



**DALHOUSIE
UNIVERSITY**

**CSCI 5408 – Data Management, Warehousing,
Analytics**

Assignment – 4

Prashit Patel - B00896717

Gitlab Repository

https://git.cs.dal.ca/pppatel/csci-5408-f2021-b00896717-prashit_patel.git

Task 1: Business Intelligence reporting using Cognos

Facts and Dimensions

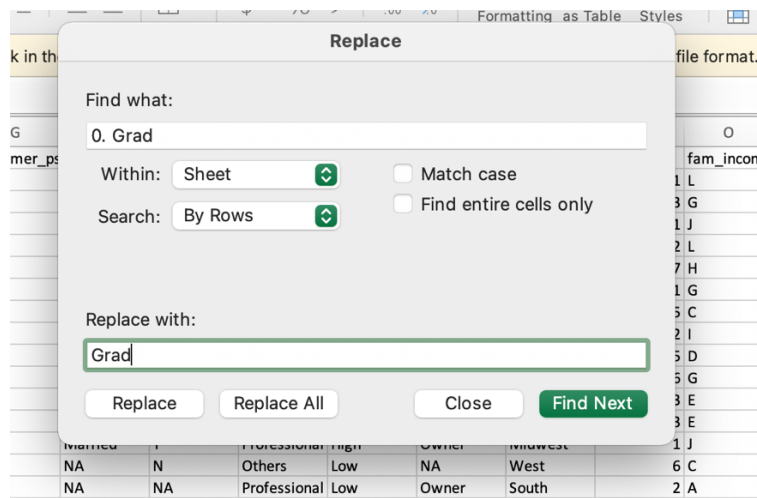
- Sale preferences can be calculated based on different dimensions such as gender, region, age, ownership of the house, occupation, education, product bought or not and online shopping experience. We can gather these sales preferences, study patterns and improve total sales based on the analysis.

Dimensions

1. **Product bought or not:** Y, N - field values describe if the customer bought the product or not. Y specifies that the customer bought the product and N specifies that the customer did not buy the product.
2. **Gender:** Male, Female - field values describe if the customer was a male or a female.
3. **Region:** Northeast, Midwest, West, South, Rest – field values describe the region of the customer.
4. **Age:** Under 25, Under 35, Under 45, Under 55, Under 65, Over 65 – fields describe the age range of the customer.
5. **Ownership of house:** Renter, Owner – describe the ownership of the house of the customer.
6. **Occupation:** Professional, Blue Collar, Sales/Service, Others, Retired, Farm – field values describe the occupation of the customer.
7. **Online shopping experience:** Y, N – field values describe if the customer had any previous experience of online shopping or not. Yes specifies that the customer had a previous online experience while N specifies no previous experience.
8. **Education:** Grad, Bach, College, HS, Less than HS – field values specifies the education of the customer.

Cleaning and Formatting



- Spreadsheet Find and Replace feature was used to clean and format data.








Fields transformed:

- gender – U was converted to NA.
- education – 0. <HS was converted to Less than HS, 1. HS was converted to HS, 2. Some College was converted to College, 3. Bach was converted to Bach, 4. Grad was converted to Grad, blanks were converted to NA.
- age – 1_Unk was converted to NA, 2_<=25 was converted to Under 25, 3_<=35 was converted to Under 35, 4_<=45 was converted to Under 45, 5_<=55 was converted to Under 55, 6_<=65 was converted to Under 65, 7_>65 was converted to Over 65.
- customer psy – U was converted to NA.
- Marriage – blanks were converted to NA.
- child – 0 was converted to N, U was converted to NA.
- mortgage – 1Low was converted to Low, 2Med was converted to Med, 3High was converted to High.
- house_owner - blanks were converted to NA.
- fam_income – U was converted to NA.

Fact Table:

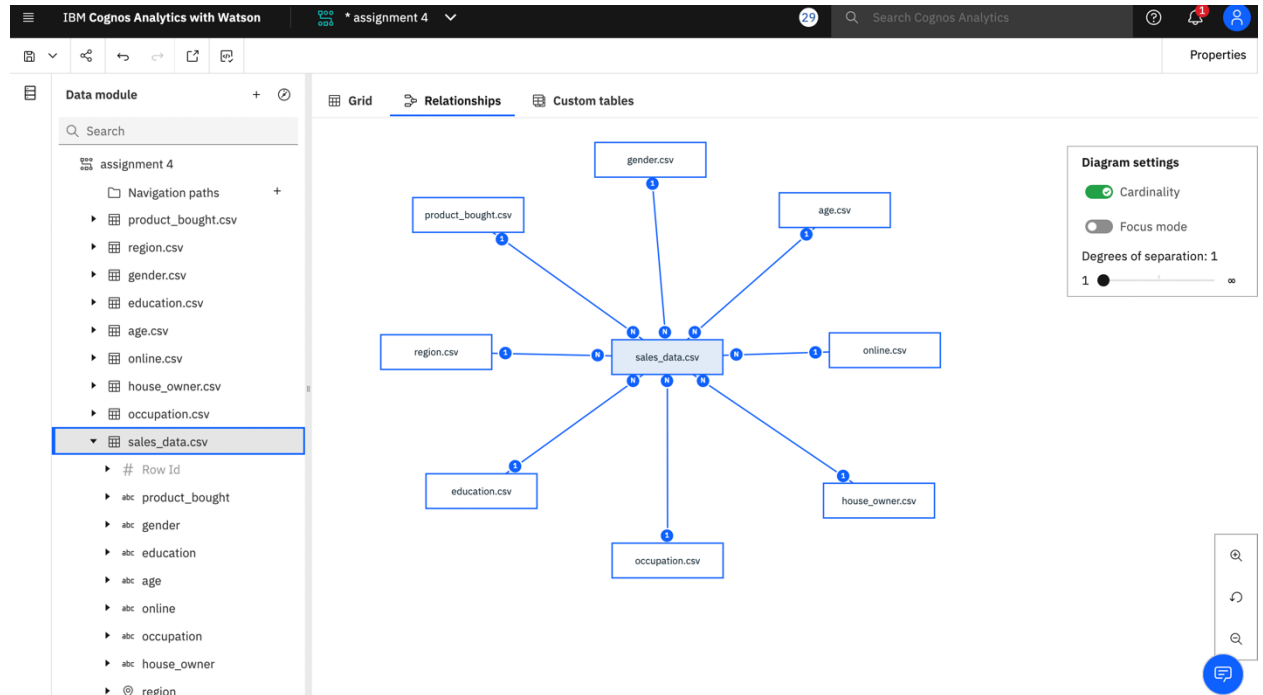
▼  sales_data.csv
▶ # Row Id
▶ abc product_bought
▶ abc gender
▶ abc education
▶ abc age
▶ abc online
▶ abc occupation
▶ abc house_owner
▶  region

Dimensions Table:

- ▶  product_bought.csv
- ▶  region.csv
- ▶  gender.csv
- ▶  education.csv
- ▶  age.csv
- ▶  online.csv
- ▶  house_owner.csv
- ▶  occupation.csv

