```
In [3]: pip install yellowbrick
```

```
Collecting yellowbrick
  Downloading yellowbrick-1.5-py3-none-any.whl (282 kB)
     ------------------------------------ 282.6/282.6 kB 758.7 kB/s eta 0:00:00
Requirement already satisfied: scikit-learn>=1.0.0 in c:\users\vaibh\anaconda3\lib\site-packages (from yellowbrick) (1.
0.2)
Requirement already satisfied: cycler>=0.10.0 in c:\users\vaibh\anaconda3\lib\site-packages (from yellowbrick) (0.11.0)
Requirement already satisfied: numpy>=1.16.0 in c:\users\vaibh\anaconda3\lib\site-packages (from yellowbrick) (1.21.5)
Requirement already satisfied: matplotlib!=3.0.0,>=2.0.2 in c:\users\vaibh\anaconda3\lib\site-packages (from yellowbric
k) (3.5.2)
Requirement already satisfied: scipy>=1.0.0 in c:\users\vaibh\anaconda3\lib\site-packages (from yellowbrick) (1.9.1)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\vaibh\anaconda3\lib\site-packages (from matplotlib!=3.0.0,
>=2.0.2->yellowbrick) (4.25.0)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\vaibh\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>
=2.0.2->yellowbrick) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\vaibh\anaconda3\lib\site-packages (from matplotlib!=3.
0.0,>=2.0.2->yellowbrick) (2.8.2)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\vaibh\anaconda3\lib\site-packages (from matplotlib!=3.0.0,
>=2.0.2->yellowbrick) (1.4.2)
Requirement already satisfied: packaging>=20.0 in c:\users\vaibh\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=
2.0.2->yellowbrick) (21.3)
Requirement already satisfied: pillow>=6.2.0 in c:\users\vaibh\anaconda3\lib\site-packages (from matplotlib!=3.0.0,>=2.
0.2->yellowbrick) (9.2.0)
Requirement already satisfied: joblib>=0.11 in c:\users\vaibh\anaconda3\lib\site-packages (from scikit-learn>=1.0.0->ye
llowbrick) (1.1.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\vaibh\anaconda3\lib\site-packages (from scikit-learn>=
1.0.0->yellowbrick) (2.2.0)
Requirement already satisfied: six>=1.5 in c:\users\vaibh\anaconda3\lib\site-packages (from python-dateutil>=2.7->matpl
otlib!=3.0.0,>=2.0.2->yellowbrick) (1.16.0)
Installing collected packages: yellowbrick
Successfully installed yellowbrick-1.5
Note: you may need to restart the kernel to use updated packages.
```

```
In [6]: import warnings
        with warnings.catch_warnings():
            warnings.simplefilter(action='ignore', category=FutureWarning)
```

```
In [4]: import pandas as pd
        import numpy as np
        from sklearn.model_selection import train_test_split
        from sklearn.linear_model import LinearRegression, LassoCV
        from sklearn.metrics import mean_squared_error, r2_score
```

```python
import matplotlib.pyplot as plt
from sklearn.model_selection import ShuffleSplit
from yellowbrick.regressor import PredictionError
from yellowbrick.model_selection import LearningCurve
from sklearn.decomposition import PCA
import seaborn as sns

np.random.seed(1)


class DataScienceProject:
    def __init__(self):
        pass

    def load_data(self, file_path):
        # Load data using pandas
        data = pd.read_csv(file_path)
        return data

    def report_missing_values(self, df ):
        # Calculate the number of missing values per column
        missing_values = df.isnull().sum()
        missing_report = pd.DataFrame(missing_values, columns=['missing_values'])
        missing_report = missing_report[missing_report['missing_values'] > 0]

        # Suggest imputation values
        imputation_values = {}
        for column in missing_report.index:
            if df[column].dtype in ['int64', 'float64']:
                skewness = df[column].skew()
                if abs(skewness) > 0.5:
                    imputation_value = df[column].median()
                    imputation_values[column] = ('median', imputation_value)
                else:
                    imputation_value = df[column].mean()
                    imputation_values[column] = ('mean', imputation_value)
            else:
                imputation_value = df[column].mode()[0]
                imputation_values[column] = ('mode', imputation_value)

        return imputation_values

    def apply_imputations(self, df, imputation_values):
        for column, (strategy, value) in imputation_values.items():
            if strategy in ['mean', 'median', 'mode']:
```

```python
                df[column].fillna(value, inplace=True)
        # print("after imp", df.isna().sum())
        return df


    def visualize_data(self, data, title_suffix=''):
        """
        Visualizes distributions of numerical and categorical features in the dataset.

        Args:
        data (DataFrame): The dataset to visualize.
        title_suffix (str): A suffix for the plot title to distinguish between original and preprocessed data.
        """
        numerical_cols = data.select_dtypes(include=['int64', 'float64']).columns
        categorical_cols = data.select_dtypes(include=['object']).columns

        # Plot for numerical features
        for col in numerical_cols:
            plt.figure(figsize=(8, 4))
            sns.histplot(data[col], kde=True)
            plt.title(f'Distribution of {col} {title_suffix}')
            plt.xlabel(col)
            plt.ylabel('Frequency')
            plt.show()

        # Plot for categorical features
        for col in categorical_cols:
            plt.figure(figsize=(8, 4))
            sns.countplot(x=col, data=data)
            plt.title(f'Distribution of {col} {title_suffix}')
            plt.xlabel(col)
            plt.ylabel('Count')
            plt.show()

    def select_data_within_iqr(self, df, iqr_factor=1.5):
        Q1 = df.quantile(0.2)
        Q3 = df.quantile(0.8)
        IQR = Q3 - Q1

        lower_bound = Q1 - (iqr_factor * IQR)
        upper_bound = Q3 + (iqr_factor * IQR)

        # Select only the rows where each column value is within the IQR bounds
        selected_data = df[~((df < lower_bound) | (df > upper_bound)).any(axis=1)]
```

```python
        # print("after iqr", selected_data.isna().sum())
        return selected_data

    def apply_log_transformation(self, df, target_column,skew_threshold=0.5):
        """
        Applies log transformation to highly skewed columns.

        Args:
        df (DataFrame): The dataframe containing the data.
        skew_threshold (float): The threshold to identify highly skewed columns.

        Returns:
        DataFrame: The dataframe with log-transformed columns.
        """
        for column in df.select_dtypes(include=['float64', 'int64']):
            if df[column].skew() > skew_threshold and column != target_column :
                df['Log_' + column] = np.log1p(df[column])

        # print(df.columns)
        # print("after log", df.isna().sum())
        return df

    def preprocess_data(self, data, target_column):
        # Report and apply imputations, and handle outliers
        imputation_values = self.report_missing_values(data)
        data = self.apply_imputations(data, imputation_values)
        data = self.apply_log_transformation(data, target_column)
        data = self.select_data_within_iqr(data)
        return data

    def train_model(self, training_data, target_column):
        # Modify to return X_train, X_test, y_train, y_test for Lasso analysis
        X = training_data.drop(target_column, axis=1)
        y = training_data[target_column]
        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)

        model = LinearRegression()
        model.fit(X_train, y_train)
        return model, X_train, X_test, y_train, y_test

    def plot_actual_vs_predicted(self, y_test, y_predict, model_type):
        """
        Plots the actual vs predicted values.

        Args:
```

```python
            y_test (array-like): The true values of the target variable.
            y_predict (array-like): The predicted values by the model.
            model_type (str): Type of the model ('Baseline' or 'Main').
        """
        plt.figure(figsize=(16, 6))
        plt.title(f"{model_type} Model: Actual vs Predicted Values")
        x_points = list(range(len(y_test)))
        plt.plot(x_points, y_test, label='Actual Values', marker='o')
        plt.plot(x_points, y_predict, label='Predicted Values', marker='x')
        plt.xlabel('Data Points')
        plt.ylabel('Target Variable')
        plt.legend()
        plt.show()


    def plot_learning_curve_and_prediction_error(self, model, X_train, X_test, y_train, y_test, model_type):
        """
        Plots the learning curve and prediction error using Yellowbrick.

        Args:
        model: The trained model.
        X_train, X_test, y_train, y_test: Training and testing data.
        model_type (str): Type of the model ('Baseline' or 'Main').
        """
        # Learning Curve
        plt.figure(figsize=(10, 6))
        lc_viz = LearningCurve(model, cv=5, scoring='r2', n_jobs=4, train_sizes=np.linspace(0.1, 1.0, 10))
        lc_viz.fit(X_train, y_train)
        lc_viz.set_title(f"{model_type} Model: Learning Curve")
        lc_viz.show()

        # Prediction Error
        plt.figure(figsize=(10, 6))
        pe_viz = PredictionError(model)
        pe_viz.fit(X_train, y_train)
        pe_viz.score(X_test, y_test)
        pe_viz.set_title(f"{model_type} Model: Prediction Error")
        pe_viz.show()


    def feature_importance_lasso(self, X_train, y_train, X_test, y_test):
        # Create and fit the LassoCV model
        lasso = LassoCV(cv=5, random_state=0)
        lasso.fit(X_train, y_train)
```

```python
        best_alpha = lasso.alpha_
        lasso_coef = lasso.coef_
        y_pred = lasso.predict(X_test)
        mse = mean_squared_error(y_test, y_pred)
        r2 = r2_score(y_test, y_pred)

        # Plotting feature importances
        plt.figure(figsize=(10, 6))
        feature_importance = pd.Series(lasso_coef, index=X_train.columns).sort_values()
        feature_importance.plot(kind='barh')
        plt.title('Feature Importances from Lasso Model')
        plt.xlabel('Importance')
        plt.ylabel('Features')
        plt.show()

        print(f"Best alpha: {best_alpha}")
        print(f"MSE: {mse}")
        print(f"R-squared: {r2:.2f}")

    def make_prediction(self, model, new_data):
        # Make predictions using the trained model
        prediction = model.predict(new_data)
        return prediction

    def train_and_evaluate(self, X_train, X_test, y_train, y_test):
        """
        Train a linear regression model and evaluate it.
        """
        model = LinearRegression()
        model.fit(X_train, y_train)
        y_pred = model.predict(X_test)

        mse = mean_squared_error(y_test, y_pred)
        r2 = r2_score(y_test, y_pred)
        return model, mse, r2

    def main_pipeline(self, file_path, target_column):
        print("===== Data Science Project Pipeline =====")
        print("[Step 1] Loading and Preprocessing Data")
        data = self.load_data(file_path)

        # Visualize original data
        self.visualize_data(data, title_suffix='(Original)')

        preprocessed_data = self.preprocess_data(data, target_column)
```

```python
        # Visualize preprocessed data
        self.visualize_data(preprocessed_data, title_suffix='(Preprocessed)')

        # Baseline Model
        print("\n[Step 2] Training and Evaluating Baseline Model")
        mse_baseline, r2_baseline, baseline_model, X_train_baseline, X_test_baseline, y_train_baseline, y_test_baseline
        y_pred_baseline = baseline_model.predict(X_test_baseline)
        self.plot_actual_vs_predicted(y_test_baseline, y_pred_baseline, "Baseline")
        self.plot_learning_curve_and_prediction_error(baseline_model, X_train_baseline, X_test_baseline, y_train_baseli

        # Main Model
        print("\n[Step 3] Training and Evaluating Main Model")
        model, X_train, X_test, y_train, y_test = self.train_model(preprocessed_data, target_column)
        _, mse_main, r2_main = self.train_and_evaluate(X_train, X_test, y_train, y_test)
        y_predict_main = model.predict(X_test)
        self.plot_actual_vs_predicted(y_test, y_predict_main, "Main")
        self.plot_learning_curve_and_prediction_error(model, X_train, X_test, y_train, y_test, "Main")

        print("\n[Step 4] Performing Lasso Feature Importance Analysis")
        self.feature_importance_lasso(X_train, y_train, X_test, y_test)

        print("\n===== Model Comparison Results =====")
        print("Baseline Model:")
        print("  - MSE: {:.3f}".format(mse_baseline))
        print("  - R-squared: {:.3f}".format(r2_baseline))
        print("Main Model:")
        print("  - MSE: {:.3f}".format(mse_main))
        print("  - R-squared: {:.3f}".format(r2_main))
        print("===================================")

        return preprocessed_data

    def compare_with_baseline(self, data, target_column):
        # Updated to return the model and train/test splits
        baseline_data = data.dropna()
        X_baseline = baseline_data.drop(target_column, axis=1)
        y_baseline = baseline_data[target_column]
        X_train_baseline, X_test_baseline, y_train_baseline, y_test_baseline = train_test_split(X_baseline, y_baseline,
        baseline_model, mse_baseline, r2_baseline = self.train_and_evaluate(X_train_baseline, X_test_baseline, y_train_
        return mse_baseline, r2_baseline, baseline_model, X_train_baseline, X_test_baseline, y_train_baseline, y_test_ba
```

