

Bivariate analysis using Pearson correlation and scatter plot of a combination of 2 columns on Melbourne dataset

For the remaining excersices where numerical columns are needed, I will use the Melbourne Housing Market dataset

https://www.kaggle.com/anthonypino/melbourne-housing-market?select=MELBOURNE_HOUSE_PRICES_LESS.csv

Some Key Details

Suburb: Suburb

Address: Address

Rooms: Number of rooms

Price: Price in Australian dollars

Method:

S - property sold;

SP - property sold prior;

PI - property passed in;

PN - sold prior not disclosed;

SN - sold not disclosed;

NB - no bid;

VB - vendor bid;

W - withdrawn prior to auction;

SA - sold after auction;

SS - sold after auction price not disclosed.

N/A - price or highest bid not available.

Type:

br - bedroom(s);

h - house,cottage,villa, semi,terrace;

u - unit, duplex;

t - townhouse;

dev site - development site;

o res - other residential.

SellerG: Real Estate Agent

Distance: Distance from CBD (Central Business District) in Kilometres

Date: Date sold

Regionname: General Region (West, North West, North, North east ...etc)

Propertycount: Number of properties that exist in the suburb.

Bedroom2 : Scraped # of Bedrooms (from different source)

Bathroom: Number of Bathrooms

Car: Number of carspots

Landsize: Land Size in Metres

BuildingArea: Building Size in Metres

YearBuilt: Year the house was built

CouncilArea: Governing council for the area

Latitude: Self explanatory

Longitude: Self explanatory

```
In [ ]: import pandas as pd
import seaborn as sns
```

```
In [ ]: df = pd.read_csv('melbourne_housing_prices.csv', sep=',')
df.head()
```

```
Out[ ]:
```

	Suburb	Address	Rooms	Type	Price	Method	SellerG	Date	Postcode	Regionna
0	Abbotsford	49 Lithgow St	3	h	1490000.0	S	Jellis	1/04/2017	3067	North Metropol
1	Abbotsford	59A Turner St	3	h	1220000.0	S	Marshall	1/04/2017	3067	North Metropol
2	Abbotsford	119B Yarra St	3	h	1420000.0	S	Nelson	1/04/2017	3067	North Metropol
3	Aberfeldie	68 Vida St	3	h	1515000.0	S	Barry	1/04/2017	3040	West Metropol
4	Airport West	92 Clydesdale Rd	2	h	670000.0	S	Nelson	1/04/2017	3042	West Metropol

```
In [ ]: houseCorrelations = df.corr(method='pearson')
houseCorrelations.style.background_gradient(cmap='coolwarm', axis=None).set_p
```

C:\Users\Stijn\AppData\Local\Temp\ipykernel_14048\2374652365.py:2: FutureWarning: this method is deprecated in favour of `Styler.format(precision=..)`

```
houseCorrelations.style.background_gradient(cmap='coolwarm', axis=None).set_
```

```
Out[ ]:
```

	Rooms	Price	Postcode	Propertycount	Distance
Rooms	1.00	0.41	0.09	-0.05	0.27
Price	0.41	1.00	0.00	-0.06	-0.25
Postcode	0.09	0.00	1.00	-0.00	0.50
Propertycount	-0.05	-0.06	-0.00	1.00	0.01
Distance	0.27	-0.25	0.50	0.01	1.00

I expected distance to the business district and price to correlate negatively, I did expect a stronger correlation

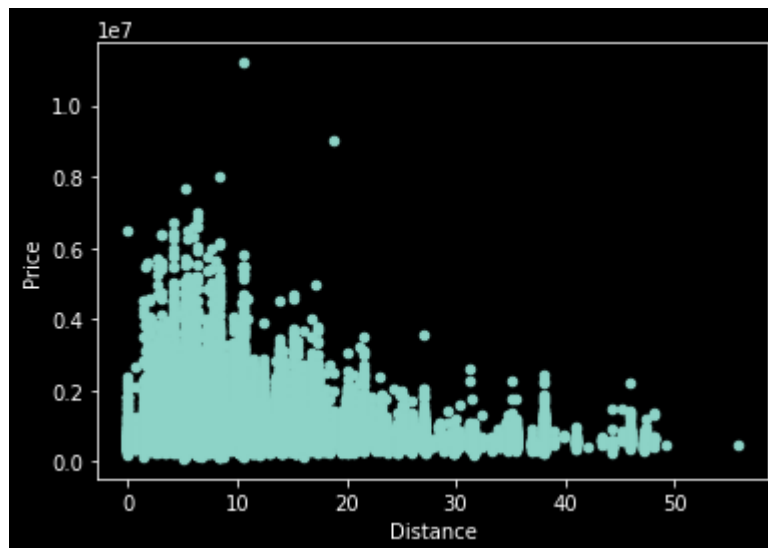
```
In [ ]: correlation = df[['Price', 'Distance']]
correlation.corr()
```

```
Out[ ]:
```

	Price	Distance
Price	1.000000	-0.253668
Distance	-0.253668	1.000000

```
In [ ]: df.plot(kind='scatter', x='Distance', y='Price')
```

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



The scatter plot shows the negative correlation very well.