



# HOUSING

# DATA ANALYSIS

Prepared by :  
**chandana T M**  
**DA/DS ONLINE**  
**batch 28**



## INTRODUCTION

The objective of this project is to perform exploratory data analysis (EDA) on a housing dataset to uncover insights and patterns that can inform decision-making. The dataset includes various attributes of housing units such as price, location, size, and features. By conducting EDA, we aim to understand the underlying structure of the data, identify key factors affecting housing prices, and provide valuable insights for stakeholders such as real estate professionals and potential buyers.

## AIM

The aim of this project is to utilize exploratory data analysis techniques to:

- Understand the key characteristics and distributions of the housing dataset.
- Identify and address data quality issues.
- Derive meaningful metrics and insights that can help in predicting housing prices.
- Highlight significant trends, correlations, and patterns in the data to guide strategic decision-making.

## **Business Problem / Problem Statement**

In the competitive real estate market, accurately determining the factors that influence housing prices is crucial for both buyers and sellers. The problem is to identify which features of the housing data—such as location, size, and amenities—most significantly impact the price of a property. Understanding these factors will assist stakeholders in making informed decisions, optimizing property value assessments, and formulating pricing strategies. The dataset provides a comprehensive view of housing units, but extracting actionable insights requires careful analysis to uncover hidden trends and relationships.

## Project Workflow

The project workflow involves several key steps:

- **Data Collection:** Gather the housing dataset from reliable sources.
- **Data Understanding:** Explore and understand the structure and contents of the data.
- **Data Cleaning:** Address missing values, outliers, and inconsistent data.
- **Bivariate Analysis:** Investigate relationships between pairs of variables.
- **Multivariate Analysis:** Explore interactions between multiple variables.
- **Insights:** Synthesize findings to provide actionable insights and recommendations.

## Data Understanding

The housing dataset comprises various attributes, including but not limited to, property prices, location details, size (in square feet), number of bedrooms and bathrooms, and additional features such as garden or garage. Initial exploration involves examining the dataset's summary statistics, distribution of key variables, and basic visualizations. Understanding the data's structure helps in identifying the types of variables (categorical, numerical) and their potential relationships. This step also includes checking for the completeness and integrity of the data.

## Data Cleaning - Missing Values Imputation, Outliers, Handling Inconsistent Values

Data cleaning involves several processes:

**Missing Values Imputation:** Identify missing values and use appropriate methods to impute them, such as mean, median, or mode for numerical values, and mode or a specific category for categorical values.

**Outliers:** Detect outliers using statistical methods (e.g., Z-score, IQR) and decide whether to remove them or transform them, depending on their impact on the analysis.

**Handling Inconsistent Values:** Correct inconsistencies in categorical values (e.g., different spellings or formats) and standardize units and formats across the dataset. This ensures uniformity and accuracy in the analysis.

Identify and address missing values in the dataset through imputation or removal, ensuring data completeness.

Detect and rectify any inconsistencies or anomalies in the data, such as erroneous entries or irregular formatting.

```
d1=df['sqft_living'].sort_values(ascending=True)
```

```
d1
```

```
3778    370.0
2919    380.0
2416    420.0
1219    430.0
4184    490.0
...
59      NaN
60      NaN
61      NaN
62      NaN
63      NaN
Name: sqft_living, Length: 4600, dtype: float64
```

```
d2=df['sqft_living'].median()
```

```
d2
```

```
1980.0
```

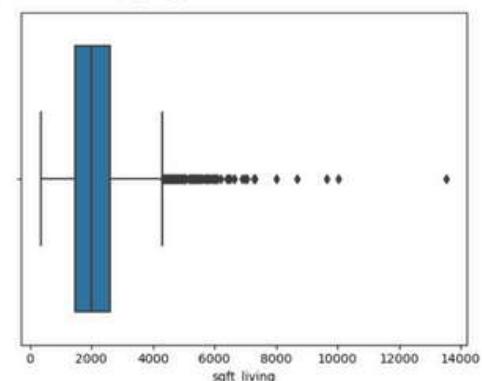
```
df['sqft_living'].fillna(d2,inplace=True)
```

```
df.isnull().sum()
```

```
date      0
price     0
bedrooms  0
bathrooms 0
sqft_living 0
...
...
```

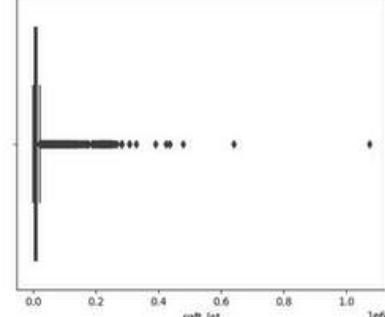
```
sns.boxplot(x=df['sqft_living'])
```

```
axes.xlabel='sqft_living'
```



```
sns.boxplot(x=df['sqft_lot'])
```

```
axes.xlabel='sqft_lot'
```



```
=df['sqft_lot'].sort_values(ascending=True)
```

```
608    638.0
15    681.0
913    704.0
604    746.0
247    747.0
...
...
NaN
0
1
2
3
Name: sqft_lot, Length: 4600, dtype: float64
```

```
1=df['sqft_lot'].median()
```

```
1
```

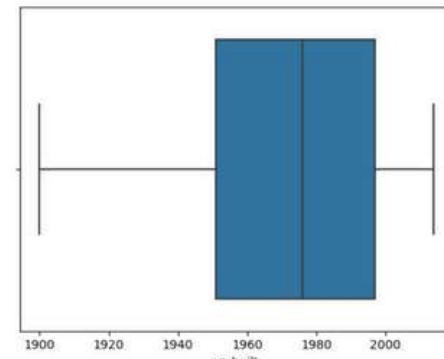
```
df['sqft_lot'].fillna(c1,inplace=True)
```

```
df.isnull().sum()
```

```
date      0
price     0
bedrooms  0
bathrooms 0
sqft_living 0
sqft_lot   0
floors     0
waterfront 0
...
...
```

```
sns.boxplot(x=df['yr_built'])
```

```
axes.xlabel='yr_built'
```



```
1931., 1941., 1947., 1948., 1949., 1950., 1951., 1952., 1953., 1954., 1955., 1956., 1957., 1958., 1959., 1960., 1961., 1962., 1963., 1964., 1965., 1966., 1967., 1968., 1969., 1970., 1971., 1972., 1973., 1974., 1975., 1976., 1977., 1978., 1979., 1980., 1981., 1982., 1983., 1984., 1985., 1986., 1987., 1988., 1989., 1990., 1991., 1992., 1993., 1994., 1995., 1996., 1997., 1998., 1999., 2000., 2001., 2002., 2003., 2004., 2005., 2006., 2007., 2008., 2009., 2010., 2011., 2012., 2013., 2014., 2015., 2016., 2017., 2018., 2019., 2020., 2021., 2022., 2023., 2024., 2025., 2026., 2027., 2028., 2029., 2030., 2031., 2032., 2033., 2034., 2035., 2036., 2037., 2038., 2039., 2040., 2041., 2042., 2043., 2044., 2045., 2046., 2047., 2048., 2049., 2050., 2051., 2052., 2053., 2054., 2055., 2056., 2057., 2058., 2059., 2060., 2061., 2062., 2063., 2064., 2065., 2066., 2067., 2068., 2069., 2070., 2071., 2072., 2073., 2074., 2075., 2076., 2077., 2078., 2079., 2080., 2081., 2082., 2083., 2084., 2085., 2086., 2087., 2088., 2089., 2090., 2091., 2092., 2093., 2094., 2095., 2096., 2097., 2098., 2099., 20100., 20101., 20102., 20103., 20104., 20105., 20106., 20107., 20108., 20109., 20110., 20111., 20112., 20113., 20114., 20115., 20116., 20117., 20118., 20119., 20120., 20121., 20122., 20123., 20124., 20125., 20126., 20127., 20128., 20129., 20130., 20131., 20132., 20133., 20134., 20135., 20136., 20137., 20138., 20139., 20140., 20141., 20142., 20143., 20144., 20145., 20146., 20147., 20148., 20149., 20150., 20151., 20152., 20153., 20154., 20155., 20156., 20157., 20158., 20159., 20160., 20161., 20162., 20163., 20164., 20165., 20166., 20167., 20168., 20169., 20170., 20171., 20172., 20173., 20174., 20175., 20176., 20177., 20178., 20179., 20180., 20181., 20182., 20183., 20184., 20185., 20186., 20187., 20188., 20189., 20190., 20191., 20192., 20193., 20194., 20195., 20196., 20197., 20198., 20199., 201000., 201001., 201002., 201003., 201004., 201005., 201006., 201007., 201008., 201009., 201010., 201011., 201012., 201013., 201014., 201015., 201016., 201017., 201018., 201019., 201020., 201021., 201022., 201023., 201024., 201025., 201026., 201027., 201028., 201029., 201030., 