

**Sri Sivasubramaniya Nadar College of Engineering, Chennai**  
(An autonomous Institution affiliated to Anna University)

Degree & Branch	B.E. Computer Science & Engineering	Semester	V
Subject Code & Name	ICS1512 & Machine Learning Algorithms Laboratory		
Academic year	2025-2026 (Odd)	Batch:2023-2028	<b>Due date: 25/07/2025</b>

**Experiment 2 : Loan Amount Prediction using Linear Regression**

**Aim:**

To predict the loan amount using Linear Regression using historical data and measure the model's performance

**Libraries used:**

- Numpy
- Pandas
- Matplotlib
- Seaborn
- sklearn

**Mathematical Description**

In this experiment, Linear Regression is used to predict the loan sanction amount based on several input features.

The mathematical model for Linear Regression is represented as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + \varepsilon$$

where:

- $y$  is the dependent variable (Loan Sanction Amount),
- $x_1, x_2, \dots, x_n$  are independent variables (features such as Age, Income, Credit Score, etc.),
- $\beta_0$  is the intercept,
- $\beta_1, \beta_2, \dots, \beta_n$  are the coefficients of the features,
- $\varepsilon$  is the error term representing noise or unexplained variance.

The coefficients  $\beta$  are estimated by minimizing the Residual Sum of Squares (RSS), defined as:

$$RSS = \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \sum_{i=1}^m \left( y_i - \left( \beta_0 + \sum_{j=1}^n \beta_j x_{ij} \right) \right)^2$$

where  $m$  is the number of observations. The model is evaluated using the following metrics:

- **Mean Absolute Error (MAE)** – average absolute difference between actual and predicted values,
- **Mean Squared Error (MSE)** – average squared difference between actual and predicted values,
- **Root Mean Squared Error (RMSE)** – square root of MSE, providing error in original units,
- **R-squared ( $R^2$ )** – proportion of variance explained by the model,
- **Adjusted  $R^2$**  – adjusted for number of predictors to prevent overfitting.

## Objectives :

Apply Linear Regression to predict the loan amount sanctioned to users using the dataset provided, visualize and interpret the results to gain insights into the model performance

## Importing Required Libraries:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, KFold, learning_curve
from sklearn.preprocessing import StandardScaler, OneHotEncoder, OrdinalEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
```

## Loading the dataset:

```
''' Load data '''
df = pd.read_csv("C:/Users/Kavi/Downloads/train.csv/train.csv")
df.head()
```

## Output:

	Customer ID	Name	Gender	Age	Income (USD)	Income Stability	Profession	Type of Employment	Location	Loan Amount Request (USD)	...	Credit Score	No. of Defaults	Haz Active Credit Card	Property ID	Property Age	Property Type	Property Location	Co-Applicant	Property Price	Loan Sanction Amount (USD)
0	C-36995	Frederica Shealy	F	56	1933.05	Low	Working	Sales staff	Semi-Urban	72809.58	...	809.44	0	NaN	746	1933.05	4	Rural	1	119933.46	54607.18
1	C-13999	America Calderone	M	32	4952.91	Low	Working	NaN	Semi-Urban	49637.47	...	780.40	0	Unprocessed	608	4652.91	2	Rural	1	247991.80	37499.38
2	C-3770	Rosetta Nema	F	65	988.19	High	Pensioner	NaN	Semi-Urban	45530.04	...	833.15	0	Unprocessed	548	988.19	2	Urban	0	72460.53	38474.43
3	C-26480	Zoe Chilly	F	65	NaN	High	Pensioner	NaN	Rural	80067.82	...	832.70	1	Unprocessed	890	NaN	2	Semi-Urban	1	32441.51	50040.54
4	C-23459	Afton Venema	F	31	2614.77	Low	Working	High skill tech staff	Semi-Urban	113858.80	...	745.55	1	Active	715	2614.77	4	Semi-Urban	1	208567.91	74008.28

5 rows × 24 columns

5 rows x 24 columns

Figure 1: Displaying first 5 rows

## Exploratory Data Analysis:

- **Descriptive statistics:**

