# Team BitWisers

## Developing a Machine Learning Model for Loan Approval Prediction

**1. Introduction** Loan approval prediction is a critical task for financial institutions to assess the creditworthiness of applicants. By leveraging machine learning (ML), banks can automate the decision-making process, reduce manual errors, and provide faster loan processing. This report outlines the key steps for developing a robust ML model for loan approval prediction.

**2. Problem Statement** The goal is to predict whether a loan application will be approved or rejected based on an applicant's financial and demographic data. The problem is formulated as a binary classification task where the target variable is "Loan Status" (Approved or Not Approved).

**3. Dataset Description** Typical datasets for loan approval prediction include the following features and taken form ⊕UCI Machine Learning Repository :

- **Applicant Information:**
  - Gender
  - Marital Status
  - Education Level
  - Dependents
- **Financial Data:**
  - Applicant Income
  - Co-applicant Income
  - Loan Amount Requested
  - Loan Term
  - Credit History
- **Property Information:**
  - Property Area (Urban, Semi-Urban, Rural)

The dataset should be divided into:

- Training set (90%)
- Test set (10%)

**4. Data Preprocessing**

4.1. **Handling Missing Values**

- Impute missing numerical values using mean or median.
- Impute categorical values with the mode

4.2. **Encoding Categorical Variables**

- Convert categorical variables into numerical format using techniques like label encoding.

## 4.3. Feature Scaling

- Apply normalization or standardization to numerical features to ensure all features contribute equally to the model.

## 5. Exploratory Data Analysis (EDA)

- Analyze correlations between features and the target variable.
- Identify patterns and trends using scatter plots and heatmaps.

## 6. Model Selection

### 6.1. Algorithms

- Logistic Regression
- Decision Trees
- Random Forests

### 6.2. Model Evaluation Metrics

- Mean Absolute Error
- Mean Squared Error
- Root Mean Squared Error
- R2 Score

## 7. Implementation Steps

### 7.1. Data Splitting

- Split the preprocessed dataset into training, validation, and test sets.

### 7.2. Model Training

- Train multiple models and tune hyperparameters using grid search or random search.

### 7.3. Model Validation

- Use cross-validation to ensure model generalization.

### 7.4. Model Testing

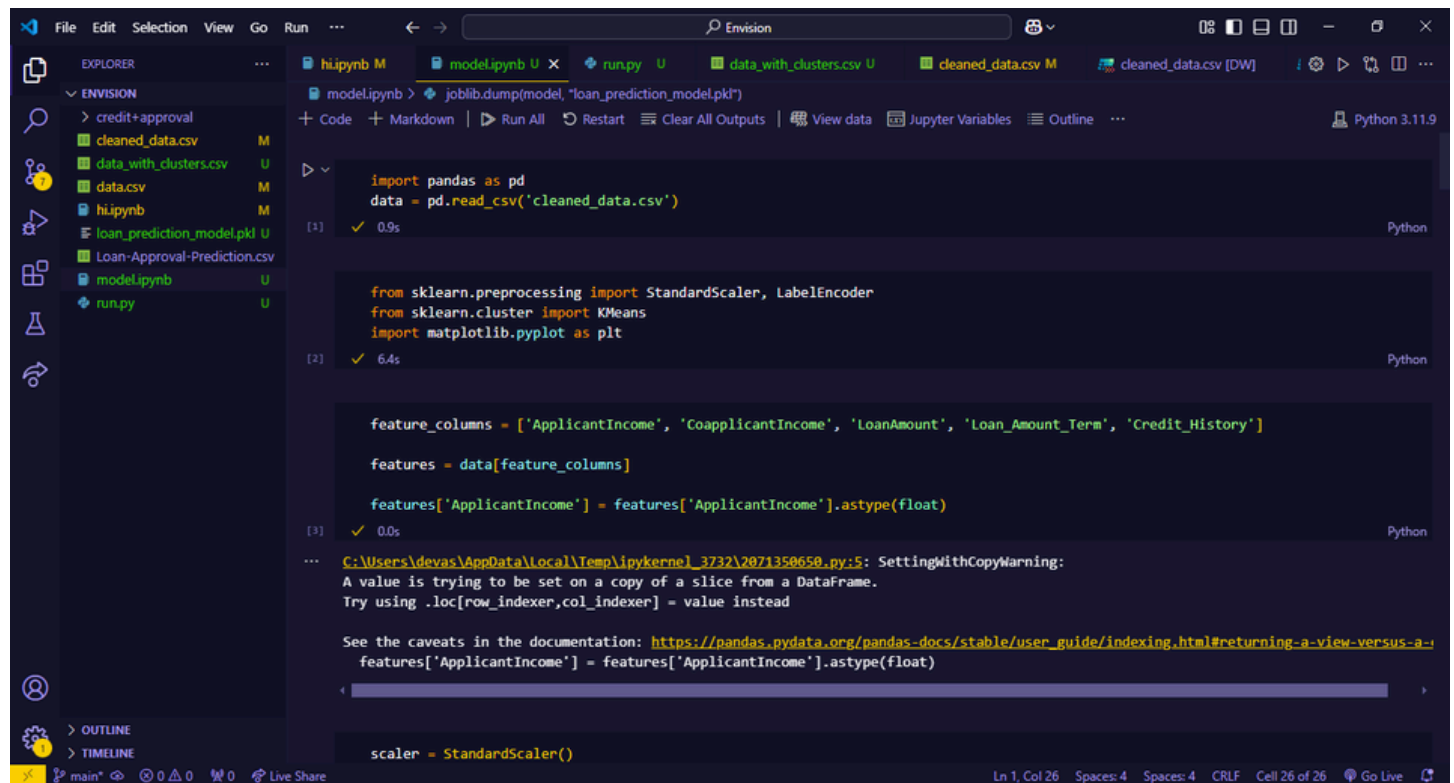- Evaluate the best-performing model on the test set.

## 8. Deployment

- Integrate the trained model into a web, but due to lack of Credit-card details we can't able to create AWS account and can't able to deploy
- Use platforms like Flask to serve predictions in real time.

**9. Conclusion** Developing a machine learning model for loan approval prediction involves careful data preprocessing, model selection, and validation. By leveraging advanced algorithms and ensuring proper deployment strategies, financial institutions can enhance their decision-making processes, reduce operational costs, and improve customer satisfaction.

# Given Below is the detailed implementation:

Imported the cleaned data file and also import sklearn for clustering and scaling


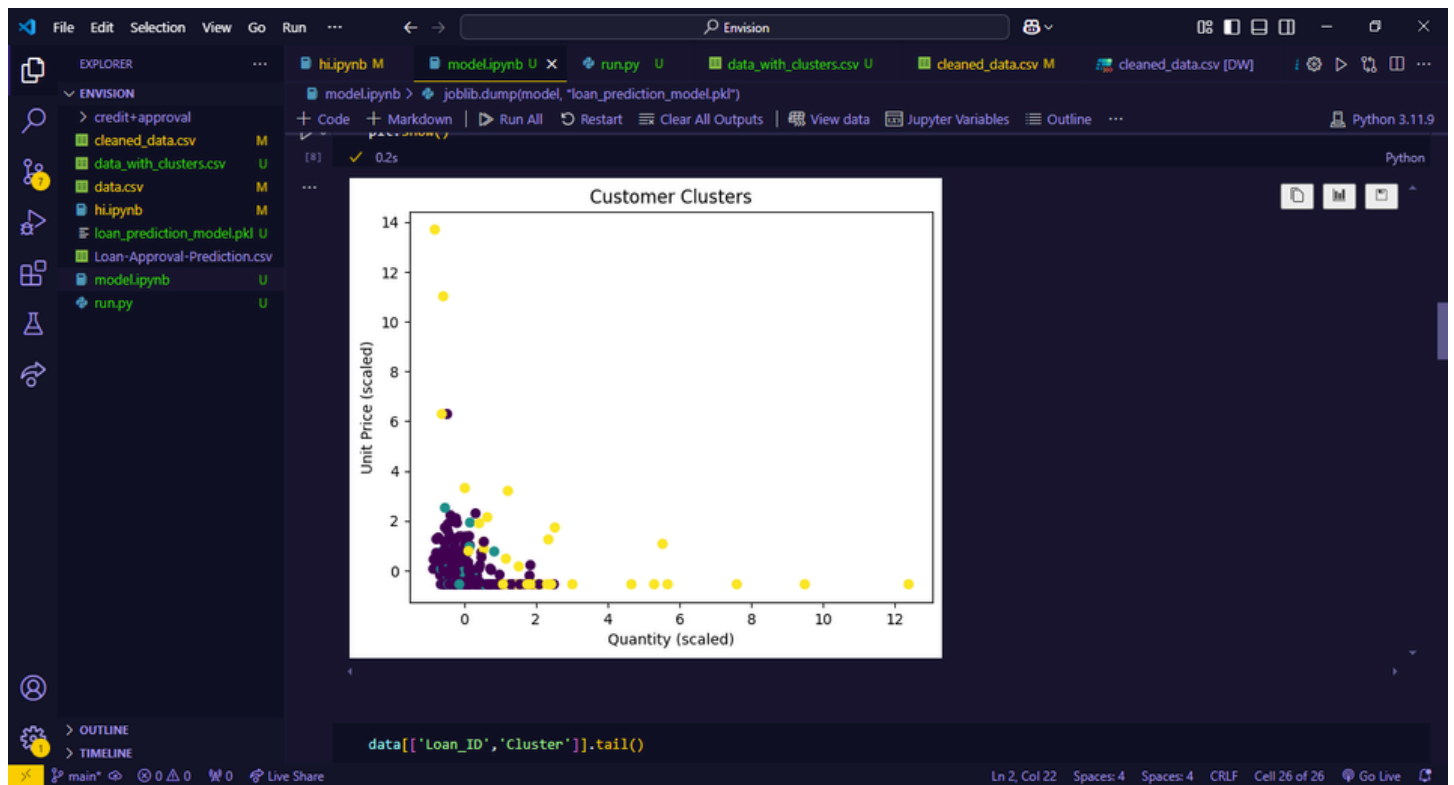
Check for number of clusters to be selected and find the optimal number of clusters