201031., 201032., 201033., 201034., 201035., 201036., 201037., 201038., 201039., 201040., 201041., 201042., 201043., 201044., 201045., 201046., 201047., 201048., 201049., 201050., 201051., 201052., 201053., 201054., 201055., 201056., 201057., 201058., 201059., 201060., 201061., 201062., 201063., 201064., 201065., 201066., 201067., 201068., 201069., 201070., 201071., 201072., 201073., 201074., 201075., 201076., 201077., 201078., 201079., 201080., 201081., 201082., 201083., 201084., 201085., 201086., 201087., 201088., 201089., 201090., 201091., 201092., 201093., 201094., 201095., 201096., 201097., 201098., 201099., 201100., 201101., 201102., 201103., 201104., 201105., 201106., 201107., 201108., 201109., 201110., 201111., 201112., 201113., 201114., 201115., 201116., 201117., 201118., 201119., 201120., 201121., 201122., 201123., 201124., 201125., 201126., 201127., 201128., 201129., 201130., 201131., 201132., 201133., 201134., 201135., 201136., 201137., 201138., 201139., 201140., 201141., 201142., 201143., 201144., 201145., 201146., 201147., 201148., 201149., 201150., 201151., 201152., 201153., 201154., 201155., 201156., 201157., 201158., 201159., 201160., 201161., 201162., 201163., 201164., 201165., 201166., 201167., 201168., 201169., 201170., 201171., 201172., 201173., 201174., 201175., 201176., 201177., 201178., 201179., 201180., 201181., 201182., 201183., 201184., 201185., 201186., 201187., 201188., 201189., 201190., 201191., 201192., 201193., 201194., 201195., 201196., 201197., 201198., 201199., 201200., 201201., 201202., 201203., 201204., 201205., 201206., 201207., 201208., 201209., 201210., 201211., 201212., 201213., 201214., 201215., 201216., 201217., 201218., 201219., 201220., 201221., 201222., 201223., 201224., 201225., 201226., 201227., 201228., 201229., 201230., 201231., 201232., 201233., 201234., 201235., 201236., 201237., 201238., 201239., 201240., 201241., 201242., 201243., 201244., 201245., 201246., 201247., 201248., 201249., 201250., 201251., 201252., 201253., 201254., 201255., 201256., 201257., 201258., 201259., 201260., 201261., 201262., 201263., 201264., 201265., 201266., 201267., 201268., 201269., 201270., 201271., 201272., 201273., 201274., 201275., 201276., 201277., 201278., 201279., 201280., 201281., 201282., 201283., 201284., 201285., 201286., 201287., 201288., 201289., 201290., 201291., 201292., 201293., 201294., 201295., 201296., 201297., 201298., 201299., 201300., 201301., 201302., 201303., 201304., 201305., 201306., 201307., 201308., 201309., 201310., 201311., 201312., 201313., 201314., 201315., 201316., 201317., 201318., 201319., 201320., 201321., 201322., 201323., 201324., 201325., 201326., 201327., 201328., 201329., 201330., 201331., 201332., 201333., 201334., 201335., 201336., 201337., 201338., 201339., 201340., 201341., 201342., 201343., 201344., 201345., 201346., 201347., 201348., 201349., 201350., 201351., 201352., 201353., 