```
print(df.info())  
print(df.describe())
```

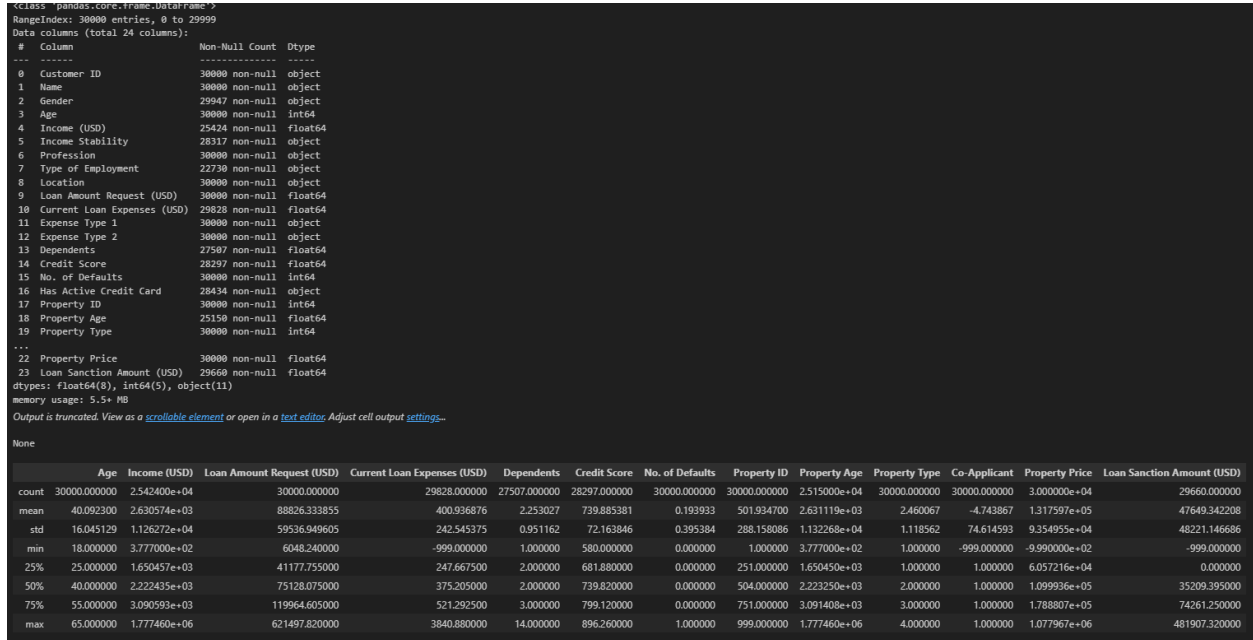


Figure 2: Descriptive Statistics

**Descriptive statistics** summarize and organize the characteristics of a dataset, providing insights into its **central tendency, dispersion, and distribution**. Key metrics such as **mean, median, mode, standard deviation, minimum, and maximum values** offer a foundational understanding of each feature's behavior. These statistics are essential for identifying data quality issues, such as **skewness, outliers, or missing values**, and guide appropriate preprocessing steps before advanced modeling.

- **Visualizing numerical features distribution:**

```
numeric_cols = df.select_dtypes(include=['int64', 'float64']).columns  
df[numeric_cols].hist(figsize=(15, 10), bins=30)  
plt.suptitle("Initial Distribution of Numerical Features (Raw Data)")  
plt.show()
```

Initial Distribution of Numerical Features (Raw Data)

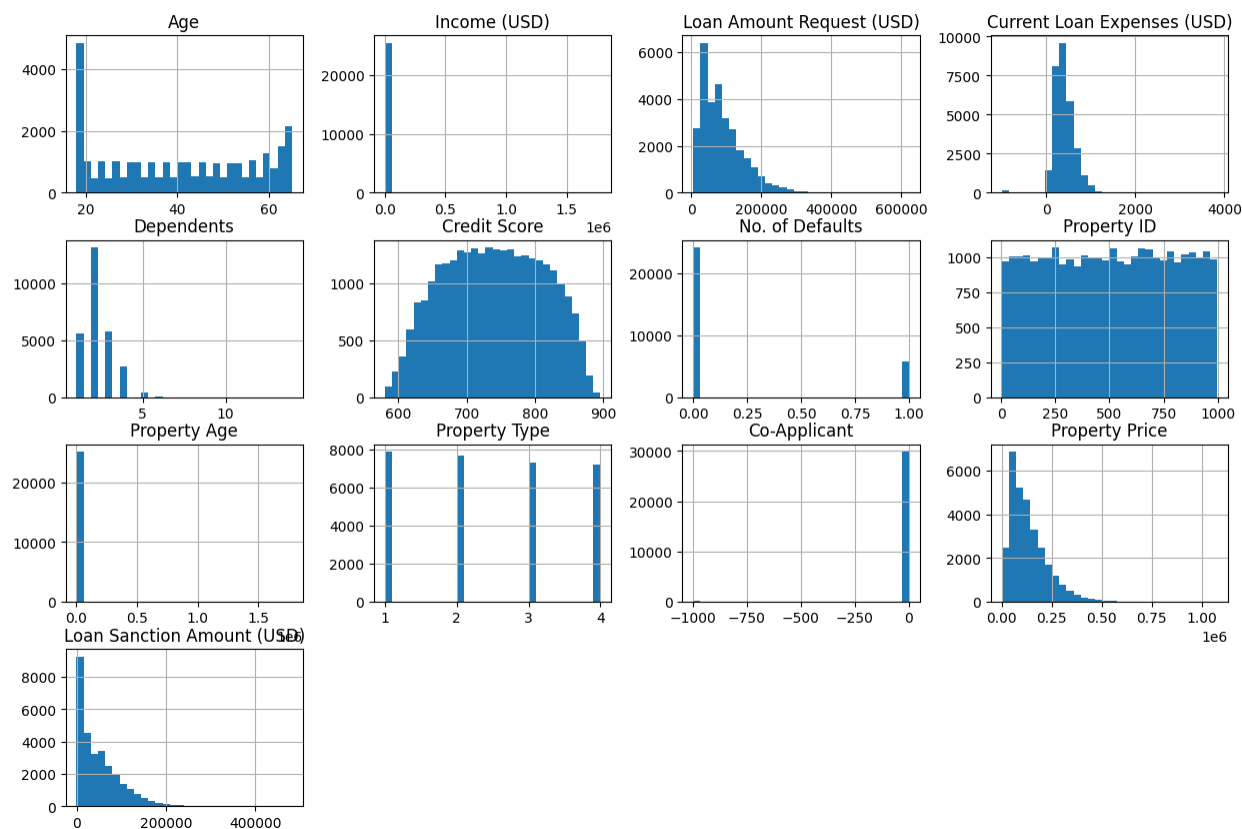


Figure 3: Numerical features distribution

### Summary:

- Age: Bimodal, possible outliers.
- Income: Strong right skew, heavy outliers.
- Loan Amount Request/Sanction: Right-skewed, lower values dominate.
- Loan Expenses: Nearly normal, mild skew.
- Dependents: Discrete, mode at 0.
- Credit Score: Bell-shaped, well-centered.
- Defaults / Co-Applicant: Binary, mostly 0s.
- Property Age: Right-skewed, most are new.
- Property Type: Categorical (4 types).
- Property Price: Right-skewed, wide range.

- **Removing unneeded columns and plotting correlation heatmap:**

```
"""Removing unneeded columns"""
```

```
df.drop(['Customer ID', 'Name', 'Property ID'], axis=1, inplace=True)
```

```
"""Correlation Heatmap"""
```

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