In [7]:
```python
# Create an instance of the DataScienceProject class
dsp = DataScienceProject()
```

```
# Run the main pipeline
prediction = dsp.main_pipeline("D:\Temp\HousingData.csv", 'MEDV')
```

===== Data Science Project Pipeline =====
[Step 1] Loading and Preprocessing Data


Distribution of CRIM (Original)

Distribution of ZN (Original)

Distribution of INDUS (Original)

Distribution of CHAS (Original)

Distribution of NOX (Original)

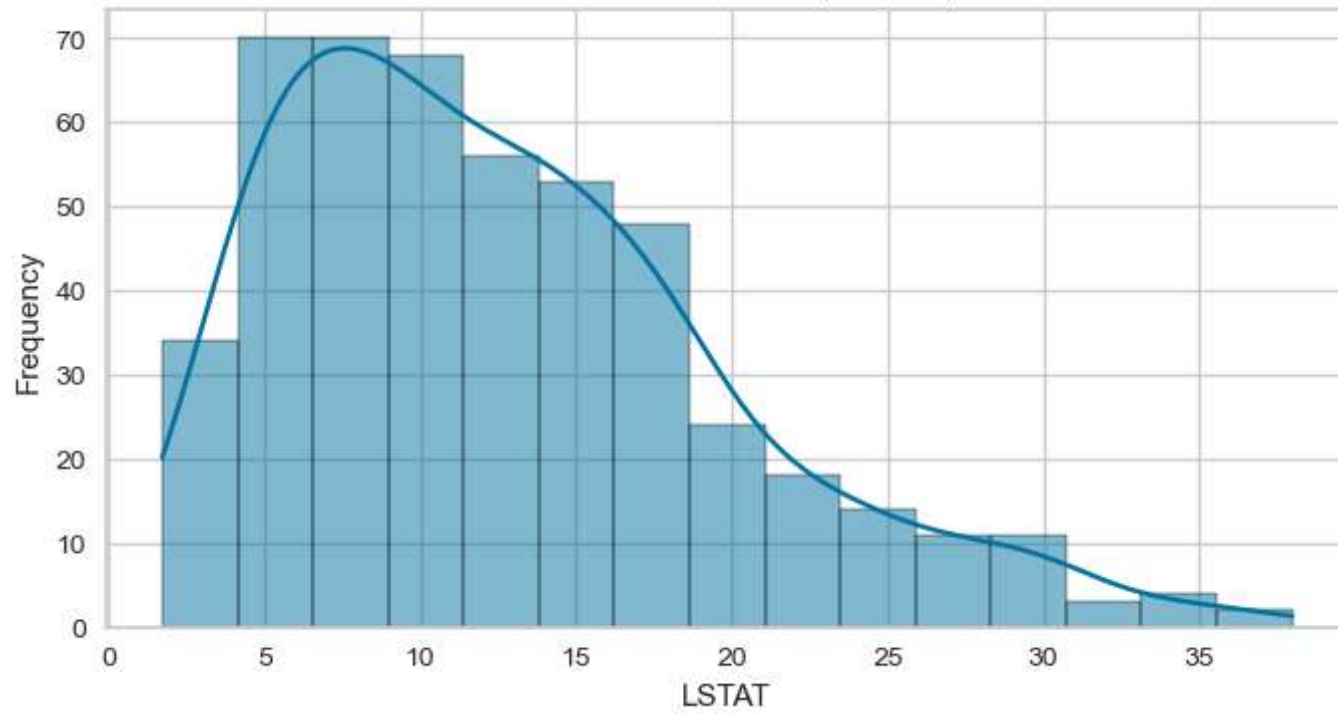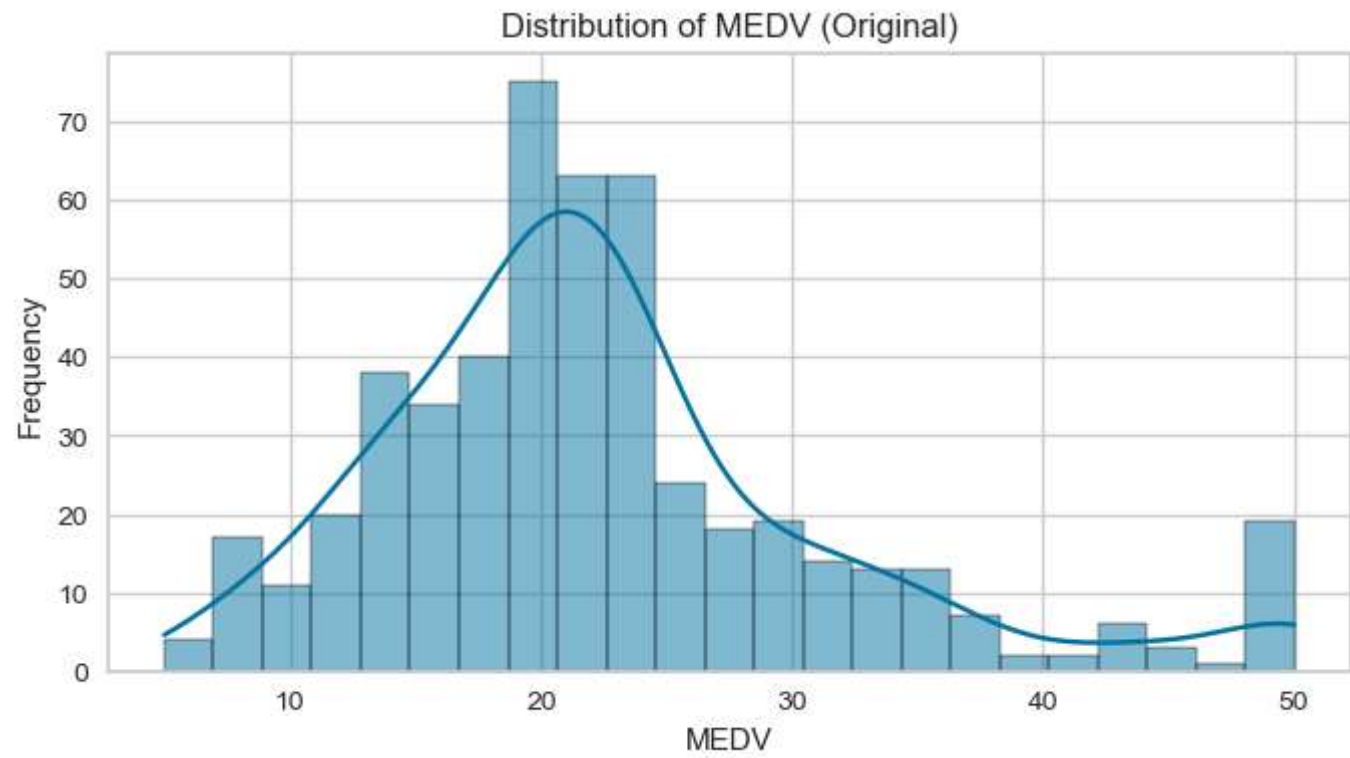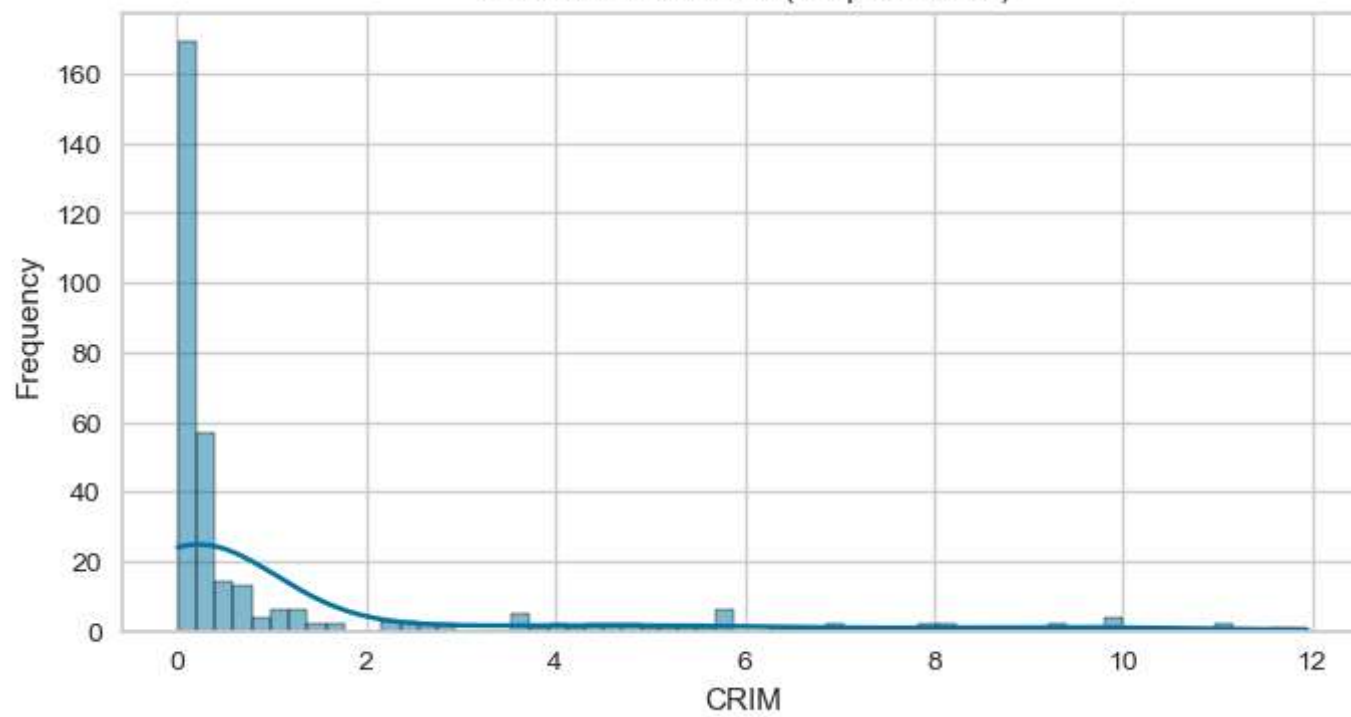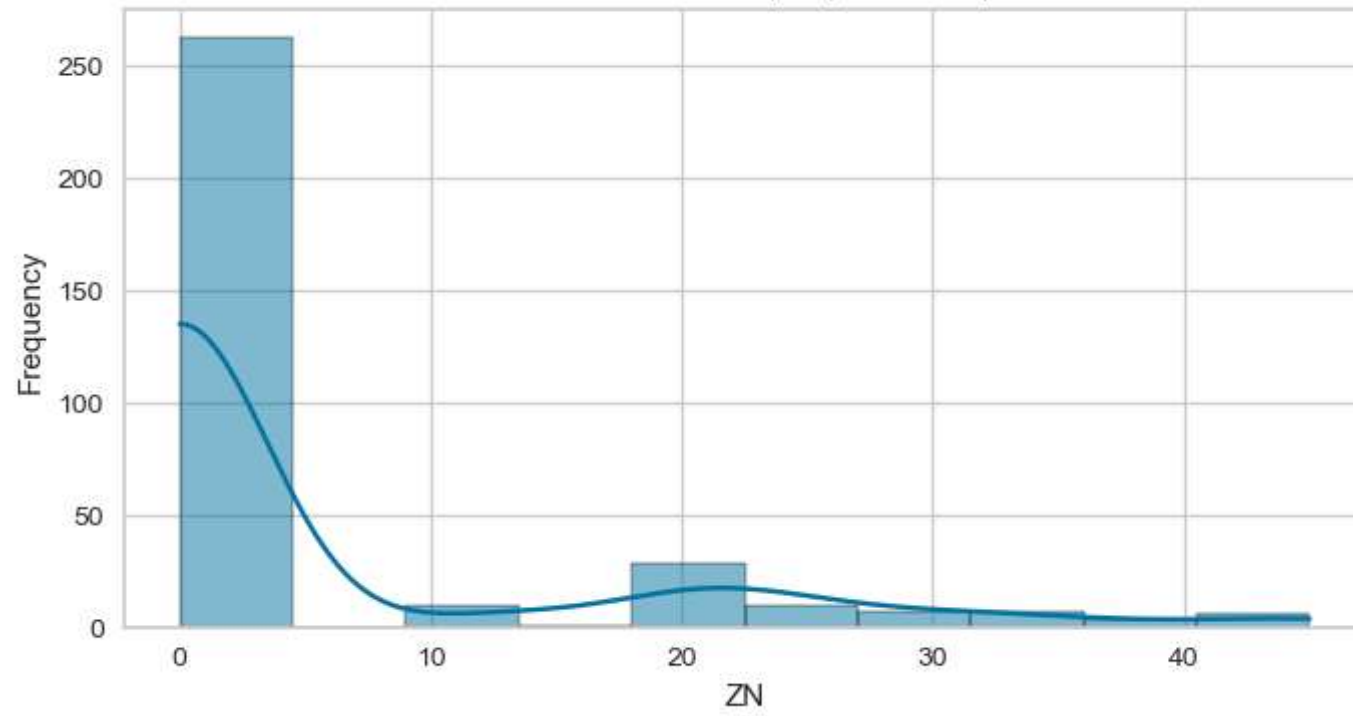Distribution of RM (Original)