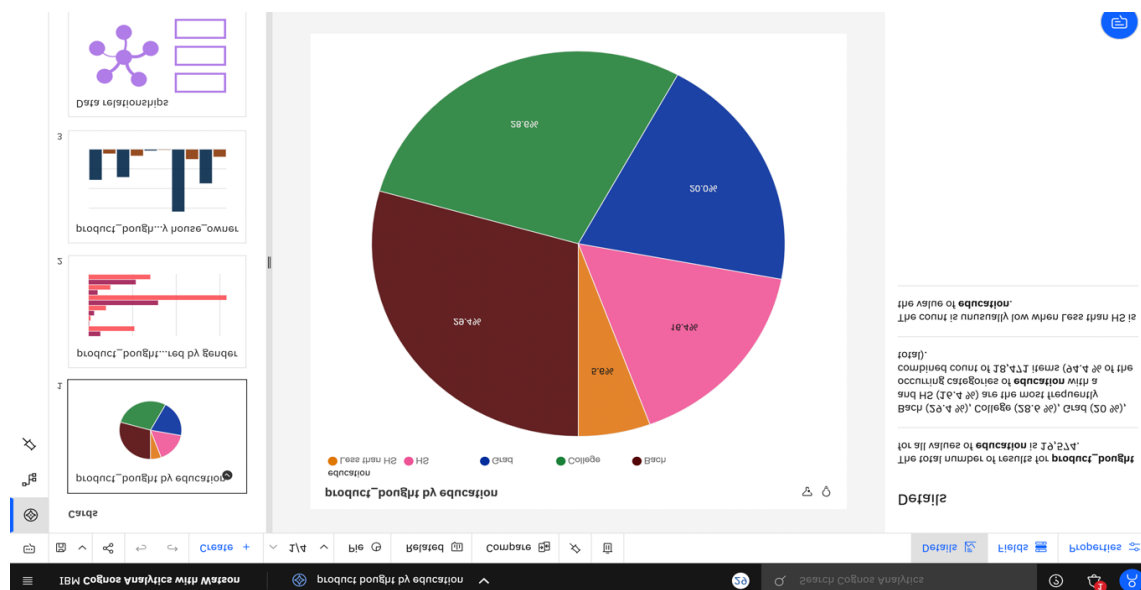
Star Schema

The above facts and dimensions were used to create a star schema. All the CSV files were imported as tables. Then, many(N) to 1 relationship was created between imported facts and tables.

