

```
In [1]: import pandas as pd
import numpy as np

In [5]: path = '/Users/User/VSCodeFiles/Projects/CaliHousing/archive'
df = pd.read_csv(path+'//housing.csv')

In [8]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   longitude       20640 non-null   float64
 1   latitude        20640 non-null   float64
 2   housing_median_age 20640 non-null   float64
 3   total_rooms     20640 non-null   float64
 4   total_bedrooms  20640 non-null   float64
 5   population      20640 non-null   float64
 6   households      20640 non-null   float64
 7   median_income    20640 non-null   float64
 8   median_house_value 20640 non-null   float64
 9   ocean_proximity 20640 non-null   object  
dtypes: float64(9), object(1)
memory usage: 1.6+ MB

In [18]: null = df.total_bedrooms.isnull().sum()
print(null)

207
207 null vs 20,000 rows. Will just remove

In [24]: df1 = df.dropna().reset_index(drop=True)
df1.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20433 entries, 0 to 20432
Data columns (total 10 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   longitude       20433 non-null   float64
 1   latitude        20433 non-null   float64
 2   housing_median_age 20433 non-null   float64
 3   total_rooms     20433 non-null   float64
 4   total_bedrooms  20433 non-null   float64
 5   population      20433 non-null   float64
 6   households      20433 non-null   float64
 7   median_income    20433 non-null   float64
 8   median_house_value 20433 non-null   float64
 9   ocean_proximity 20433 non-null   object  
dtypes: float64(9), object(1)
memory usage: 1.6+ MB

In [28]: duplicates = df1.duplicated()
print(duplicates.sum())

0
No duplicates

In [38]: ## Need to one hot encode categorical data
pd.unique(df1['ocean_proximity'])

Out[38]: array(['NEAR BAY', '<1H OCEAN', 'INLAND', 'NEAR OCEAN', 'ISLAND'],
       dtype=object)

In [36]: df_dummies = pd.get_dummies(df1,columns=['ocean_proximity'])
df_dummies

Out[36]:
   longitude latitude housing_median_age total_rooms total_bedrooms population households median_income median_house_value
0 -122.23 37.88 41.0 880.0 129.0 322.0 126.0 8.3252 452600.0
1 -122.22 37.86 21.0 7099.0 1106.0 2401.0 1138.0 8.3014 358500.0
2 -122.24 37.85 52.0 1467.0 190.0 496.0 177.0 7.2574 352100.0
3 -122.25 37.85 52.0 1274.0 235.0 558.0 219.0 5.6431 341300.0
4 -122.25 37.85 52.0 1627.0 280.0 565.0 259.0 3.8462 342200.0
...
20428 -121.09 39.48 25.0 1665.0 374.0 845.0 330.0 1.5603 78100.0
20429 -121.21 39.49 18.0 697.0 150.0 356.0 114.0 2.5568 77100.0
20430 -121.22 39.43 17.0 2254.0 485.0 1007.0 433.0 1.7000 92300.0
20431 -121.32 39.43 18.0 1860.0 409.0 741.0 349.0 1.8672 84700.0
20432 -121.24 39.37 16.0 2785.0 616.0 1387.0 530.0 2.3886 89400.0
20433 rows × 14 columns

In [37]: import matplotlib.pyplot as plt
import seaborn as sns

In [39]: sns.pairplot(df_dummies)

Out[39]:
<seaborn.axisgrid.PairGrid at 0x1d8f9a292b0>

<img alt="A 14x14 grid of scatter plots showing the relationships between all pairs of numerical features in the dataset. The features include longitude, latitude, housing_median_age, total_rooms, total_bedrooms, population, households, median_income, median_house_value, and five binary columns for ocean_proximity (NEAR BAY, &lt;1H OCEAN, INLAND, NEAR OCEAN, ISLAND). The diagonal shows histograms of each individual feature. The off-diagonal shows scatter plots for every pair of features, revealing various correlations and patterns in the data." data-bbox="100 500 980 750>

Only obvious possible linear relationship with median house value is median income. Linear regression is not a viable option. Will try other classifiers. Going to use random forest classifiers, no need to preprocess data

In [43]: Y = df_dummies['median_house_value']
X = df_dummies.drop(labels='median_house_value', axis = 1)

In [47]: ## Splitting up data into training, validation, and testing data. 60/20/20
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.4)

In [52]: X_val, X_test, y_val, y_test = train_test_split(X_test, y_test, test_size=0.5)

In [61]: ## First model will be random forest regressor
## First use gridsearch to find best parameters

from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import max_error, r2_score

In [63]: params_rf = {
    'n_estimators':[50,75,125],
    'max_depth':[3,5,None],
    'min_samples_split':[2,3],
    'max_features':[4,6]
}

scoring = ('max_error', 'r2')

In [65]: rf = RandomForestRegressor(random_state=5)
rf_cv = GridSearchCV(rf, param_grid=params_rf, cv = 5, scoring=scoring, refit='r2')

In [69]: rf_cv.fit(X_train, y_train)

Out[69]:
* GridSearchCV
  * estimator: RandomForestRegressor
    * RandomForestRegressor
      * ...

