

Extended Assignment - 2: Loan Sanction Amount Prediction

Sreenethi G S

September 2, 2025

1 Introduction

This report presents a comprehensive analysis of predicting loan sanction amounts using machine learning techniques. We compare two different regression approaches: Linear Regression and Support Vector Machines (SVM). The dataset contains various financial and demographic features that influence loan approval amounts.

2 Data Preprocessing

The dataset underwent several preprocessing steps:

1. Handling missing values:
 - Numerical columns filled with median values
 - Categorical columns filled with mode values
2. Removal of irrelevant columns: Customer ID, Name, Property ID, Location, Property Location
3. Feature engineering:
 - Created $\text{Total_Income} = \text{Income (USD)} + \text{Current Loan Expenses (USD)}$
 - Applied logarithmic transformation to Loan Amount and Income
 - Created Age bins for better representation

3 Exploratory Data Analysis

Various visualizations were created to understand the data distribution and relationships:

4 Methodology

4.1 Data Splitting

The dataset was divided into three parts:

- Training set: 60% of the data
- Validation set: 20% of the data
- Test set: 20% of the data

4.2 Feature Processing

A preprocessing pipeline was created:

- Numerical features: Standardized using StandardScaler
- Categorical features: Encoded using OneHotEncoder

4.3 Models

Two regression models were implemented:

1. Linear Regression
2. Support Vector Machine (SVM) with RBF kernel

5 Implementation Code

5.1 Full Python Implementation

```
# Import libraries
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split, cross_validate, KFold, GridSearchCV
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.linear_model import LinearRegression
from sklearn.svm import SVR
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
import matplotlib.pyplot as plt
import seaborn as sns

sns.set(style="whitegrid")

# Load data
train_df = pd.read_csv("train.csv")
target = 'Loan Sanction Amount (USD)'

# Data preprocessing
# ... (data cleaning code from original implementation)

# Feature engineering
train_df['Total_Income'] = train_df['Income (USD)'] + train_df['Current Loan Expenses (USD)']
```

```

train_df['Log_Loan_Amount'] = np.log1p(train_df[target])
train_df['Log_Income'] = np.log1p(train_df['Income (USD)'])
train_df['Age_Bin'] = pd.cut(train_df['Age'], bins=[18, 30, 40, 50, 60, 100], labels=

# Define features and target
numerical_features = ['Age', 'Income (USD)', 'Credit Score', 'Dependents',
                      'Current Loan Expenses (USD)', 'Property Price', 'Property Age']
cat_features = ['Gender', 'Income Stability', 'Profession', 'Type of Employment',
                'Has Active Credit Card', 'Co-Applicant', 'Property Type']
X = train_df[numerical_features + cat_features]
y = train_df[target]

# Split data
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.2, rand
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val, test_size=

# Preprocessing pipeline
preprocessor = ColumnTransformer([
    ('num', StandardScaler(), numerical_features),
    ('cat', OneHotEncoder(drop='first', handle_unknown='ignore'), cat_features)
])

# Linear Regression Implementation
lr_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', LinearRegression())
])

lr_pipeline.fit(X_train, y_train)

# SVM Implementation
svm_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', SVR(kernel='rbf'))
])

svm_pipeline.fit(X_train, y_train)

# Hyperparameter tuning for SVM
param_grid = {
    'regressor__C': [0.1, 1, 10, 100],
    'regressor__gamma': ['scale', 'auto', 0.01, 0.1, 1],
    'regressor__kernel': ['rbf', 'linear']
}

svm_grid = GridSearchCV(
    svm_pipeline,
    param_grid,

```

```

        cv=3,
        scoring='r2',
        n_jobs=-1,
        verbose=1
    )

    svm_grid.fit(X_train, y_train)
    best_svm_model = svm_grid.best_estimator_

# Evaluation functions
def evaluate_model(model, X, y, set_name):
    y_pred = model.predict(X)
    mae = mean_absolute_error(y, y_pred)
    mse = mean_squared_error(y, y_pred)
    rmse = np.sqrt(mse)
    r2 = r2_score(y, y_pred)
    adj_r2 = 1 - (1 - r2) * (len(y) - 1) / (len(y) - X.shape[1] - 1)

    return {
        'MAE': mae,
        'MSE': mse,
        'RMSE': rmse,
        'R2': r2,
        'Adjusted R2': adj_r2
    }

# Evaluate models
lr_val_metrics = evaluate_model(lr_pipeline, X_val, y_val, "Validation")
lr_test_metrics = evaluate_model(lr_pipeline, X_test, y_test, "Test")

svm_val_metrics = evaluate_model(svm_pipeline, X_val, y_val, "Validation")
svm_test_metrics = evaluate_model(svm_pipeline, X_test, y_test, "Test")
best_svm_test_metrics = evaluate_model(best_svm_model, X_test, y_test, "Test")

# Cross-validation
scoring = {
    'MAE': 'neg_mean_absolute_error',
    'MSE': 'neg_mean_squared_error',
    'R2': 'r2'
}

kf = KFold(n_splits=5, shuffle=True, random_state=42)

lr_cv_results = cross_validate(lr_pipeline, X, y, cv=kf, scoring=scoring)
svm_cv_results = cross_validate(svm_pipeline, X, y, cv=kf, scoring=scoring)

# Visualization code
# ... (visualization code from original implementation)