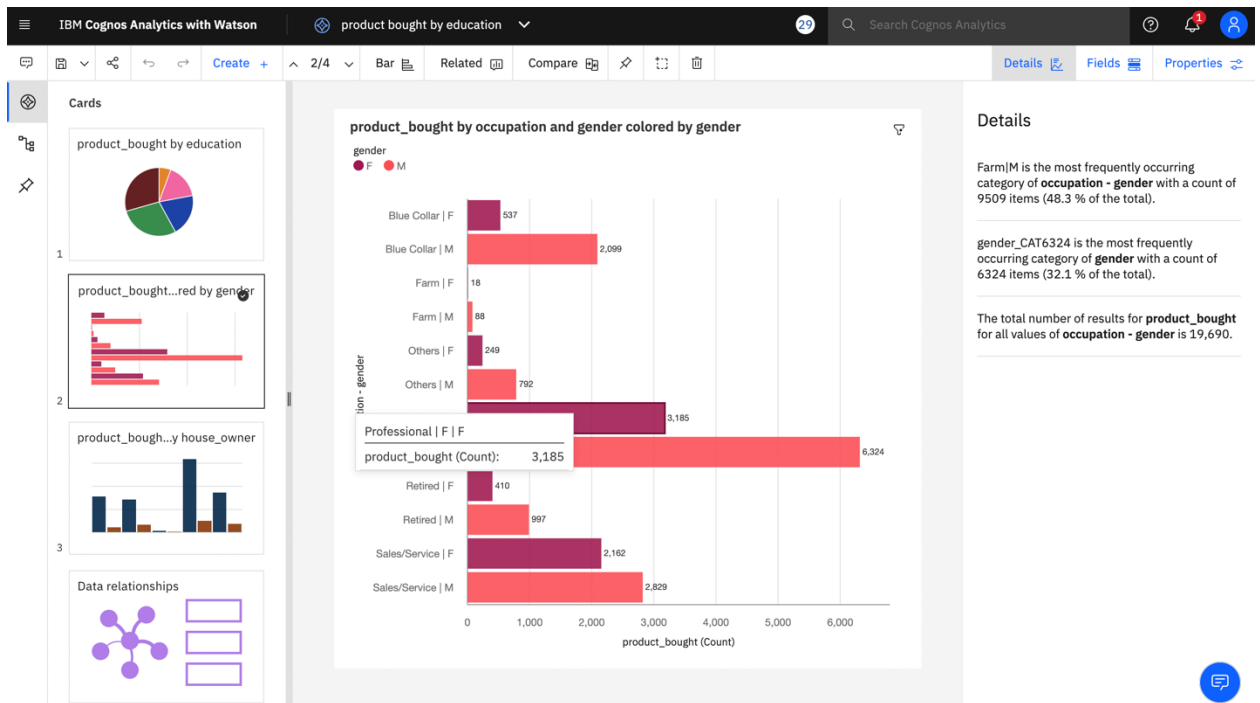


Visual Analysis

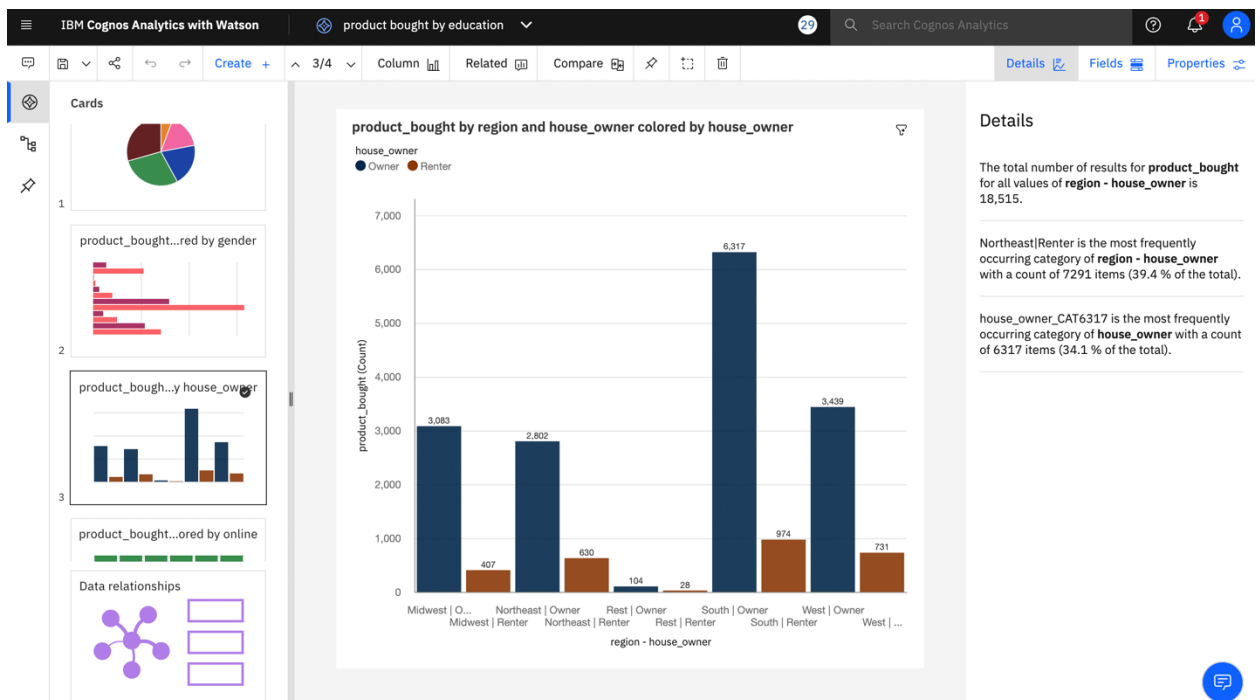
1. product bought by education: Percentages can be obtained for the type of education background of the customers that bought the product.



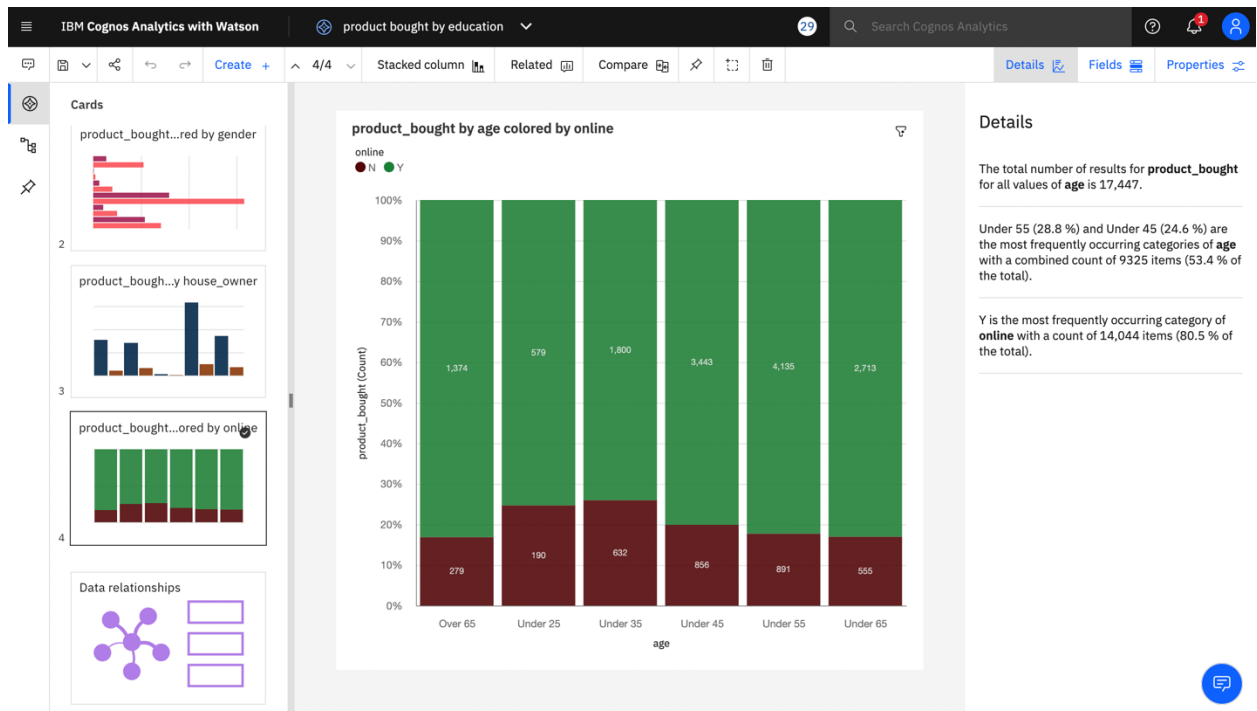
- product bought by occupation and gender: The total number of customers that bought the product based on the occupation and gender.



- product bought by region and house_owner: Total number of products bought according to regions and house ownership type of the customer.

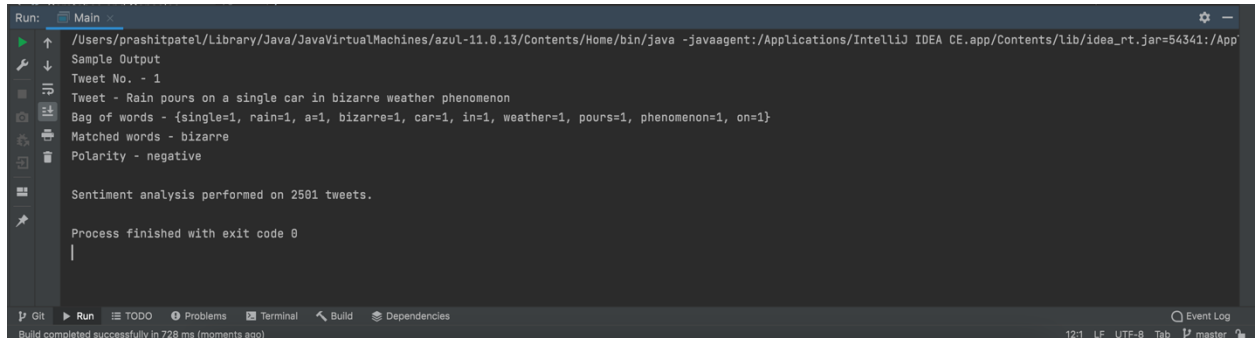


4. product bought by age and online shopping experience: Percentage and the total number of customers that bought the product based on age group and previous online experience amongst them.



