

Assignment #4

CSCI 5408 (Data Management, Warehousing, Analytics)
Faculty of Computer Science, Dalhousie University

Date Given: Mar 26, 2021

Due Date: Apr 7, 2021 at 11:59 pm

Late Submissions are not accepted and will result in deductions of 10%/day

Disclaimer: This assignment requires students to work on BI Framework, and sentiment/semantic analysis. Submissions related to this assignment will not be used for commercial purposes.

Objective:

- The objective of this assignment is to understand BI framework, creating star/snowflake schema, and concept of sentiment and semantic analysis.

Plagiarism Policy:

- This assignment is an individual task. Collaboration of any type amounts to a violation of the academic integrity policy and will be reported to the AIO.
- Content cannot be copied verbatim from any source(s). Please understand the concept and write in your own words. In addition, cite the actual source. Failing to do so will be considered as plagiarism and/or cheating.
- The Dalhousie Academic Integrity policy applies to all material submitted as part of this course. Please understand the policy, which is available at:
https://www.dal.ca/dept/university_secretariat/academic-integrity.html

Assignment Rubric

	Excellent (25%)	Proficient (15%)	Marginal (5%)	Unacceptable (0%)	Problem # where applied
Completeness including Citation	All required tasks are completed	Submission highlights tasks completion. However, missed some tasks in between, which created a disconnection	Some tasks are completed, which are disjoint in nature.	Incorrect and irrelevant	Problem #3
Correctness	All parts of the given tasks are correct	Most of the given tasks are correct. However, some portions need	Most of the given tasks are incorrect. The submission	Incorrect and unacceptable	Problem #2

		minor modifications	requires major modifications.		
Novelty	The submission contains novel contribution in key segments, which is a clear indication of application knowledge	The submission lacks novel contributions. There are some evidences of novelty, however, it is not significant	The submission does not contain novel contributions. However, there is an evidence of some effort	There is no novelty	Problem #1
Clarity	The written or graphical materials, and developed applications provide a clear picture of the concept, and highlights the clarity	The written or graphical materials and developed applications do not show clear picture of the concept. There is room for improvement	The written or graphical materials, and developed applications fail to prove the clarity. Background knowledge is needed	Failed to prove the clarity. Need proper background knowledge to perform the tasks	Problem #1

Citation:

McKinney, B. (2018). The impact of program-wide discussion board grading rubrics on students' and faculty satisfaction. *Online Learning*, 22(2), 289-299.

Tasks

- This assignment requires you to submit programming codes on gitLab, and a single PDF file on Brightspace.

Problem #1

Business Intelligence Reporting using Cognos

1. Download the weather dataset available on <https://www.kaggle.com/PROPPG-PPG/hourly-weather-surface-brazil-southeast-region?select=sudeste.csv>
2. Explore the dataset and identify data field(s) that could be measured by certain factors or dimensions. (Follow recorded lecture #18, and synchronous session #18)

Example: In a Sales dataset, you may find a measurable field "total sales", which could be analyzed by other factors such as, "products", "time", "location" etc. These factors are known as dimensions. Depending on the data, you may also find possibilities of slice and dice, i.e. analysis could be possible in more granular level; From **total sales by city** to **total sales by store**

3. Write ½ page explanation on how did you select the measurable field, i.e. fact and what are the possible dimensions. Include this part in your PDF file.

4. Clean the dataset, if required perform formatting. You can perform the cleaning and formatting using spreadsheet operation or programming script. If you use program add that in GitLab, if you use other methods, write the steps in the PDF file.
5. Create Cognos account and import your dataset. Identify the dimensions, and create/import the dimension tables.
6. Based on your understanding of the domain (please read the information/metadata available on the dataset source, i.e. Kaggle), create star schema or snowflake schema. Provide justification of your model creation in the PDF file.
7. In addition to justification, attach screenshot of the model (star schema or snowflake schema) in the PDF file.
8. Display visual analysis of the data in a suitable format, e.g. bar graph showing temperature change in terms of a suitable dimension. Add the screenshot of the analysis on the pdf or add a screen recording of the analysis on your .zip folder.

Problem #2

Sentiment Analysis

9. To perform this task, you need to consider the processed tweets (“messages” or “texts” only, ignore metadata) that you obtained in previous assignment.
10. Write a script to remove URL and/or any special characters. (If not done in Assignment 3. Otherwise ignore this step)
11. Write a script to create bag-of-words for each tweet. (code from online or other sources are not accepted)
e.g. tweet1 = “hey i m happy in Canada not the room”
bow1 = {“hey”:1, “i”:1, “m”:1, “happy”:1, “in”:1, “Canada”:1, “not”:1, “the”:1, “room”:1}
12. Compare each bag-of-words with a list of positive and negative words. You can download list of positive and negative words from online source(s).
13. Tag each tweet as “positive”, “negative”, or “neutral”. You can add an additional column to present your finding.

Tweet	Message/tweets	match	polarity
1	hey i m happy in Canada not the room	Happy, not	neutral

Problem #3

Semantic Analysis

14. For this task, consider the processed tweet collection that you created in Assignment
15. Use the following steps to compute TF-IDF (term frequency-inverse document frequency)
- Suppose, you have 500 tweets (messages only) that are stored in 500 JSON arrays. You need to consider these data points as the total number of **documents (N)**. In this case $N=500$
Now, use the search query “flu”, “snow”, “cold”, and search in how many documents these words have appeared.

Total Documents	500		
Search Query	Document containing term(df)	Total Documents(N)/ number of documents term appeared (df)	$\text{Log}_{10}(N/df)$
flu	20	500/20	1.39
snow	40	500/40	1.09
cold	10	500/10	1.69

- Once you build the above table, you need to find which document has the highest occurrence of the word “cold”. You can find this by performing frequency count of the word per document.

Term	Canada	
Canada appeared in 20 documents	Total Words (m)	Frequency (f)
Article #1	6	2
Article #2	10	1
:	:	:
Article #20	8	1

- You should print the news article, which has the highest relative frequency. You can find this by computing (f/m) .

Assignment 4 Submission Format:

- 1) Compress all your reports/files into a single .zip file and give it a meaningful name.

You are free to choose any meaningful file name, preferably - **BannerId_Lastname_firstname_5408_A4** but avoid generic names like assignment-4.

- 2) Submit your reports only in PDF format.

Please avoid submitting .doc/.docx and submit only the PDF version. You can merge all the reports into a single PDF or keep them separate. **You should also include output (if any) and test cases (if any) in the PDF file.**

- 3) Your executable code/script needs to be submitted on <https://git.cs.dal.ca/>