

There are two projects contributing 50% to the total assessments of this unit. The projects are to be submitted to csubmit during the semester.

- Marking scheme and more details of the projects will be available here as the semester progresses. The overall objectives of the projects are to build a data warehouse from real-world datasets, and to carry out basic data mining activities including association rule mining, classification and clustering.

Project 1 contributes 25% to the final grade of this unit, and requires submission to [cssubmit-cits3401](#). Project 1 is an individual effort on data warehouse design and implementation, due on **Friday 23:59 pm 2nd April** ([cssubmit-cits3401](#)). In this project, we will use the World Bank COVID-19 dataset as the source of data for the data warehouse. A copy of the World Bank COVID-19 dataset is [here](#).

The overall objectives of this project are to build a data warehouse from the given data, and to **answer the following 4 business queries**.

- You may follow Kimball's four steps to designing a data warehouse. To realise the four steps, you can start by drawing and refining a StarNet.

1. Observe the data carefully. You can find all the information needed to answer the 4 business queries from the data provided. Draw a StarNet with the aim to identify the dimensions and concept hierarchies for each dimension. This should be based on the lowest level of information you have access to.
2. Use the StarNet footprints to illustrate how the 4 business queries can be answered with your design. Refine the StarNet if the business queries cannot be answered, for example, by adding more dimensions or concept hierarchies.
3. Once the StarNet diagram is completed, draw it using software such as Microsoft Visio (free to download under the [Azure Education](#)) or an online drawing service (i.e. draw.io) or a drawing program of your own choice. Paste it to your report.
4. Implement a suitable schema (star/snowflake) using SQL Server Management Studio (SSMS). Paste the database diagram generated by SSMS onto your report.
5. Implement Data Cleaning, Integration and ETL processes (you can use any program languages or/and any software you prefer) and load the World Bank COVID-19 dataset to populate the tables in SQL Server. You may need to create separate data files for your dimension tables.
6. Use SQL Server Data Tools to build a multi-dimensional analysis service solution, with a cube designed to answer your business queries. Make sure the concept hierarchies match your StarNet design.
7. Use Power BI to visualise the data returned from your business queries, and paste the visualisation results to your report.
8. List all the OLAP operations used when answering the 4 business queries. Briefly explain how the OLAP operations are used in answering the queries.

Try to complete as many of the above tasks as possible. You should submit a report in .pdf format. Your report should include the StarNet and query foot-prints, the star/snowflake/galaxy schema and the visualisation results to answer the business queries, as well as the explanation on how you perform data cleaning/pre-processing.

The followings are the **files needed for Project 1 submission**.

- A PDF report consists of the design, implementation and usage of the data warehouse to answer the queries, the StarNet and query foot-prints, the Star/Snowflake/galaxy Schema, the description of the data cleaning/preprocessing/ETL process for data transformation, and the visualisation results to answer the business queries.
- The SQL Script file and the CSV files for building and populating the tables of your data warehouse.
- The Visual Studio solution project file of the analysis service multi-dimensional project.
- The Power BI file (.pbix).

- [5 marks x 2] Data cleaning/pre-processing/ETL process for data transformation with code or screenshots or explanation
- [5 marks x 2] 4 compulsory business queries that the StarNet can answer and Power BI visualisation corresponding to the 4 business queries
- [5 marks] Concept hierarchies and corresponding StarNet
- [5 marks] Star/Snowflake schema for data warehouse design
- [5 marks] SQL Script file for building and loading the database
- [5 marks] Coherence between the design and implementation, quality and complexity of the solution, reproducibility of the solution
- [5 marks] The correctness of the answer of task 8

Data warehousing exercises are often open-ended. In other words, there is almost always a better solution. **You can interpret the scale of marks as:**

- 5 - Exemplary (comprehensive solution demonstrating professional application of the knowledge taught in the class with initiative beyond just meeting the project requirement. I.e. a highly automated and highly fault tolerant solution for data cleaning/preprocessing and/or other data operations, a deep understanding in dataset with excellent schema design, a very clear, pretty, and convincing powerBI visualisation design, a formal style and well-written report.)
- 4 - Proficient (correct application of the taught concepts)
- 3 - Satisfactory (managed to meet most of the project requirement)
- 2 - Developing (some skills are demonstrated but need revision)
- 1 - Not yet Satisfactory (minimal effort)
- 0 - Not attempted.