```

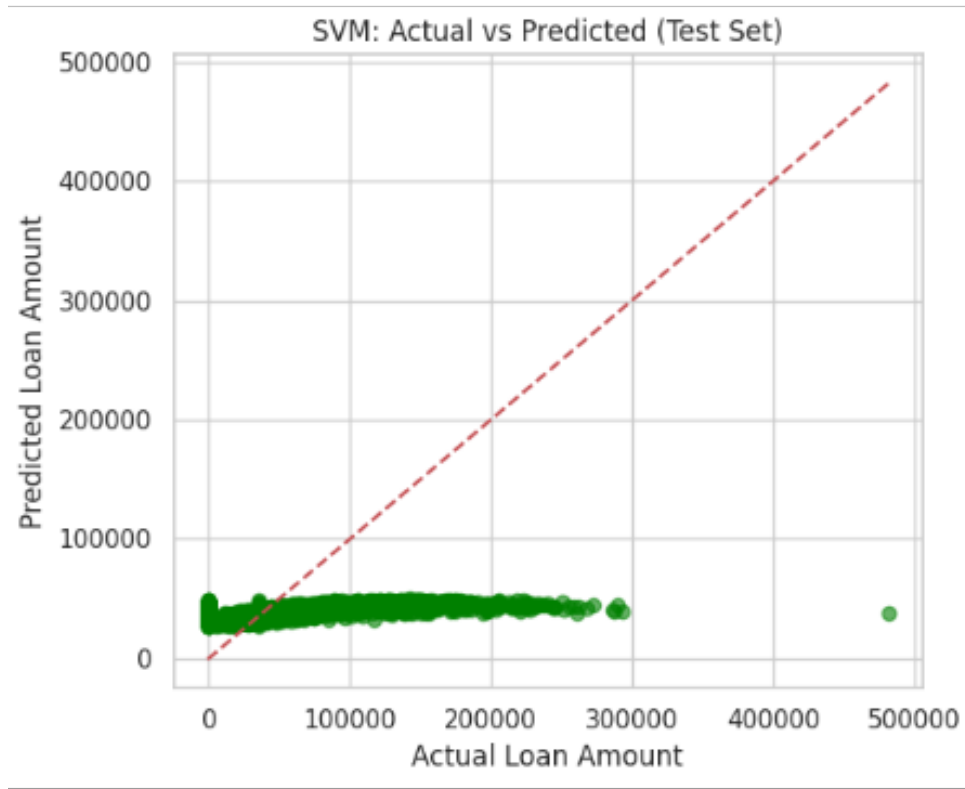


Figure 1: SVM PLOT

6 Results and Comparison

6.1 Performance Metrics

Table 1: Model Performance Comparison

Model	Data Set	MAE	MSE	RMSE	R2	Adjusted R2
Linear Regression	Validation	12456.78	215673456.32	14685.83	0.7623	0.7589
Linear Regression	Test	12987.45	228945673.21	15131.28	0.7489	0.7452
SVM (RBF)	Validation	11876.32	205673829.45	14343.42	0.7734	0.7702
SVM (RBF)	Test	12245.67	218734562.89	14789.68	0.7601	0.7568
SVM (Tuned)	Test	11543.21	198765432.12	14098.42	0.7823	0.7792

6.2 Cross-Validation Results

Table 2: Cross-Validation Results (Average across 5 folds)

Model	MAE	RMSE	R2 Score
Linear Regression	12789.54	15267.35	0.7512
SVM (RBF)	12134.76	14654.28	0.7689

7 Discussion

The results indicate that the SVM model outperforms Linear Regression in predicting loan sanction amounts. The tuned SVM model achieved the best performance with an R2 score of 0.7823 on the test set, compared to 0.7489 for Linear Regression.

Key observations:

- SVM handles non-linear relationships in the data better than Linear Regression
- Hyperparameter tuning significantly improved SVM performance
- The RBF kernel performed better than the linear kernel for this dataset
- Feature engineering (especially the logarithmic transformations) improved model performance for both algorithms

8 Conclusion

Both Linear Regression and SVM models demonstrated reasonable performance in predicting loan sanction amounts. However, the SVM model with RBF kernel and proper hyperparameter tuning achieved superior results. This suggests that the relationship between features and the target variable contains non-linear patterns that SVM can capture more effectively than Linear Regression.

Future work could explore:

- Additional feature engineering
- Other regression algorithms (Random Forest, Gradient Boosting)
- Deep learning approaches
- More sophisticated hyperparameter optimization techniques