determined which variables have the strongest predictive power
- **Residual Analysis:** Examined model residuals to validate assumptions

## 7. Linear Regression & Correlation

Linear regression analysis was performed to understand relationships between variables:

- **Correlation Analysis:** Examined linear relationships between numerical variables
- **Regression Coefficients:** Quantified the impact of independent variables on the dependent variable
- **R-squared Value:** Measured the proportion of variance explained by the model
- **P-values:** Assessed statistical significance of predictor variables
- **Multicollinearity Check:** Verified independence of predictor variables
- **Model Assumptions:** Validated linearity, homoscedasticity, and normality of residuals
- **Predictive Accuracy:** Evaluated model performance on test data

## 8. Conclusion

This comprehensive data analysis project successfully achieved its objectives:

- **Data Quality:** Thorough preprocessing ensured clean, reliable data for analysis
- **Insights Generated:** Statistical analysis revealed meaningful patterns and relationships
- **Visualization Effectiveness:** Charts and plots effectively communicated findings
- **Model Development:** Predictive models were built and validated successfully
- **Normality Assessment:** Q-Q plot analysis provided insights into data distribution
- **Actionable Recommendations:** Findings can inform decision-making processes
- **Future Work:** Additional feature engineering and advanced modeling techniques could further improve results