

Problem statement

predicting the house price in USA.To create a model to help him estimate of what the house would sell for.

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: df=pd.read_csv(r"C:\Users\user\Desktop\Ash\Datasets\USA_Housing.csv")
df
```

Out[2]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.458574	5.682861	7.009188	4.09	23086.800503	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanielstown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386
...
4995	60567.944140	7.830362	6.137356	3.46	22837.361035	1.060194e+06	USNS Williams\nFPO AP 30153- 7653
4996	78491.275435	6.999135	6.576763	4.02	25616.115489	1.482618e+06	PSC 9258, Box 8489\nAPO AA 42991-3352
4997	63390.686886	7.250591	4.805081	2.13	33266.145490	1.030730e+06	4215 Tracy Garden Suite 076\nJoshualand, VA 01...
4998	68001.331235	5.534388	7.130144	5.44	42625.620156	1.198657e+06	USS Wallace\nFPO AE 73316
4999	65510.581804	5.992305	6.792336	4.07	46501.283803	1.298950e+06	37778 George Ridges Apt. 509\nEast Holly, NV 2...

5000 rows × 7 columns

To display top 10 rows

```
In [3]: df.head(10)
```

```
Out[3]:
```

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.458574	5.682861	7.009188	4.09	23086.800503	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanielstown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386
5	80175.754159	4.988408	6.104512	4.04	26748.428425	1.068138e+06	06039 Jennifer Islands Apt. 443\nTracyport, KS...
6	64698.463428	6.025336	8.147760	3.41	60828.249085	1.502056e+06	4759 Daniel Shoals Suite 442\nNguyenburgh, CO ...
7	78394.339278	6.989780	6.620478	2.42	36516.358972	1.573937e+06	972 Joyce Viaduct\nLake William, TN 17778-6483
8	59927.660813	5.362126	6.393121	2.30	29387.396003	7.988695e+05	USS Gilbert\nFPO AA 20957
9	81885.927184	4.423672	8.167688	6.10	40149.965749	1.545155e+06	Unit 9446 Box 0958\nDPO AE 97025

Data Cleaning And Pre-Processing

In [4]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 7 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Avg. Area Income                      5000 non-null   float64
1   Avg. Area House Age                   5000 non-null   float64
2   Avg. Area Number of Rooms             5000 non-null   float64
3   Avg. Area Number of Bedrooms          5000 non-null   float64
4   Area Population                       5000 non-null   float64
5   Price                                 5000 non-null   float64
6   Address                               5000 non-null   object
dtypes: float64(6), object(1)
memory usage: 273.6+ KB
```

In [5]: *# Display the statistical summary*
df.describe()

Out[5]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price
count	5000.000000	5000.000000	5000.000000	5000.000000	5000.000000	5.000000e+03
mean	68583.108984	5.977222	6.987792	3.981330	36163.516039	1.232073e+06
std	10657.991214	0.991456	1.005833	1.234137	9925.650114	3.531176e+05
min	17796.631190	2.644304	3.236194	2.000000	172.610686	1.593866e+04
25%	61480.562388	5.322283	6.299250	3.140000	29403.928702	9.975771e+05
50%	68804.286404	5.970429	7.002902	4.050000	36199.406689	1.232669e+06
75%	75783.338666	6.650808	7.665871	4.490000	42861.290769	1.471210e+06
max	107701.748378	9.519088	10.759588	6.500000	69621.713378	2.469066e+06

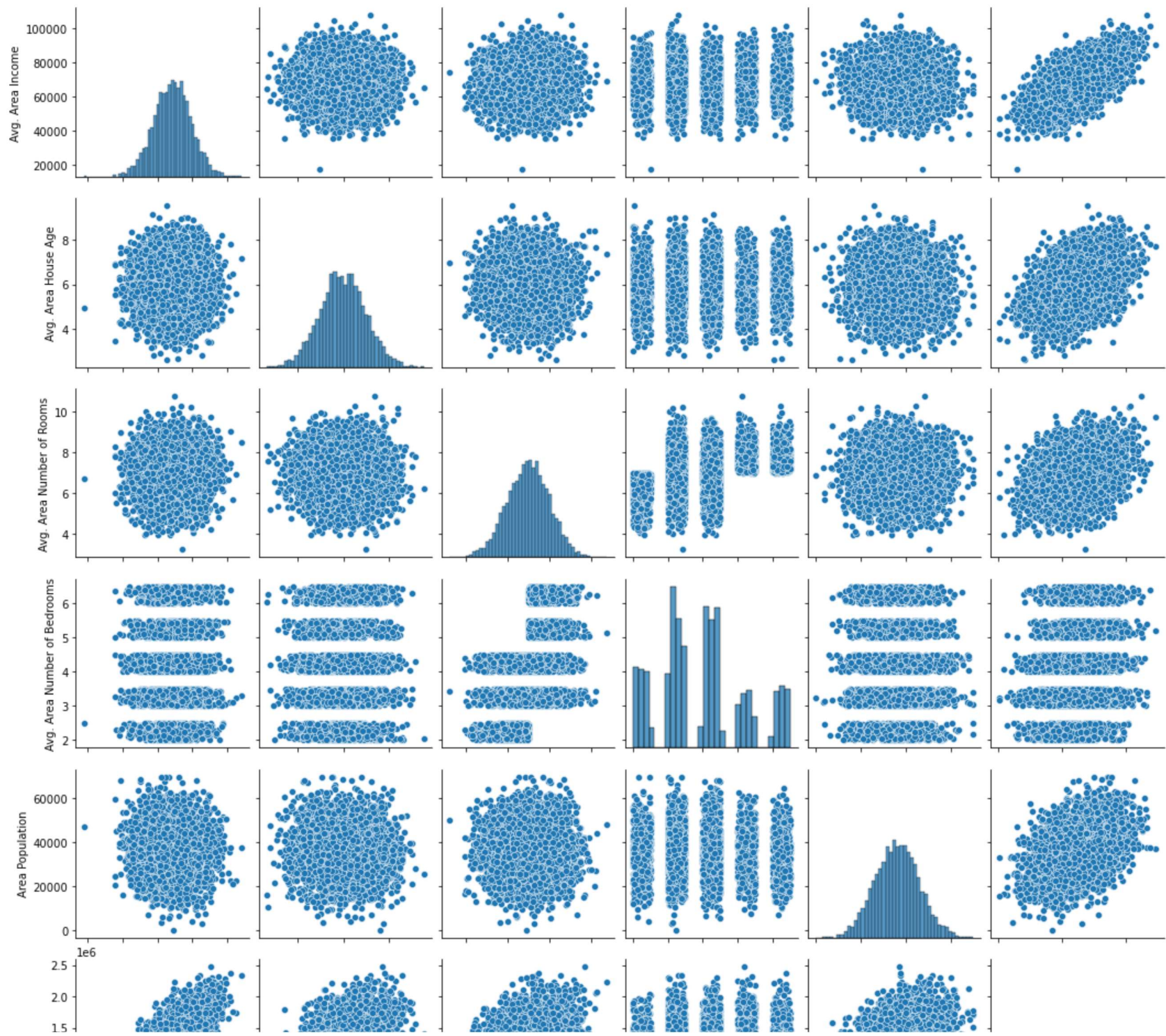
```
In [6]: # To display the col headings  
df.columns
```

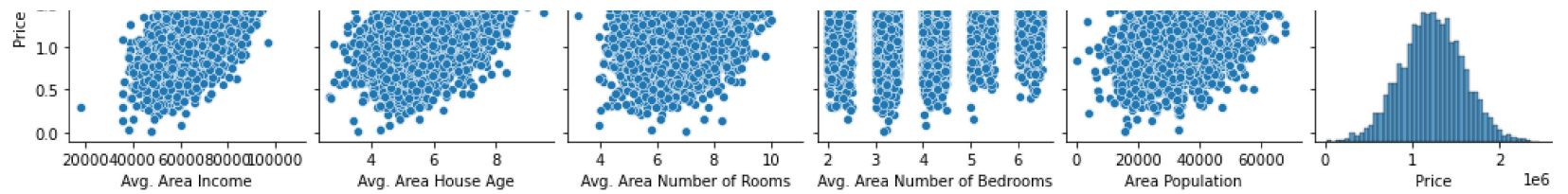
```
Out[6]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',  
              'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],  
              dtype='object')
```

EDA and Visualization

```
In [7]: sns.pairplot(df)
```

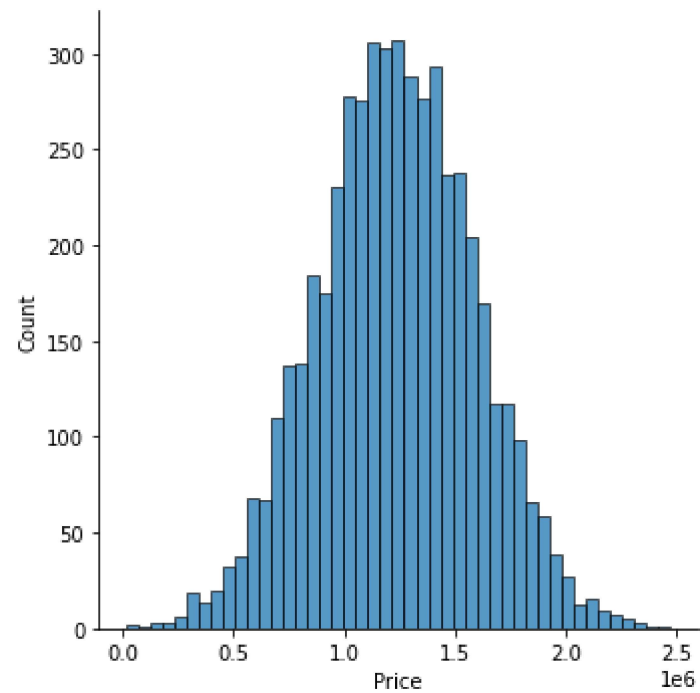
```
Out[7]: <seaborn.axisgrid.PairGrid at 0x23c14bef6d0>
```



```
In [8]: sns.displot(df['Price'])
```

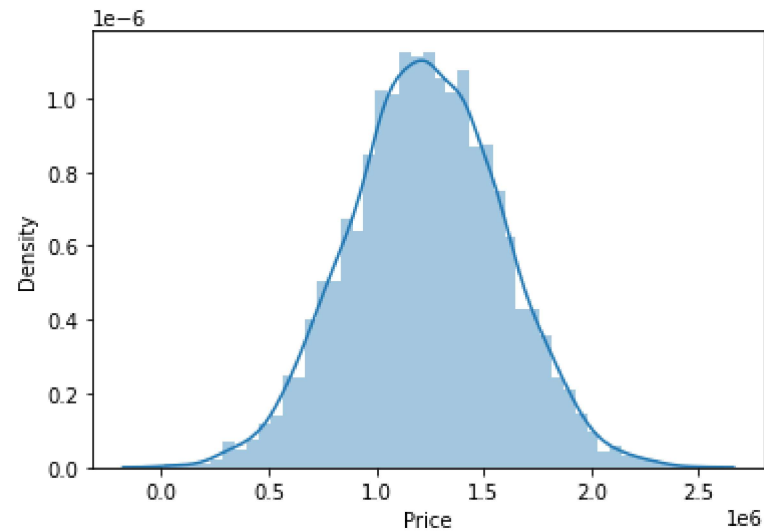
```
Out[8]: <seaborn.axisgrid.FacetGrid at 0x23c16653160>
```



```
In [9]: # We use displot in older version we get distplot use displot  
sns.distplot(df['Price'])
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

```
Out[9]: <AxesSubplot:xlabel='Price', ylabel='Density'>
```



```
In [10]: df1=df[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
               'Avg. Area Number of Bedrooms', 'Area Population', 'Price']]
df
```

Out[10]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.458574	5.682861	7.009188	4.09	23086.800503	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanielstown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386
...
4995	60567.944140	7.830362	6.137356	3.46	22837.361035	1.060194e+06	USNS Williams\nFPO AP 30153- 7653
4996	78491.275435	6.999135	6.576763	4.02	25616.115489	1.482618e+06	PSC 9258, Box 8489\nAPO AA 42991-3352
4997	63390.686886	7.250591	4.805081	2.13	33266.145490	1.030730e+06	4215 Tracy Garden Suite 076\nJoshualand, VA 01...
4998	68001.331235	5.534388	7.130144	5.44	42625.620156	1.198657e+06	USS Wallace\nFPO AE 73316
4999	65510.581804	5.992305	6.792336	4.07	46501.283803	1.298950e+06	37778 George Ridges Apt. 509\nEast Holly, NV 2...

5000 rows × 7 columns

```
In [11]: sns.heatmap(df1.corr())
```

```
Out[11]: <AxesSubplot:>
```



To train the model - MODEL BUILD

Going to train linear regression model; We split our data into 2 variables x and y where x is independent var(input) and y is dependent on x(output), we could ignore address col as it is not required for our model

```
In [12]: x=df1[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',  
               'Avg. Area Number of Bedrooms', 'Area Population']]  
y=df1[['Price']]
```

To split the dataset into test data

```
In [13]: # importing lib for splitting test data  
from sklearn.model_selection import train_test_split
```

```
In [22]: x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3)
```

```
In [23]: from sklearn.linear_model import LinearRegression  
  
lr=LinearRegression()  
lr.fit(x_train,y_train)
```

```
Out[23]: LinearRegression()
```

```
In [24]: print(lr.intercept_)  
  
[-2652763.62949875]
```

```
In [25]: print(lr.score(x_test,y_test))  
  
0.9146315175365524
```

```
In [26]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=["Co-efficient"])  
coeff
```

```

-----
ValueError                                Traceback (most recent call last)
C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\internals\managers.py in create_block_manager_from_blocks(blocks, axes)
    1674         blocks = [
-> 1675             make_block(
    1676                 values=blocks[0], placement=slice(0, len(axes[0])), ndim=2

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\internals\blocks.py in make_block(values, placement, klass, ndim, dtype)
    2741
-> 2742     return klass(values, ndim=ndim, placement=placement)
    2743

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\internals\blocks.py in __init__(self, values, placement, ndim)
    141         if self._validate_ndim and self.ndim and len(self.mgr_locs) != len(self.values):
--> 142             raise ValueError(
    143                 f"Wrong number of items passed {len(self.values)}, "

```

ValueError: Wrong number of items passed 5, placement implies 1

During handling of the above exception, another exception occurred:

```

ValueError                                Traceback (most recent call last)
<ipython-input-26-97bc44621db8> in <module>
----> 1 coeff=pd.DataFrame(lr.coef_,x.columns,columns=["Co-efficient"])
      2 coeff

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\frame.py in __init__(self, data, index, columns, dtype, copy)
    556         mgr = init_dict({data.name: data}, index, columns, dtype=dtype)
    557     else:
-> 558         mgr = init_ndarray(data, index, columns, dtype=dtype, copy=copy)
    559
    560         # For data is list-like, or Iterable (will consume into list)

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\internals\construction.py in init_ndarray(values, index, columns, dtype, copy)
    236         block_values = [values]
    237
-> 238     return create_block_manager_from_blocks(block_values, [columns, index])
    239

```

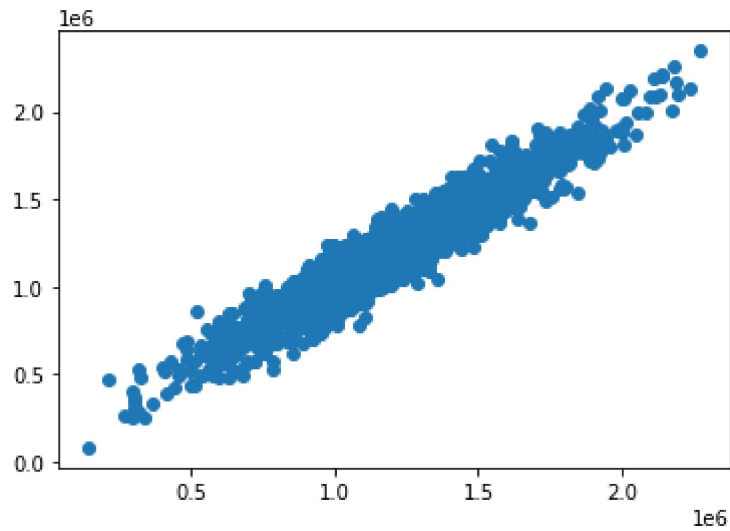
240

```
C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\internals\managers.py in create_block_manager_from_blocks(blocks, axes)
    1685         blocks = [getattr(b, "values", b) for b in blocks]
    1686         tot_items = sum(b.shape[0] for b in blocks)
-> 1687         raise construction_error(tot_items, blocks[0].shape[1:], axes, e)
    1688
    1689
```

ValueError: Shape of passed values is (1, 5), indices imply (5, 1)

```
In [27]: pred = lr.predict(x_test)
        plt.scatter(y_test, pred)
```

Out[27]: <matplotlib.collections.PathCollection at 0x23c188b3bb0>



```
In [28]: len(x.index)
```

Out[28]: 5000

In [29]: `len(y.index)`

Out[29]: 5000

In []: