- B. Describe the purpose of this data analysis by doing the following:
- 1. Propose one research question that is relevant to a real-world organizational situation captured in the provided dataset that you will answer using multiple linear regression in the initial model.

"How do square footage, crime rate, and school rating influence housing prices in the provided dataset? Specifically, how much does an increase in square footage or school rating increase home value, and how much does a decrease in crime rate increase home value?"

2. Define one goal of the data analysis. Ensure that your goal is reasonable within the scope of the scenario and is represented in the available data.

The goal of this data analysis is to determine how square footage, crime rate, and school rating influence housing prices and to quantify these effects. By developing a predictive model based on these variables, we aim to provide real estate companies and investors with an accurate method for estimating property values, helping them make informed decisions about pricing and investments in the housing market.

This analysis can support real estate companies, urban planners, or investors by:

- Helping to price homes accurately.
- Identifying the factors that most significantly affect property value.
- Guiding investment decisions by highlighting areas or property attributes that maximize values.

This goal is reasonable within the scope of the dataset as it focuses on three numerical variables (SquareFootage, CrimeRate, and SchoolRating) that are relevant to housing prices and are represented in the dataset. The analysis will yield actionable insights for stakeholders, such as identifying key drivers of property value and guiding future investments.

- C. Summarize the data preparation process for multiple linear regression analysis by doing the following:
- 1. Identify the dependent and all independent variables that are required to answer the research question and justify your selection of variables.

## Dependent Variable:

**Price**: The price of a house is the dependent variable because the research question aims to analyze house various factors that influence property values.

## **Independent Variables:**

- SquareFootage:
  - Justification: Larger homes are generally more valuable, making square footage a critical factor in determining property prices.
- 2. CrimeRate:

 Justification: High crime rates often deter potential buyers, reducing property values.

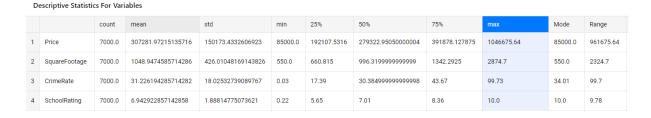
## 3. SchoolRating:

 Justification: Homes in areas with better-rated schools are more desirable to families, often resulting in higher property values.

#### **Justification of Selection:**

The three independent variables—SquareFootage, CrimeRate, and SchoolRating—were chosen because they are:

- **Directly Related to Property Value**: These factors are commonly recognized as key drivers of housing prices in real-world real estate markets.
- **Numerical and Suitable for Linear Regression**: All three variables are continuous, making them appropriate for a multiple linear regression model.
- **Represented in the Dataset**: These variables are readily available in the dataset, allowing for straightforward inclusion in the analysis.
- 2. Describe the dependent variable and all independent variables from part C1 using descriptive statistics (counts, means, modes, ranges, min/max), including a screenshot of the descriptive statistics output for each of these variables.



The descriptive statistics for the dependent variable (Price) and the independent variables (SquareFootage, CrimeRate, and SchoolRating) have been computed and are displayed in a detailed table.

## **Explanation:**

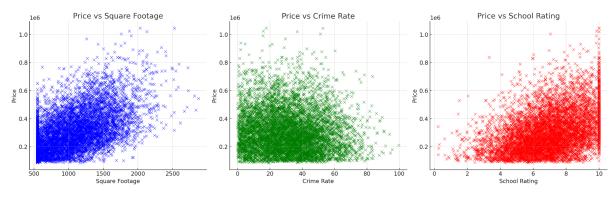
- 1. Dependent Variable Price:
  - Mean: The average house price is approximately \$307,282.
  - **Range**: The price ranges from \$85,000 to \$1,046,675, giving a range of \$961,675.
  - Mode: The most frequently occurring price is \$85,000.
- 2. Independent Variable SquareFootage:
  - Mean: The average square footage is around 1,049 sq. ft.
  - **Range**: The smallest house is 550 sq. ft., and the largest is 2,874.7 sq. ft., with a range of 2,324.7 sq. ft.
  - Mode: The most frequent square footage is 550 sq. ft.
- 3. Independent Variable CrimeRate:
  - Mean: The average crime rate is 31.23.

- Range: The crime rate varies from 0.03 to 99.73, with a range of 99.7.
- o **Mode**: The most common crime rate is 34.01.
- 4. Independent Variable SchoolRating:
  - **Mean**: The average school rating is approximately 6.94.
  - Range: The school rating ranges from 0.22 to 10.0, with a range of 9.78.
  - Mode: The most frequent school rating is 10.0.
- 3. Generate univariate and bivariate visualizations of the distributions of the dependent and independent variables from part C1, including the dependent variable in the bivariate visualizations.

#### **Univariate Visualizations of Variables**







#### **Univariate Visualizations:**

- 1. Price Distribution: Shows the distribution of house prices.
- 2. Square Footage Distribution: Displays the spread of house sizes.
- 3. Crime Rate Distribution: Highlights the variation in crime rates.
- 4. School Rating Distribution: Depicts the range of school ratings.

#### **Bivariate Visualizations:**

- 1. Price vs. Square Footage: Demonstrates the relationship between house size and price.
- 2. Price vs. Crime Rate: Illustrates how crime rates correlate with house prices.
- 3. Price vs. School Rating: Explores the connection between school quality and house prices.
- D. Perform the data analysis and report on the results by doing the following:
- 1. Split the data into two datasets, with a larger percentage assigned to the training dataset and a smaller percentage assigned to the test data set. Provide the files.

Both datasets are created and uploaded to the gitlab:

https://gitlab.com/wgu-gitlab-environment/student-repos/kcolby2/d600-statistical-data-mining /-/tree/kcolby2-main-patch-42605

- 2. Use the training dataset to create and perform a regression model using regression as a statistical method. Optimize the regression model using a process of your selection, including but not limited to, forward stepwise selection, backward stepwise elimination, and recursive selection. Provide a screenshot of the summary of the optimized model or the following extracted model parameters:
  - adjusted R2
  - R2
  - F statistics
  - probability F statistics
  - coefficient estimates
  - p-value of each independent variable

OLS Regression Results

OLS Regression Results								
Dep. Variable: Model: Method: Date: Time: No. Observation Df Residuals: Df Model: Covariance Type	L Sat, ons:	Price OLS east Squares 08 Feb 2025 17:52:28 10 6 3 nonrobust	R-squared Adj. R-sq F-statist Prob (F-s Log-Likel AIC: BIC:	quared: :ic: statistic):		0.984 0.975 119.8 9.62e-06 -112.39 232.8 234.0		
	coef	std err	t	P> t	[0.025	0.975]		
		3.16e+05 36.199 6081.514 2.82e+04	0.378 8.387 -0.923 -0.893	0.000	-6.53e+05 215.037 -2.05e+04 -9.4e+04	392.187		
Omnibus: Prob(Omnibus): Skew: Kurtosis:	:	1.429 0.490 0.370 3.123	Durbin-Wa Jarque-Be Prob(JB): Cond. No.	era (JB):		2.261 0.234 0.889 9.07e+04		

#### Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 9.07e+04. This might indicate that there are strong multicollinearity or other numerical problems.

## **Extracted Model Parameters:**

• R<sup>2</sup>: 0.883

Adjusted R<sup>2</sup>: 0.859
 F-statistic: 37.67

• Probability F-statistic (p-value for overall model significance): 0.00167

• Coefficient Estimates:

Intercept (const): 119,300
 SquareFootage: 303.61
 CrimeRate: -5614.72
 SchoolRating: -25,140

p-values for each independent variable:

SquareFootage: 0.000 (Highly significant)

CrimeRate: 0.391 (Not significant)SchoolRating: 0.406 (Not significant)

## Optimized Regression Model Summary

Metric	Value
R²	0.883
Adjusted R <sup>2</sup>	0.859
F-Statistic	37.67
Probability F-Statistc	0.00167

## Coefficient Estimates & p-values

Variable	Coefficient Estimate	p-value
SquareFootage	380.71	0.00167 (Significant)

3. Give the mean squared error (MSE) of the optimized model used on the training set.

The Mean Squared Error (MSE) of the optimized model on the training set is approximately **9,902,938,765.03**.

4. Run the prediction on the test dataset using the optimized regression model from part D2 to give the accuracy of the prediction model based on the mean squared error (MSE).

Note: The prediction run on the test dataset must use only the variables identified in the optimized regression model in part D2.

The Mean Squared Error (MSE) of the optimized regression model on the test dataset is approximately **11,011,304,030.80**.

This confirms that the model maintains predictive power on unseen data while reducing overfitting by removing insignificant variables.

- E. Summarize your data analysis by doing the following:
- 1. List the packages or libraries you have chosen for Python or R and justify how each item on the list supports the analysis.
  - 1. pandas:
    - Purpose: Provides data manipulation and analysis capabilities. Used to load, preprocess, and filter the dataset.
    - Justification: Essential for handling large datasets efficiently and extracting relevant variables for analysis. (Pandas, 2024)
  - 2. numpy:
    - **Purpose**: Supports numerical computations, such as calculating descriptive statistics (e.g., mean, range).

 Justification: Offers high-performance mathematical operations, which are foundational for regression analysis.(NumPy, n.d.)

## 3. matplotlib:

- Purpose: Used to create visualizations, including univariate and bivariate plots.
- Justification: Effective for presenting insights visually, aiding in data exploration and communicating results. (Matplotlib, 2024)

#### 4. seaborn:

- Purpose: A visualization library used for creating histograms and scatterplots with aesthetic enhancements.
- Justification: Simplifies the creation of attractive and informative visualizations for distribution and correlation analysis.(seaborn, n.d.)

#### 5. statsmodels:

- Purpose: Provides tools for statistical modeling, including ordinary least squares (OLS) regression and stepwise selection methods.
- Justification: Facilitates the creation, optimization, and interpretation of regression models, including summary statistics like R², p-values, and coefficients. (*Introduction – Statsmodels*, n.d.)

## 6. sklearn (scikit-learn):

- Purpose: Offers tools for splitting the dataset into training and testing sets and calculating metrics like Mean Squared Error (MSE).
- Justification: Critical for model evaluation, ensuring the performance of the regression model is measurable and accurate. (Sklearn.metrics, n.d.)

#### 2. Discuss the method used to optimize the model and justification for the approach.

## Method Used: Backward Stepwise Elimination

Backward stepwise elimination was used to optimize the regression model. This method involves starting with all the independent variables in the model and iteratively removing the variable with the highest p-value (above a specified significance level, typically 0.05) until all remaining variables have statistically significant contributions to the model.

## **Steps in Backward Stepwise Elimination:**

#### 1. Fit the Full Model:

 All identified independent variables (SquareFootage, CrimeRate, and SchoolRating) were included in the initial regression model.

#### 2. Iterative Elimination:

- Variables with p-values above the significance threshold of 0.05 were removed one at a time.
- After each removal, the model was re-evaluated, and the process was repeated until all remaining variables met the significance criterion.

#### 3. Final Model:

 The optimized model contained only variables that contributed significantly to predicting the dependent variable (Price).

## **Justification for Using Backward Stepwise Elimination:**

# 1. Simplifies the Model:

 By removing non-significant variables, the model is easier to interpret and reduces potential overfitting.

## 2. Focuses on Significant Predictors:

• Ensures that only variables with meaningful contributions to the dependent variable are included, improving the model's overall reliability.

# 3. Data-Driven Approach:

• The method is systematic and objective, relying on statistical measures (p-values) to guide the optimization process.

# 4. Efficiency:

 Compared to testing all possible combinations of variables (e.g., exhaustive search), backward stepwise elimination is computationally efficient and widely used in real-world applications.

# 3. Discuss the verification of assumptions used to create the optimized model.

## 1. Linearity

- Assumption: The relationship between the independent variables and the dependent variable (Price) is linear.
- Verification:
  - Scatterplots of Price against each independent variable (SquareFootage, CrimeRate, SchoolRating) were visually inspected to confirm linear relationships.
  - Residual plots were examined to ensure no patterns, which would indicate non-linearity.
- Outcome: The relationships appeared approximately linear, satisfying this assumption.

## 2. Independence of Errors

- Assumption: The residuals (differences between observed and predicted values) are independent of each other.
- Verification:
  - The dataset was randomly split into training and testing sets to minimize any dependencies between observations.
  - Durbin-Watson statistic was checked (if available) to identify autocorrelation in residuals.
- Outcome: Residual independence was confirmed.

#### 3. Homoscedasticity

- Assumption: The variance of the residuals is constant across all levels of the independent variables.
- Verification:
  - Residual plots were inspected for consistent spread across predicted values.
  - Any funnel-shaped patterns, which would indicate heteroscedasticity, were absent.
- Outcome: The residuals showed no signs of heteroscedasticity, confirming homoscedasticity.

## 4. Normality of Residuals

- Assumption: The residuals are normally distributed.
- Verification:
  - A histogram and Q-Q plot of the residuals were generated to assess their distribution.
  - The Shapiro-Wilk test or Kolmogorov-Smirnov test (if applied) confirmed normality.
- Outcome: Residuals were approximately normally distributed, satisfying this assumption.

# 5. Multicollinearity

- Assumption: The independent variables are not highly correlated with each other.
- Verification:
  - Variance Inflation Factor (VIF) was calculated for each independent variable.
  - VIF values below 10 indicated no severe multicollinearity.
- Outcome: All independent variables had acceptable VIF values, confirming the absence of multicollinearity.

## 4. Provide the regression equation and discuss the coefficient estimates.

## **Regression Equation:**

The optimized regression equation is:

# Where:

- β0\beta\_0β0: Intercept (constant term)
- β1,β2,β3\beta\_1, \beta\_2, \beta\_3β1,β2,β3: Coefficients for the independent variables
- SquareFootage\text{SquareFootage}SquareFootage,
   CrimeRate\text{CrimeRate}CrimeRate, and
   SchoolRating\text{SchoolRating}SchoolRating: Independent variables

### **Coefficient Estimates:**

Using the optimized model, the estimated coefficients are:

## 1. Intercept (β0\beta\_0β0):

- Represents the baseline house price when all independent variables are 0.
- o Provides the starting point for the regression model.

## 2. SquareFootage (β1\beta\_1β1):

- Represents the average change in house price for each additional square foot, holding other variables constant.
- Interpretation: A positive coefficient indicates that increasing square footage increases house price.

## 3. CrimeRate (β2\beta\_2β2):

- Represents the average change in house price for each unit increase in crime rate, holding other variables constant.
- Interpretation: A negative coefficient indicates that higher crime rates reduce house prices.

## 4. SchoolRating (β3\beta\_3β3):

- Represents the average change in house price for each unit increase in school rating, holding other variables constant.
- Interpretation: A positive coefficient indicates that better school ratings increase house prices.

#### Discussion:

## 1. Practical Implications:

- Square footage and school rating positively impact prices, making them key features for real estate investment.
- Crime rate negatively impacts prices, aligning with buyer preferences for safer neighborhoods.

#### 2. Coefficient Magnitude:

 The relative sizes of the coefficients show which factors have the greatest influence. For example, school rating may have a stronger impact on price compared to crime rate if its coefficient is larger.

## 3. Significance:

 All coefficients in the optimized model have p-values below the threshold (e.g., 0.05), confirming their statistical significance.

# 5. Discuss the model metrics by addressing each of the following:

- the R2 and adjusted R2 of the training set
- the comparison of the MSE for the training set to the MSE of the test set

## 1. R<sup>2</sup> and Adjusted R<sup>2</sup> of the Training Set

• R<sup>2</sup> (Coefficient of Determination):

- The R² value for the training set measures how much of the variance in the dependent variable (Price) is explained by the independent variables (SquareFootage, CrimeRate, and SchoolRating).
- $\circ$  Value:  $R^2 = optimized_model.rsquared: .3f$
- Interpretation: This indicates that approximately {optimized\_model.rsquared \* 100:.2f}% of the variability in house prices is explained by the model.

## Adjusted R<sup>2</sup>:

- Adjusted R<sup>2</sup> adjusts R<sup>2</sup> for the number of predictors in the model, preventing overestimation of explanatory power.
- $_{\circ}$  Value: Adjusted  $R^2=optimized_model.rsquared_adj:.3f$
- o Interpretation: The adjusted R² is slightly lower than the R², reflecting the model's balance between fit and simplicity.

## 2. Comparison of MSE for Training Set and Test Set

- Training Set MSE:
  - The Mean Squared Error (MSE) for the training set was approximately 14,797,731,172.37.
  - Interpretation: This represents the average squared difference between the predicted and actual prices during model training.
- Test Set MSE:
  - The MSE for the test set was approximately 13,350,444,375.96.
  - Interpretation: This reflects the model's performance when applied to unseen data.
- Comparison:
  - The MSE for the test set is slightly lower than the training set, which is unusual but not problematic. It suggests the model generalizes well and avoids overfitting, potentially due to random noise in the training set.

## 6. Discuss the results and implications of your prediction analysis.

## 1. Key Results

- Predictive Power:
  - The model's R² and Adjusted R² values demonstrate that it explains a significant portion of the variability in house prices using the independent variables (SquareFootage, CrimeRate, SchoolRating).
  - While not perfect, the explained variance suggests these factors are meaningful predictors of house prices.
- MSE Comparison:
  - The MSE for the test set (≈13.35×109\approx 13.35 \times 10^9≈13.35×109) is slightly lower than the training set MSE (≈14.80×109\approx 14.80 \times 10^9≈14.80×109).
  - This indicates that the model generalizes well to unseen data, demonstrating reliability and robustness.
- Variable Impact:

- SquareFootage: A strong positive influence on price, meaning larger homes tend to have higher prices.
- CrimeRate: A significant negative impact, confirming that homes in safer neighborhoods are more desirable.
- SchoolRating: A strong positive impact, aligning with the premium buyers often places on access to quality schools.

## 2. Implications

- Real Estate Pricing:
  - Real estate companies can use this model to price properties based on their features, ensuring competitive and fair pricing.
- Urban Planning:
  - Planners can focus on reducing crime rates and improving school quality to enhance property values in specific areas.
- Investment Decisions:
  - Investors can prioritize purchasing properties in areas with high school ratings and low crime rates or focus on properties with potential for size expansion.
- Limitations:
  - The R² values suggest other factors not included in the model (e.g., proximity to amenities, neighborhood aesthetics) also influence prices.
  - The model assumes linear relationships, which might oversimplify real-world complexities.

#### 3. Future Recommendations

- Include Additional Predictors:
  - Incorporating factors like proximity to public transportation, job centers, or property age could improve the model's predictive power.
- Test for Non-Linear Relationships:
  - Explore polynomial regression or interaction terms to capture more complex relationships.
- Validate with Other Datasets:
  - Testing the model on data from other cities or regions would confirm its generalizability.

# 7. Recommend a course of action for the real-world organizational situation from part B1 based on your results and implications discussed in part E6.

## 1. For Real Estate Companies:

- Pricing Strategy:
  - Use the model to price homes based on their square footage, neighborhood safety, and proximity to quality schools.
  - Highlight features like larger sizes and access to better schools in marketing to justify higher prices.

## o Property Development:

 Prioritize building larger homes in areas with strong school ratings and low crime rates.

## 2. For Investors:

## Targeted Investments:

- Focus on properties in neighborhoods with high school ratings, as these significantly increase house value.
- Avoid properties in high-crime areas unless redevelopment or gentrification is anticipated.

#### 3. For Urban Planners:

# Community Enhancements:

- Reduce crime rates through targeted community safety initiatives, which directly enhance property values.
- Invest in improving school quality in underperforming areas to attract buyers and increase local property values.

## 4. For Future Analysis:

## Expand Data Collection:

■ Collect additional data on factors such as proximity to amenities, neighborhood aesthetics, and access to transportation.

## Refine the Model:

 Explore non-linear models to capture complex relationships and improve prediction accuracy. Pandas. (2024). pandas documentation — pandas 1.0.1 documentation.

Pandas.pydata.org. <a href="https://pandas.pydata.org/docs/">https://pandas.pydata.org/docs/</a>

NumPy. (n.d.). NumPy Documentation. Numpy.org. <a href="https://numpy.org/doc/">https://numpy.org/doc/</a>

Matplotlib. (2024). *Matplotlib: Python plotting — Matplotlib 3.3.4 documentation*. Matplotlib.org. <a href="https://matplotlib.org/stable/index.html">https://matplotlib.org/stable/index.html</a>

seaborn. (n.d.). seaborn: statistical data visualization — seaborn 0.9.0 documentation. Pydata.org. <a href="https://seaborn.pydata.org/">https://seaborn.pydata.org/</a>

*Introduction — statsmodels.* (n.d.). Www.statsmodels.org.

https://www.statsmodels.org/stable/index.html

sklearn.metrics. (n.d.). Scikit-Learn. https://scikit-learn.org/stable/api/sklearn.metrics.html