

labset-1

February 17, 2025

```
[1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.datasets import fetch_california_housing
```

```
[2]: data = fetch_california_housing()
print(data)
```

```
{'data': array([[ 8.3252      , 41.          , 6.98412698, ...,
2.55555556,
37.88      , -122.23      ],
[ 8.3014      , 21.          , 6.23813708, ..., 2.10984183,
37.86      , -122.22      ],
[ 7.2574      , 52.          , 8.28813559, ..., 2.80225989,
37.85      , -122.24      ],
...,
[ 1.7         , 17.          , 5.20554273, ..., 2.3256351 ,
39.43      , -121.22      ],
[ 1.8672      , 18.          , 5.32951289, ..., 2.12320917,
39.43      , -121.32      ],
[ 2.3886      , 16.          , 5.25471698, ..., 2.61698113,
39.37      , -121.24      ]]), 'target': array([4.526, 3.585, 3.521,
..., 0.923, 0.847, 0.894]), 'frame': None, 'target_names': ['MedHouseVal'],
'feature_names': ['MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population',
'AveOccup', 'Latitude', 'Longitude'], 'DESCR': '..
_california_housing_dataset:\n\nCalifornia Housing
dataset\n-----\n\n**Data Set Characteristics:**\n\nNumber
of Instances: 20640\n\nNumber of Attributes: 8 numeric, predictive attributes
and the target\n\nAttribute Information:\n    - MedInc          median income in
block group\n    - HouseAge      median house age in block group\n    - AveRooms
average number of rooms per household\n    - AveBedrms        average number of
bedrooms per household\n    - Population  block group population\n    -
AveOccup            average number of household members\n    - Latitude      block
group latitude\n    - Longitude    block group longitude\n\nMissing Attribute
Values: None\n\nThis dataset was obtained from the StatLib
repository.\nhttps://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html\n\nThe
target variable is the median house value for California districts,\nexpressed
```

in hundreds of thousands of dollars (\$100,000).\n\nThis dataset was derived from the 1990 U.S. census, using one row per census\nblock group. A block group is the smallest geographical unit for which the U.S.\nCensus Bureau publishes sample data (a block group typically has a population\nof 600 to 3,000 people).\n\nA household is a group of people residing within a home. Since the average\nnumber of rooms and bedrooms in this dataset are provided per household, these\ncolumns may take surprisingly large values for block groups with few households\nand many empty houses, such as vacation resorts.\n\nIt can be downloaded/loaded using the\nfunc:`sklearn.datasets.fetch_california_housing` function.\n\n.. topic:: References\n\n- Pace, R. Kelley and Ronald Barry, Sparse Spatial Autoregressions,\nStatistics and Probability Letters, 33 (1997) 291-297\n']

```
[3]: # Step 1: Load the California Housing dataset
data = fetch_california_housing(as_frame=True)
housing_df = data.frame
print(housing_df)
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	\
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	
...	
20635	1.5603	25.0	5.045455	1.133333	845.0	2.560606	39.48	
20636	2.5568	18.0	6.114035	1.315789	356.0	3.122807	39.49	
20637	1.7000	17.0	5.205543	1.120092	1007.0	2.325635	39.43	
20638	1.8672	18.0	5.329513	1.171920	741.0	2.123209	39.43	
20639	2.3886	16.0	5.254717	1.162264	1387.0	2.616981	39.37	

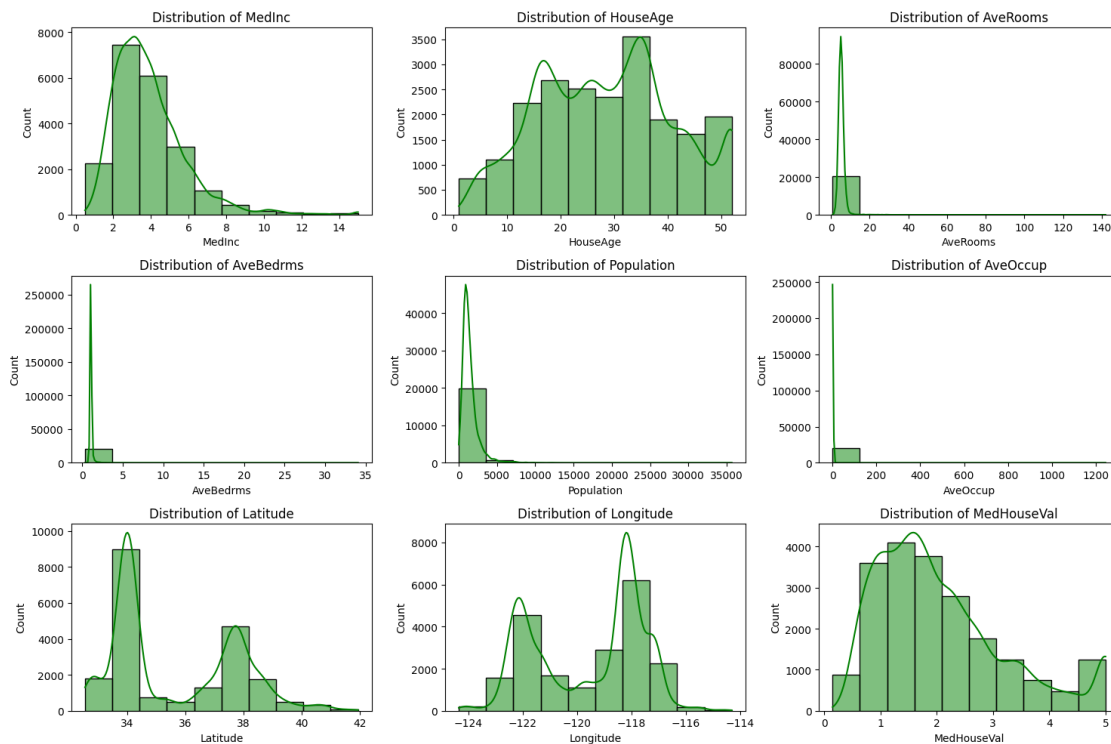
	Longitude	MedHouseVal
0	-122.23	4.526
1	-122.22	3.585
2	-122.24	3.521
3	-122.25	3.413
4	-122.25	3.422
...
20635	-121.09	0.781
20636	-121.21	0.771
20637	-121.22	0.923
20638	-121.32	0.847
20639	-121.24	0.894

[20640 rows x 9 columns]

```
[5]: # Step 2: Create histograms for numerical features
numerical_features = housing_df.select_dtypes(include = [np.number]).columns
print(numerical_features)
```

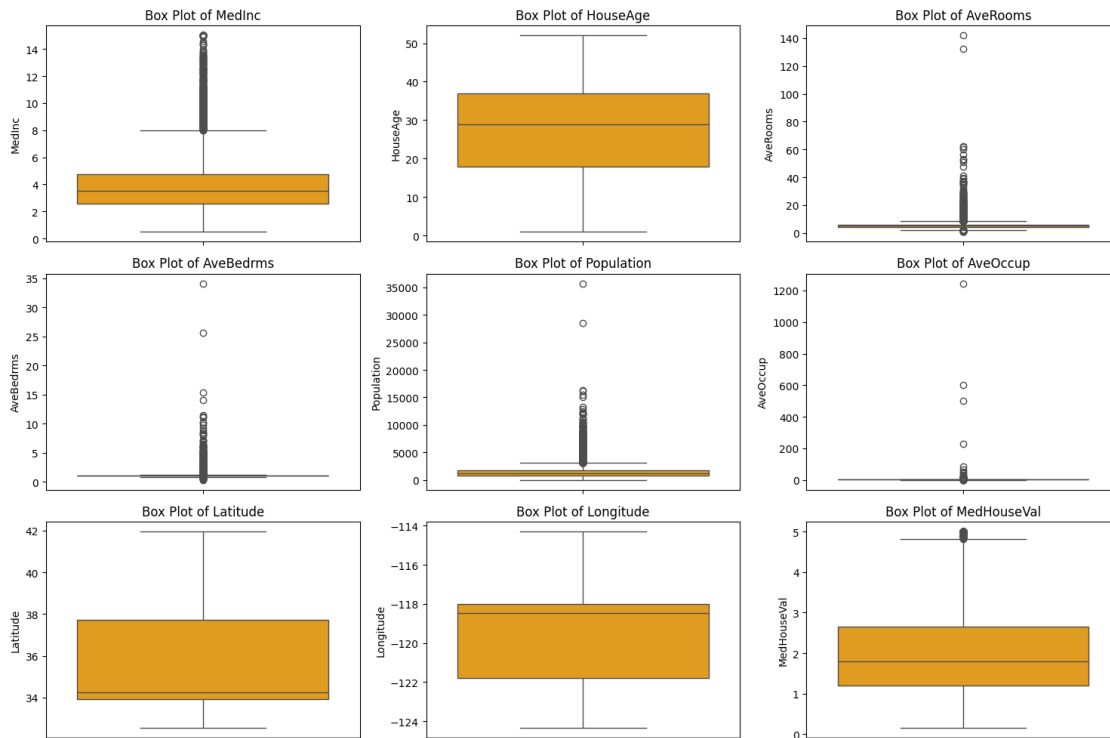
```
Index(['MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population', 'AveOccup',
      'Latitude', 'Longitude', 'MedHouseVal'],
      dtype='object')
```

```
[6]: # Plot histograms
plt.figure(figsize=(15, 10))
for i, feature in enumerate(numerical_features):
    plt.subplot(3, 3, i + 1)
    sns.histplot(housing_df[feature], kde = True, bins=10, color='green')
    plt.title(f'Distribution of {feature}')
plt.tight_layout()
```



```
[7]: # Step 3: Generate box plots for numerical features
plt.figure(figsize=(15, 10))
for i, feature in enumerate(numerical_features):
    plt.subplot(3, 3, i + 1)
    sns.boxplot(housing_df[feature], color='orange')
    plt.title(f'Box Plot of {feature}')
plt.tight_layout()
```

```
plt.show()
```



```
[8]: # Step 4: Identify outliers using the IQR method
print("Outliers Detection:")
outliers_summary = {}
for feature in numerical_features:
    Q1 = housing_df[feature].quantile(0.25)
    Q3 = housing_df[feature].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = housing_df[(housing_df[feature] < lower_bound) |
    ↪ (housing_df[feature] > upper_bound)]
    outliers_summary[feature] = len(outliers)
    print(f"{feature}: {len(outliers)} outliers")
```

Outliers Detection:
MedInc: 681 outliers
HouseAge: 0 outliers
AveRooms: 511 outliers
AveBedrms: 1424 outliers
Population: 1196 outliers
AveOccup: 711 outliers
Latitude: 0 outliers

Longitude: 0 outliers
MedHouseVal: 1071 outliers

```
[9]: # Optional: Print a summary of the dataset
print("\nDataset Summary:")
print(housing_df.describe())
```

Dataset Summary:

	MedInc	HouseAge	AveRooms	AveBedrms	Population \
count	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000
mean	3.870671	28.639486	5.429000	1.096675	1425.476744
std	1.899822	12.585558	2.474173	0.473911	1132.462122
min	0.499900	1.000000	0.846154	0.333333	3.000000
25%	2.563400	18.000000	4.440716	1.006079	787.000000
50%	3.534800	29.000000	5.229129	1.048780	1166.000000
75%	4.743250	37.000000	6.052381	1.099526	1725.000000
max	15.000100	52.000000	141.909091	34.066667	35682.000000

	AveOccup	Latitude	Longitude	MedHouseVal
count	20640.000000	20640.000000	20640.000000	20640.000000
mean	3.070655	35.631861	-119.569704	2.068558
std	10.386050	2.135952	2.003532	1.153956
min	0.692308	32.540000	-124.350000	0.149990
25%	2.429741	33.930000	-121.800000	1.196000
50%	2.818116	34.260000	-118.490000	1.797000
75%	3.282261	37.710000	-118.010000	2.647250
max	1243.333333	41.950000	-114.310000	5.000010

```
[ ]:
```