

Report: Comparison of Linear Regression, Lasso, and Ridge Regularization Techniques

1. Methodology:

The goal of this analysis was to evaluate the impact of regularization techniques, specifically Lasso and Ridge regression, on the performance of a predictive model for housing prices using the California Housing dataset. The following steps were as follows:

a. Data Preprocessing:

- i. One-hot encoding was applied to the categorical feature `ocean_proximity`.
- ii. The dataset was split into training and testing sets.
- iii. Features were normalized using `StandardScaler` to ensure uniform scaling.

b. Models Implemented:

- i. Linear Regression: Used as the baseline model.
- ii. Lasso Regression: Applied to perform feature selection by shrinking coefficients of less important features to zero.
- iii. Ridge Regression: Applied to reduce the impact of multicollinearity by shrinking coefficients uniformly without eliminating any features.

c. Evaluation Metrics:

- i. Mean Squared Error (MSE): Measures the average squared difference between predicted and actual values.
- ii. R-squared: Measures the proportion of variance explained by the model.

2. Findings:

	MSE	R2
Linear Regression	4.945e+09	0.638382
Lasso Regression	4.948e+09	0.638140
Ridge Regression	4.953e+09	0.637819

a. Linear Regression:

- i. Achieved the lowest MSE and highest R-squared among the three models.

- ii. However, it lacks robustness to multicollinearity and is prone to overfitting if irrelevant features are present.

b. Lasso Regression:

- i. The MSE was slightly higher, and R-squared was marginally lower compared to Linear Regression.
- ii. The regularization term α caused the coefficients of less significant features to shrink to zero, effectively performing feature selection.
- iii. Despite this, the performance was very close to Linear Regression, indicating minimal overfitting in the baseline model.

c. Ridge Regression:

- i. The MSE was slightly higher than both Linear and Lasso regressions.
- ii. Unlike Lasso, Ridge does not eliminate features but shrinks all coefficients closer to zero, leading to a more balanced model.
- iii. The performance was nearly identical to Lasso and Linear Regression.

3. Impact of Regularization Term (Alpha)

a. Lasso Regression:

- i. The alpha parameter determines the strength of regularization. At $\alpha = 1.0$, Lasso slightly reduced the impact of less significant features without over-penalizing the coefficients.
- ii. A higher alpha would likely result in more aggressive feature elimination, but it could also reduce the model's ability to generalize.

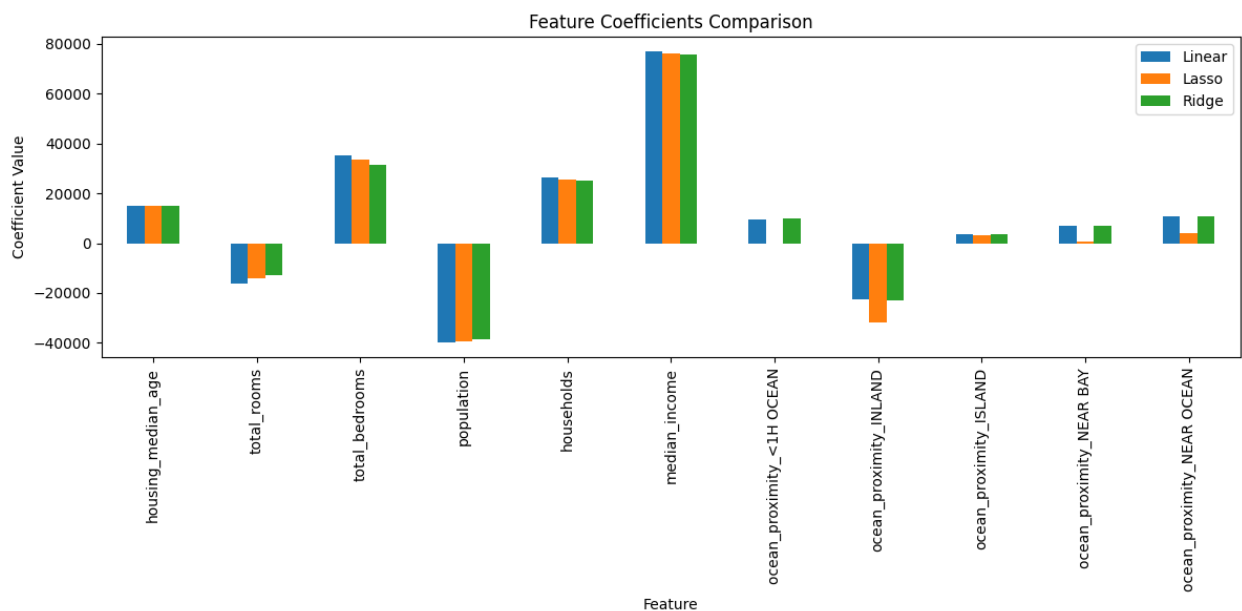
b. Ridge Regression:

- i. The alpha parameter shrinks coefficients uniformly. At $\alpha = 1.0$, the model maintained all features but reduced the magnitude of coefficients slightly, preventing overfitting.
- ii. A higher alpha would increase the shrinkage, potentially leading to underfitting.

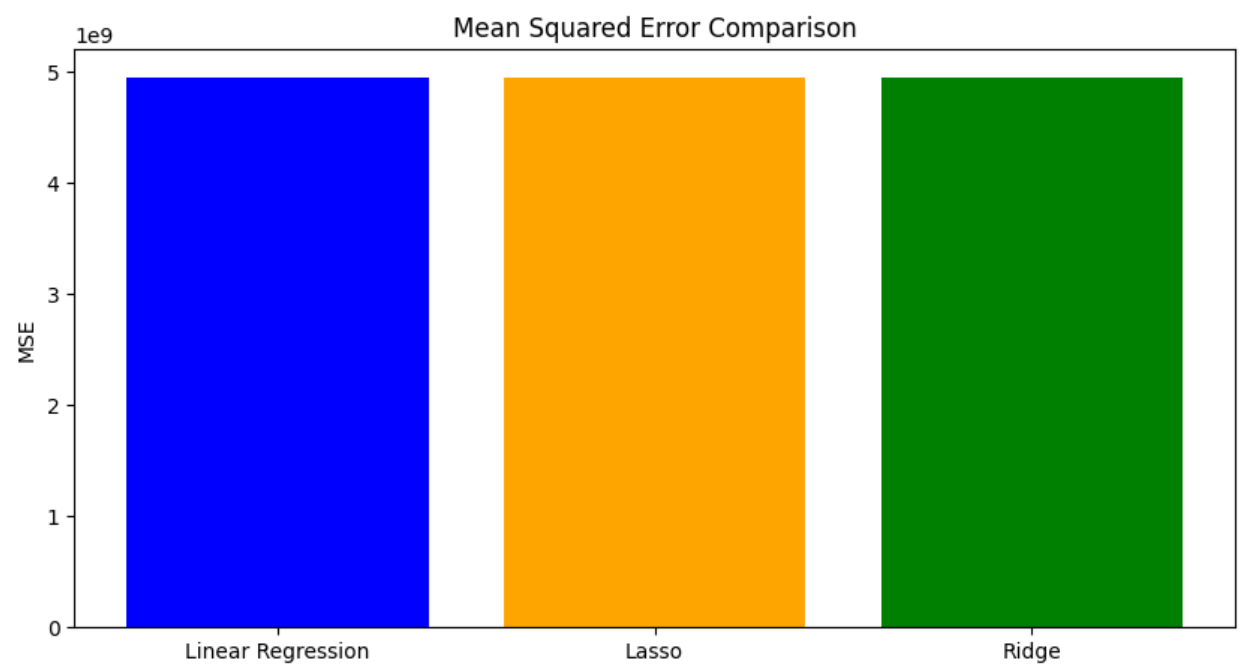
4. Recommendation:

For this use-case, **Linear Regression** is sufficient given its marginally better performance. However, for more complex datasets, **Ridge Regression** is recommended for its robustness.

Feature Coefficients Comparison



Mean Squared Error Comparison:



R-squared Comparison:

