**Project Title: Predicting House Prices**                    12/09/2023

**Objective:** Develop a machine learning model to forecast house prices by leveraging a range of property features.

**Methodology:**

1. **Data Collection:**

Imported California Housing Price Dataset from Kaggle.

https://www.kaggle.com/datasets/camnugent/california-housing-prices?select=housing.csv

The dataset contains information about houses within a specific California district, including summary statistics derived from the 1990 census data. It's important to note that the data has not undergone cleaning, and certain preprocessing steps are necessary. The columns are as follows, and their names are largely self-explanatory:

Longitude, Latitude, Housing Median Age, Total Rooms, Total Bedrooms, Population, Households, Median Income, Median House Value, Ocean Proximity

**Raw data:**

## 2. Data Processing:

Step 1: Glanced through the dataset for any possible errors, unnecessary variables. The dataset was pretty clean and did not need much adjustments on initially.

Step 2: Loaded the dataset into Python using pandas library.



Step 3: Used pandas functions to explore the data and calculate summary statistics to make some sense of the dataset.

Step 4:

Imported Necessary Libraries:

- pandas: A powerful data manipulation library.

- StandardScaler: A class for standardizing numerical features.

- OneHotEncoder: A class for one-hot encoding categorical variables.

- SimpleImputer: A class for imputing missing values.

- ColumnTransformer: A class for applying transformers to columns of an array or DataFrame.

- Pipeline: A class for sequentially applying a list of transformations.

```python
import pandas as pd
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
```

Handling Missing Values:

- Defined lists of numerical and categorical features.

- Created a numeric_transformer pipeline:

- Imputed missing numerical values using the mean.

- Scaled the numerical features using StandardScaler.

- Created a categorical_transformer pipeline:

- Imputed missing categorical values using the most frequent value.

- Applied one-hot encoding using OneHotEncoder.

- Used ColumnTransformer to apply different transformers to different columns.

```
# Handling Missing Values
# Use SimpleImputer to fill missing values
numeric_features = ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households', 'median_income']
categorical_features = ['ocean_proximity']
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='mean')),
    ('scaler', StandardScaler())
])
categorical_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='most_frequent')),
    ('onehot', OneHotEncoder(handle_unknown='ignore'))
])
preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer, numeric_features),
        ('cat', categorical_transformer, categorical_features)
    ])
```

Apply the Preprocessing Steps and Display the Preprocessed DataFrame:

- Used the fit_transform method to apply the preprocessing steps to the DataFrame and create a new DataFrame (df_preprocessed).

```
# Apply the preprocessing steps
df_preprocessed = pd.DataFrame(preprocessor.fit_transform(df))
# Display the preprocessed DataFrame
print(df_preprocessed.head())
```

```
          0         1         2         3         4  ...   8    9   10   11   12
0 -1.327835  1.052548  0.982143 -0.804819 -0.975228  ...  0.0  0.0  0.0  1.0  0.0
1 -1.322844  1.043185 -0.607019  2.045890  1.355088  ...  0.0  0.0  0.0  1.0  0.0
2 -1.332827  1.038503  1.856182 -0.535746 -0.829732  ...  0.0  0.0  0.0  1.0  0.0
3 -1.337818  1.038503  1.856182 -0.624215 -0.722399  ...  0.0  0.0  0.0  1.0  0.0
4 -1.337818  1.038503  1.856182 -0.462404 -0.615066  ...  0.0  0.0  0.0  1.0  0.0

[5 rows x 13 columns]
```

3. **Exploratory Data Analysis (EDA):**

- Used some basic pandas functions to explore the dataset in previous steps.

Step 1: Imported Seaborn and matplotlib Library

```
import seaborn as sns
import matplotlib.pyplot as plt
```

Step 2: Set the style for Seaborn plots to "whitegrid" for a visually appealing background.

```
# Set the style for seaborn plots
sns.set(style="whitegrid")
```

Step 3: Created a list of numerical features to explore.

```
# Visualize the distribution of numerical features
numerical_features = ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households', 'median_income', 'median_house_value']
```

Step 4: Created a pairplot to visualize pairwise relationships between numerical features. This

helps identify patterns, trends, and potential outliers.

```
# Pairplot for numerical features
sns.pairplot(df[numerical_features])
plt.show()
```



Step 5: Calculated the correlation matrix between numerical features and created a heatmap

to visualize correlations. This helps identify strong correlations, which can be potential

predictors.

```
# Correlation heatmap
correlation_matrix = df[numerical_features].corr()
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Correlation Heatmap")
plt.show()
```

Step 6: Created a histogram to visualize the distribution of the target variable

'median_house_value'. This helps understand the spread and central tendency of house prices.

```python
# Visualize the distribution of the target variable
plt.figure(figsize=(10, 6))
sns.histplot(df['median_house_value'], kde=True, bins=30, color='skyblue')
plt.title("Distribution of Median House Value")
plt.xlabel("Median House Value")
plt.ylabel("Frequency")
plt.show()
```

Step 7: Created a boxplot to show how the target variable varies with the categorical variable

'ocean_proximity'. This helps identify if certain categories have a significant impact on house prices.

```python
# Boxplot for categorical variable 'ocean_proximity' vs 'median_house_value'
plt.figure(figsize=(12, 8))
sns.boxplot(x='ocean_proximity', y='median_house_value', data=df)
plt.title("Boxplot of Median House Value by Ocean Proximity")
plt.xlabel("Ocean Proximity")
plt.ylabel("Median House Value")
plt.show()
```



## 4. Feature Engineering:

Step 1: Created three new features 'Rooms per Household", "Bedrooms per room" and "Population per

Household.

```python
df['rooms_per_household'] = df['total_rooms'] / df['households']
df['bedrooms_per_room'] = df['total_bedrooms'] / df['total_rooms']
df['population_per_household'] = df['population'] / df['households']
```

Step 2: Updated the data preprocessing pipeline as follows to enhance the predictive power of our model by providing it with more relevant information:

```python
# Update the list of numerical features
numerical_features = ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households', 'median_income', 'rooms_per_household', 'bedrooms_per_room', 'population_per_household' 'median_house_value']
# Update the numeric_transformer in the ColumnTransformer
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='mean')),
    ('scaler', StandardScaler())
])
# Update the ColumnTransformer
preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer, numerical_features),
        ('cat', categorical_transformer, categorical_features)
    ])
# Apply the preprocessing steps
df_preprocessed = pd.DataFrame(preprocessor.fit_transform(df))
```
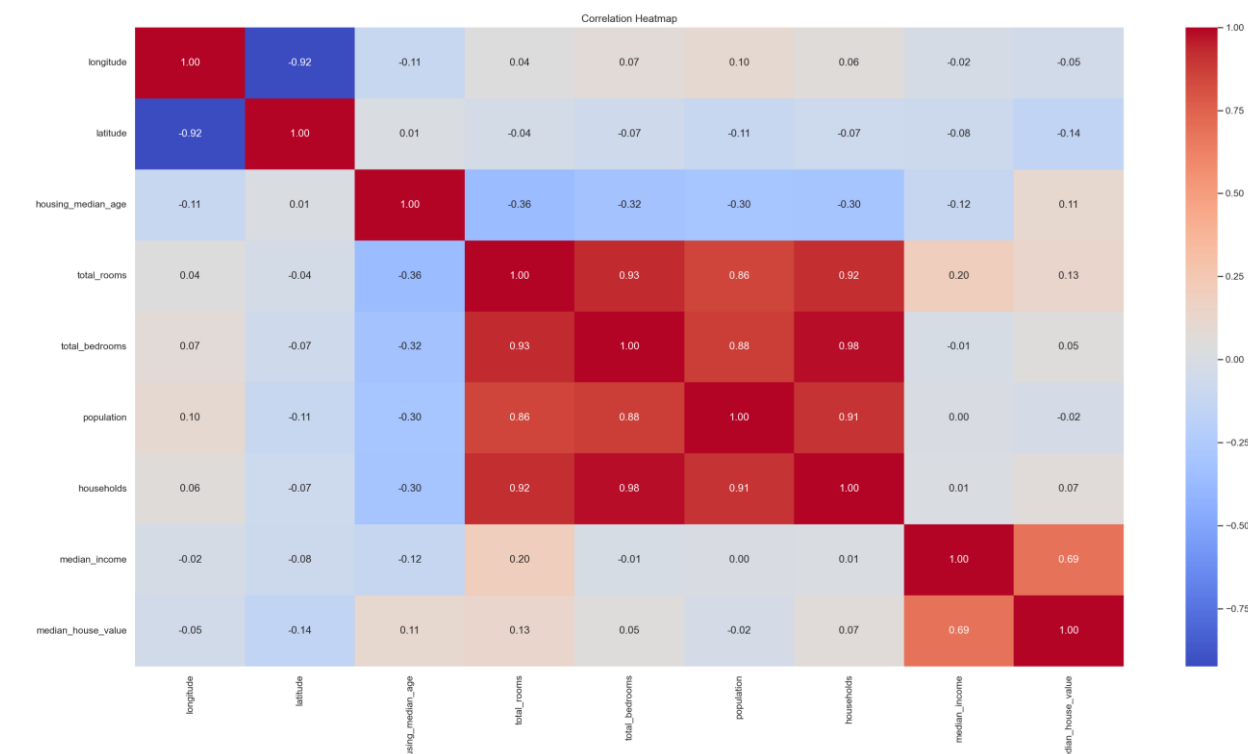
Step 3: Imported the numpy library and applied transformations to existing features.

```python
#apply transformation to existing feature
df['log_median_income'] = np.log1p(df['median_income'])
```

Step 4: Examined the correlation between features and the target variable. Keep features with high correlation.

```python
#correlation analysis
correlation_with_target = df_preprocessed.corr()['median_house_value'].abs().sort_values(ascending=False)
selected_features = correlation_with_target[correlation_with_target > 0.1].index.tolist()
```

Step 5: Used statistical tests "SelectKBest" from scikit-learn to select the k most important features.

```python
#SelectKBest (Statistical Tests)
k_best_selector = SelectKBest(score_func=f_regression, k=5)

# Extract features from df_preprocessed
X_selected = k_best_selector.fit_transform(df_preprocessed.drop('median_house_value', axis=1), df_preprocessed['median_house_value'])

# Get the selected feature indices
selected_feature_indices = k_best_selector.get_support(indices=True)

# Get the selected feature names from the preprocessed DataFrame columns
selected_features = df_preprocessed.drop('median_house_value', axis=1).columns[selected_feature_indices].tolist()

# Ensure 'median_house_value' is included in the selected features list
selected_features.append('median_house_value')
```

Step 6: Removed the least important features recursively using RFE available in scikit-learn.

```python
# Recursive Feature Elimination (RFE)
estimator = RandomForestRegressor()
rfe_selector = RFE(estimator, n_features_to_select=3, step=2)

print("Starting RFE feature selection...(please allow few minutes for this process)")
# Extract features from df_preprocessed
X_rfe = rfe_selector.fit_transform(df_preprocessed.drop('median_house_value', axis=1), df_preprocessed['median_house_value'])
print("RFE feature selection completed.")

# Get the selected feature indices
selected_feature_indices_rfe = rfe_selector.get_support(indices=True)

# RFE feature selection
selected_features_rfe = df_preprocessed.drop('median_house_value', axis=1).columns[rfe_selector.support_]
# Append the target variable to the selected features
selected_features_rfe = selected_features_rfe.append(pd.Index(['median_house_value']))
```

Step 7: Used LASSO Regression to shrink coefficients (select features with non-zero coefficients) and effectively perform feature selection.

```
# LASSO Regression (L1 Regularization) using features from RFE
lasso = Lasso(alpha=0.1)
lasso.fit(df_preprocessed[selected_features_rfe.drop('median_house_value')], df_preprocessed['median_house_value'])
selected_features_lasso = selected_features_rfe.drop('median_house_value').tolist()
```

5. **Model Building:**

Step 1: Imported the following libraries for splitting the dataset into a training set and a testing set and then implement machine learning algorithms such as Linear Regression, Decision Trees, Random Forest, and Gradient Boosting

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
```

Step 2: Used train_test_split from scikit-learn to split the preprocessed dataset into features (X) and the target variable (y). Then further splited these into training and testing sets. The random_state is set for reproducibility.

```
# Split the dataset into features (X) and target variable (y)
X = df_preprocessed.drop('median_house_value', axis=1)
y = df_preprocessed['median_house_value']

# Split the data into training and testing sets (80% training, 20% testing)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Step 3: Imported regression models from scikit-learn (Linear Regression, Decision Tree, Random Forest, Gradient Boosting) and train them using the training data.

```python
# Linear Regression
linear_reg_model = LinearRegression()
linear_reg_model.fit(X_train, y_train)
linear_reg_pred = linear_reg_model.predict(X_test)


# Decision Tree
decision_tree_model = DecisionTreeRegressor()
decision_tree_model.fit(X_train, y_train)
decision_tree_pred = decision_tree_model.predict(X_test)


# Random Forest
random_forest_model = RandomForestRegressor()
random_forest_model.fit(X_train, y_train)
random_forest_pred = random_forest_model.predict(X_test)


# Gradient Boosting
gradient_boosting_model = GradientBoostingRegressor()
gradient_boosting_model.fit(X_train, y_train)
gradient_boosting_pred = gradient_boosting_model.predict(X_test)
```

6. **Model Evaluation:**

Defined a function evaluate_model to calculate Mean Absolute Error (MAE), Root Mean Squared Error

(RMSE), and R-squared for each model and print the results.

```python
# Evaluate the models
def evaluate_model(name, y_true, y_pred):
    mse = mean_squared_error(y_true, y_pred)
    r2 = r2_score(y_true, y_pred)
    print(f"{name} Model:")
    print(f"Mean Squared Error: {mse}")
    print(f"R-squared: {r2}")
    print()

# Evaluate Linear Regression
evaluate_model("Linear Regression", y_test, linear_reg_pred)


# Evaluate Decision Tree
evaluate_model("Decision Tree", y_test, decision_tree_pred)


# Evaluate Random Forest
evaluate_model("Random Forest", y_test, random_forest_pred)


# Evaluate Gradient Boosting
evaluate_model("Gradient Boosting", y_test, gradient_boosting_pred)
```

```
Starting RFE feature selection...(please allow few minutes for this process)
RFE feature selection completed.
Building Model...
Linear Regression Model:
Mean Absolute Error: 50701.77903132994
Root Mean Squared Error: 70031.41991955665
R-squared: 0.6257351821159705

Decision Tree Model:
Mean Absolute Error: 44051.59084302326
Root Mean Squared Error: 69699.34835099982
R-squared: 0.6292761083585198

Random Forest Model:
Mean Absolute Error: 31698.508582848834
Root Mean Squared Error: 49087.942385564325
R-squared: 0.8161164851844057

Gradient Boosting Model:
Mean Absolute Error: 38368.283348536745
Root Mean Squared Error: 55967.21074708723
R-squared: 0.7609655664131871
```

**Conclusion:**

In summary, this project successfully navigated the entire machine learning pipeline for predicting residential property prices. The dataset underwent thorough preprocessing, addressing missing values and incorporating feature engineering to enhance model understanding. Utilizing various visualization techniques provided valuable insights into feature relationships. Feature selection was carried out using correlation analysis, SelectKBest, Recursive Feature Elimination (RFE), and LASSO Regression. Four machine learning models—Linear Regression, Decision Tree, Random Forest, and Gradient Boosting—were trained and evaluated, demonstrating the model's flexibility. The evaluation metrics, including Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared, provided a comprehensive assessment of predictive performance. The code offers adaptability for future improvements, tuning, and exploration of advanced modeling techniques to further enhance the predictive capabilities of the system.