

Assignment 1: Exploratory data analysis



```
1 begin
2     using PlutoUI
3     using PlutoTest
4     using Random, Statistics, LinearAlgebra, Dates
5     using Distributions, StatsBase
6     using DataFrames , XLSX
7     using CairoMakie
8     using Chain
9     using TidierData
10    using FileIO
11    using Polynomials
12    using Match
13    using Latexify
14    const colors = Makie.wong_colors()
15 end;
```

Table of Contents

Assignment 1: Exploratory data analysis

- Load the Melbourne Housing data set

 - Load into a dataframe

 - Generate a summary table of house types

- Create dfPlot , a subset of data for plotting

- Geographical distribution and density of housing

- House price patterns by main city districts

- Relationship between price and proximity to the CBD

Load the Melbourne Housing data set

This dataset was downloaded from Kaggle

(<https://www.kaggle.com/datasets/anthonypino/melbourne-housing-market>). The data was scraped from publicly available results posted every week from Domain.com.au. The dataset includes Address, Type of Real estate, Suburb, Method of Selling, Rooms, Price, Real Estate Agent, Date of Sale and distance from C.B.D.

Load into a dataframe

df =

	Suburb	Address	Rooms	Type	Method	SellerG	Date
1	"Abbotsford"	"68 Studley St"	2	"h"	"SS"	"Jellis"	2016-
2	"Airport West"	"154 Halsey Rd"	3	"t"	"PI"	"Nelson"	2016-
3	"Albert Park"	"105 Kerferd Rd"	2	"h"	"S"	"hockingstuart"	2016-
4	"Albert Park"	"85 Richardson St"	2	"h"	"S"	"Thomson"	2016-
5	"Alphington"	"30 Austin St"	3	"h"	"SN"	"McGrath"	2016-
6	"Alphington"	"6 Smith St"	4	"h"	"S"	"Brace"	2016-
7	"Alphington"	"5/6 Yarralea St"	3	"h"	"S"	"Jellis"	2016-
8	"Altona"	"158 Queen St"	3	"h"	"VB"	"Greg"	2016-
9	"Altona North"	"1 Beuron Rd"	3	"h"	"SP"	"Williams"	2016-
10	"Altona North"	"45 Hearn St"	5	"h"	"S"	"FN"	2016-
more							
34857	"Westmeadows"	"42 Pascoe St"	4	"h"	"S"	"Barry"	2017-

```
1 df = DataFrame(XLSX.readtable("Melbourne_housing.xlsx", 1))
```

(1) Calculate the number of missing observations for each df column

Store the result in `dfMissing`. Use column names `Column` and `MissingValues`. It's important to use these names to facilitate grading. See the test below.

Tip

Use `count(ismissing, df[!, col])` to return the number of missing values in the column `col`.

Your code here

```

1 begin
2   #New data frame for columns and their missing values
3   dfMissing = DataFrame(Column = String[], MissingValues = Int[])
4   #Loop through each column of original dataframe
5   for col in names(df)
6     missing_count = count(ismissing, df[:, col])
7     #Add a new row when missing content is found
8     push!(dfMissing, (col, missing_count))
9   end
10 end

```

```
sum(dfMissing.MissingValues) == 100954
```

```
1 @test sum(dfMissing.MissingValues) == 100954
```

Generate a summary table of house types

(2) Re-Calculate the table of house types

With fewer missing observations excluded, the overall number of records should now be larger.

Tip

Use the `@chain` macro to extract and summarise the data from `df`.

Visit the `TidierData.jl` GitHub repo to see the macros for the various data operations you can perform (e.g., `@group_by`, `@mutate`, `@summarize`).

Your code here

houseType =

	Type	n
1	"h"	23980
2	"t"	3580
3	"u"	7297

```

1 houseType = @chain df begin
2   #Select Type column
3   @select(Type)
4   @group_by(Type)
5   #counting number of occurrences
6   @summarize(n = n())
7 end

```

```
sum(houseType.n) == 34857
```

```
1 @test sum(houseType.n) == 34857
```

(3) Create a dataframe for subsequent plots

Include only the variables you will need from the original df. You may also have to create one or two new ones to help with the plots. Note that the `Colour` column is used to provide a numerical code for Region for graphing purposes.

Include: Regionname, Price, Rooms, Latitude, Longitude, Distance, Type.

Create: Region, PriceK, TypeFull, Colour

Select records for just the five metropolitan regions.

Tip

Be sure to name the new dataframe `dfPlot` to ensure the test below is passed.

Create `dfPlot`, a subset of data for plotting

Your code here

	Region	PriceK	TypeFull	Colour	Rooms	Latitude	Longitude	Distance
1	"WM"	840.0	"Townhouse"	"2"	3	-37.718	144.878	13.5
2	"SM"	1275.0	"House"	"3"	2	-37.8459	144.956	3.3
3	"SM"	1455.0	"House"	"3"	2	-37.845	144.954	3.3
4	"NM"	2000.0	"House"	"1"	4	-37.7707	145.032	6.4
5	"NM"	1110.0	"House"	"1"	3	-37.7854	145.032	6.4
6	"WM"	520.0	"House"	"2"	3	-37.87	144.825	13.8
7	"WM"	1085.0	"House"	"2"	5	-37.8388	144.857	11.1
8	"WM"	781.0	"House"	"2"	3	-37.8291	144.835	11.1
9	"SM"	599.0	"Unit"	"3"	2	-37.8543	145.026	6.3
10	"WM"	455.0	"Unit"	"2"	2	-37.7756	144.917	5.9
more								
20622	"NM"	791.0	"House"	"1"	4	-37.6763	144.894	16.5

```

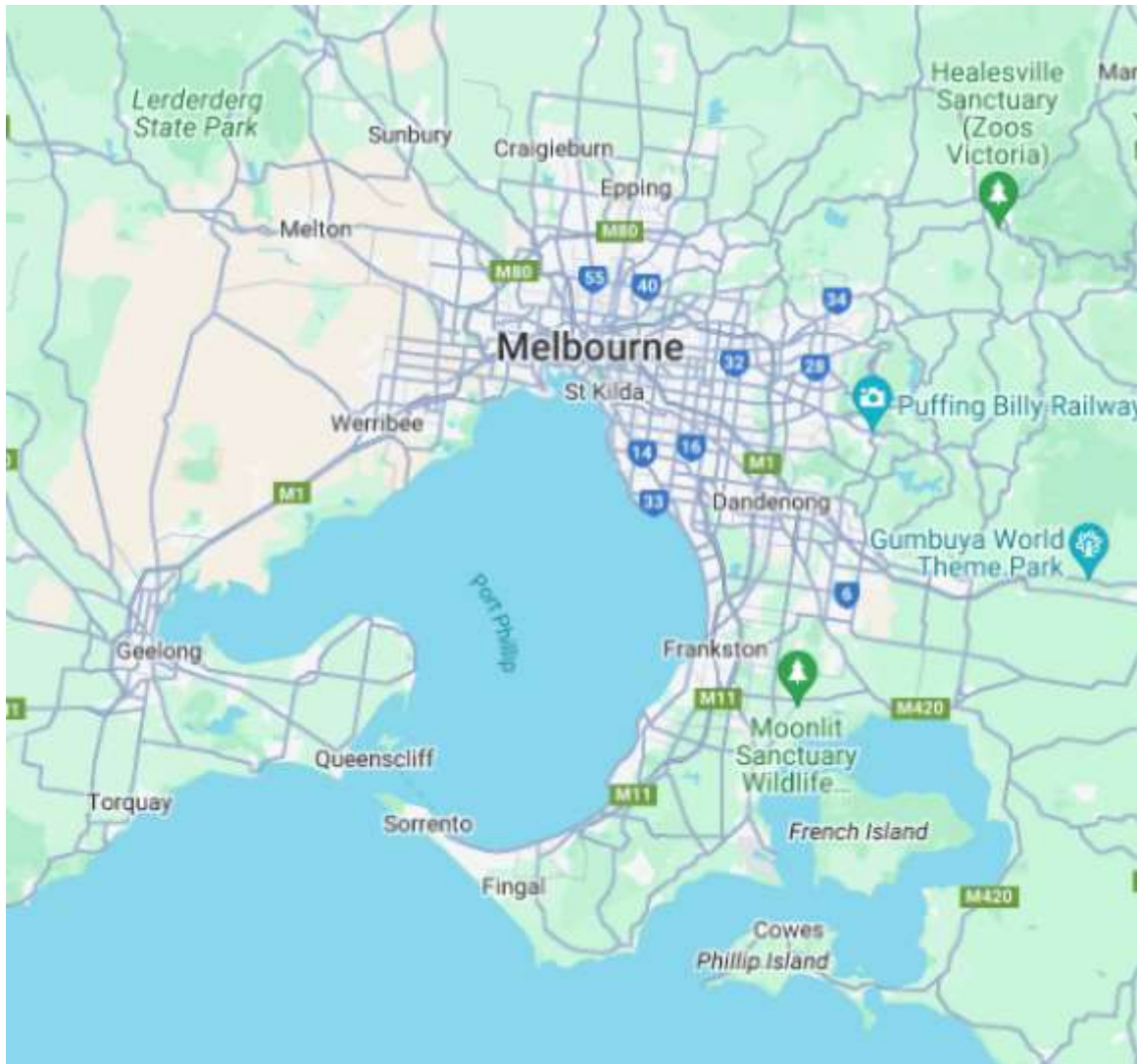
1 begin
2     dfPlot = @chain df begin
3         @select(Regionname, Price, Rooms, Latitude, Longitude, Distance, Type)
4         @filter(Regionname in ["Northern Metropolitan", "Western Metropolitan",
5             "Southern Metropolitan", "Eastern Metropolitan", "South-Eastern
6             Metropolitan"])
7         @drop_missing()
8         #Create new columns to replace old
9         @mutate(
10             Region = replace(Regionname,
11                 "Northern Metropolitan" => "NM",
12                 "Western Metropolitan" => "WM",
13                 "Southern Metropolitan" => "SM",
14                 "Eastern Metropolitan" => "EM",
15                 "South-Eastern Metropolitan" => "SEM"),
16             PriceK = Price/1000,
17             TypeFull= replace(Type,
18                 "h" => "House",
19                 "t" => "Townhouse",
20                 "u" => "Unit"),
21             Colour = replace(Regionname,
22                 "Northern Metropolitan" => 1,
23                 "Western Metropolitan" => 2,
24                 "Southern Metropolitan" => 3,
25                 "Eastern Metropolitan" => 4,
26                 "South-Eastern Metropolitan" => 5))
27         #Select new columns for output
28         @select(Region, PriceK, TypeFull, Colour, Rooms, Latitude, Longitude, Distance)
29     end
30 end

```

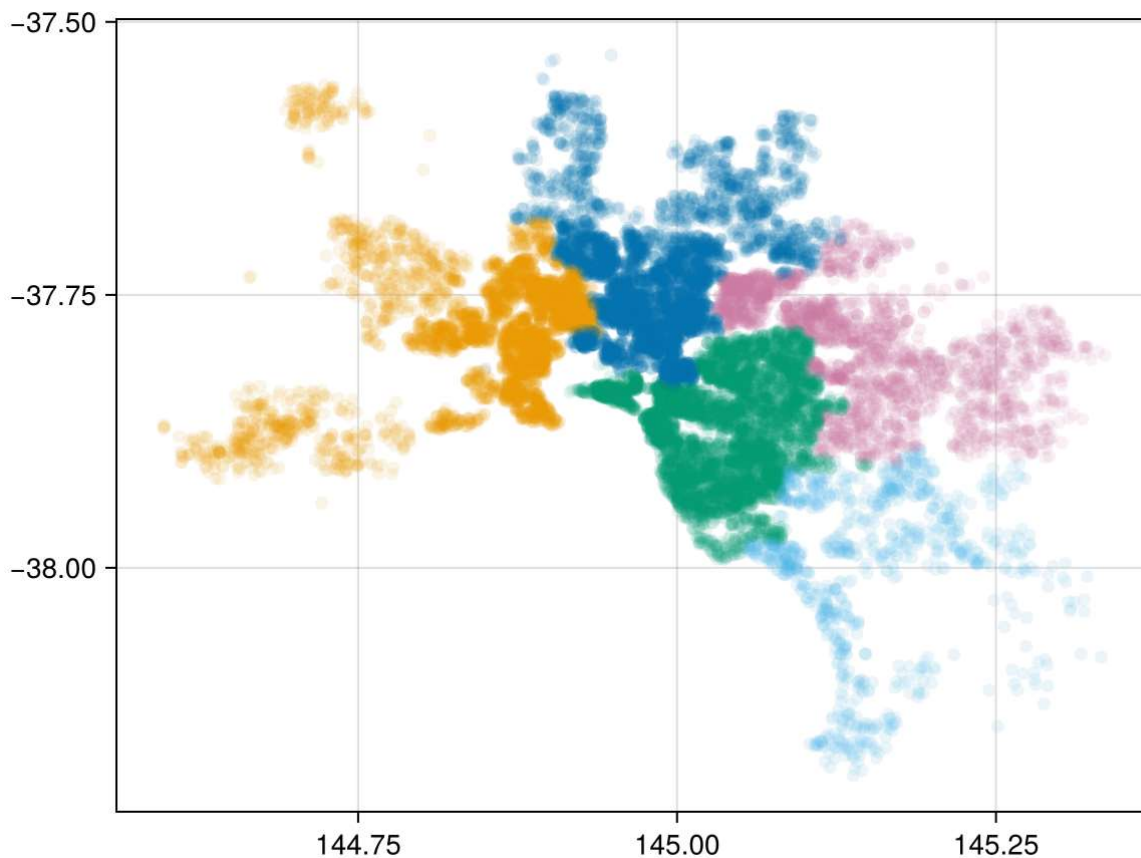
```
nrow(dfPlot) == 20622
```

```
1 @test nrow(dfPlot) == 20622
```

Geographical distribution and density of housing



This error will go once you define `dfPlot`.



```
1 let
2   xs = dfPlot.Longtitude
3   ys = dfPlot.Latitude
4   points = Point2f.(xs, ys)
5   scatter(points, alpha=0.1,
6           color = colors[parse.(Int, dfPlot.Colour)])
7 end
```

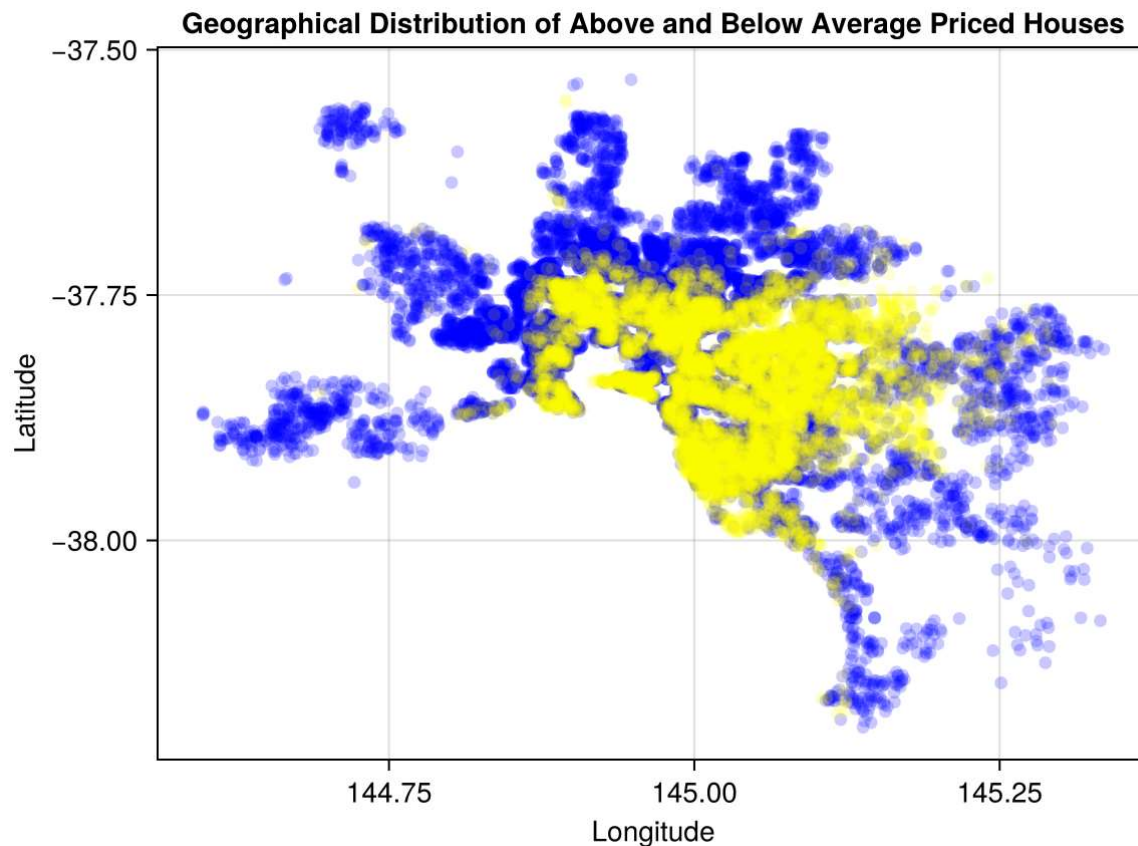
(4) Make an additional geographic plot of above- and below-average priced houses

Discuss the patterns you find.

Tip

You can use an expression like `avgPrice = mean(dfPlot.PriceK)` to calculate the average house price in the dataframe.

Your code here



```

1 begin
2   #Calculate average
3   avgPrice = mean(dfPlot.PriceK)
4   #Create new subsets for those above average and below
5   aboveAvg = dfPlot[dfPlot.PriceK .> avgPrice, :]
6   belowAvg = dfPlot[dfPlot.PriceK .<= avgPrice, :]
7
8   #Create points for subsets by their longitude and latitude
9   points_above = Point2f.(aboveAvg.Longitude, aboveAvg.Latitude)
10  points_below = Point2f.(belowAvg.Longitude, belowAvg.Latitude)
11
12 let
13   fig = Figure()
14   ax = Axis(fig[1, 1], title="Geographical Distribution of Above and Below Average
    Priced Houses")
15
16   #Plot points for those below
17   scatter!(ax, points_below, label="Below Average Price", color=:blue, alpha=0.2)
18   #Plot points for those above
19   scatter!(ax, points_above, label="Above Average Price", color=:yellow,
    alpha=0.155)
20
21
22   ax.xlabel = "Longitude"
23   ax.ylabel = "Latitude"
24
25   #Return the figure
26   fig
27 end

```

28 end

From the graph I noticed above-averaged priced homes clustered in areas closer to the city center and along the waterfront. As we move away from these two areas, the distribution shifts with below-averaged priced homes clustering farther out. Even though this clustering of below averaged homes starts away from the city center, there are still a few below average priced scattered in the more city areas

(5) Suggest why the South Metro distribution tends to be bimodal

Support your conjecture with data and/or graphs.

From the Melbourne House Prices by Region graph below, we can see that the Southern Metro area exhibits a bimodal distribution, with two distinct peaks: one at the lower end of the price spectrum (~500) and another at the higher end (~1500).

The Geographical Distribution and Density Housing graph shows that the Southern Metro area is situated between the city center and the coastline. Additionally, the Above and Below Average Cost graph shows that most of the housing in this area is on the above average side, yet there still remains a significant number of below-average priced options.

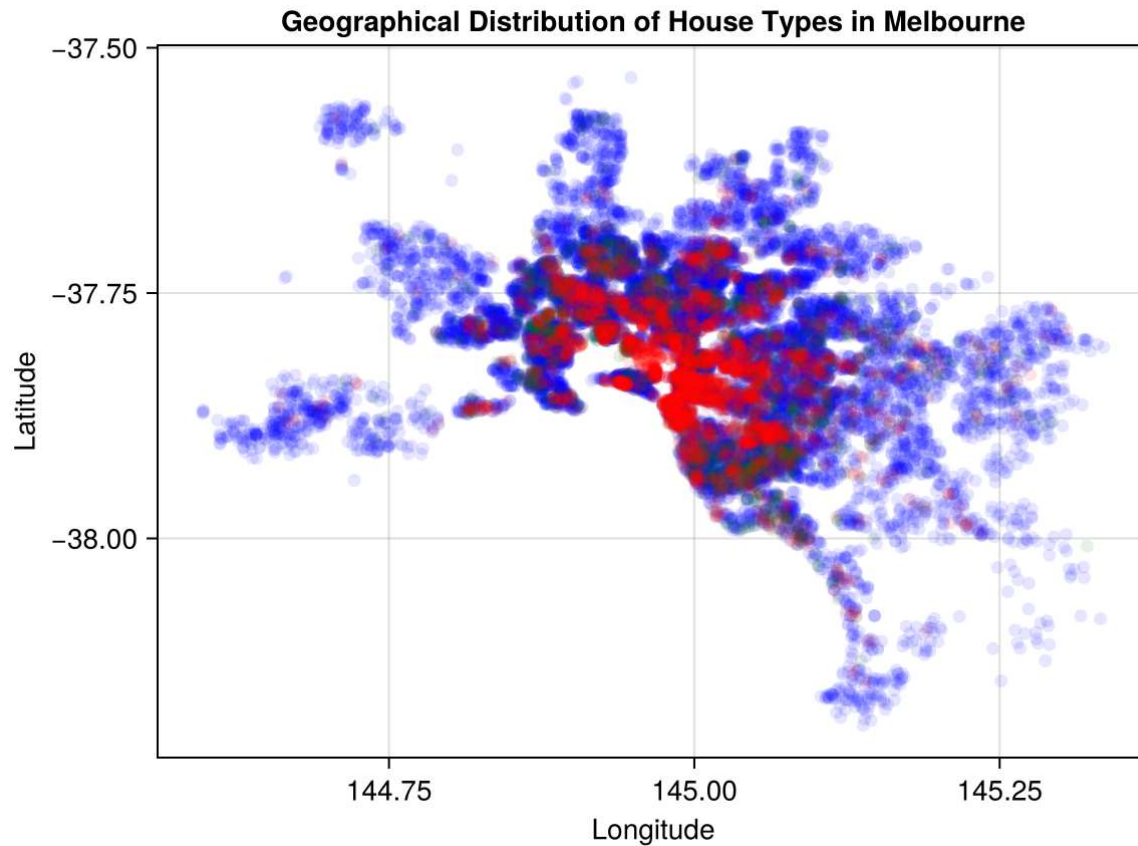
My hypothesis for this phenomenon is influenced by the fact that the models mentioned above don't take into consideration the different house types. Like I said before this area is located right by the city which typically tends to have more apartments (units) and townhouses, that tend to cost less due to their limited space.

To see if this hypothesis was right I organized the Southern Metro area by Type below:

	TypeFull	Count	AveragePriceK
1	"Townhouse"	55	889.368
2	"House"	896	933.828
3	"Unit"	78	551.971

From this updated dataset we can see that units cost significantly less than houses and townhouses.

To explore this further, I created the graph below, which shows the geographical distribution of different housing types across all areas of Melbourne (red = units, blue = houses, green = townhouses).

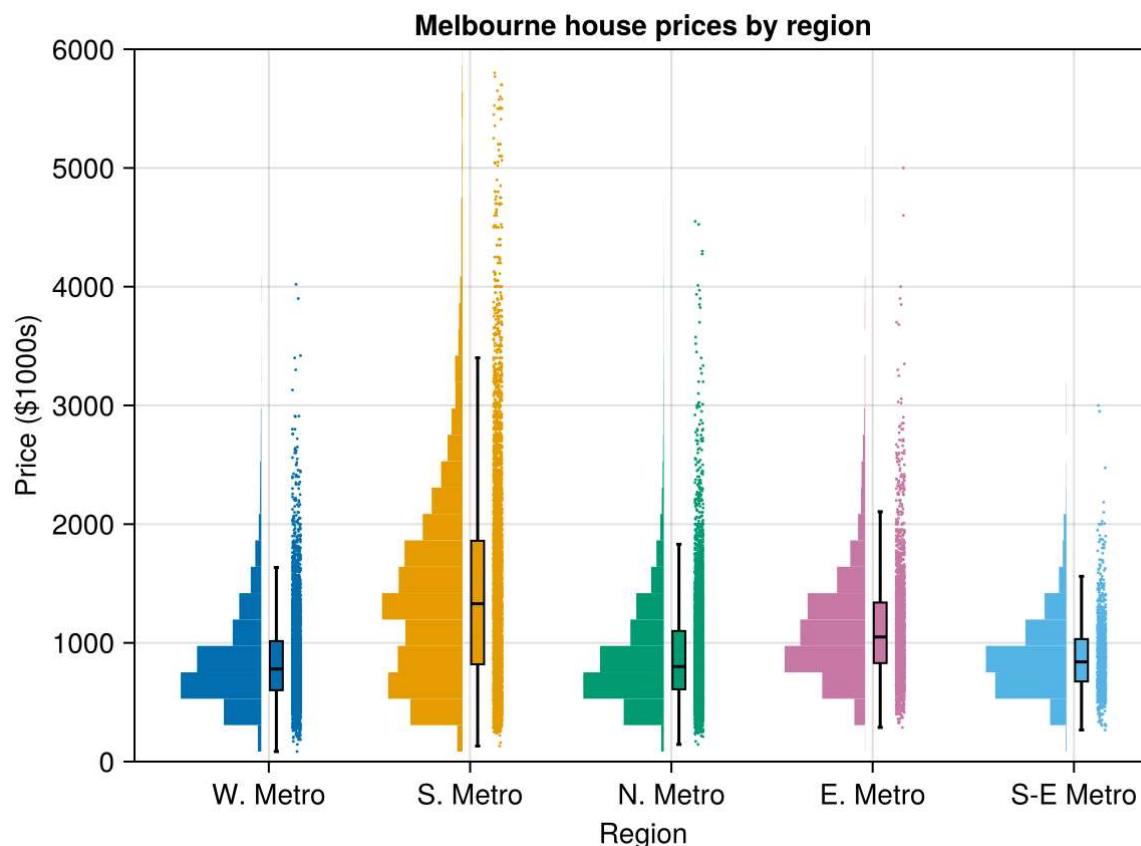


Most of these units are located in the part of the Southern Metro area closest to the city center.

In conclusion the SM's closeness to the city center causes a dip towards the cheaper end of housing because of its trend of apartments. Conversely, the higher-end prices are likely due to the presence of more expensive housing near the shoreline.

House price patterns by main city districts

This error will go once you define `dfPlot`.



```

1 let
2   category_labels = dfPlot.Region
3   data_array = dfPlot.PriceK
4
5   raincloud(category_labels, data_array;
6     axis = (; xlabel = "Region",
7       ylabel = "Price (\$1000s)",
8       limits = (nothing, (0,6000)),
9       title = "Melbourne house prices by region",
10      xticks= (1:5, ["W. Metro", "S. Metro", "N. Metro", "E. Metro", "S-E
      Metro"])),
11     plot_boxplots = true, cloud_width=0.5, clouds=hist, hist_bins=50,
12     color = colors[indexin(category_labels, unique(category_labels))])
13
14 end

```

Relationship between price and proximity to the CBD

(6) Create a dataframe of CBD distance and average price

You should use `trunc()` to return the nearest integral value of the Distance measure in `dfPlot`, then group and summarise price data within each of the resulting kilometer blocks.

Be sure to name the new dataframe `distPrice` and the variables `Dist`, `n`, and `AvgPrice` for test purposes.

Your code here

	Dist	n	AvgPrice
19	16	589	1038.82
20	20	456	694.659
21	34	57	604.648
22	24	122	717.068
23	22	90	836.426
24	18	445	731.071
25	19	196	739.292
26	38	111	773.184
27	32	9	713.722
28	21	209	1066.74
29	35	31	747.597
30	31	134	576.498
31	27	107	745.518
32	28	39	657.256
33	0	22	586.682
34	26	16	745.562
35	33	6	586.333

```

1 begin
2     distPrice = @chain dfPlot begin
3         #return nearest integral value
4         @mutate(Dist = trunc.(Distance))
5         @mutate(Dist = Int.(Dist))
6         @group_by(Dist)
7         @summarize(
8             n = n(),
9             AvgPrice = mean(PriceK)
10        )
11     end
12 end

```

(7) Graph and fit a polynomial curve to these data

Use the `Polynomials.jl` library. Provide correlation and goodness-of-fit measures (R^2 and RMSE) on the graph.

Bind a slider to the variable controlling the polynomial order to show different fits.

Discuss what polynomial order best describes the relationship. Note that this does not necessarily mean gives the best R^2 .

Tip

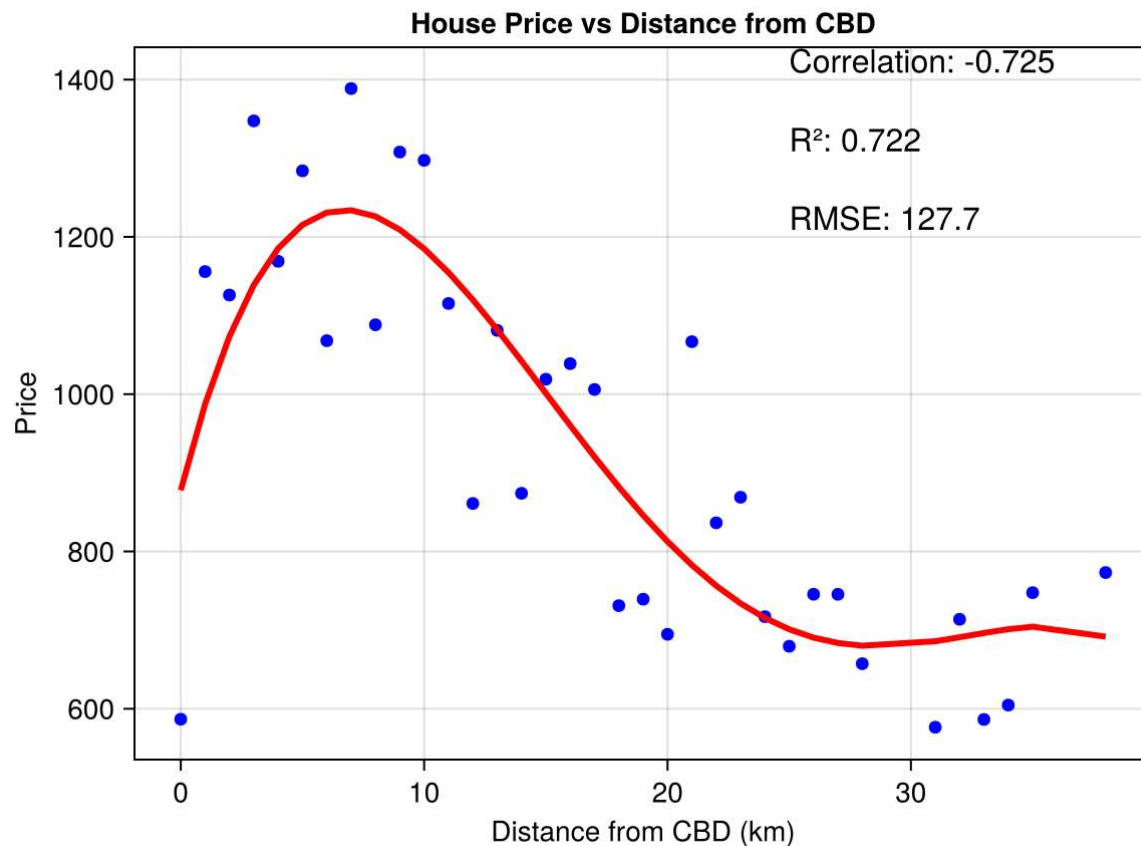
Use `corr()` from the `Statistics` library to calculate the correlation between distance and average price.

The R^2 goodness of fit measure is calculated using: $R^2 = 1 - \frac{\sum (y_{\text{actual}} - y_{\text{pred}})^2}{\sum (y_{\text{actual}} - \bar{y})^2}$

The RMSE measure is calculated using: $\text{RMSE} = \sqrt{\frac{\sum (y_{\text{actual}} - y_{\text{pred}})^2}{N}}$

Your code here

```
1 poly_order = 2;
```



```

1
2 let
3   xs = distPrice.Dist
4   ys = distPrice.AvgPrice
5
6   function fit_polynomial(xs, ys, degree)
7     #Creating the fitted line
8     p = Polynomials.fit(xs, ys, degree)
9     #Predict values using fitted polynomial
10    ys_pred = p.(xs)
11
12    # Calculate R², RMSE, and correlation
13    R2 = 1 - sum((ys .- ys_pred).^2) / sum((ys .- mean(ys)).^2)
14    rmse = sqrt(mean((ys .- ys_pred).^2))
15    correlation = cor(xs, ys)
16
17    fig = Figure()
18    ax = Axis(fig[1, 1],
19      title="House Price vs Distance from CBD",
20      xlabel="Distance from CBD (km)",
21      ylabel="Price"
22    )
23
24    #Original data points (distance & price)
25    scatter!(ax, xs, ys, color=:blue)
26
27    # Sort xs and ys_pred for proper line plotting
28    sorted_indices = sortperm(xs)
29    sorted_xs = xs[sorted_indices]

```



```

30 sorted_ys_pred = ys_pred[sorted_indices]
31
32 #Plot fitted line
33 lines!(ax, sorted_xs, sorted_ys_pred, color=:red, linewidth=3)
34
35 # Set custom y-ticks
36 y_ticks = 600:200:1400
37 ax.yticks = y_ticks
38
39 text!(ax, 25, 1400, text="Correlation: $(round(correlation, digits=3))",
fontSize=16)
40 text!(ax, 25, 1300, text="R²: $(round(R2, digits=3))", fontsize=16)
41 text!(ax, 25, 1200, text="RMSE: $(round(rmse, digits=1))", fontsize=16)
42
43 return fig
44 end
45
46 fit_polynomial(xs, ys, poly_order) #function call
47 end

```



```
1 @bind poly_order PlutoUI.Slider(1:1:5, default=3)
```

```
round(cor(distPrice.Dist, distPrice.AvgPrice), digits = 3) == -0.725
```

```
1 @test round(cor(distPrice.Dist, distPrice.AvgPrice), digits=3) == -0.725
```

The degree 3 polynomial best captures the underlying relationship between house prices and distance from the CBD. While the degree 5 polynomial provides the best statistical fit—reflected by a higher R^2 and lower RMSE—it runs the risk of overfitting. Overfitting occurs when a model is too finely tuned to the specific dataset it was trained on, leading to poor performance on new data.

In contrast, the degree 3 polynomial provides a better balance. It accurately reflects the general trend in the data without being overly sensitive to small variations, making it more likely to generalize well to additional data.