```
sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm')
```

```
plt.title('Correlation Heatmap')
```

```
plt.show()
```

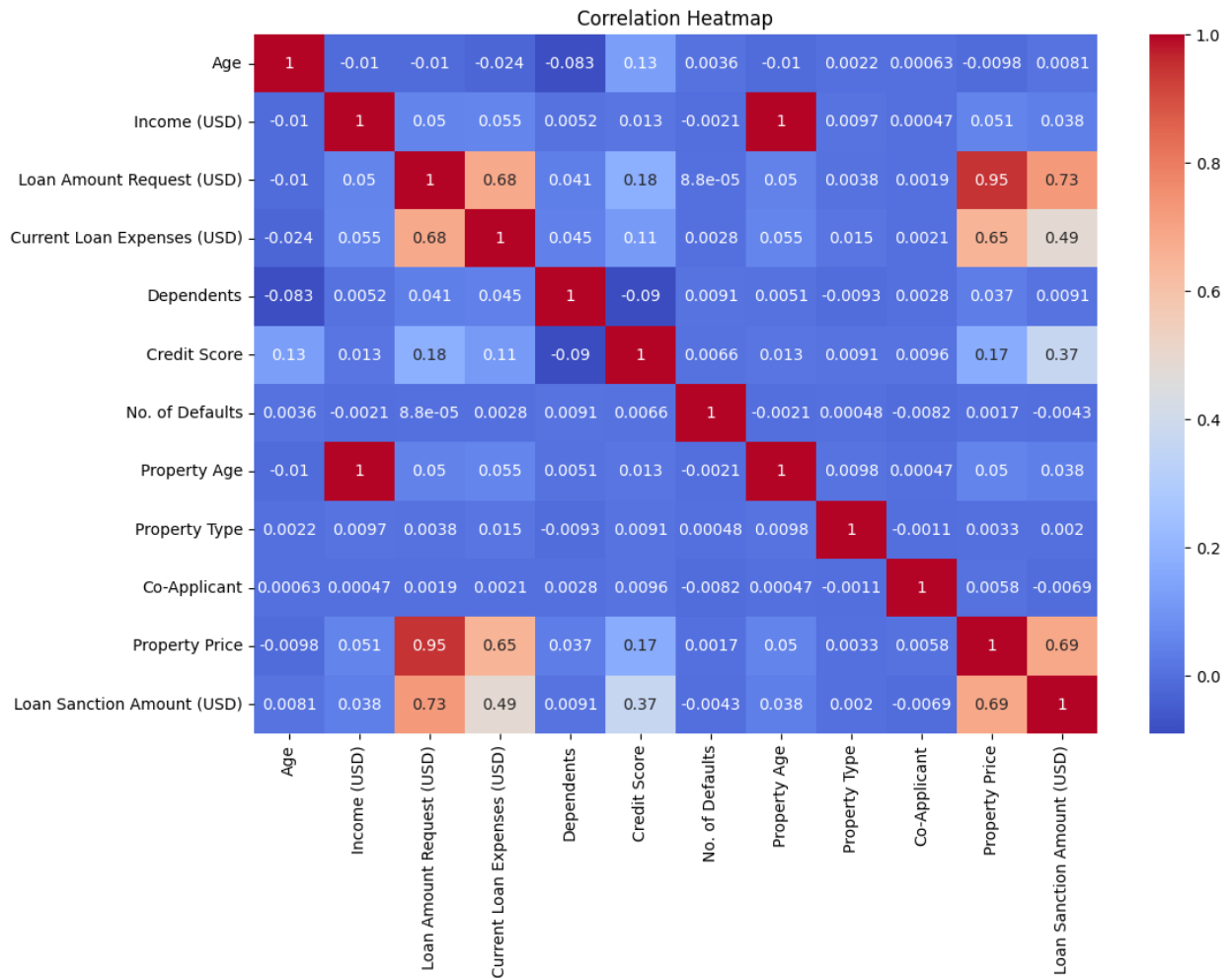


Figure 4: Correlation Heatmap

**Highly Correlated features with Target:** Loan Amount Request, Property price, Loan Expenses

**Moderate Correlation:** Current Loan Expenses with Credit Score

**Negligibly Correlated features with Target :** Age, Income, Dependents, Property Age, Co-Applicant, No.of Defaults

- **Outlier Detection and Removal**

```
''' Boxplots for outlier detection '''
num_cols = df.select_dtypes(include=['int64','float64']).columns
for col in num_cols:
    plt.figure(figsize=(6, 2))
    sns.boxplot(data=df, x=col)
    plt.title(f"Boxplot for col")
    plt.show()

''' Outlier removal using IQR method '''
Q1 = df[num_cols].quantile(0.25)
Q3 = df[num_cols].quantile(0.75)
IQR = Q3 - Q1
df = df[((df[num_cols] < (Q1 - 1.5 * IQR)) | (df[num_cols] > (Q3 + 1.5 * IQR))).any(axis =
1)]
```

- **Data Cleaning and Preprocessing:**

```
''' Recompute column types after dropping '''
num_cols = df.select_dtypes(include=['int64', 'float64']).columns.tolist()
cat_cols = df.select_dtypes(include='object').columns.tolist()
ordinal_cols = ['Income Stability', 'Location', 'Has Active Credit Card', 'Expense Type
1', 'Expense Type 2']
cat_cols = [col for col in cat_cols if col not in ordinal_cols]

num_cols.remove('Loan Amount Request (USD)')

X = df.drop('Loan Amount Request (USD)', axis=1)
y = df['Loan Amount Request (USD)']

numeric_transformer = Pipeline(steps=[('imputer', SimpleImputer(strategy='mean')), ('scaler',
StandardScaler())])

categorical_transformer = Pipeline(steps=[('imputer', SimpleImputer(strategy='most_frequent')), ('encoder',
OneHotEncoder(handle_unknown='ignore'))])

ordinal_transformer = Pipeline(steps=[('imputer', SimpleImputer(strategy='most_frequent')),
('ordinal', OrdinalEncoder())])

''' Update categorical transformer'''
preprocessor = ColumnTransformer(transformers=[('num', numeric_transformer, num_cols), ('ord',
ordinal_transformer, ordinal_cols), ('cat', categorical_transformer, list(set(cat_cols) - set(ordinal_cols)))])
```

- **Learning Curve** : Shows model performance with increasing data

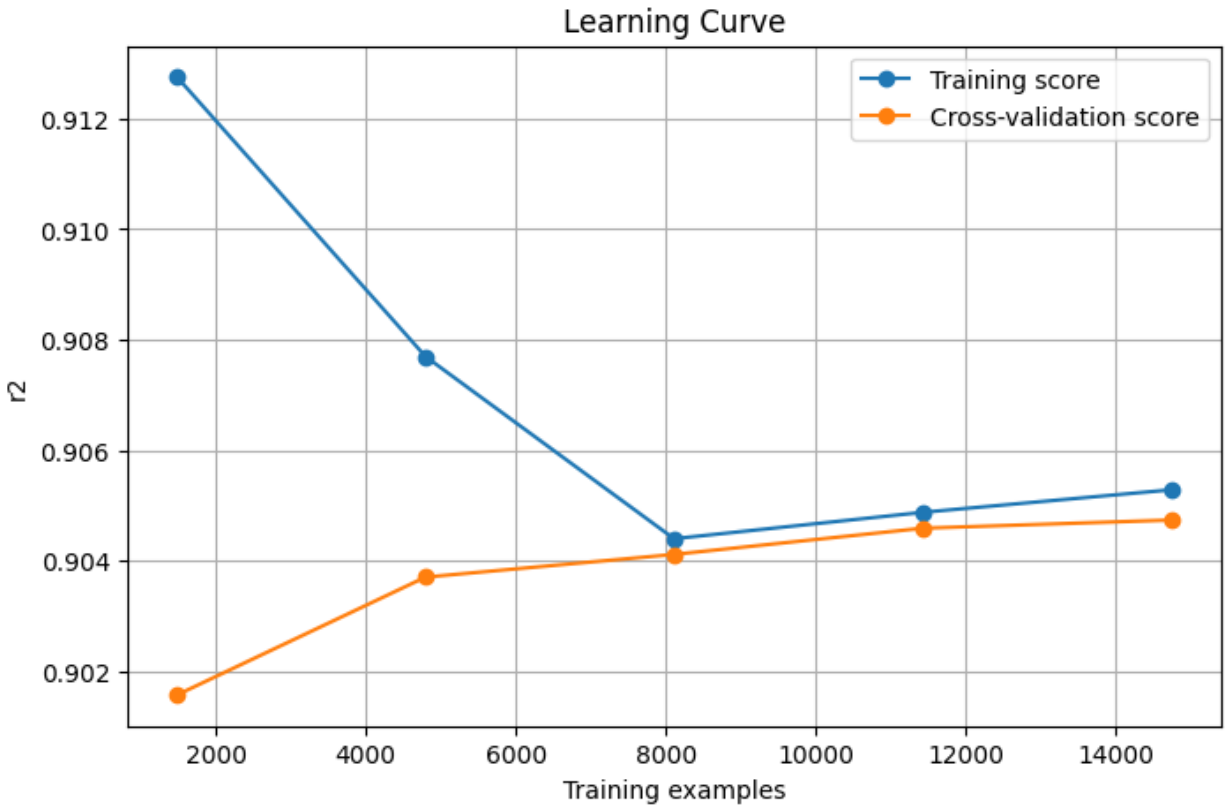


Figure 5: Training vs Cross Validation

**Low Bias** : Small gap b/w training and validation = not much underfitting

**Low Variance** : Training and validation scores shrink as data increases = Not overly sensitive to training data

**Model Generalization** : Both lines converge around 0.904  $R^2$  = Good Generalization

- **K-Fold Cross Validation:**

