

Cherrybrook Real Estate Data Analysis Project

By: Nipunn Khurana

1. Project Overview

This project focuses on analysing the housing market in Cherrybrook, NSW, using real-world property data. The goal was to uncover insights into pricing trends, land size distribution, school catchment impacts, and sales agent performance. The analysis followed the data analytics pipeline: data extraction, cleaning, transformation, and visualization.

2. Data Extraction using Selenium

To begin with, I extracted raw property listing data from Domain.com.au using Selenium, a Python automation tool that simulates human interaction with web pages. I built a custom script that:

1. Navigated through multiple listing pages for Cherrybrook.
2. Captured relevant data fields such as address, price, sold date, link, property type, number of bedrooms, bathrooms, car spaces, and land area.
3. Stored the data in a structured format (CSV) for further analysis.

```
14 import time
15
16 class DomainListings:
17     def __init__(self, driver_path):
18         self.driver_path = driver_path
19         self.driver = None
20
21     def start_browser(self):
22         chrome_options = Options()
23         chrome_options.add_argument("--start-maximized")
24
25         # ✅ Block images to speed up page load
26         prefs = {"profile.managed_default_content_settings.images": 2}
27         chrome_options.add_experimental_option("prefs", prefs)
28
29         service = Service(self.driver_path)
30         self.driver = webdriver.Chrome(service=service, options=chrome_options)
31         print("Browser started successfully (with images blocked).")
32
33     def navigate_to_url(self, url):
34         if self.driver is None:
35             print("Error: Browser is not started. Call 'start_browser' first.")
36             return
37         self.driver.get(url)
38         print(f"Navigated to {url}")
39
40     def click_sold_button_and_apply_filter(self):
41         try:
42             wait = WebDriverWait(self.driver, 5)
43
44             # Step 1: Click the 'Sold' button
45             sold_button = wait.until(EC.element_to_be_clickable((By.CSS_SELECTOR, 'button[data-testid="sold-navigation"]')))
46             sold_button.click()
47             print("✅ Clicked the 'Sold' button.")
48
49             # Step 2: Click the 'Filters' button
50             filters_button = wait.until(EC.element_to_be_clickable((By.CSS_SELECTOR, 'button[data-testid="search-filters-button-desktop"]')))
51             filters_button.click()
52             print("✅ Clicked the 'Filters' button.")
53         except Exception as e:
54             print(f"Error: {e}")
```

This method enabled me to gather up-to-date, location-specific real estate information which would otherwise not be readily available in public datasets.

3. Data Cleaning and Manipulation in SQL

After collecting the raw dataset, I imported the CSV file into a SQL database to perform data cleaning and preprocessing. Key steps included:

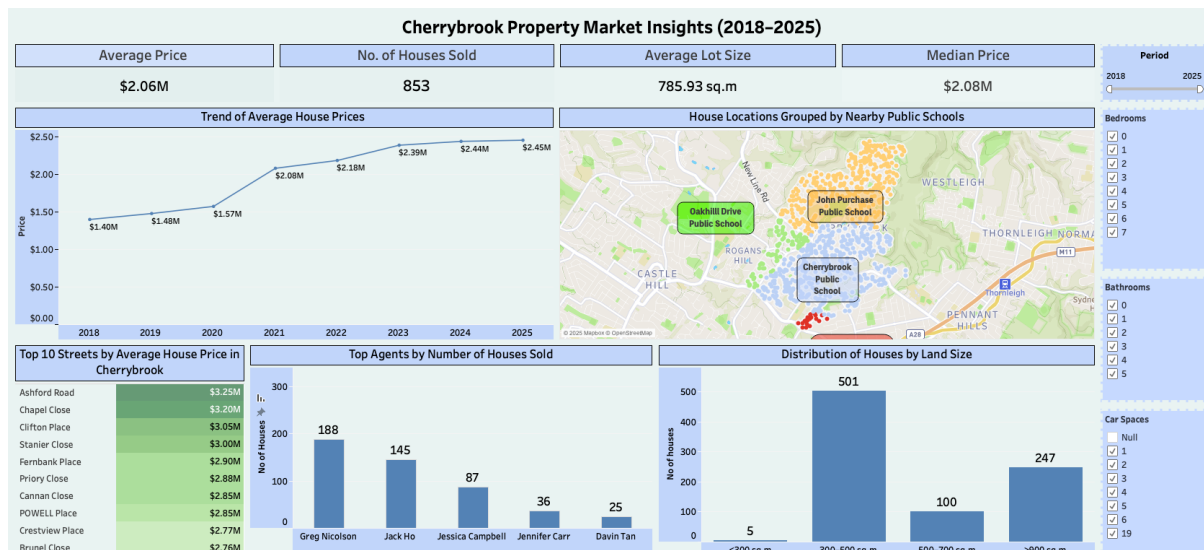
1. Parsing and standardizing address components.
2. Extracting numeric values from text fields such as price and area.
3. Handling missing values, such as estimating area where missing or filtering out incomplete records.
4. Formatting dates and calculating useful derived fields like price per square metre.
5. Using SQL for data wrangling ensured accuracy, efficiency, and reusability for multiple queries and filters.

```
--  
28 -- Task 1.  
29 ALTER TABLE cherrybrook  
30 DROP COLUMN Page, Listing_Number  
31  
32 -- Task 2.  
33  
34 UPDATE cherrybrook  
35 SET  
36     Car_Spaces = NULLIF(Car_Spaces, 'N/A'),  
37     Area = NULLIF(Area, 'N/A')  
38 WHERE Car_Spaces = 'N/A' OR Area = 'N/A';  
39  
40 -- Car_Spaces has 3 N/A values, Area has 207 N/A values, Agent_Name has 131 missing values,  
41 -- For Car_spaces, we will replace 'N/A' values with NULL  
42 -- For Area, replace 'N/A' by NULL  
43 -- For Agent_Name, dont replace 'N/A'  
44  
45 -- Task 3.  
46  
47 WITH CTE AS (  
48     SELECT *, ROW_NUMBER() OVER (PARTITION BY Link, Address, Price ORDER BY (SELECT NULL)) AS rn  
49     FROM cherrybrook  
50 )  
51 DELETE FROM CTE  
52 WHERE rn > 1;  
53
```

4. Visualization in Tableau

Once the dataset was cleaned and transformed, I imported it into Tableau Public to design an interactive dashboard. The dashboard was structured to answer key stakeholder questions and includes the following components:

1. Trend line of average property prices from 2018 to 2025
2. Geographical map of houses categorized by public school catchment
3. Bar charts showing agent-wise and area-wise property sales
4. Top 10 streets in Cherrybrook based on average price
5. KPIs like average price, median price, number of houses sold, and average land size



You can view the live dashboard here:

[Tableau Public – Cherrybrook Insights Dashboard](#)

5. Key Insights and Findings

Based on the dashboard visualizations, several important trends were identified:

- 🏠 The average house price in Cherrybrook increased significantly from \$1.4M in 2018 to \$2.45M in 2025, showing a clear upward trend in property value.
- 🗺️ School catchment areas had a noticeable effect on property distribution and clustering, especially around Cherrybrook Public, John Purchase, and Oakhill Drive Public School.
- 🇮🇹 Most houses were sold in the 300–500 sq.m range, followed by >900 sq.m, suggesting a strong demand for medium-sized lots.
- 🧑 Greg Nicolson and Jack Ho were the leading agents in Cherrybrook by volume of houses sold.

📌 Top-priced streets include Ashford Road, Chapel Close, and Clifton Place, with average prices exceeding \$3 million.

6. Conclusion

This end-to-end project demonstrates how real estate data can be scraped, structured, cleaned, and visualized to extract meaningful business insights. The combination of Selenium, SQL, and Tableau created a robust pipeline for analysis, and the findings can support homebuyers, real estate agents, or investors in making informed decisions.