# Data Exploration

1.  Descriptive Statistics: Exploring basic statistics like mean, median, standard deviation, min/max values for each column. This gives a quick overview of the data's distribution and scale.
2.  Data Visualization:
    *   Histograms: Plotted histograms for each feature to understand their distribution. This can reveal skewness or outliers.
    *   Scatter Plots: Useful for examining relationships between features, especially with respect to the target variable ('median_house_value').
    *   Correlation Matrix: Generated a correlation matrix and visualize it using a heatmap to identify potential multicollinearity or strong correlations with the target variable.
3.  Missing Values: Checked for missing or NaN values in each column. The presence of missing data can significantly impact model performance.
4.  Categorical Data Analysis: If there are any categorical features, analyzing the distribution of categories and their relationship with the target variable.

# Preprocessing

1.  Handling Missing Data:
    *   Imputation: If missing values are found, we may consider imputing them using statistical methods (mean, median) or more complex techniques like k-NN imputation.
    *   Removal: In cases where imputation is not suitable, we may remove rows/columns with a high percentage of missing values.
2.  Outlier Detection and Treatment:
    *   Detect outliers using statistical methods (e.g., Z-scores, IQR).
    *   Decide whether to remove them, cap them, or use transformations to reduce their impact.
3.  Feature Engineering:
    *   New Features: Create new features that might be relevant for prediction (e.g., 'rooms_per_household').
    *   Transformation: Apply transformations (log, square root, etc.) to features that are not normally distributed.
4.  Encoding Categorical Variables:
    *   If there are categorical variables, use techniques like one-hot encoding or label encoding to convert them into a format suitable for machine learning algorithms.
5.  Dimensionality Reduction (Optional):
    *   Techniques like PCA (Principal Component Analysis) can be explored to reduce the number of features while retaining most of the information.

By thoroughly exploring and preprocessing the data, we not only gain valuable insights but also prepare the dataset for more effective modelling.

# Feature Selection

To determine which features are most useful for predicting 'median_house_value', we'll consider the following steps:

1. Correlation Analysis: We'll start by analyzing the correlation of each feature with the target variable. Features with a higher correlation may be more useful for our prediction.
2. Feature Importance using a Model: We can use a machine learning model like a Random Forest to determine the importance of each feature.
3. Multicollinearity Check: It's important to check for multicollinearity because highly correlated predictors can lead to unreliable and unstable estimates of regression coefficients.
4. Domain Knowledge Consideration: We should also consider domain knowledge. For example, features like 'latitude' and 'longitude' might be critical in real estate price prediction.

Based on the outputs from the feature selection techniques used, analyzing which features may be most useful in predicting 'median_house_value' and consider whether a subset of the features might be more appropriate:

## Correlation Analysis

- High Correlation: 'median_income' has a strong positive correlation (0.688075) with 'median_house_value'. This suggests it's likely a significant predictor.
- Low Correlation: Features like 'population', 'longitude', and 'latitude' show low to moderate negative correlations. While 'longitude' and 'latitude' may capture geographical trends, their impact might be less direct compared to other features.

## Feature Importances from Random Forest

- Most Important: 'median_income' is again the most important feature, followed by 'longitude' and 'latitude'.
- Less Important: 'total_rooms' and 'households' appear to be less important.

## Variance Inflation Factor (VIF) Scores

- High Multicollinearity: 'total_bedrooms', 'households', 'latitude', and 'longitude' have very high VIF scores, indicating significant multicollinearity.
- Acceptable VIF: 'median_income' and 'housing_median_age' have moderately high VIF scores but are still within an acceptable range.

## Analysis and Decision

- Key Feature: 'median_income' stands out as a key feature due to its high correlation with the target variable and its significant feature importance score.

- Geographic Features: Despite high multicollinearity, 'latitude' and 'longitude' might be capturing essential geographic trends. However, their direct inclusion might not be as effective due to multicollinearity issues.
- Reducing Multicollinearity: Features with high VIF scores, especially 'total_bedrooms' and 'households', might need to be dropped or combined to reduce multicollinearity. Alternatively, you could transform these features (e.g., combining 'total_rooms', 'total_bedrooms', and 'households' into more meaningful metrics like 'rooms_per_household').
- Subset Selection: Considering the correlation, feature importance, and VIF scores, a subset of features ('median_income', 'housing_median_age', 'latitude', 'longitude', and potentially transformed features) might be more effective.

**Techniques Utilized**

- Correlation Analysis: To understand linear relationships with the target variable.
- Random Forest Feature Importance: For empirical evidence of feature relevance.
- Variance Inflation Factor (VIF): To identify and address multicollinearity among predictors.

In conclusion, not all features in the dataset are equally useful for predicting 'median_house_value'. A subset of carefully selected and potentially transformed features, informed by the analysis, would likely yield a more effective and efficient predictive model.

# Feature Normalization

After selecting the relevant features, we need to normalize them. Normalization ensures that all features contribute equally to the prediction and improves the convergence of most algorithms.

**Selected Features**

Given this information, our feature selection prioritizes:

- 'median_income': Due to its high correlation and feature importance.
- 'housing_median_age': Due to its moderate correlation and importance, and relatively lower multicollinearity.

We will exclude 'latitude' and 'longitude' despite their importance due to extreme multicollinearity. Similarly, 'total_rooms', 'total_bedrooms', 'population', and 'households' are excluded due to their high multicollinearity and lower importance.

# Train, Validate, Test Split

To perform the train, validate, and test split, we will follow a common approach used in machine learning. The rationale for splitting the data is to ensure that we have separate datasets for training the model, tuning hyperparameters (validation), and finally, evaluating the model's performance on unseen data (testing). This helps in assessing the generalization capability of the model.

A typical split ratio is 60% for training, 20% for validation, and 20% for testing. However, these ratios can be adjusted based on the size of the dataset and the specific requirements of the project.

# Machine Learning Model 1

**Model Choice: Linear Regression**

**Why**: Linear Regression is a straightforward and effective baseline model for regression tasks. It's particularly suitable for datasets with linear relationships between features and the target variable.

**Building the Model**:

Step 1: Starting by creating and training a Linear Regression model using the training dataset. Then, you'll use the validation dataset to evaluate its performance.

Step 2: Hyper-parameter Tuning, for Linear Regression: In its basic form, Linear Regression doesn't have hyperparameters to tune.

# Machine Learning Model 2

**Model Choice: Decision Tree Regressor**

**Why**:

- Non-Linear Modeling: Unlike Linear Regression, Decision Trees can model non-linear relationships, which might be present in real estate data.
- Feature Interactions: Decision Trees naturally consider interactions between different features, which can be crucial in predicting house values.
- Interpretability: Despite being more complex than Linear Regression, Decision Trees are still relatively interpretable, as the decision process can be visualized and understood.

**Building the Model:**

Step 1: Building a Decision Tree Regressor model and evaluating it on the validation set.

Step 2: Hyperparameter Tuning, tuning hyperparameters such as max_depth, min_samples_split, and min_samples_leaf to improve the model. This can be done using GridSearchCV or RandomizedSearchCV.

# Machine Learning Model 3

**Model Choice: Random Forest Regressor**

For the third machine learning model in predicting 'median_house_value', considering the performance of previous models, a good choice would be the Random Forest Regressor. This model is a powerful ensemble learning method that often yields high accuracy and handles complex datasets with interrelated features effectively.

**Why**:

- Ensemble Learning: Random Forest is an ensemble of Decision Trees, which generally leads to better performance and robustness than a single Decision Tree. It reduces the risk of overfitting, which is often a concern with single Decision Trees.
- Handling Complex Interactions and Non-linearity: Random Forest can capture complex interactions between features and is adept at handling non-linear relationships, which are common in real estate data.
- Feature Importance Insights: Like Decision Trees, Random Forest provides insights into feature importance, helping in understanding which factors most significantly impact house values.
- Flexibility and Generalization: Random Forest works well with both numerical and categorical data and tends to generalize well to new, unseen data.

**Building the Model:**

Step 1: Building a Random Forest Regressor model and evaluating it on the validation set.

Step 2: Hyperparameter Tuning, for tuning, we'll adjust parameters like n_estimators, max_depth, and min_samples_split.