```

''' Define the full pipeline '''
pipeline = Pipeline(steps=[('preprocessor', preprocessor), ('regressor', LinearRegression())])

''' K-Fold CV '''
kfold = KFold(n_splits=5, shuffle=True, random_state=42)
scores = []
mse_scores = []
mae_scores = []
rmse_scores = []
for train_idx, val_idx in kfold.split(X):
    X_train_fold, X_val_fold = X.iloc[train_idx], X.iloc[val_idx]
    y_train_fold, y_val_fold = y.iloc[train_idx], y.iloc[val_idx]

    pipeline.fit(X_train_fold, y_train_fold)
    y_pred = pipeline.predict(X_val_fold)

    mae = mean_absolute_error(y_val_fold, y_pred)
    mse = mean_squared_error(y_val_fold, y_pred)
    rmse = np.sqrt(mse)
    r2 = r2_score(y_val_fold, y_pred)
    print(f'MAE : {mae:.2f} , MSE : {mse:.2f} , RMSE : {rmse:.2f} , R2 : {r2:.2f}')
    scores.append(r2)
    mae_scores.append(mae)
    mse_scores.append(mse)
    rmse_scores.append(rmse)

print(f'Average R2 across 5 folds: np.mean(scores):.4f')
print(f'Average MAE across 5 folds: np.mean(mae_scores):.4f')
print(f'Average MSE across 5 folds: np.mean(mse_scores):.4f')
print(f'Average RMSE across 5 folds: np.mean(rmse_scores):.4f')

```

Table 1: Cross-Validation Results ( $K = 5$ )

Fold	MAE	MSE	RMSE	R <sup>2</sup> Score
Fold 1	10326.97	247284558.37	15725.28	0.90
Fold 2	10317.09	242315282.95	15566.48	0.90
Fold 3	10270.01	222268480.96	14908.67	0.90
Fold 4	9921.38	220742234.20	14857.40	0.91
Fold 5	9955.12	203990745.53	14282.53	0.92
<b>Average</b>	10158.1119	227320260.4028	15068.0729	0.9047



- Train/Validate/Test:

```
# Step 1: Train/Test Split (80% train, 20% test)
X_temp, X_test, y_temp, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Step 2: Further split training into Train/Validation (75% train, 25% val of the 80%)
X_train, X_val, y_train, y_val = train_test_split(X_temp, y_temp, test_size=0.25, random_state=42)
# Final Split: 60% train, 20% val, 20% test

# Step 3: Fit the pipeline on training data
pipeline.fit(X_train, y_train)

# Step 4: Predict on validation data
y_val_pred = pipeline.predict(X_val)

# Step 5: Validation metrics
val_mse = mean_squared_error(y_val, y_val_pred)
val_mae = mean_absolute_error(y_val, y_val_pred)
val_rmse = np.sqrt(val_mse)
val_r2 = r2_score(y_val, y_val_pred)

print("----- Validation Metrics -----")
print(f"Validation MAE: {val_mae:.2f}")
print(f"Validation RMSE: {val_rmse:.2f}")
print(f"Validation R2: {val_r2:.4f}")

# Step 6: Final training on combined train+val for final test eval (optional)
# pipeline.fit(X_temp, y_temp) # Uncomment if you want final model retrained on more data

# Step 7: Predict on test data
y_test_pred = pipeline.predict(X_test)

# Step 8: Test metrics
test_mse = mean_squared_error(y_test, y_test_pred)
test_mae = mean_absolute_error(y_test, y_test_pred)
test_rmse = np.sqrt(test_mse)
test_r2 = r2_score(y_test, y_test_pred)

print("\n----- Test Metrics -----")
print(f"Test MAE: {test_mae:.2f}")
print(f"Test RMSE: {test_rmse:.2f}")
print(f"Test R2: {test_r2:.4f}")
```

Figure 6: Training, Validating and Testing

Phase	MAE	MSE	RMSE	R <sup>2</sup> Score
Validation	10127.45	229494187.86	15149.07	0.9035
Testing	10333.42	248208063.75	15754.62	0.8977

Table 2: Validation/ Test Metrics

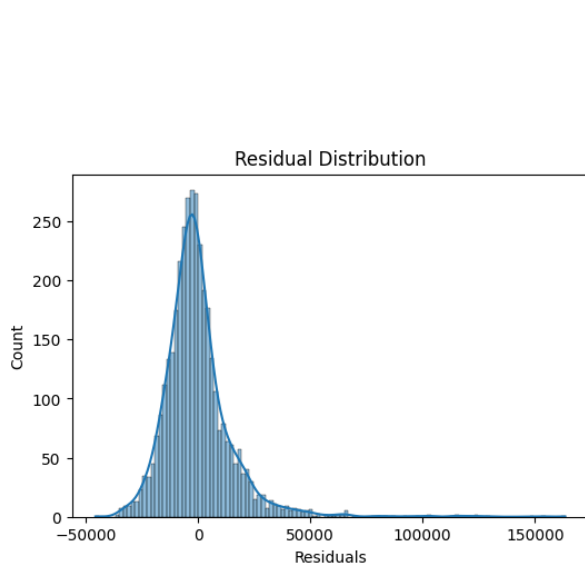
- **Residual Plot and Model Performance assessment:**

```
''' Residual Plot '''
```

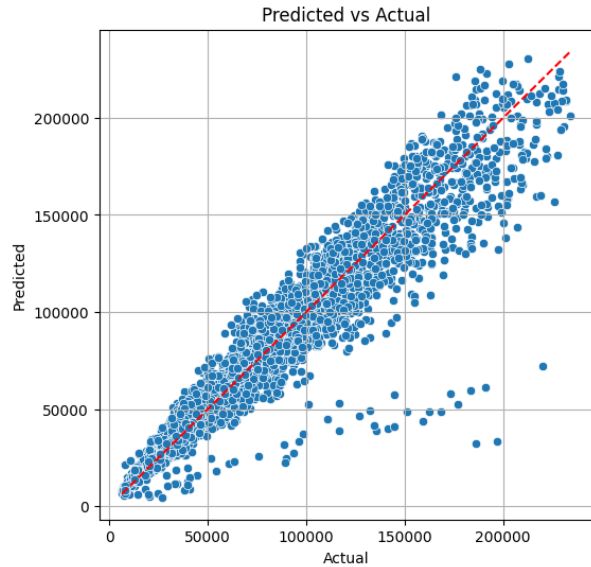
```
residuals = y_test - y_test_pred
plt.figure(figsize=(6,4))
sns.histplot(residuals, kde=True)
plt.title("Residual Distribution")
plt.xlabel("Residuals")
plt.show()
```