Distribution of AGE (Original)

Distribution of DIS (Original)

Distribution of RAD (Original)

Distribution of TAX (Original)

Distribution of PTRATIO (Original)
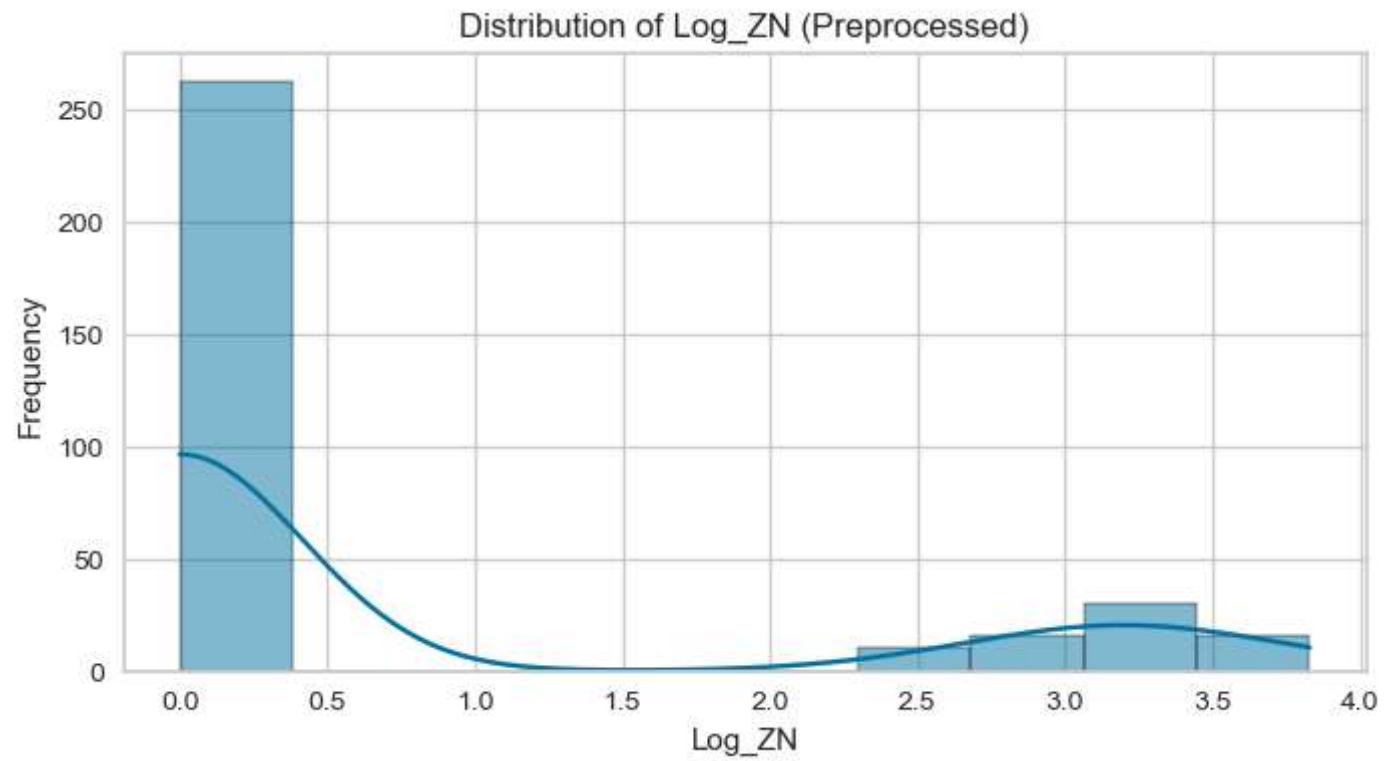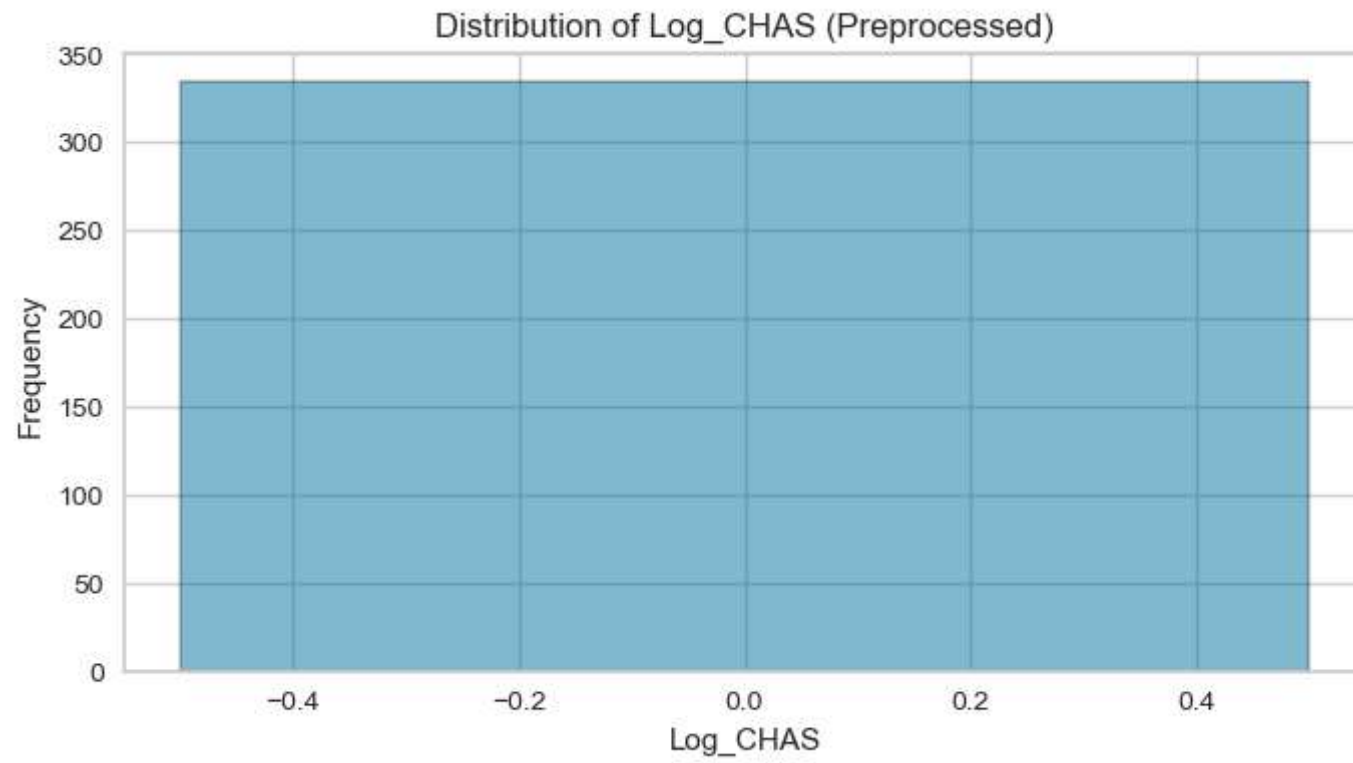
Distribution of B (Original)

Distribution of LSTAT (Original)

Distribution of MEDV (Original)
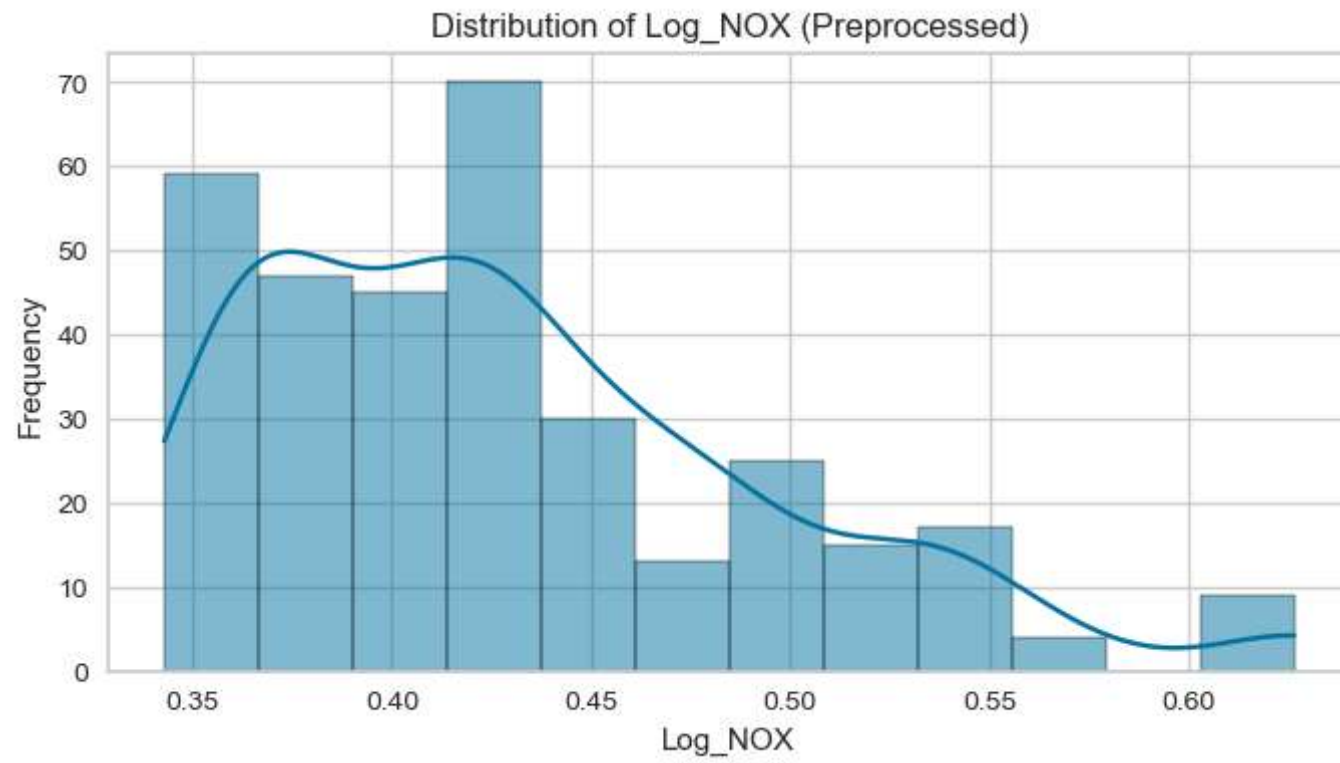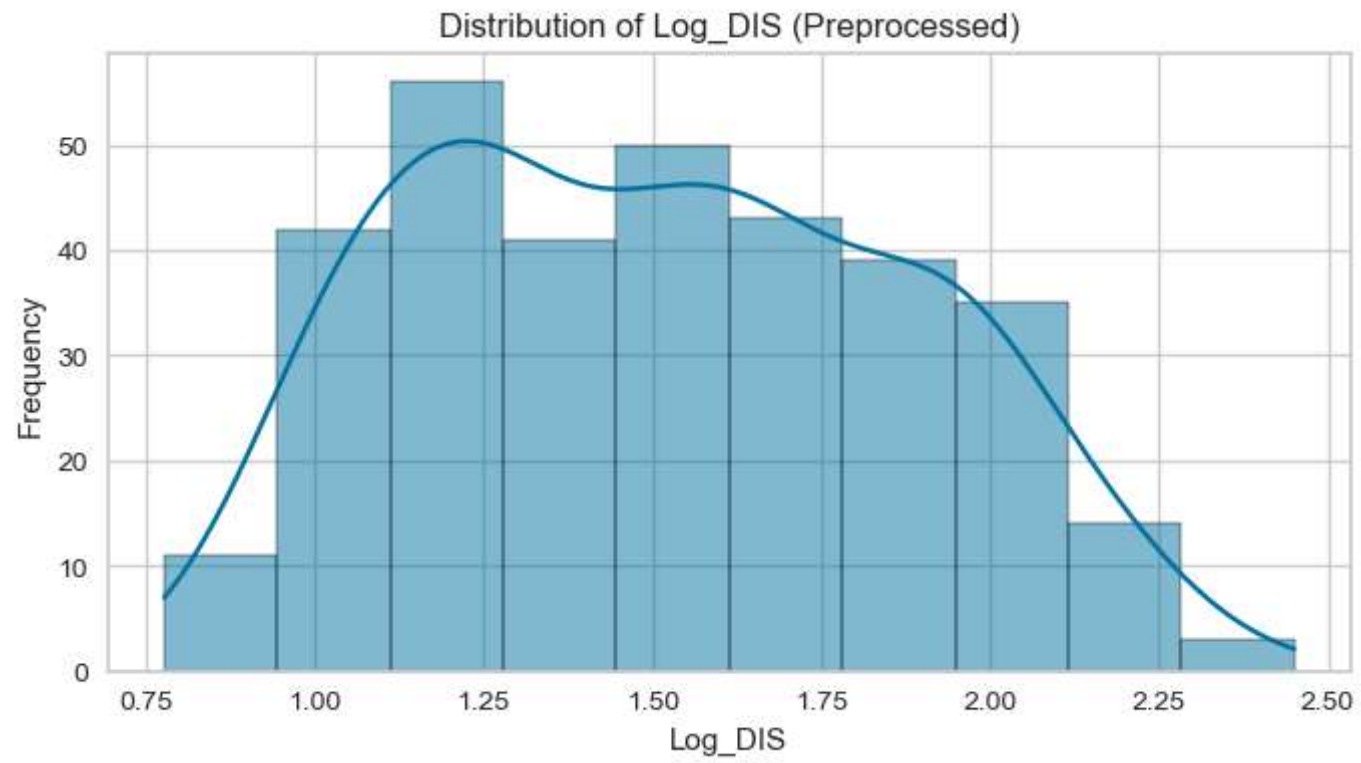
Distribution of CRIM (Preprocessed)