In [70]: rf_cv.best_estimator_
* RandomForestRegressor
RandomForestRegressor(max_features=6, n_estimators=125, random_state=5)

In [71]: rf_cv.best_score_
0.809573320643534

In [72]: y_pred = rf_cv.predict(X_val)

In [78]: sns.scatterplot(x=y_pred, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 750 980 900>

Only obvious possible linear relationship with median house value is median income. Linear regression is not a viable option. Will try other classifiers. Going to use random forest classifiers, no need to preprocess data

In [43]: Y = df_dummies['median_house_value']
X = df_dummies.drop(labels='median_house_value', axis = 1)

In [47]: ## Splitting up data into training, validation, and testing data. 60/20/20
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.4)

In [52]: X_val, X_test, y_val, y_test = train_test_split(X_test, y_test, test_size=0.5)

In [61]: ## First model will be random forest regressor
## First use gridsearch to find best parameters

from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import max_error, r2_score

In [63]: params_rf = {
    'n_estimators':[50,75,125],
    'max_depth':[3,5,None],
    'min_samples_split':[2,3],
    'max_features':[4,6]
}

scoring = ('max_error', 'r2')

In [65]: rf = RandomForestRegressor(random_state=5)
rf_cv = GridSearchCV(rf, param_grid=params_rf, cv = 5, scoring=scoring, refit='r2')

In [69]: rf_cv.fit(X_train, y_train)

Out[69]:
* GridSearchCV
  * estimator: RandomForestRegressor
    * RandomForestRegressor
      * ...

In [70]: rf_cv.best_estimator_
* RandomForestRegressor
RandomForestRegressor(max_features=6, n_estimators=125, random_state=5)

In [71]: rf_cv.best_score_
0.809573320643534

In [72]: y_pred = rf_cv.predict(X_val)

In [78]: sns.scatterplot(x=y_pred, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 900 980 1050>

In [73]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [74]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

SVM Regressor not effective

In [75]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [76]: gb_cv.fit(X_train, y_train)

Out[76]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [77]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [78]: gb_cv.best_score_
0.8099260579842472

In [79]: y_pred_gb = gb_cv.predict(X_val)

In [80]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1050 980 1200>

In [81]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [82]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [83]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [84]: gb_cv.fit(X_train, y_train)

Out[84]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [85]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [86]: gb_cv.best_score_
0.8099260579842472

In [87]: y_pred_gb = gb_cv.predict(X_val)

In [88]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1200 980 1350>

In [89]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [90]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [91]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [92]: gb_cv.fit(X_train, y_train)

Out[92]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [93]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [94]: gb_cv.best_score_
0.8099260579842472

In [95]: y_pred_gb = gb_cv.predict(X_val)

In [96]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1350 980 1500>

In [97]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [98]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [99]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [100]: gb_cv.fit(X_train, y_train)

Out[100]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [101]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [102]: gb_cv.best_score_
0.8099260579842472

In [103]: y_pred_gb = gb_cv.predict(X_val)

In [104]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1500 980 1650>

In [105]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [106]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [107]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [108]: gb_cv.fit(X_train, y_train)

Out[108]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [109]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [110]: gb_cv.best_score_
0.8099260579842472

In [111]: y_pred_gb = gb_cv.predict(X_val)

In [112]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1650 980 1800]

In [113]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [114]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [115]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [116]: gb_cv.fit(X_train, y_train)

Out[116]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [117]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [118]: gb_cv.best_score_
0.8099260579842472

In [119]: y_pred_gb = gb_cv.predict(X_val)

In [120]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1800 980 1950]

In [121]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [122]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [123]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [124]: gb_cv.fit(X_train, y_train)

Out[124]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [125]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [126]: gb_cv.best_score_
0.8099260579842472

In [127]: y_pred_gb = gb_cv.predict(X_val)

In [128]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 1950 980 2100]

In [129]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [130]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [131]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [132]: gb_cv.fit(X_train, y_train)

Out[132]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [133]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [134]: gb_cv.best_score_
0.8099260579842472

In [135]: y_pred_gb = gb_cv.predict(X_val)

In [136]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 2100 980 2250]

In [137]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [138]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [139]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [140]: gb_cv.fit(X_train, y_train)

Out[140]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [141]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.0}

In [142]: gb_cv.best_score_
0.8099260579842472

In [143]: y_pred_gb = gb_cv.predict(X_val)

In [144]: sns.scatterplot(x=y_pred_gb, y=y_val)
plt.ylabel('Actual Values')
plt.xlabel('Predicted Values')
plt.title('Actual vs Predicted Values')
plt.show()

<img alt="A scatter plot titled 'Actual vs Predicted Values'. The x-axis is labeled 'Predicted Values' and ranges from 0 to 500,000. The y-axis is labeled 'Actual Values' and ranges from 0 to 500,000. The data points are blue dots, forming a dense cloud that follows the 1:1 line, indicating that the predicted values are generally close to the actual values." data-bbox="100 2250 980 2400]

In [145]: r2_gb = r2_score(y_val,y_pred_gb)
gb_regressor = {'Model':'Gradient Boosting Regressor','r2': r2_gb}
validation_table = validation_table.append(gb_regressor, ignore_index=True)

In [146]: validation_table
Model      r2
0  RandomForestRegressor  0.810541
1  SVMRegressor  0.076776
2  Gradient Boosting Regressor  0.806554

Gradient Boosting Regressor not effective

In [147]: from sklearn.ensemble import GradientBoostingRegressor

gb_params = {
    'learning_rate':[0.1,0.3,0.5],
    'n_estimators':[75,100,200],
    'subsample':[0.5,0.75,1.0],
    'min_samples_split':[2,3,5],
}
gb_cv = GridSearchCV(gb, param_grid=gb_params, scoring=scoring, cv=5, refit='r2')

In [148]: gb_cv.fit(X_train, y_train)

Out[148]:
* GridSearchCV
  * estimator: GradientBoostingRegressor
    * GradientBoostingRegressor
      * ...

In [149]: gb_cv.best_params_
{'learning_rate': 0.3,
 'min_samples_split': 3,
 'n_estimators': 200,
 'subsample': 1.
```

