



JSPM's
**RAJARSHI SHAHU COLLEGE OF
ENGINEERING
TATHAWADE, PUNE-33**



DEPARTMENT OF COMPUTER ENGINEERING

Case Study on House Price Prediction using Linear, Ridge and Polynomial Regression.

Submitted by

Name: Rushikesh Gajanan Bobade
Roll No: CS3146
PRN: RBT23CS049
Class & Division: TY - A
Course: Machine Learning

1. Abstract

This case study aims to predict house prices using multiple regression techniques—**Linear Regression**, **Ridge Regression**, and **Polynomial Regression**—and to analyze how regularization and non-linear feature expansion affect model performance.

Using the **Ames Housing dataset**, the study compares models based on **RMSE (Root Mean Squared Error)**, **MAE (Mean Absolute Error)**, and **R² score**.

Results show that **Ridge Regression** improves generalization over the basic linear model, while **Polynomial Regression** captures non-linear patterns but risks overfitting if not regularized.

2. Problem Statement

Real estate pricing depends on numerous features like location, area, number of rooms, year built, and materials.

The challenge is to build an accurate regression model that estimates house prices based on these attributes.

3. Dataset Description

Dataset Used: Ames Housing Dataset

Source: Kaggle - House Prices: Advanced Regression Techniques

Basic Info:

- **Target Variable:** SalePrice
 - **Features:** 80+ (numeric and categorical)
 - **Rows:** ~1460
 - **Missing Values:** Present in several features (LotFrontage, GarageYrBlt, etc.)
-

4. Data Preprocessing

Steps:

1. **Handle missing values** – median imputation for numeric, mode for categorical features.
2. **Encode categorical features** – one-hot encoding using OneHotEncoder.

-
3. **Scale numeric features** – standardization for models sensitive to feature scales.
 4. **Train-test split** – 80/20 division with fixed random seed.
-

5. Model Development

We use three regression models:

1. **Linear Regression:** Baseline model without regularization.
 2. **Ridge Regression:** Adds L2 penalty to reduce overfitting.
 3. **Polynomial Regression:** Introduces non-linear relationships via polynomial feature expansion (degree=2).
-

6. Model Evaluation Metrics

Metric	Description
RMSE	Root Mean Squared Error (lower is better)
MAE	Mean Absolute Error
R ²	Coefficient of Determination (closer to 1 = better fit)

7. Full Functional Code

```
# =====  
  
# Case Study 2: House Price Prediction using Linear, Ridge & Polynomial Regression  
  
# =====  
  
# -----  
  
# 1. Import Libraries  
  
# -----  
  
import pandas as pd  
  
import numpy as np  
  
from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score  
  
from sklearn.linear_model import LinearRegression, Ridge  
  
from sklearn.preprocessing import StandardScaler, OneHotEncoder,  
PolynomialFeatures  
  
from sklearn.compose import ColumnTransformer  
  
from sklearn.pipeline import Pipeline  
  
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score  
  
from sklearn.impute import SimpleImputer  
  
import matplotlib.pyplot as plt  
  
import seaborn as sns  
  
import joblib
```

```
RANDOM_STATE = 42
```

```
# -----
```

```
# 2. Load Dataset
```

```
# -----
```

```
# Download dataset from Kaggle (train.csv)
```

```
df = pd.read_csv("train.csv")
```

```
print("Dataset loaded successfully.")
```

```
print(df.head())
```

```
# -----
```

```
# 3. Select Features and Target
```

```
# -----
```

```
y = df['SalePrice']
```

```
X = df.drop(columns=['SalePrice', 'Id']) # Remove ID & target
```

```
# Identify numeric and categorical columns
```

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

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

```
# -----
```

4. Preprocessing Pipelines

```
# -----
```

```
numeric_transformer = Pipeline(steps=[
```

```
    ('imputer', SimpleImputer(strategy='median')),
```

```
    ('scaler', StandardScaler())
```

```
])
```

```
categorical_transformer = Pipeline(steps=[
```

```
    ('imputer', SimpleImputer(strategy='most_frequent')),
```

```
    ('encoder', OneHotEncoder(handle_unknown='ignore'))
```

```
])
```

```
preprocessor = ColumnTransformer(transformers=[
```

```
    ('num', numeric_transformer, num_cols),
```

```
    ('cat', categorical_transformer, cat_cols)
```

```
])
```

```
# -----  
  
# 5. Define Models  
  
# -----  
  
pipe_linear = Pipeline(steps=[  
  
    ('preprocessor', preprocessor),  
  
    ('regressor', LinearRegression())  
  
])  
  
  
  
pipe_ridge = Pipeline(steps=[  
  
    ('preprocessor', preprocessor),  
  
    ('regressor', Ridge(random_state=RANDOM_STATE))  
  
])  
  
  
  
pipe_poly = Pipeline(steps=[  
  
    ('preprocessor', preprocessor),  
  
    ('poly', PolynomialFeatures(degree=2, include_bias=False)),  
  
    ('scaler', StandardScaler()),  
  
    ('regressor', Ridge(random_state=RANDOM_STATE))
```

)

6. Train-Test Split

X_train, X_test, y_train, y_test = train_test_split(

X, y, test_size=0.2, random_state=RANDOM_STATE

)

7. Train and Evaluate Models

def evaluate_model(model, X_train, X_test, y_train, y_test):

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

rmse = np.sqrt(mean_squared_error(y_test, y_pred))

mae = mean_absolute_error(y_test, y_pred)

r2 = r2_score(y_test, y_pred)

return rmse, mae, r2

```
results = []

models = {

    "Linear Regression": pipe_linear,

    "Ridge Regression": pipe_ridge,

    "Polynomial Regression": pipe_poly

}

for name, model in models.items():

    print(f"\nTraining {name}...")

    rmse, mae, r2 = evaluate_model(model, X_train, X_test, y_train, y_test)

    results.append((name, rmse, mae, r2))

    print(f"RMSE: {rmse:.2f}, MAE: {mae:.2f}, R2: {r2:.3f}")

results_df = pd.DataFrame(results, columns=["Model", "RMSE", "MAE", "R2"])

print("\nModel Performance Summary:")

print(results_df)

# -----
```

```
# 8. Hyperparameter Tuning for Ridge
```

```
# -----
```

```
param_grid = {'regressor__alpha': [0.1, 1.0, 10.0, 100.0]}
```

```
grid = GridSearchCV(pipe_ridge, param_grid, cv=5,  
scoring='neg_root_mean_squared_error', n_jobs=-1)
```

```
grid.fit(X_train, y_train)
```

```
print(f"\nBest Ridge Alpha: {grid.best_params_['regressor__alpha']}")
```

```
print(f"Best CV RMSE: {-grid.best_score_:.2f}")
```

```
# -----
```

```
# 9. Visualizations
```

```
# -----
```

```
plt.figure(figsize=(8, 5))
```

```
sns.barplot(x='Model', y='RMSE', data=results_df)
```

```
plt.title('Model Comparison (RMSE)')
```

```
plt.ylabel('Root Mean Squared Error')
```

```
plt.show()
```

```
# Predicted vs Actual for best model
```

```
best_model_name = results_df.sort_values('RMSE').iloc[0]['Model']
```

```
best_model = models[best_model_name]

best_model.fit(X_train, y_train)

y_pred_best = best_model.predict(X_test)

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

sns.scatterplot(x=y_test, y=y_pred_best)

plt.xlabel("Actual SalePrice")

plt.ylabel("Predicted SalePrice")

plt.title(f'{best_model_name} – Predicted vs Actual')

plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', ls='--')

plt.show()
```

```
# -----
```

```
# 10. Save Best Model
```

```
# -----
```

```
joblib.dump(best_model, 'best_house_price_model.joblib')
```

```
print(f"\n ✅ Best model ({best_model_name}) saved as  
'best_house_price_model.joblib'")
```

8. Results Summary

Model	RMSE	MAE	R ²
Linear Regression	~33000	~23000	0.85
Ridge Regression	~31000	~21500	0.87
Polynomial Regression	~29500	~21000	0.89

(Results vary slightly depending on dataset version and preprocessing)

9. Discussion

- **Linear Regression** provides a simple, interpretable baseline but may underfit.
 - **Ridge Regression** reduces overfitting and stabilizes coefficients using L2 regularization.
 - **Polynomial Regression** captures non-linear relationships but risks overfitting unless regularized (hence Ridge used).
 - The **best performing model is Polynomial Ridge Regression** with degree=2.
 -
-

10. Conclusion

Ridge regularization and polynomial features significantly improved predictive performance over standard linear regression.

Further enhancements could include:

- Feature selection using Recursive Feature Elimination (RFE)
 - Outlier treatment and log-transforming SalePrice
 - Using ensemble models (Random Forest, Gradient Boosting)
-

11. Appendix

Dependencies:

pip install pandas numpy scikit-learn matplotlib seaborn joblib

Model Saved As:

best_house_price_model.joblib

Dataset Required: train.csv (from Kaggle House Prices competition)