Distribution of ZN (Preprocessed)

Distribution of INDUS (Preprocessed)

# Distribution of CHAS (Preprocessed)

Distribution of NOX (Preprocessed)

Distribution of RM (Preprocessed)

Distribution of AGE (Preprocessed)

Distribution of DIS (Preprocessed)

Distribution of RAD (Preprocessed)

# Distribution of TAX (Preprocessed)

Distribution of PTRATIO (Preprocessed)

Distribution of B (Preprocessed)

Distribution of LSTAT (Preprocessed)

Distribution of MEDV (Preprocessed)

Distribution of Log_CRIM (Preprocessed)

Distribution of Log_ZN (Preprocessed)

Distribution of Log_CHAS (Preprocessed)

Distribution of Log_NOX (Preprocessed)

Distribution of Log_DIS (Preprocessed)

Distribution of Log_RAD (Preprocessed)

Distribution of Log_TAX (Preprocessed)
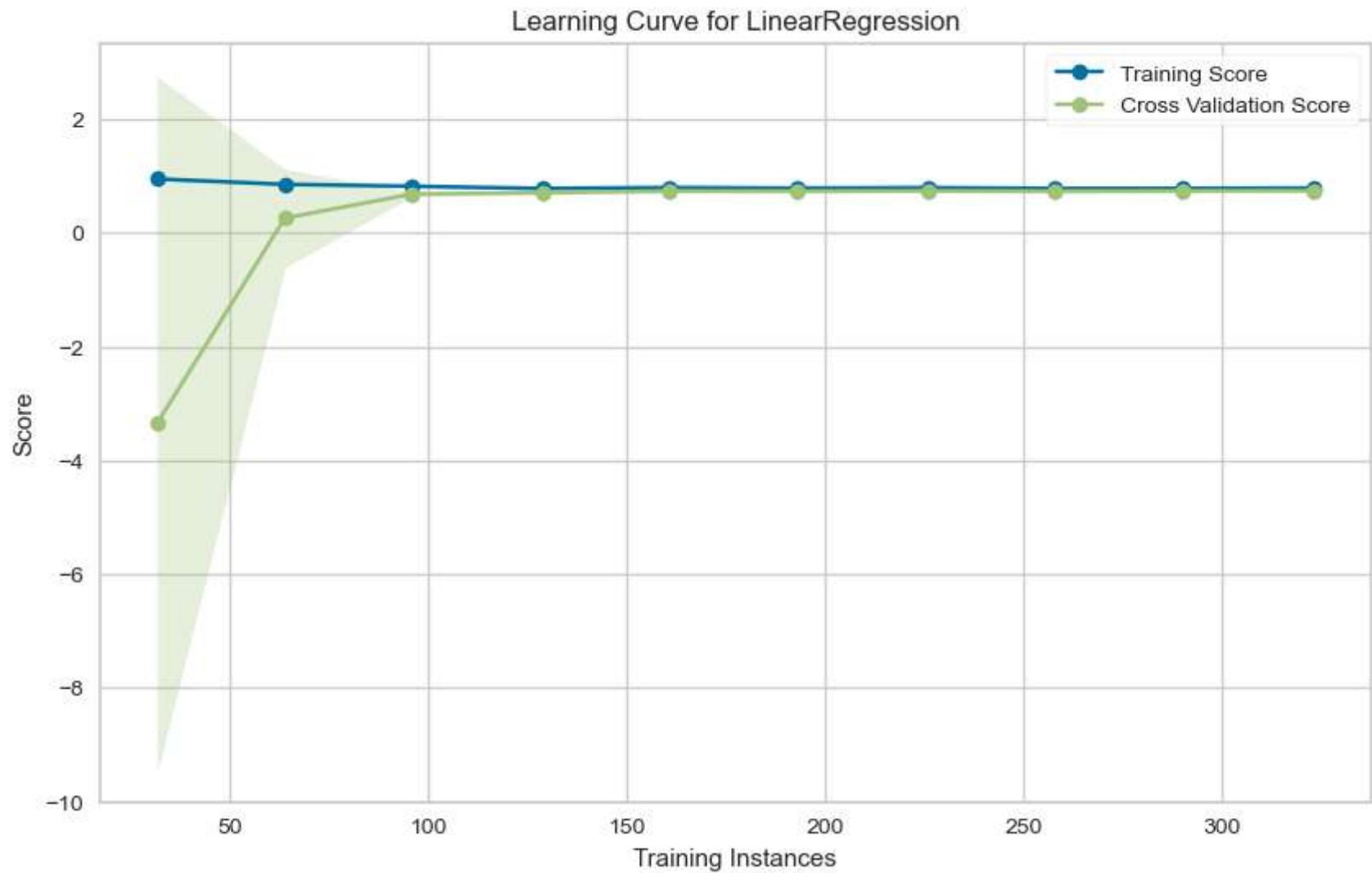
Distribution of Log_LSTAT (Preprocessed)

[Step 2] Training and Evaluating Baseline Model
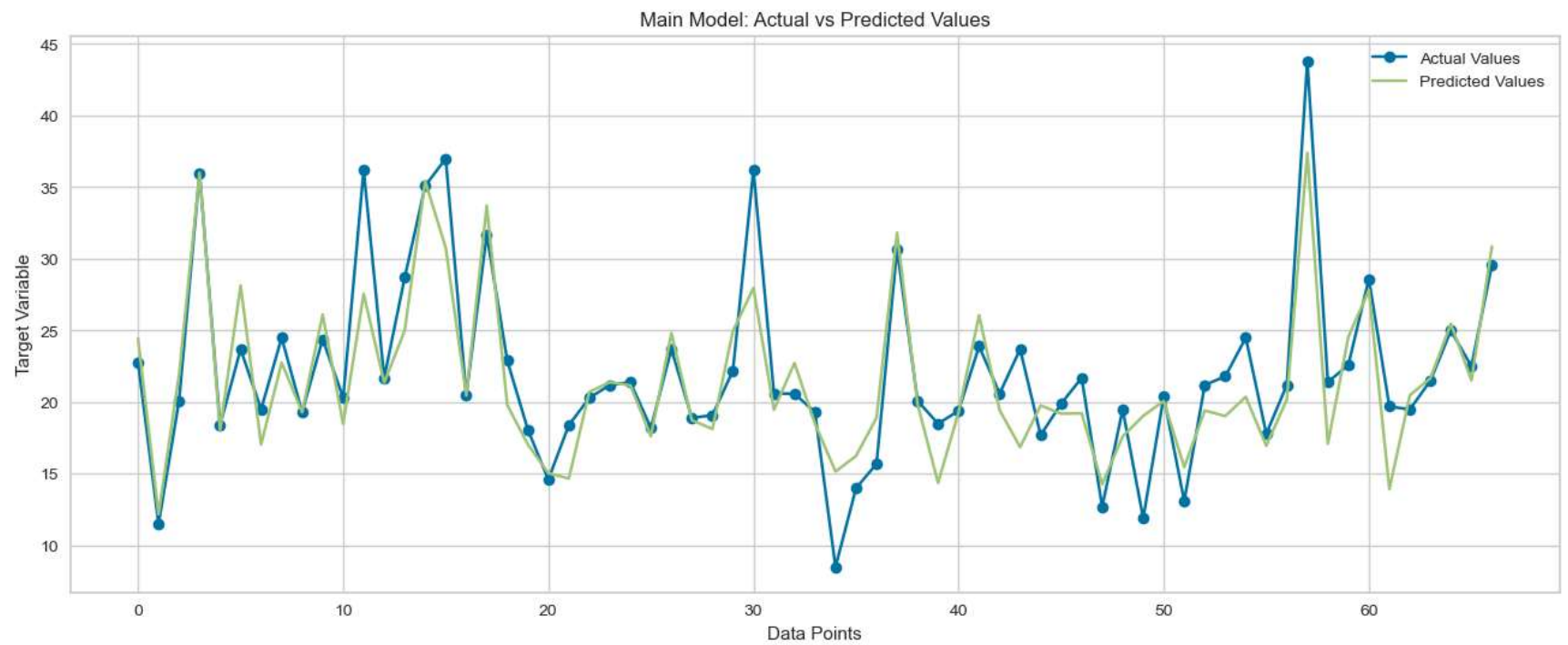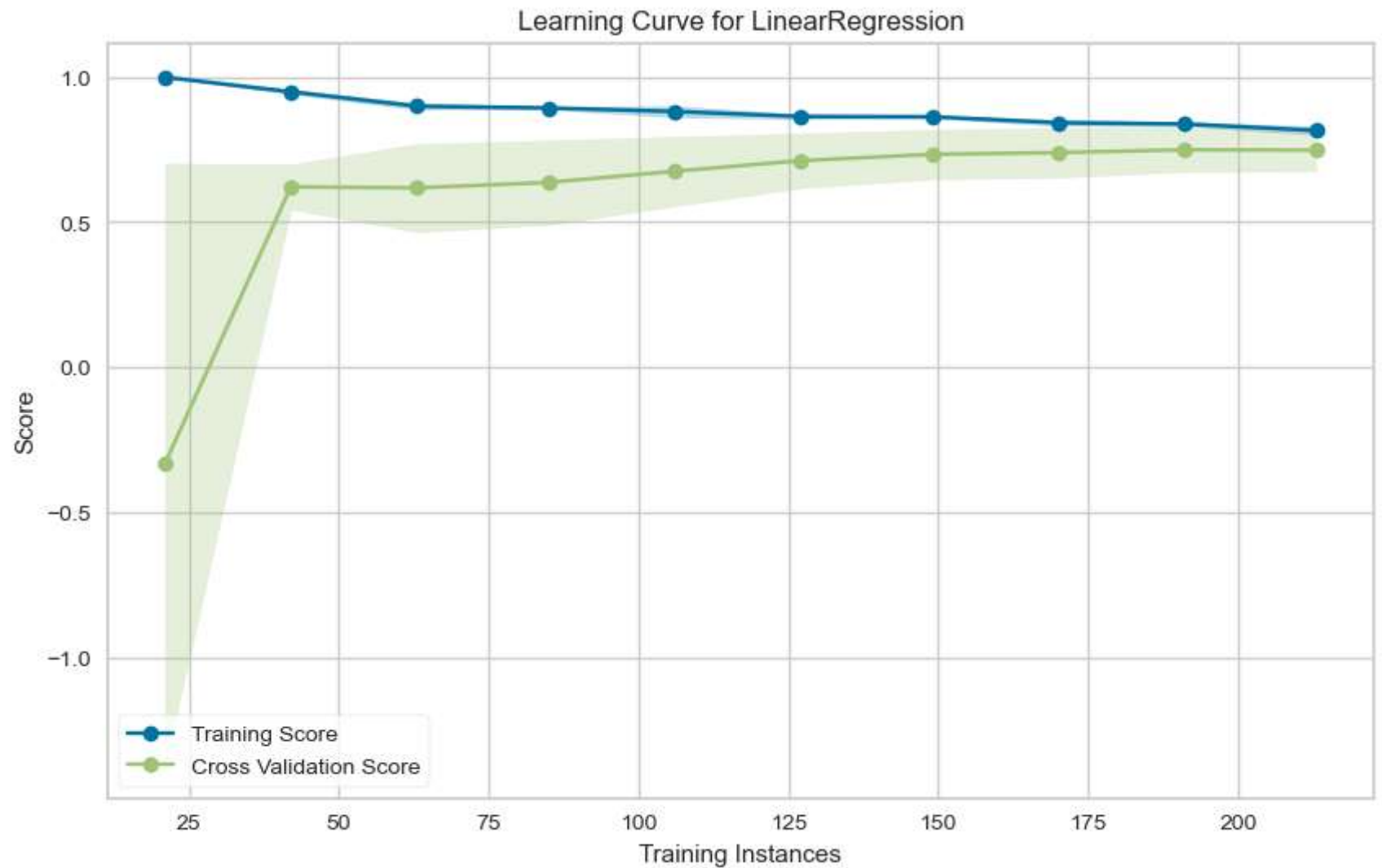
Baseline Model: Actual vs Predicted Values

Learning Curve for LinearRegression

```
C:\Users\vaibh\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but L
inearRegression was fitted with feature names
  warnings.warn(
```

Prediction Error for LinearRegression

[Step 3] Training and Evaluating Main Model

Main Model: Actual vs Predicted Values

Learning Curve for LinearRegression

C:\Users\vaibh\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names
  warnings.warn(

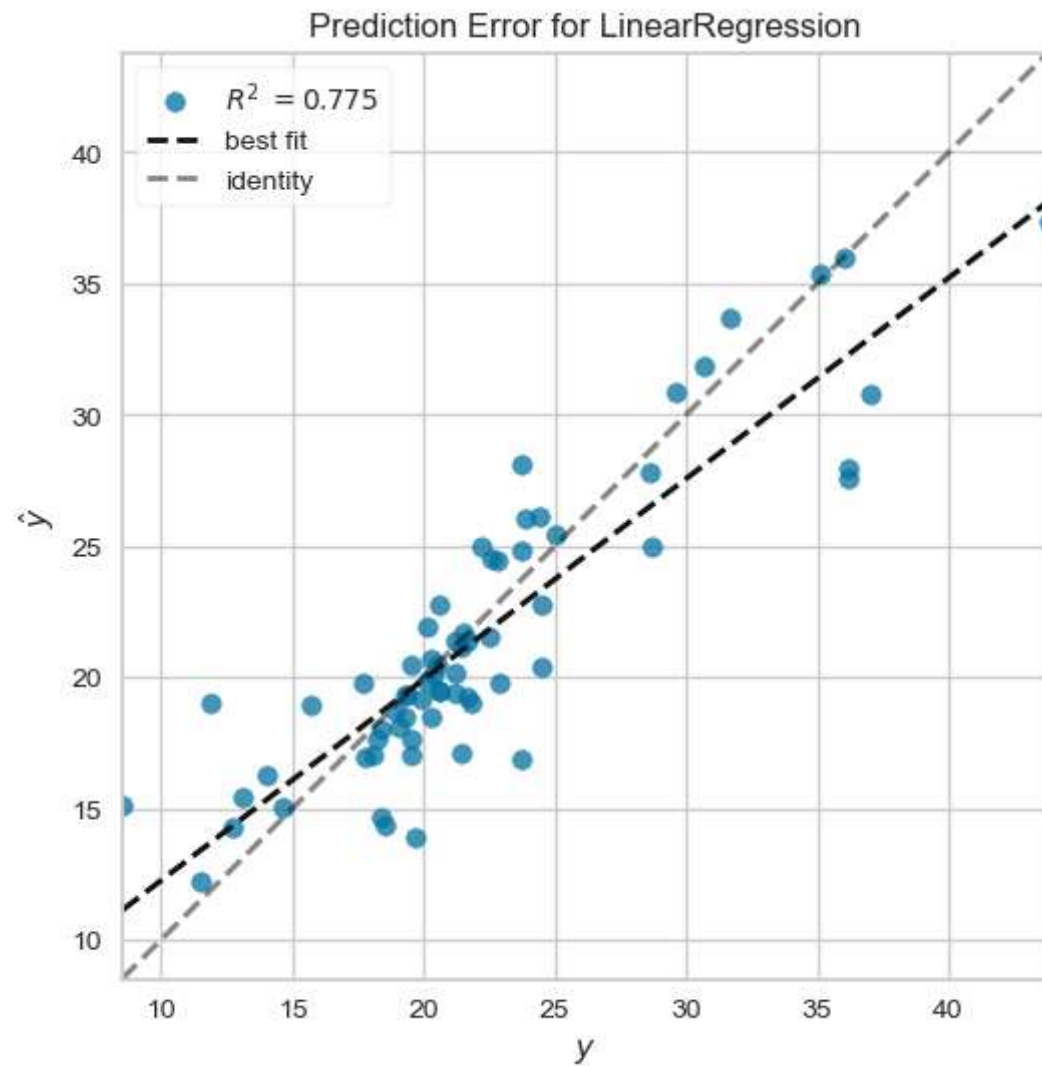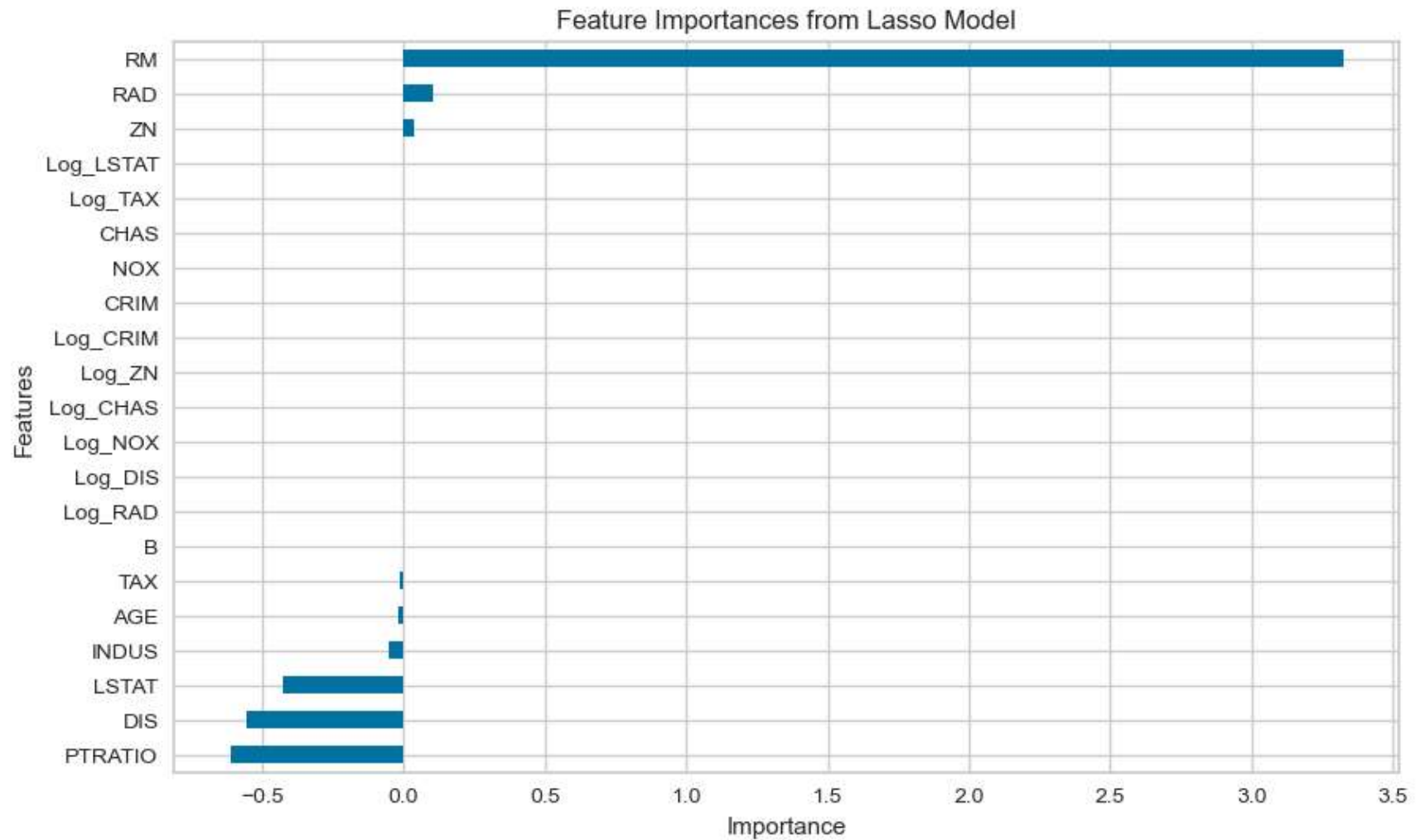Prediction Error for LinearRegression

[Step 4] Performing Lasso Feature Importance Analysis

Feature Importances from Lasso Model

Best alpha: 0.47841076323135406
MSE: 13.048323531370661
R-squared: 0.68

===== Model Comparison Results =====
Baseline Model:
  - MSE: 18.325
  - R-squared: 0.815
Main Model:
  - MSE: 9.213
  - R-squared: 0.775
====================================

In [ ]: