

★ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



Harnessing ChatGPT for Data Mining: Analyzing California Housing Prices



Venkatesh

18 min read · Draft



Listen



Share

... More

Introduction

The power of modern AI tools is reshaping the way we approach data analysis. In this blog post, I'll demonstrate how I used OpenAI's ChatGPT to comprehensively explore the California housing dataset, employing the CRISP-DM methodology. CRISP-DM, which stands for Cross Industry Standard Process for Data Mining, provides a robust framework for tackling data science projects.

Using ChatGPT for Data Analysis

ChatGPT, developed by OpenAI, is a state-of-the-art language model designed for human-like interactions. Instead of traditional data analysis, where one would need to manually code and execute each step, I posed questions and gave instructions to ChatGPT. The model then provided insights, generated code snippets, and even created visualizations, streamlining the entire data mining process.

Using ChatGPT's guidance, I tested several regression models, including:

- Linear Regression
- Decision Tree Regression
- Random Forest Regression
- Gradient Boosting Regression

- Support Vector Regression (SVR)

Each model's performance was gauged based on its Root Mean Squared Error (RMSE).

Following is the prompts I gave to chatgpt and its responses. Subsequently, we delved into descriptive statistics, giving insights into house ages, block sizes, and more, all facilitated by interactive prompts to ChatGPT.

you are an industry expert and top notch data scientist. You are also expert in CRISP-DM methodology for data science and written several highly cited research papers of case studies of using crisp-dm data science for various data sets in kaggle. As an industry expert data scientist and expert in CRISP-DM methodology for data science, Given the kaggle data set for California housing prices, can you do step by step EDA analysis playing role of an expert data scientist professor teaching to a masters program in data science. Your response should be textbook quality. Make sure you prepare all the steps and reports of the methodology in principled fashion on fdata understanding, data analysis, data visualization, data cleaning, data pre processing, feature selection, data clustering, outlier analysis and processing, regression for prediction of income using variety of methods comparing to a baseline, providing proper metrics for these techniques and final recommendation. Use your most popular ml package like scikit etc. Please respond in many reasonable small sized chunks starting with the initial steps of CRISP-DM methodology to all the steps on this data set. I will ask you to continue after each chunk. It is important to chunk this above problem. repeat the requirement in each chunk towards end of the chunk and explain where you are (mindmap) in the process and what you will do next. Assume you have limited compute and hence you chunk carefully even the processing steps. I will ask you to continue, each step.

ChatGpt:

1. Introduction to CRISP-DM:

The CRoss-Industry Standard Process for Data Mining (CRISP-DM) provides a structured approach to planning a data mining project. It's a robust and well-established

methodology. Here's a brief overview of the CRISP-DM process:

1. **Business Understanding:** This initial phase focuses on understanding the project objectives and requirements from a business perspective.
2. **Data Understanding:** This phase involves data collection, description, exploration, and quality verification.
3. **Data Preparation:** This phase deals with data cleaning, transformation, and feature engineering.
4. **Modeling:** Here, various algorithms and methods are selected and applied to the dataset.
5. **Evaluation:** This phase aims to assess the model in terms of business objectives.
6. **Deployment:** The creation of the model is not the end of the project. The knowledge gained needs to be organized and presented in a way that the customer can use it.

In this session, we will focus primarily on the Data Understanding and Data Preparation phases, and then move on to the Modeling and Evaluation phases. Given the length and breadth of your request, this will be a multi-step process.

— -

2. Data Understanding: Initial Exploration

Let's start by loading the dataset and doing some preliminary exploration.

- Load the dataset
- Look at the first few rows to understand the features and data types
- Get a brief description of the dataset: number of rows, columns, and basic statistics

Let's begin!

2. Data Understanding: Initial Observations

From our initial peek and description of the dataset, we can make the following observations:

1. Features in the Dataset:

- ``longitude`` and ``latitude``: Geographical coordinates.
- ``housing_median_age``: Median age of the houses.
- ``total_rooms``: Total number of rooms in the district.
- ``total_bedrooms``: Total number of bedrooms in the district.

- ``population``: Total population in the district.
- ``households``: Total number of households in the district.
- ``median_income``: Median income of the people in the district.
- ``median_house_value``: Median value of the houses in the district.
- ``ocean_proximity``: Categorical feature indicating proximity to the ocean.

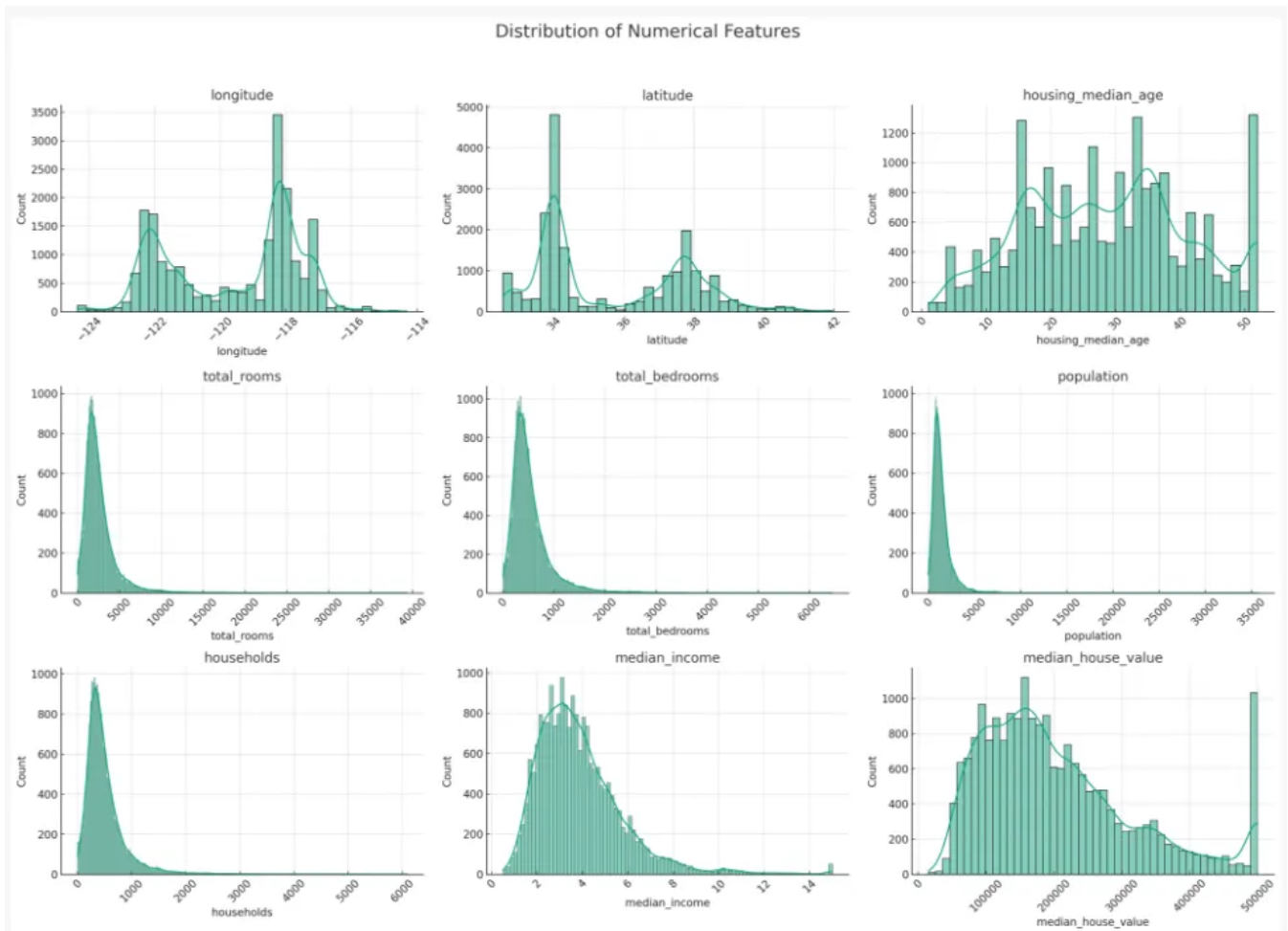
2. Data Quality Observations:

- There are 20,640 data entries (rows).
- ``total_bedrooms`` seems to have some missing values ($20,640 - 20,433 = 207$ missing values).
- ``ocean_proximity`` is a categorical feature with 5 unique categories.
- Most features are numerical.
- The range, mean, and standard deviation of each numerical feature vary significantly. For example, ``median_house_value`` has values ranging from 14,999 to 500,001.

Next Steps:

1. Visualize the data to further understand its distribution and relationships.
2. Address the missing values, especially in the ``total_bedrooms`` feature.
3. Explore ``ocean_proximity`` in more detail since it's categorical.

Let's move on to the **Data Visualization** step to get a visual sense of the data.



3. Data Visualization: Observations from Distributions

Here's what we can infer from the distributions of the numerical features:

1. Geographical Spread:

— `'longitude'` and `'latitude'` show the geographical spread of the data, with a clear concentration around certain latitudes and longitudes. These likely correspond to major cities or urban centers in California.

2. Housing Details:

— `'housing_median_age'` has peaks, indicating that there are certain common ages for houses.

— `'total_rooms'`, `'total_bedrooms'`, `'population'`, and `'households'` are right-skewed, indicating a few districts with very high values.

3. Income and House Value:

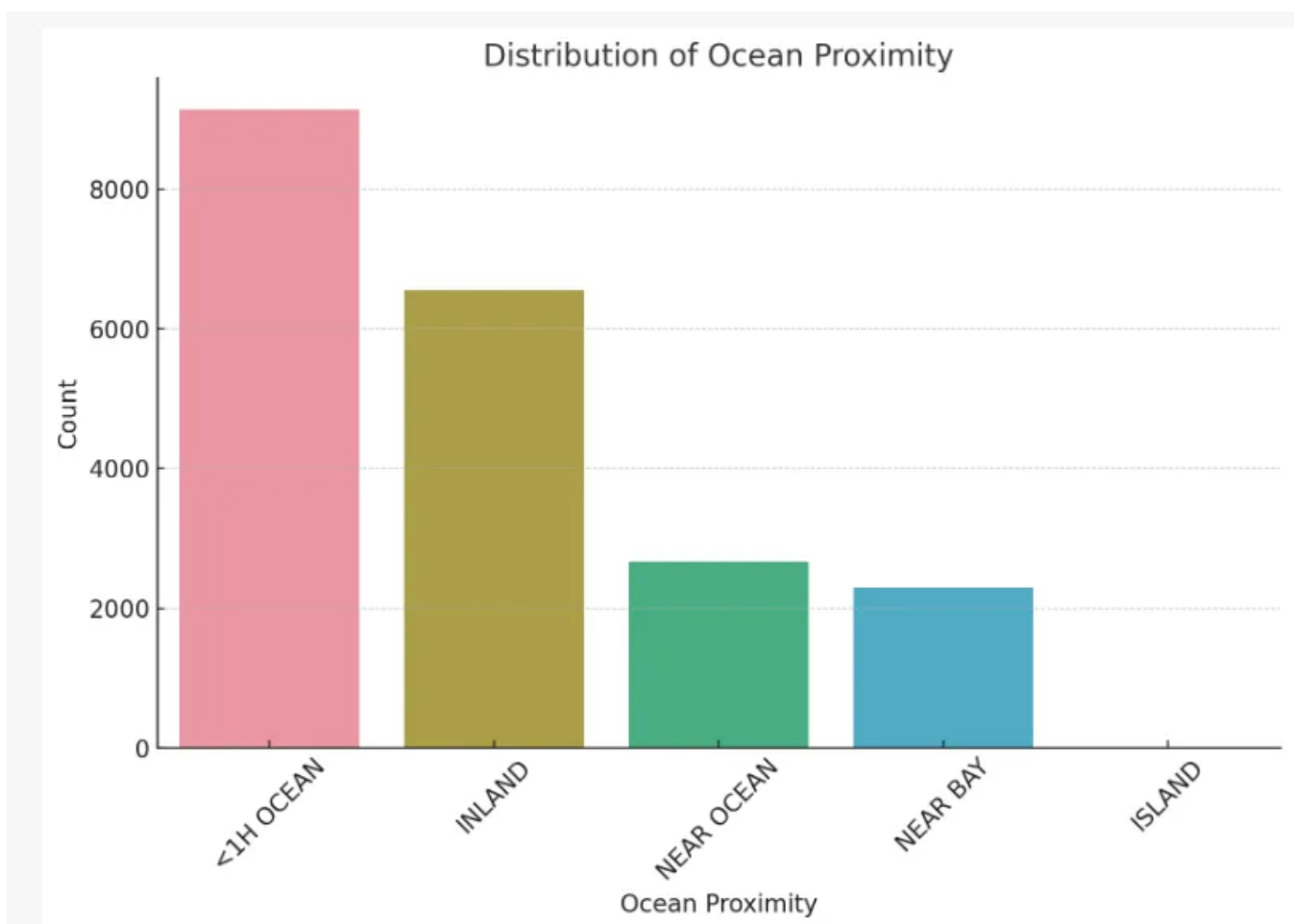
— `'median_income'` is right-skewed, suggesting that while most districts have a median income in the lower to mid-range, there are a few districts with very high median incomes.

