



**CE/CZ4123/SC4023 BIG DATA MANAGEMENT**

**SEMESTER GROUP PROJECT**

**COLLEGE OF COMPUTING AND DATA SCIENCE  
NANYANG TECHNOLOGICAL UNIVERSITY**

## 1 ASSIGNMENT DESCRIPTION

The goal of this semester's project is to conduct a simple analysis on the resale flats to have a flavor of the data management process. You will be provided with a series of transaction records concerning the resale of HDB flats over the last 10 years (2014 to 2023) in Singapore. These transactions present comprehensive information including approval date, location, flat model, price, etc, enabling us to execute various queries, such as the average resale flat price on a particular street or the price trends for specific types of HDB flats. To facilitate a more manageable workload and evaluation process, we propose to select several specific queries for you to implement.

From the perspective of a potential flat buyer, one may be interested in accessing the statistics of the resale HDB flats whose area meets some requirements in a certain location over several months. Therefore, your program is required to compute the minimum price, average price, standard deviation of price, and minimum price per square meter of resale HDB flats whose area is not smaller than 80 square meters ( $\geq 80m^2$ ) within a specific town over two consecutive months in a given year.

You are expected to write a program to manage the data in a **column-oriented** manner, including data storage and processing. Your program should first receive queries, scan the data columns to find matched lines, and compute the results according to associated query content. To be specific, a query is composed of a target time (YYYY-MM to YYYY-(MM+1)), a matched town, and a query content. These factors are determined by your matriculation number as follows: a) The last digit of the year of the target time (YYYY) equals the last digit of the matriculation number; b) the commencing month (MM) equals the second last digit of the matriculation number (note that "0" represents October); c) the matched town depends on the third last digit of the matriculation number as Table 1 presents; d) there are four query contents in total that are listed in Table 2, and the area requirement ( $\geq 80m^2$ ) is applicable to all these contents.

Digit	0	1	2	3	4
Town	BEDOK	BUKIT PANJANG	CLEMENTI	CHOA CHU KANG	HOUGANG
Digit	5	6	7	8	9
Town	JURONG WEST	PASIR RIS	TAMPINES	WOODLANDS	YISHUN

**Table 1:** List of towns corresponding to the third last digit in matriculation number.

Please use the matriculation numbers of **anyone** of your group members to generate four queries encompassing **all** four query contents in arbitrary sequence and record the results, respectively. Please round the results to the **hundredth place** in your report.

Minimum Price	Standard Deviation of Price
Average Price	Minimum Price per Square Meter

**Table 2:** List of query contents for testing your program.

**Example:** For querying the average price, a student with matriculation number **A1234567B** should scan the resale HDB flats in **JURONG WEST** from **Jun. 2017** to **Jul. 2017** to compute the associated average price of the matched flats.

In this case, an example query for the average price of the resale HDB flats not smaller than 80 square meters in JURONG WEST during Jun. 2017 to Jul. 2017 should be equivalent to the following SQL query:

#### Task in SQL

```
1 WITH Tab1 AS (  
2     SELECT *  
3     FROM ResalePricesSingapore  
4     WHERE (YEAR(Month) = 2017)  
5           AND (MONTH(Month) >= 6)  
6           AND (MONTH(Month) <= 7)  
7           AND (Town = 'JURONG WEST')  
8           AND (Area >= 80)  
9 )  
10 SELECT AVG(Resale_Price) FROM Tab1
```

## 2 INPUT FORMAT

The input file `ResalePricesSingapore.csv` is the historical transaction records of the resale HDB flat in Singapore during the past 10 years (Jan. 2014 – Dec. 2023). The data is extracted from an open access dataset published on Singapore's national open data collection website\* maintained by Data.gov.sg team.

The input data is given in .csv format. You can download the data via NTU Learn. The first row is the title row. Each following row contains a line of transaction information as listed, which are separated by a comma “,”.

- **Month:** approval date of the resale, in the format YYYY-MM.
- **Town:** the town of the associated HDB flat.
- **Block:** the block of the associated HDB flat.
- **Street\_Name:** the street of the associated HDB flat.
- **Flat\_Type:** the type of flat of the associated HDB flat. In Singapore, there are 1-room flats up to 5-room flats, as well as executive flats.
- **Flat\_Model:** it implies the approximate size and the number of rooms for the HDB flat, categorized into types such as Standard, Improved, New Generation, etc.
- **Storey\_Range:** In this dataset, the storey range is given in a range of 3 (e.g. 10 to 12, which means the flat is based on the 10th to 12th storey).
- **Floor\_Area:** the floor area of the associated HDB flat in square meters.
- **Lease\_Commence\_Date:** the commence date of the flat lease in (months and) years.
- **Resale\_Price:** the resale price of the associated HDB flat.

Please note that you should store the entire table and may focus on particular information of interest for your task.

## 3 OUTPUT FORMAT

Your output file `ScanResult_<MatricNum>.csv` should contain all query results associated with the chosen matriculation number. The first row is the title row and the following

---

\*Data from: [Data.gov.sg](https://data.gov.sg).

rows present the query information and results as listed below which are separated by commas “,”. If there is no qualified data in your target range, please take "No result" as the query result.

- **Year:** the year in the query, in the format of YYYY.
- **Month:** the month you **start** to collect statistics, in the format of MM.
- **Town:** the town where the queried HDB flats locate.
- **Category:** type of statistic associated with the query. Please use the expression provided in Table 2 in output, such as “Average Price” or “Standard Deviation of Price”.
- **Value:** the value of the query outcome.

**Example:** Suppose two consecutive queries access the average price and standard deviation of price for HDB flats whose area is larger or equal to 80 square meters in CHOA CHU KANG during January 2023 to February 2023, respectively, which are 530000.05 and 4100.18. Then the corresponding result rows in the output file should be:

ScanResult_A1234313B.csv (example)					
1	Year,Month,town,Category,Value				
2	2023,01,CHOA CHU KANG,Average Price,530000.05				
3	2023,01,CHOA CHU KANG,Standard Deviation of Price,4100.18				
4	...				

We would like to remind you that **all of the 4 types of query contents presented in Table 2** should be tested with your program.

## 4 SUBMISSION

**Time:** During Week 14 (By April 25 unless otherwise specified)

**Method:** Via NTULearn

The required files include the output file, the source code of your program, and an assignment report. They should be compressed and submitted in a .zip file. Name the .zip file with your **group number**. The requirements of each files are as follows:

- Output Files ScanResult\_<MatricNum>.csv: the scan results following the requirements in **Output Format** Section. Do not include any raw or intermediate data files.
- Source Code source: the file or folder containing the source codes that input the file ResalePricesSingapore.csv and the matriculation number, and output the corresponding ScanResult\_<MatricNum>.csv. Source codes should be well-commented and contains essential documentations to help understand the functionalities.
- Report Report.pdf: the report exported in .pdf format. Your report sections and contents should follow the requirements in **Report Format** in Appendix. The report should be at most 5 pages (single column, font size 11pt, excluding cover page and contribution form). If your really want to place big figures such as screenshots (or

design figures) that you cannot squeeze into the 5 main pages after your best trial, you may optionally add an Appendix within two pages. This section is fully optional which should only contain figures and corresponding captions, and the assessment would still be mainly based on the 5-page main content.

## 5 FORMING GROUPS

The expected group size is 3. We encourage you to form groups autonomously by editing the [Online Form](#) before **February 6** (Week 4 Thursday). Students who are not involved in any group will be randomly assigned to a group by the TA.

## 6 ASSESSMENT

This is a **group project**. Your submission will be evaluated in multiple aspects, including design sophistication (e.g., whether the program meets basic requirements and whether there are additional optimizations), output accuracy, code quality (e.g., whether you can reuse some functions for conciseness), and report quality. Late submission will be penalized. The evaluation of an individual is based on the contribution form.

## 7 GENERAL GUIDELINES

1. If you are not familiar with the `.csv` format input file, you can regard it as a plain text file (just like `.txt` format).
2. Please note assuming that the `Month` column is monotonically increasing may lead to inaccurate query outcomes. Additionally, the resale HDB flats locating in the same town may not be strictly clustered together in the `Town` column.
3. While we recommend Java, you are free to choose any programming language in case you are not familiar with Java.
4. Ensure that your program is implemented in the column-store manner. Avoid high-level data tools when storing and processing the data. Example 1: Python pandas is not column-oriented. Example 2: simple SQL implementation is not column-oriented.
5. Please validate the accuracy of your query results with supplementary tools, such as Microsoft Excel or Google Sheets, and include supporting evidence in your report. For instance, you can attach the associated Office Scripts in Excel (or Excel formulas) and screenshots of the results to compare with your outputs.
6. The computation time of your program will not be evaluated as the working environments and hardware configurations vary widely. However, we still encourage you to implement optimizations to improve efficiency and present the enhancements in your report. You may consider dealing with various scenarios, such as handling data too large to fit in main memory, speeding up data scanning with additional index or specialized data layouts, and boosting computational efficiency through data compression techniques, etc. The design sophistication will be counted into assessment.
7. The code and report should be developed on your own, and using AI tools to generate codes and reports is not allowed.

# **REPORT FORMAT**

Name and Matriculation Number

## **1 Data Storage**

In this section, explain how your program handles and stores the data. You may present your design and experience (whether success or failure) related to:

- How to store the data in the column-store approach<sup>‡</sup>;
- How to design data columns for efficient processing<sup>‡</sup>;
- How to read and write the input/output files;
- How to handle possible exceptions such as empty qualified entries in the code.

## **2 Data Processing**

In this section, explain how your program scans the data and finds the values. You may present contents related to:

- How to scan columns according to task conditions;
- How to decide and record the statistics (e.g. the minimum values);
- How to improve the efficiency in scanning columns<sup>‡</sup>;

## **3 Experiment Result**

In this section, present the experimental results that your program successfully complete the tasks. The following contents are compulsory:

- Screenshots that your program executes and outputs results successfully;
- Evaluations that the output results are correct (You may put screenshots in the report comparing your output results with the correct results output by other tools such as Excel).

---

<sup>‡</sup>Exploration and improvements on these aspects are encouraged.

# CONTRIBUTION FORM

Group Number

Name	Matriculation Number	Detailed Individual Contribution	Percentage (100% in total)

**Name and Signature from all group members:**

Name and Signature of Member 1

Name and Signature of Member 2

Name and Signature of Member 3

Name and Signature of Member 4