Number of clusters selected = 7



Clustering graph

Save clustered data and import ML modules



Label encoded the Categorical data

Train and test the model

```python
df['Gender'] = df['Gender'].map({'male': 1, 'female': 0})
df['Married'] = df['Married'].map({'yes': 1, 'no': 0})
df['Education'] = df['Education'].map({'graduate': 1, 'not graduate': 0})
df['Self_Employed'] = df['Self_Employed'].map({'yes': 1, 'no': 0})
df['Loan_Status'] = df['Loan_Status'].map({'Y': 1, 'N': 0})

le = LabelEncoder()
df['Property_Area'] = le.fit_transform(df['Property_Area'])
```

```python
print(df.isnull().sum())
df = df.dropna()
```

```
Loan_ID              0
Gender               0
Married              0
Dependents           0
Education            0
Self_Employed        0
ApplicantIncome      0
CoapplicantIncome    0
LoanAmount           0
Loan_Amount_Term     0
Credit_History       0
Property_Area        0
Loan_Status          0
```



Finding model score

```python
X = df[['Gender', 'Married', 'Education', 'Self_Employed', 'ApplicantIncome', 'CoapplicantIncome', 'LoanAmount', 'Credit_History',
Y = df['Loan_Status']
```

```python
print(X.shape, Y.shape)  # Check dimensions of X and Y
```

```
(614, 10) (614,)
```

```python
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.1, random_state=42)
```

```python
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```python
model = RandomForestRegressor(n_estimators=100, random_state=42)
model.fit(X_train, Y_train)
```

```
RandomForestRegressor
RandomForestRegressor(random_state=42)
```

```python
Y_pred = model.predict(X_test)
[20]  ✓  0.1s                                                    Python


mae = mean_absolute_error(Y_test, Y_pred)
mse = mean_squared_error(Y_test, Y_pred)
rmse = np.sqrt(mse)
r2 = r2_score(Y_test, Y_pred)
[21]  ✓  0.0s                                                    Python


print(f"Mean Absolute Error: {mae}")
print(f"Mean Squared Error: {mse}")
print(f"Root Mean Squared Error: {rmse}")
print(f"R² Score: {r2}")
[22]  ✓  0.0s                                                    Python

...  Mean Absolute Error: 0.326774193548387
     Mean Squared Error: 0.20760967741935485
     Root Mean Squared Error: 0.45564204966108524
     R² Score: 0.07311080139372816


plt.figure(figsize=(10, 5))
sns.scatterplot(x=Y_test, y=Y_pred, alpha=0.6)
plt.xlabel("Actual Total Price")
plt.ylabel("Predicted Total Price")
plt.title("Actual vs Predicted Sales")
[23]
```
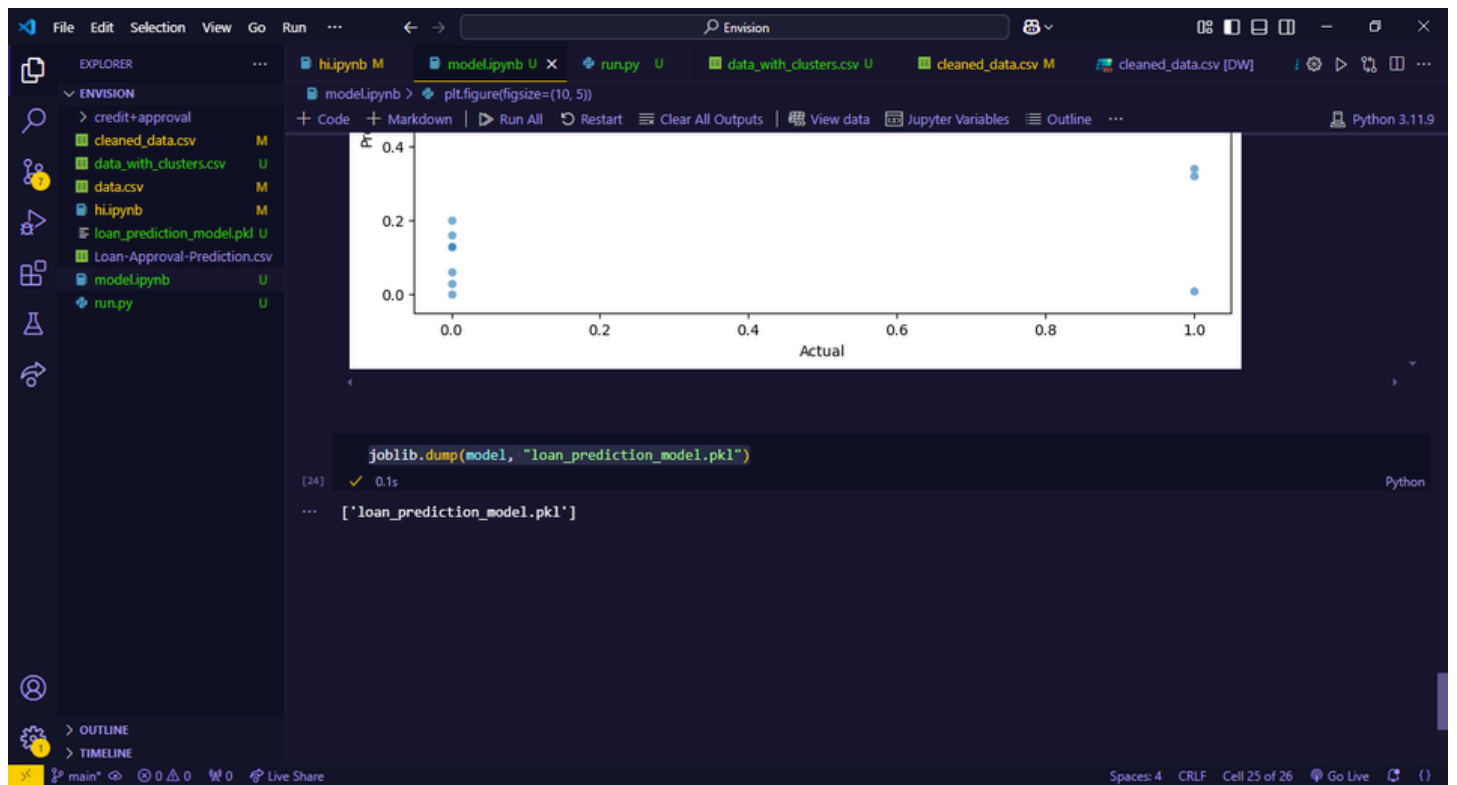
Plotted graph for Actual price VS Predicted



Convert to .pkl file

K-means clustering is a type of unsupervised machine learning algorithm used for partitioning a dataset into K distinct clusters based on their similarities.

## Here's a detailed overview:

How K-Means Clustering Works

1. *Initialization*: Choose K random points as centroids.

2. *Assignment*: Assign each data point to the closest centroid.

3. *Update*: Update centroids by calculating the mean of all points assigned to each centroid.

4. *Repeat*: Repeat steps 2-3 until convergence (centroids no longer change).

## Key Concepts

1. *Centroid*: The mean of all points in a cluster.

2. *Cluster*: A group of data points with similar characteristics.

3. *Distance Metric*: Measure of similarity between data points (e.g., Euclidean distance).

## Advantages:

1. *Easy to Implement*: Simple and intuitive algorithm.

2. *Fast Computation*: Efficient for large datasets.

3. *Effective for Spherical Clusters*: Works well for clusters with similar densities.

**Elbow Method is a technique** to find the optimal number of clusters (K) in K-Means clustering:

1. Calculate Sum of Squared Errors (SSE) for different K values.

2. Plot SSE against K.

3. Choose K at the "elbow point" where SSE decreases more gradually.

This method helps identify the optimal K value for clustering.

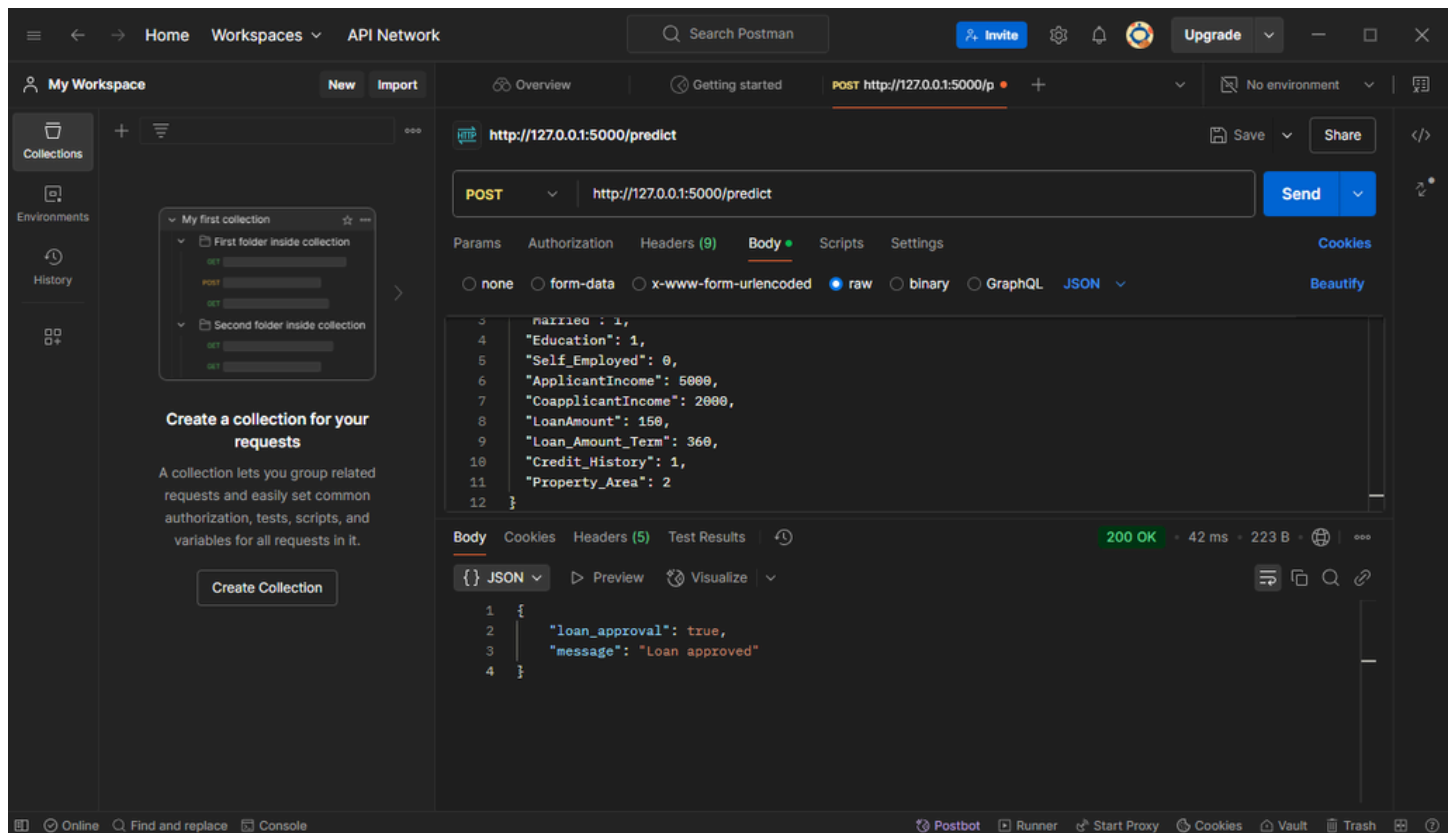# Sales Prediction Model :

1. Collect historical sales data

2. Preprocess data

3. Choose ML algorithm (e.g., Linear Regression, Random Forest)

4. Train and evaluate model

5. Deploy model for future sales predictions



Deployed locally by PostMan application and tested the prediction model