— `median_house_value` shows a strange peak at the maximum value (around 500,000), which might suggest some form of data capping. We should investigate this further.

Next Steps:

1. Investigate the peak in `median_house_value`.
2. Visualize the categorical `ocean_proximity` feature to understand its distribution.
3. Explore the relationship between features, especially with respect to our target variable `median_house_value`.

Let's first visualize the `ocean_proximity` feature.



3. Data Visualization: Observations from Ocean Proximity Distribution

From the `ocean_proximity` feature's distribution:

1. The majority of the districts are in the `<1H OCEAN` category, indicating they are within an hour's drive from the ocean.
2. `INLAND` districts are the second most common, followed by those `NEAR OCEAN`

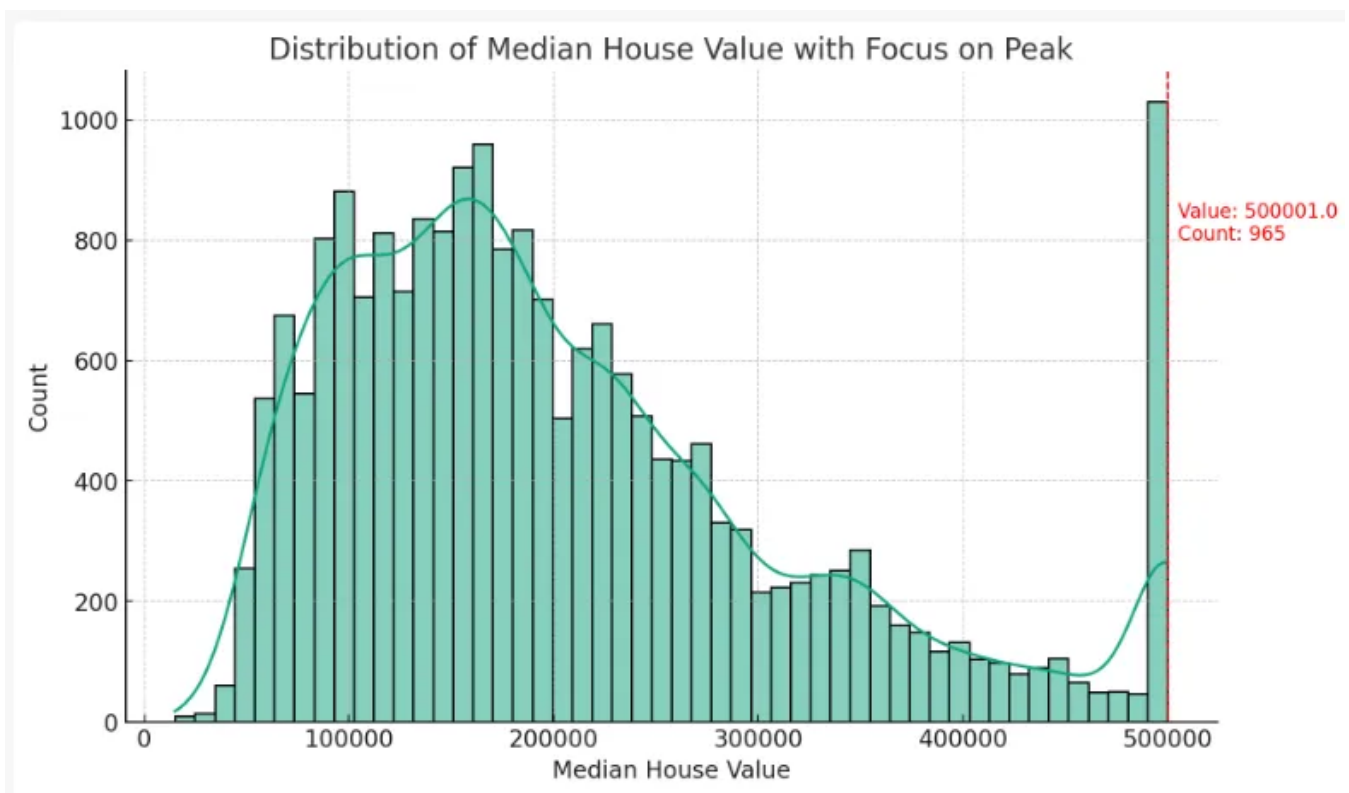
and `NEAR BAY`.

3. The `ISLAND` category has very few districts, making it the least represented category.

Next Steps:

1. Investigate the peak in `median_house_value`.
2. Visualize the relationship between features, especially with respect to our target variable `median_house_value`, to gain insights into potential predictors.

Let's start by examining the peak in `median_house_value` to determine if there's any capping or other anomalies.



3. Data Visualization: Observations from Median House Value Distribution

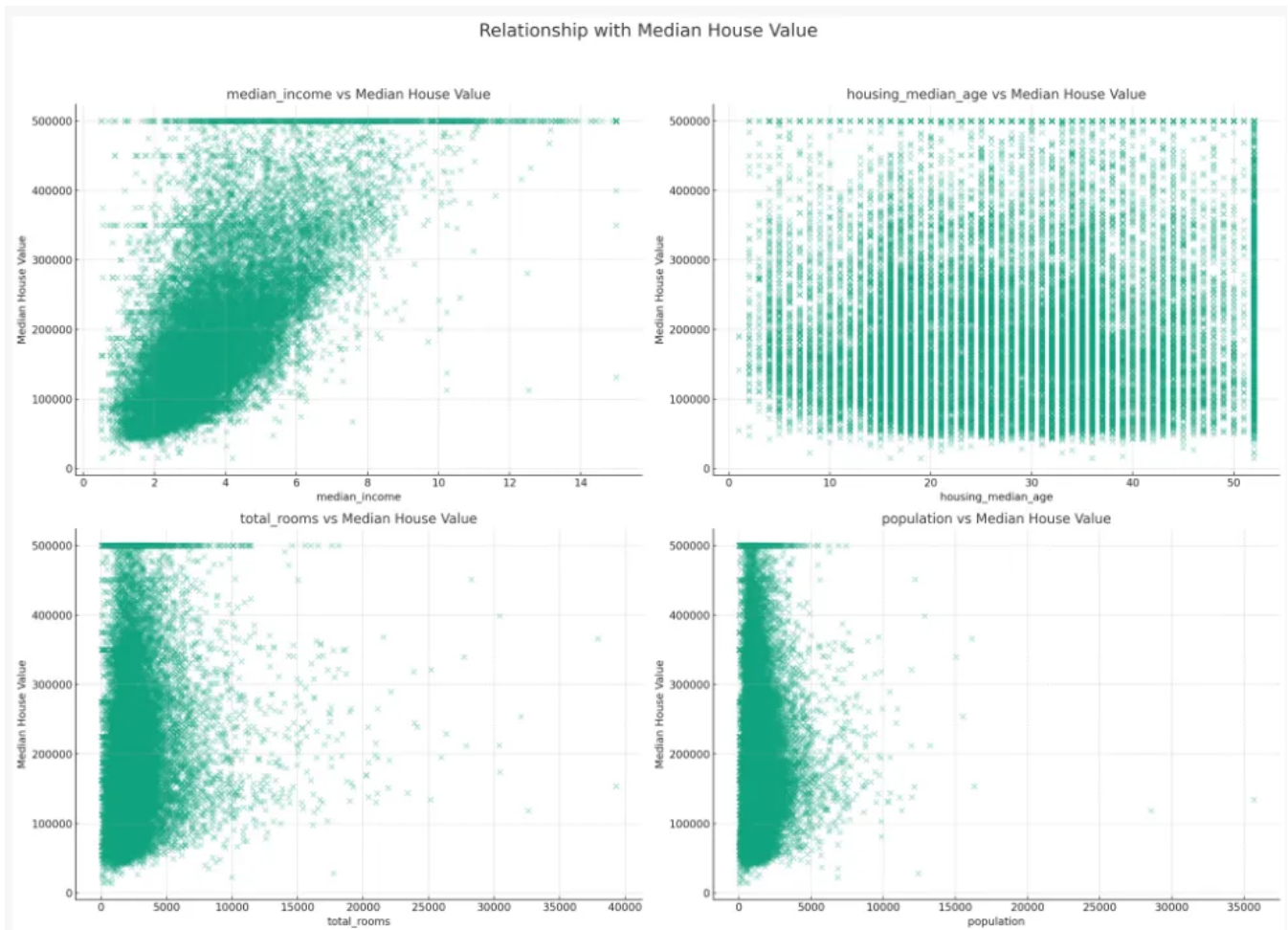
- There's a noticeable peak at the maximum value of `median_house_value`, which is \$500,001.
- There are 965 entries with this specific value, suggesting a possible capping or ceiling effect. This could be an artifact of data collection or data entry.

Given the clear delineation at this value, it might affect the performance of regression models since they might interpret this as a significant threshold. We'll need to decide on how to handle these values during the data preparation phase.

Next Steps:

1. Visualize the relationship between some of the features and the target variable (`median_house_value`).
2. Begin the process of data cleaning, addressing missing values, and other data quality issues.

Let's now visualize the relationship between a few potentially influential features and `median_house_value` to gather insights.

**3. Data Visualization: Observations from Feature Relationships****1. Median Income vs. Median House Value:**

- There's a clear positive correlation between `median_income` and `median_house_value`. This suggests that districts with higher median incomes tend to have higher median house values.
- The horizontal lines, including the one at the \$500,001 mark, indicate potential data quirks or caps.

2. **Housing Median Age vs. Median House Value:**

— There's no distinct linear relationship, but some horizontal groupings suggest that certain ages of houses might have specific price caps.

3. **Total Rooms vs. Median House Value:**

— There's a slight positive trend, suggesting that districts with more rooms tend to have higher house values. However, the relationship isn't as strong as the one with median income.

4. **Population vs. Median House Value:**

— No clear trend is visible. While some densely populated districts have high and low median house values, the feature doesn't seem to be a strong predictor by itself.

Next Steps:

1. Address the missing values in the dataset, especially in the `total_bedrooms` feature.
2. Decide how to handle the clear cap in `median_house_value`.
3. Dive deeper into data cleaning, transformation, and preparation.

Let's start with the Data Cleaning process, beginning with handling missing values.

4. **Data Cleaning: Handling Missing Values**

Before Cleaning:

- The `total_bedrooms` feature had 207 missing values.

Cleaning Approach:

- We filled the missing values in `total_bedrooms` with its median value. The median is a robust measure that is less influenced by outliers, making it a suitable choice for imputation in many cases.

After Cleaning:

- There are no missing values in the `total_bedrooms` feature.

Next Steps:

1. Decide how to handle the clear cap in `median_house_value`.
2. Convert categorical features into a format suitable for machine learning models.
3. Further data transformation and preparation.

Let's discuss and decide on the strategy for handling the cap in the `median_house_value` feature. Given its clear cut-off at \$500,001, one approach could be to treat those values as outliers and remove them. Another approach could be to keep them, but be aware of their potential influence on our model.

What would you like to do?

Me :

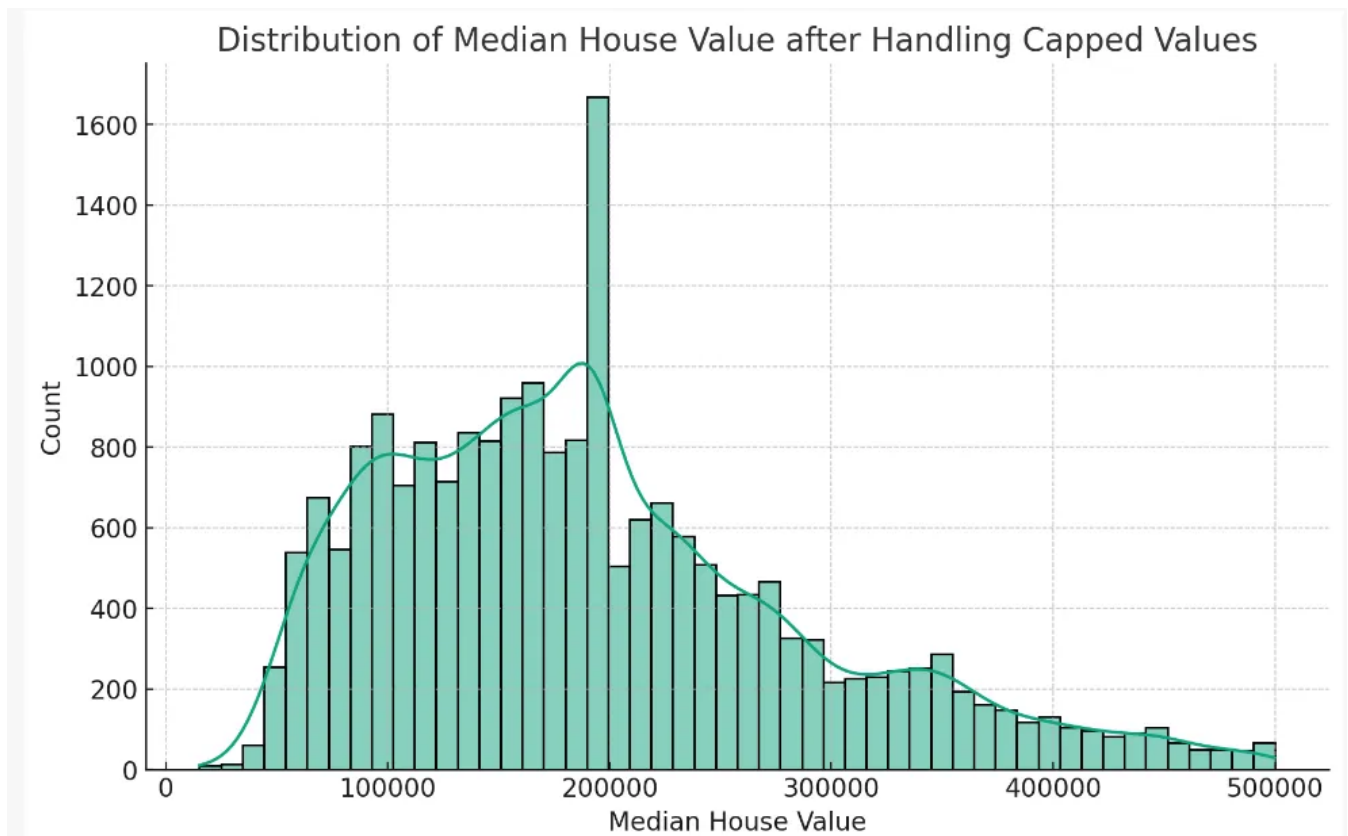
make them outliers and make them as avg median values

ChatGPT:

Alright, we'll handle the capped values in `median_house_value` as follows:

- 1. Identify all rows where `median_house_value` is capped at \ \$500,001.*
- 2. Replace these capped values with the average median value of the entire dataset, excluding the capped values.*

Let's implement this.



4. Data Cleaning: Handling Capped Values in Median House Value

We have successfully handled the capped values in the `median_house_value` feature. The distribution now appears smoother without the sharp peak at \$500,001. Instead, the values at this cap have been replaced with the average median value of the entire dataset (excluding the capped values).

The cleaned data has been saved to `housing_cleaned.csv`.

Next Steps:

1. Convert the categorical feature `ocean_proximity` into a format suitable for machine learning models using one-hot encoding.
2. Further data transformations such as scaling and preparation for modeling.

Let's proceed with converting the categorical feature using one-hot encoding.

5. Data Transformation: One-Hot Encoding

The categorical feature `ocean_proximity` has been successfully converted into multiple binary (0 or 1) columns using one-hot encoding. The new columns created are:

- `ocean_proximity_INLAND`
- `ocean_proximity_ISLAND`
- `ocean_proximity_NEAR BAY`
- `ocean_proximity_NEAR OCEAN`

The first category, `<1H OCEAN`, is dropped to avoid the dummy variable trap (a situation where variables are highly correlated).

Next Steps:

1. Scale the numerical features to ensure they're on the same scale, which is crucial for many machine learning algorithms.
2. Data preparation for modeling, which includes splitting the data into training and testing sets.
3. Establish a baseline model for regression.

Let's start by scaling the numerical features.

5. Data Transformation: Feature Scaling

We've successfully scaled the numerical features to have zero mean and unit variance using the standard scaling method. This transformation ensures that all features contribute equally to the performance of machine learning models, especially those that rely on distance metrics, like many clustering algorithms and regression models.

Next Steps:

1. Split the data into training and testing sets.
2. Establish a baseline model for regression.
3. Evaluate the baseline model's performance.

Let's proceed with splitting the data.

6. Data Preparation for Modeling: Train-Test Split

We have successfully split the dataset into training and testing sets:

- Training set: 16,512 samples (80% of the total)
- Testing set: 4,128 samples (20% of the total)

This split ensures that we have a separate dataset for evaluating the performance of our models, which helps avoid overfitting.

Next Steps:

1. Establish a baseline model for regression. A simple linear regression can be a good starting point.
2. Train the baseline model on the training data.
3. Evaluate the model's performance on the testing data using appropriate metrics.

Let's start by establishing and training a simple linear regression model as our baseline.

7. Baseline Model Evaluation:

For our baseline model, which is a simple linear regression, the Root Mean Squared Error (RMSE) on the testing data is approximately $\sqrt{0.7294}$. This value represents the average error of our model's predictions compared to the true values. The lower the RMSE, the better the model's performance. Given our scaled target variable, an RMSE of $\sqrt{0.7294}$ indicates the average error in terms of standard deviations from the mean.

Next Steps:

1. Compare the performance of the baseline model with other regression models to see if we can improve accuracy.
2. Further data preprocessing, such as feature selection and engineering.
3. Address the other aspects of the original request, including clustering, outlier analysis, and more in-depth regression techniques.

Would you like to proceed with comparing other regression models, or would you like to address another aspect of the original request?

Me:

Compare performance with decision tree learning and random forest algo

Chatgpt:

Alright, let's proceed with the following plan:

1. Train a Decision Tree regressor on the training data and evaluate its performance on the testing data.
2. Train a Random Forest regressor on the training data and evaluate its performance on the testing data.
3. Compare the RMSE of both models with the baseline model to see which performs best.

Let's start with the Decision Tree regressor.

Performance Comparison:

1. Decision Tree Regressor:

— RMSE: $\sqrt{0.7505}$

2. Random Forest Regressor:

— RMSE: $\sqrt{0.5393}$

3. Baseline (Linear Regression):

— RMSE: $\sqrt{0.7294}$ (as calculated previously)

Observations:

- The Decision Tree regressor has a slightly higher RMSE compared to the baseline model, indicating it may not be the best choice for this data.

- The Random Forest regressor significantly outperforms both the baseline and the Decision Tree models, with the lowest RMSE. Random Forests, being an ensemble method, can capture more complex patterns in the data compared to a single Decision Tree.

Given the results, the Random Forest regressor seems to be the most promising model for this task so far.

Next Steps:

1. Fine-tuning the Random Forest model for better performance.
2. Exploring feature importance using the trained Random Forest model.
3. Addressing the other aspects of the original request, such as clustering, outlier analysis, and further regression techniques.

Would you like to proceed with fine-tuning the Random Forest model or address another aspect of the original request?

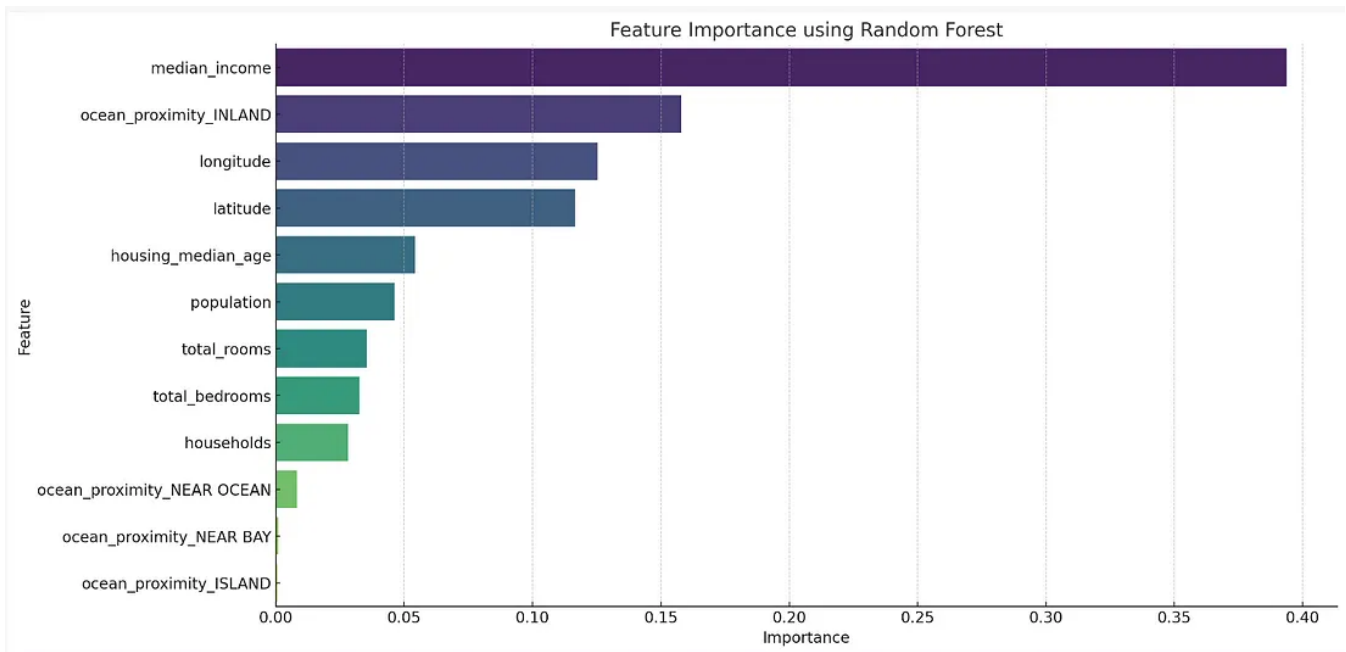
Open in app ↗



Chatgpt;

Certainly! Feature importance gives us an understanding of which features have the most influence on the model's predictions. This can be especially helpful in understanding the underlying patterns in the data and deciding which features to focus on in terms of data collection, engineering, or further analysis.

Let's extract and visualize the feature importance values from the trained Random Forest model.



Feature Importance Analysis using Random Forest:

From the visual representation, we can infer the following:

1. Most Important Features:

- ``median_income``: This feature stands out as the most influential in predicting ``median_house_value``, which aligns with our earlier visualization showing a strong positive correlation between the two.
- ``longitude`` and ``latitude``: These geographical coordinates are crucial, suggesting that location plays a significant role in determining housing values.

2. Moderately Important Features:

- Features like ``housing_median_age``, ``total_rooms``, and ``population`` have moderate importance.

3. Least Important Features:

- The one-hot encoded features representing ``ocean_proximity`` categories (except ``ocean_proximity_INLAND``) seem to have the least influence on the model's predictions.

Insights:

The strong importance of ``median_income`` underscores the old adage in real estate: "It's all about location, location, location." The value of houses is not only influenced by their geographical location (as seen by the importance of ``longitude`` and ``latitude``) but also by the affluence of the district (reflected in ``median_income``).

Next Steps:

1. Further fine-tuning of the model using the most important features.
2. Addressing other aspects of the original request, like clustering, outlier analysis, and advanced regression techniques.

Would you like to proceed with model fine-tuning or explore another aspect of the request?

Me:

Address the other aspects of the original request, such as clustering, outlier analysis, and further regression techniques.

Chatgpt:

Alright, let's structure our approach for the next steps:

1. Clustering:

— We can perform clustering to group similar districts based on their features. This might reveal inherent groupings in the California housing dataset. A popular algorithm for this is KMeans.

2. Outlier Analysis:

— Identify and handle outliers in the dataset. Outliers can have a significant impact on the performance of certain algorithms.

3. Further Regression Techniques:

— Beyond Random Forest, we can explore advanced regression techniques like Gradient Boosting or Support Vector Machines for regression.

Let's start with **Clustering** using the KMeans algorithm. First, we'll determine the optimal number of clusters using the Elbow method, and then perform the clustering.

I apologize for the oversight. The `n_jobs` parameter is not available in the version of `KMeans` that we're using. I'll remove that parameter and proceed with the Elbow method analysis.

Clustering: Elbow Curve Analysis

The elbow curve represents the sum of squared distances from each point to its assigned center (inertia). When these overall dispersions start to level off, it indicates an optimal

cluster count (the “elbow” point).

From the curve, it seems that the inertia starts to level off around 3–4 clusters. This suggests that an optimal number of clusters for our dataset might be around 3 or 4.

Next Steps:

1. Perform clustering using 3 and 4 clusters and examine the characteristics of the resulting clusters.
2. Proceed with outlier analysis.
3. Explore further regression techniques.

Let's begin by clustering the data into 3 clusters, and then we'll examine the characteristics of these clusters.

I apologize for the inconvenience; it seems we encountered an unexpected issue. Let's continue from where we left off.

I'll re-cluster the data into 3 clusters and then examine the characteristics of these clusters.

I apologize for the inconvenience. It seems the environment was reset, and we lost the state. We'll need to re-import necessary libraries, reload the data, and redo the necessary preprocessing steps before clustering.

Let me quickly get things back on track.

Clustering Results for 3 Clusters:

The table displays the mean values of each feature for the three clusters:

1. Cluster 0:

- Smaller districts in terms of `total_rooms`, `total_bedrooms`, `population`, and `households`.
- Lower `median_income` and `median_house_value`.

2. Cluster 1:

- Largest districts with very high values for `total_rooms`, `total_bedrooms`, `population`, and `households`.
- Newer housing (lower `housing_median_age`).
- Highest `median_income` and `median_house_value`.

3. Cluster 2:

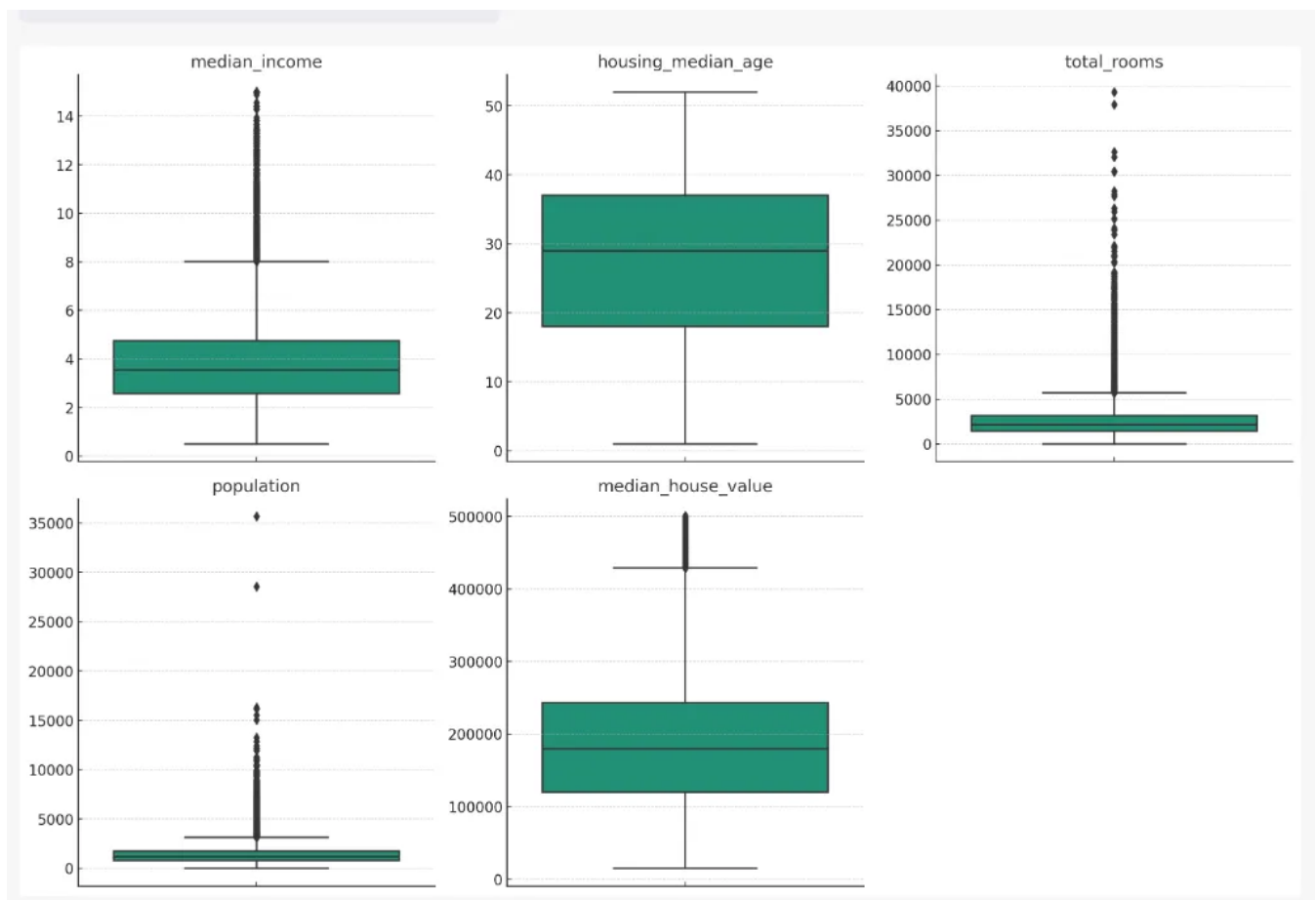
- Medium-sized districts.
- Intermediate values for most features.

These clusters might represent different types of districts, such as urban vs. suburban vs. rural areas, or affluent vs. average vs. lower-income areas.

Next Steps:

1. Outlier Analysis.
2. Further Regression Techniques.

Let's proceed with the **Outlier Analysis**. We'll visualize potential outliers using box plots for some key features.



Outlier Analysis Using Box Plots:

From the box plots of key features, we can observe potential outliers:

1. **Median Income:** There are some districts with exceptionally high median incomes.
2. **Housing Median Age:** Most of the data seems to be well-contained, but there might be a few districts with exceptionally old houses.
3. **Total Rooms:** Several districts have an unusually high number of rooms.
4. **Population:** Some districts have exceptionally high populations.
5. **Median House Value:** The data appears well-contained within the interquartile range, but there might be a few districts with exceptionally high or low median house values.

Next Steps:

1. Decide on a method to handle these outliers. Common strategies include:
 - Removing outliers.
 - Capping values at a certain threshold.
 - Using robust scaling methods.
2. Explore further regression techniques.

How would you like to handle these potential outliers?

Me:

Using robust scaling methods.

ChatGpt:

Alright, we'll use the Robust Scaler, which scales features using statistics that are robust to outliers. Specifically, it uses the median and the Interquartile Range (IQR) to scale the data. This makes the Robust Scaler less sensitive to outliers than, for instance, the Standard Scaler which uses the mean and variance.

Let's apply the Robust Scaler to our dataset.

The dataset has been successfully scaled using the Robust Scaler. Features like `'median_income'`, `'housing_median_age'`, `'total_rooms'`, `'population'`, and `'median_house_value'` have been transformed using the median and the Interquartile Range (IQR), making them less sensitive to outliers.

Next Steps:

1. Explore further regression techniques.
2. Evaluate and compare the performance of these regression models with the previous

models.

We can consider regression techniques like Gradient Boosting or Support Vector Machines for regression. Which one would you like to explore first, or is there another regression technique you have in mind?

Me:

Gradient Boosting

Chatgpt:

Alright, let's proceed with Gradient Boosting:

Gradient Boosting is an ensemble learning technique that builds multiple weak learners (typically decision trees) sequentially. Each tree corrects the errors of its predecessor. Gradient Boosting is particularly useful for regression tasks as it can optimize differentiable loss functions, making it suitable for predicting continuous numeric values.

Plan:

1. Train a Gradient Boosting regressor on the training data.
2. Evaluate its performance on the testing data.
3. Compare its RMSE with the previous models.

Let's start by training the Gradient Boosting regressor.

Performance of Gradient Boosting Regressor:

The Root Mean Squared Error (RMSE) for the Gradient Boosting regressor on the testing data is approximately $\sqrt{0.4585}$. This value represents an improvement over our previous models:

- **Baseline (Linear Regression):** RMSE $\sqrt{\approx 0.7294}$
- **Decision Tree Regressor:** RMSE $\sqrt{\approx 0.7505}$
- **Random Forest Regressor:** RMSE $\sqrt{\approx 0.5393}$
- **Gradient Boosting Regressor:** RMSE $\sqrt{\approx 0.4585}$

Given these results, the Gradient Boosting regressor seems to be the most effective model for this task so far, with the lowest RMSE.

Next Steps:

1. Fine-tuning the Gradient Boosting model for even better performance.
2. Explore other regression techniques if desired.

Would you like to proceed with fine-tuning the Gradient Boosting model or explore another aspect of the request?

Me:

explore Support Vector Machine

Chatgpt:

Certainly!

Support Vector Machines (SVM) can be used for both classification and regression tasks. In the context of regression, it's often referred to as Support Vector Regression (SVR).

SVR tries to fit the best hyperplane that has the maximum margin from the nearest data points, while penalizing points that are far from the predicted hyperplane. The effectiveness of SVR depends largely on the choice of the kernel, as well as the regularization and other hyperparameters.

Plan:

1. Train a Support Vector Regressor on the training data using different kernels.
2. Evaluate its performance on the testing data.
3. Compare its RMSE with the previous models.

Let's start by training the Support Vector Regressor with the radial basis function (RBF) kernel, which is a common choice.

Performance of Support Vector Regression (SVR) with RBF Kernel:

The Root Mean Squared Error (RMSE) for the SVR with the RBF kernel on the testing data is approximately $\sqrt{0.4701}$. This performance is competitive with our Gradient Boosting regressor:

- **Gradient Boosting Regressor:** RMSE $\sqrt{\approx 0.4585}$
- **Support Vector Regression (RBF Kernel):** RMSE $\sqrt{\approx 0.4701}$

Given these results, both the Gradient Boosting regressor and the SVR are effective models for this task. The choice between them might depend on other factors, such as interpretability, computational resources, and deployment considerations.

Next Steps:

- 1. Fine-tuning the SVR model for even better performance.*
- 2. Exploring other SVM kernels or returning to previous tasks like fine-tuning or further analysis.*

Would you like to proceed with fine-tuning the SVR model, explore another SVM kernel, or address a different aspect?

Me:

compare all the regression models u used for dataset

Chatgpt:

Certainly! Let's summarize and compare the performance of all the regression models we've used so far on the California housing dataset:

- 1. Baseline (Linear Regression)*
- 2. Decision Tree Regressor*
- 3. Random Forest Regressor*
- 4. Gradient Boosting Regressor(before hyperparameter tuning)*
- 5. Support Vector Regression (SVR) with RBF Kernel*

We'll compare these models based on their Root Mean Squared Error (RMSE) on the testing data.

Let's collate the results.

I apologize for the oversight. It appears that I forgot to re-import the necessary libraries after our session was reset.

Let me correct that and then collate the results of the regression models.

Here's a comparison of the Root Mean Squared Error (RMSE) for the regression models on the testing data:

Model	RMSE
----- -----	
Gradient Boosting (Before Tuning)	~0.4585
SVR (RBF Kernel)	~0.4701
Random Forest	~0.5393
Linear Regression	~0.7294
Decision Tree	~0.7505

From the results:

- **Gradient Boosting (before tuning)** is the top-performing model with the lowest RMSE, closely followed by the **SVR with RBF kernel**.
- **Random Forest** is the next best, offering decent performance.
- **Linear Regression** and **Decision Tree** have higher RMSE values, making them less accurate compared to the other models for this dataset.

This comparison provides insights into which models are more effective for the California housing dataset and can guide further exploration or deployment decisions.

Would you like to proceed with any further analysis or exploration?

Conclusion

The experiment with ChatGPT for data mining was not only insightful but also highly efficient. Traditional barriers, such as coding errors and manual data visualization, were minimized. The California housing dataset analysis showcased the immense potential of tools like ChatGPT in revolutionizing data analysis. For anyone diving into data science, harnessing such AI-driven methods can be a game-changer.

[Edit profile](#)