201354., 201355., 201356., 201357., 201358., 201359., 201360., 201361., 201362., 201363., 201364., 201365., 201366., 201367., 201368., 201369., 201370., 201371., 201372., 201373., 201374., 201375., 201376., 201377., 201378., 201379., 201380., 201381., 201382., 201383., 201384., 201385., 201386., 201387., 201388., 201389., 201390., 201391., 201392., 201393., 201394., 201395., 201396., 201397., 201398., 201399., 201400., 201401., 201402., 201403., 201404., 201405., 201406., 201407., 201408., 201409., 201410., 201411., 201412., 201413., 201414., 201415., 201416., 201417., 201418., 201419., 201420., 201421., 201422., 201423., 201424., 201425., 201426., 201427., 201428., 201429., 201430., 201431., 201432., 201433., 201434., 201435., 201436., 201437., 201438., 201439., 201440., 201441., 201442., 201443., 201444., 201445., 201446., 201447., 201448., 201449., 201450., 201451., 201452., 201453., 201454., 201455., 201456., 201457., 201458., 201459., 201460., 201461., 201462., 201463., 201464., 201465., 201466., 201467., 201468., 201469., 201470., 201471., 201472., 201473., 201474., 201475., 201476., 201477., 201478., 201479., 201480., 201481., 201482., 201483., 201484., 201485., 201486., 201487., 201488., 201489., 201490., 201491., 201492., 201493., 201494., 201495., 201496., 201497., 201498., 201499., 201500., 201501., 201502., 201503., 201504., 201505., 201506., 201507., 201508., 201509., 201510., 201511., 201512., 201513., 201514., 201515., 201516., 201517., 201518., 201519., 201520., 201521., 201522., 201523., 201524., 201525., 201526., 201527., 201528., 201529., 201530., 201531., 201532., 201533., 201534., 201535., 201536., 201537., 201538., 201539., 201540., 201541., 201542., 201543., 201544., 201545., 201546., 201547., 201548., 201549., 201550., 201551., 201552., 201553., 201554., 201555., 201556., 201557., 201558., 201559., 201560., 201561., 201562., 201563., 201564., 201565., 201566., 201567., 201568., 201569., 201570., 201571., 201572., 201573., 201574., 201575., 201576., 201577., 201578., 201579., 201580., 201581., 201582., 201583., 201584., 201585., 201586., 201587., 201588., 201589., 201590., 201591., 201592., 201593., 201594., 201595., 201596., 201597., 201598., 201599., 201600., 201601., 201602., 201603., 201604., 201605., 201606., 201607., 201608., 201609., 201610., 201611., 201612., 201613., 201614., 201615., 201616., 201617., 201618., 201619., 201620., 201621., 201622., 201623., 201624., 201625., 201626., 201627., 201628., 201629., 201630., 201631., 201632., 201633., 201634., 201635., 201636., 201637., 201638., 201639., 201640., 201641., 201642., 201643., 201644., 201645., 201646., 201647., 201648., 201649., 201650., 201651., 201652., 201653., 201654., 201655., 201656., 201657., 201658., 201659., 201660., 201661., 201662., 201663., 201664., 201665., 201666., 201667., 201668., 201669., 201670., 201671., 201672., 201673., 201674., 201675., 201676., 201677., 201678., 201679., 201680., 201681., 201682., 201683., 201684., 201685., 201686., 201687., 201688., 201689., 201690., 201691., 201692., 201693., 201694., 201695., 201696., 201697., 201698., 201699., 201700., 201701., 201702., 201703., 201704., 201705., 201706., 201707., 201708., 201709., 201710., 201711., 201712., 201713., 201714., 201715., 201716., 201717., 201718., 201719., 201720., 201721., 201722., 201723., 201724., 201725., 201726., 201727., 201728., 201729., 201730., 201731., 201732., 201733., 201734., 201735., 201736., 201737., 201738., 201739., 201740., 201741., 201742., 201743., 201744., 201745., 201746., 201747., 201748., 201749., 201750., 201751., 201752., 201753., 201754., 201755., 201756., 201757., 201758., 201759., 201760., 201761., 201762., 201763., 201764., 201765., 201766., 201767., 201768., 201769., 201770., 201771., 201772., 201773., 201774., 201775., 201776., 201777., 201778., 201779., 201780., 201781., 201782., 201783., 201784., 201785., 201786., 201787., 201788., 201789., 201790., 201791., 201792., 201793., 201794., 201795., 201796., 201797., 201798., 201799., 201800., 201801., 201802., 201803., 201804., 201805., 201806., 201807., 201808., 201809., 201810., 201811., 201812., 201813., 201814., 201815., 201816., 201817., 201818., 201819., 201820., 201821., 201822., 201823., 201824., 201825., 201826., 201827., 201828., 201829., 201830., 201831., 201832., 201833., 201834., 201835., 201836., 201837., 201838., 201839., 201840., 201841., 201842., 201843., 201844., 201845., 201846., 201847., 201848., 201849., 201850., 201851., 201852., 201853., 201854., 201855., 201856., 201857., 201858., 201859., 201860., 201861., 201862., 201863., 201864., 201865., 201866., 201867., 201868., 201869., 201870., 201871., 201872., 201873., 201874., 201875., 201876., 201877., 201878., 201879., 201880., 201881., 201882., 201883., 201884., 201885., 201886., 201887., 201888., 201889., 201890., 201891., 201892., 201893., 201894., 201895., 201896., 201897., 201898., 201899., 201900., 201901., 201902., 201903., 201904., 201905., 201906., 201907., 201908., 201909., 201910., 201911., 201912., 201913., 201914., 201915., 201916., 201917., 201918., 201919., 201920., 201921., 201922., 201923., 201924., 201925., 201926., 201927., 201928., 201929., 201930., 201931., 201932., 201933., 201934., 201935., 201936., 201937., 201938., 201939., 201940., 201941., 201942., 201943., 201944., 201945., 201946., 201947., 201948., 201949., 201950., 201951., 201952., 201953., 201954., 201955., 201956., 201957., 201958., 201959., 201960., 201961., 201962., 201963., 201964., 201965., 201966., 201967., 201968., 201969., 201970., 201971., 201972., 201973., 201974., 201975., 201976., 201977., 201978., 201979., 201980., 201981., 201982., 201983., 201984., 201985., 201986., 201987., 201988., 201989., 201990., 201991., 201992., 201993., 201994., 201995., 201996., 201997., 201998., 201999., 
```

## **Bivariate Analysis**

Bivariate analysis explores relationships between two variables.

Techniques include:

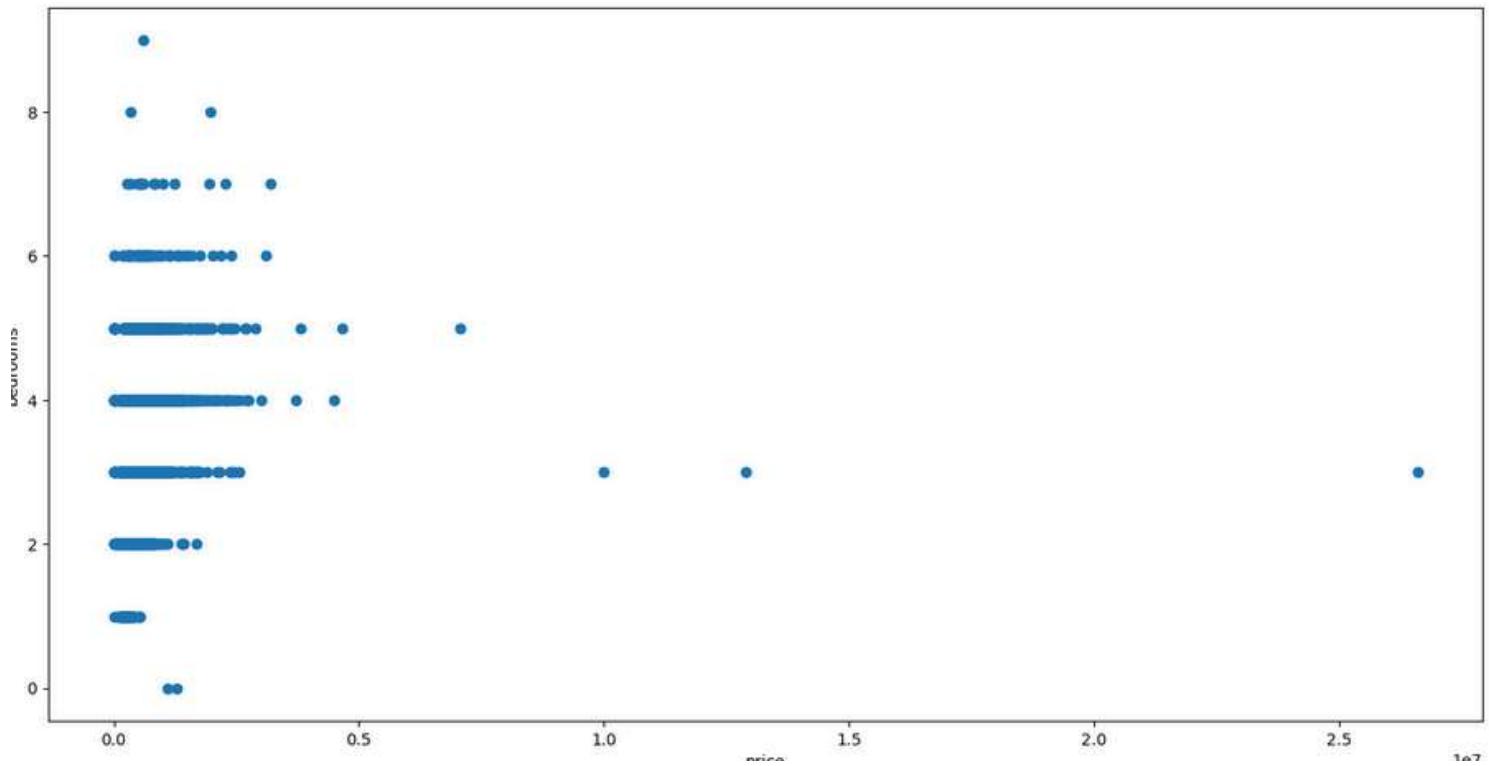
Scatter Plots: To visualize relationships between numerical variables like price and size.

Correlation Analysis: Calculating correlation coefficients to assess the strength and direction of relationships.

Box Plots: To compare distributions of numerical variables across different categories.

Bivariate analysis helps in identifying trends, correlations, and potential causations between variables.

```
fig, ax = plt.subplots(figsize=(16,8))
ax.scatter(df['price'], df['bedrooms'])
ax.set_xlabel('price')
ax.set_ylabel('bedrooms')
plt.show()
```



## Multivariate Analysis

### **correlation**

Correlation is used to find the relationship between two variables which is important in real life because we can predict the value of one variable with the help of other variables, who is being correlated with it. It is a type of Bivariate statistics since two variables are involved here.

Correlation is a statistical measure that expresses the extent to which two variables are linearly related.

```
q1sdata.quantile(0.25)
q1
price      322500.00
bedrooms     3.00
bathrooms    1.75
sqft_living   1470.00
sqft_lot      5001.00
floors       1.00
waterfront    0.00
view         0.00
condition     3.00
sqft_above    1190.00
sqft_basement  0.00
yr_built      1951.00
yr_renovated   0.00
Name: 0.25, dtype: float64

q3sdata.quantile(0.75)
q3
price      655000.0
bedrooms     4.0
bathrooms    2.5
sqft_living   2610.0
sqft_lot      11000.0
floors       2.0
waterfront    0.0
view         0.0
condition     4.0
sqft_above    2300.0
...  
price        241
bedrooms     210
bathrooms    141
sqft_living  134
sqft_lot      945
floors       8
waterfront    33
view        460
condition     6
sqft_above    121
sqft_basement  82
yr_built      0
yr_renovated   0
dtype: int64

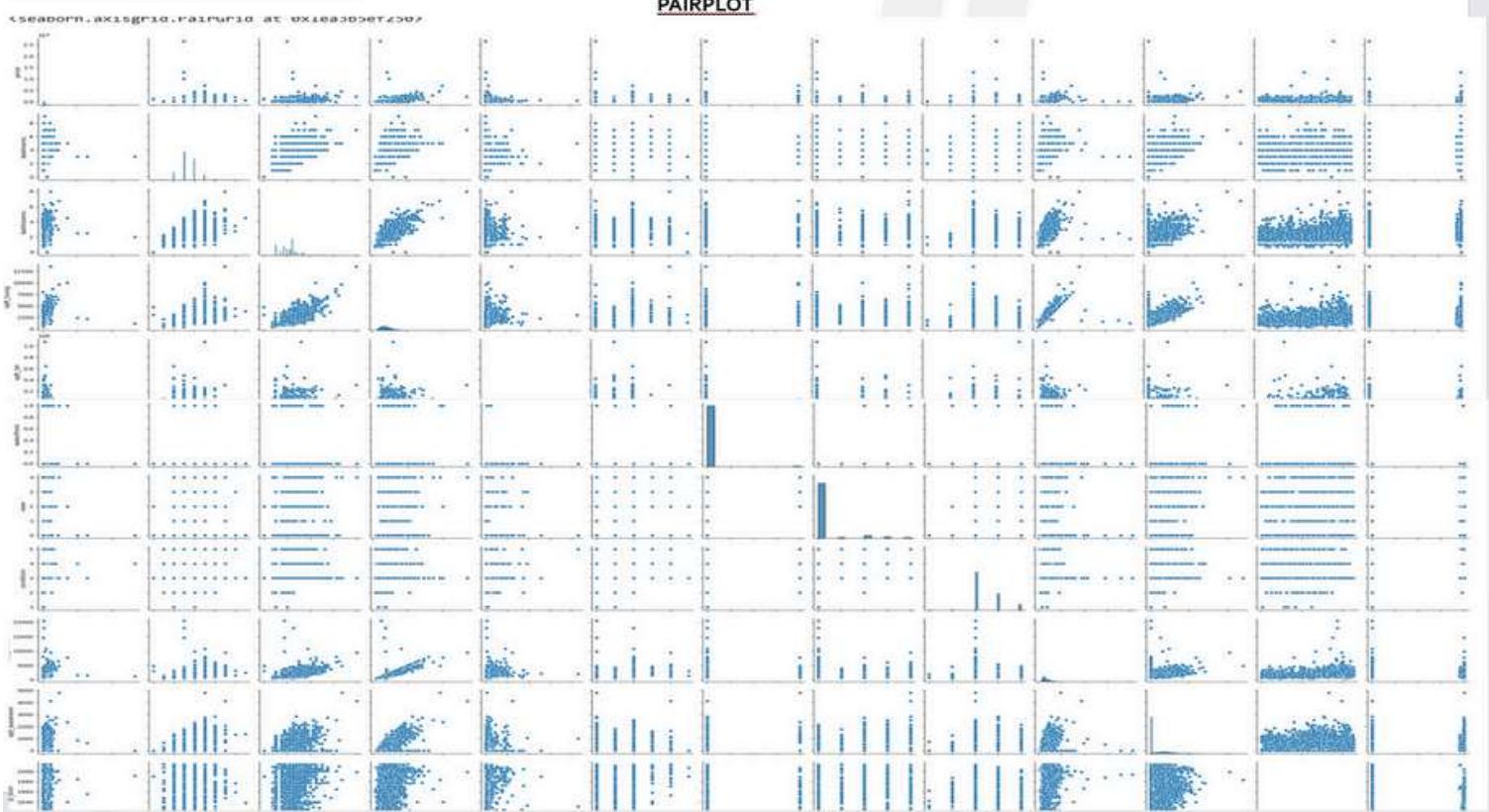
df.shape
(4600, 18)

data_clean = data[~((data < (q1 - 1.5 * IQR)) | (data > (q3 + 1.5 * IQR))).any(axis=1)]
data_clean.shape
```

```
data.corr()
```

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	sqft_above	sqft_basement	yr_built	yr_renovated
<b>price</b>	1.000000	0.192793	0.319311	0.419000	0.048293	0.143775	0.131387	0.220758	0.036660	0.311732	0.208591	0.023030	-0.031068
<b>bedrooms</b>	0.192793	1.000000	0.545920	0.593302	0.068796	0.177895	-0.003483	0.111028	0.025080	0.429573	0.334165	0.141070	-0.061082
<b>bathrooms</b>	0.319311	0.545920	1.000000	0.760145	0.107251	0.486428	0.076232	0.211960	-0.119994	0.612828	0.298020	0.461237	-0.215886
<b>sqft_living</b>	0.419000	0.593302	0.760145	1.000000	0.209090	0.343607	0.118055	0.310315	-0.061593	0.772572	0.446905	0.284979	-0.123501
<b>sqft_lot</b>	0.048293	0.068796	0.107251	0.209090	1.000000	0.003832	0.017283	0.074024	0.000406	0.200370	0.034200	0.050154	-0.022048
<b>floors</b>	0.143775	0.177895	0.486428	0.343607	0.003832	1.000000	0.022024	0.031211	-0.275013	0.460266	-0.255510	0.465032	-0.233996
<b>waterfront</b>	0.131387	-0.003483	0.076232	0.118055	0.017283	0.022024	1.000000	0.360935	0.000352	0.068900	0.097501	-0.024839	0.008625
<b>view</b>	0.220758	0.111028	0.211960	0.310315	0.074024	0.031211	0.360935	1.000000	0.063077	0.150515	0.321602	-0.065374	0.022967
<b>condition</b>	0.036660	0.025080	-0.119994	-0.061593	0.000406	-0.275013	0.000352	0.063077	1.000000	-0.168267	0.200632	-0.398093	-0.186818
<b>sqft_above</b>	0.311732	0.429573	0.612828	0.772572	0.200370	0.460266	0.068900	0.150515	-0.168267	1.000000	-0.034381	0.369645	-0.143899
<b>sqft_basement</b>	0.208591	0.334165	0.298020	0.446905	0.034200	-0.255510	0.097501	0.321602	0.200632	-0.034381	1.000000	-0.160952	0.043125
<b>yr_built</b>	0.023030	0.141070	0.461237	0.284979	0.050154	0.465032	-0.024839	-0.065374	-0.398093	0.369645	-0.160952	1.000000	-0.320952
<b>yr_renovated</b>	-0.031068	-0.061082	-0.215886	-0.123501	-0.022048	-0.233996	0.008625	0.022967	-0.186818	-0.143899	0.043125	-0.320952	1.000000

```
sns.pairplot(data)
```



## **Overall Insights Obtained from Analysis**

Based on the EDA, several key insights are typically obtained:

Price Drivers: Identifying which variables (e.g., location, size, number of bedrooms) most significantly impact housing prices.

Trends and Patterns: Uncovering trends such as price fluctuations over time or differences across neighborhoods.

Segment Characteristics: Understanding the characteristics of different property segments or clusters.

Recommendations: Providing actionable recommendations for stakeholders, such as pricing strategies, investment opportunities, or market positioning.

These insights help in making data-driven decisions and understanding market dynamics.

## **Conclusion**

In conclusion, the exploratory data analysis of the housing dataset has provided valuable insights into the factors affecting housing prices and the overall market trends. By addressing data quality issues, deriving meaningful metrics, and performing detailed analyses, we have gained a deeper understanding of the housing market. These findings offer practical implications for real estate professionals and potential buyers, enabling more informed decisions and strategic planning.