```
''' Predicted vs Actual '''
```

```
plt.figure(figsize=(6,6))
sns.scatterplot(x=y_test, y=y_test_pred)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r-')
plt.title("Predicted vs Actual")
plt.xlabel("Actual")
plt.ylabel("Predicted")
plt.grid()
plt.show()
```



(a) Caption for image 1



(b) Caption for image 2

Figure 7: Model Performance Assessment

#### INTERPRETATION:

- **Residual Distribution:** Looks reasonably tight and Gaussian-ish(though slightly skewed) = **Errors are small and almost normally distributed**
- **Predicted vs Actual:** Points are closely scattered around the best fit line. Visible spread at high values but general prediction trend is tight = **Good generalization with minor underfitting** at high targets

Table 3: Summary of Results for Loan Amount Prediction

Description	Result
Dataset Size (after preprocessing)	18,438
Train/Validation/Test Split Ratio	60/20/20
Feature(s) Used for Prediction	Gender, Age, Income (USD), Income Stability, Location, Loan Amount Request (USD), Current Loan Expenses (USD), Expense Type 1, Expense Type 2, Dependents, Credit Score, No. of Defaults, Has Active Credit Card, Property Age, Property Type, Property Location, Co-Applicant, Property Price, Loan Sanction Amount (USD)
Model Used	Linear Regression
Cross-Validation Used?	Yes
If Yes, Number of Folds (K)	5
Reference to CV Results Table	Table 1
Mean Absolute Error (MAE) on Test Set	10,333.42
Mean Squared Error (MSE) on Test Set	248,208,063.75
Root Mean Squared Error (RMSE) on Test Set	15754.62
R <sup>2</sup> Score on Test Set	0.8977
Adjusted R <sup>2</sup> Score on Test Set	0.9057
Most Influential Feature(s)	Income (USD), Loan Amount Request (USD), Property price, Loan Expenses, Credit Score
Observations from Residual Plot	Residuals appear tight and Gaussian-like (centered at 0, slight right tail), suggesting errors are small and normally distributed.
Interpretation of Predicted vs Actual Plot	Data points are tightly clustered around the best-fit line, indicating good generalization with slight underfitting.
Any Overfitting or Underfitting Observed?	Minor underfitting at higher loan amounts.
Justification	Spread in residuals increases at higher target values.

## Best Practices

- **Data Preprocessing:** Handle missing values carefully (drop or impute), and remove irrelevant columns such as IDs and names that do not contribute to prediction and removing the outliers.
- **Feature Engineering:** Create new meaningful features (e.g., total income) and apply transformations (e.g., log transformation for skewed data) to improve model performance.
- **Scaling and Encoding:** Use scaling (e.g., `StandardScaler`) for numerical features and one-hot encoding for categorical features to prepare data for linear regression.
- **Train-Validate-Test Split:** Use proper splits (e.g., 60/20/20) and consider cross-validation to ensure the model generalizes well and to prevent overfitting.
- **Model Evaluation:** Use multiple metrics (MAE, MSE, RMSE,  $R^2$ ) to assess different aspects of model performance.
- **Residual Analysis:** Analyze residual plots to detect model bias or heteroscedasticity and decide if further feature engineering or alternative models are needed.

## Learning Outcomes

- Able to identify and handling missing values, outliers and further scaling operations
- Able to interpret the variable distributions and relationships between the target and remaining features
- Able to train and evaluate linear regression model using CV and multiple performance metrics
- Able to assess model behaviour using residual and predicted vs actual plots to detect overfitting/underfitting patterns