Task 2: Sentiment Analysis

Screenshot of sample output:

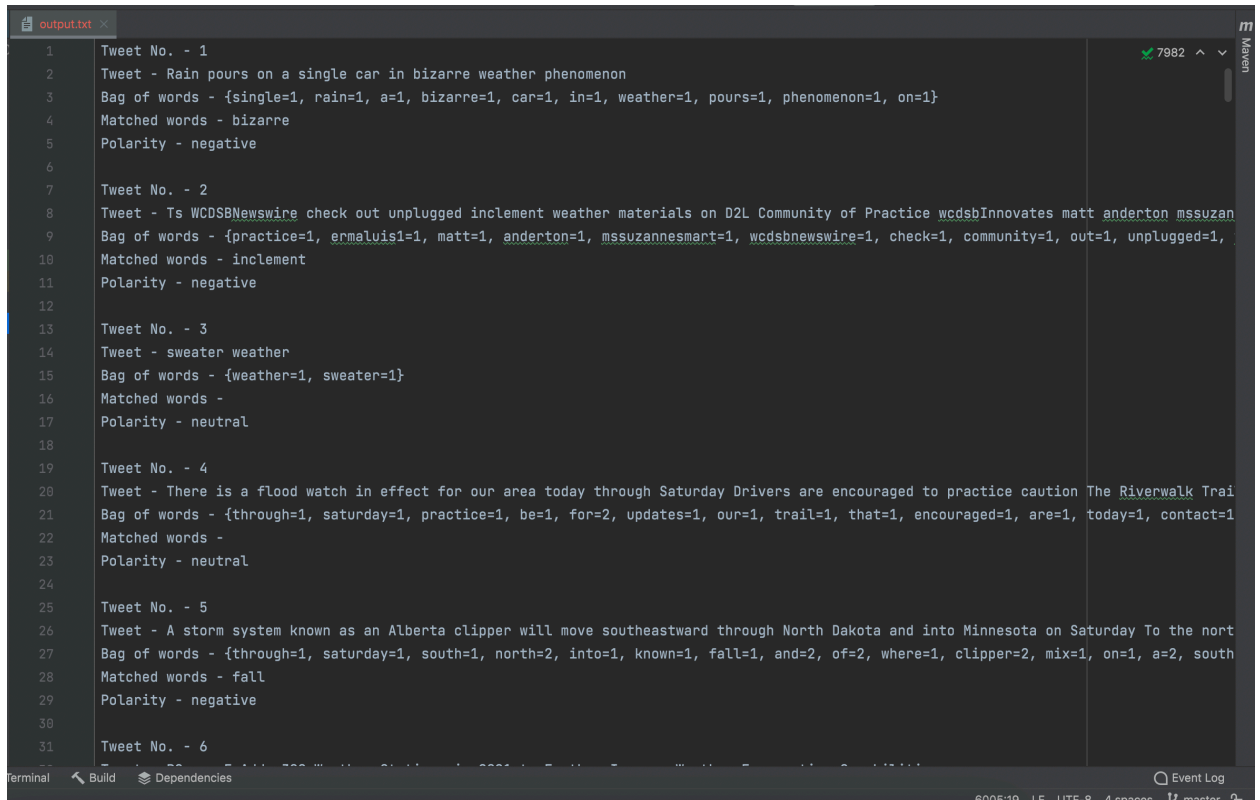


```
Run: Main
/Users/prashitpatel/Library/Java/JavaVirtualMachines/azul-11.0.13/Contents/Home/bin/java -javaagent:/Applications/IntelliJ IDEA CE.app/Contents/lib/idea_rt.jar=54341:App
Sample Output
Tweet No. - 1
Tweet - Rain pours on a single car in bizarre weather phenomenon
Bag of words - {single=1, rain=1, a=1, bizarre=1, car=1, in=1, weather=1, pours=1, phenomenon=1, on=1}
Matched words - bizarre
Polarity - negative

Sentiment analysis performed on 2501 tweets.

Process finished with exit code 0
```

The complete output of the program can be found in the output file in program directory on GitLab.



```
output.txt
1 Tweet No. - 1
2 Tweet - Rain pours on a single car in bizarre weather phenomenon
3 Bag of words - {single=1, rain=1, a=1, bizarre=1, car=1, in=1, weather=1, pours=1, phenomenon=1, on=1}
4 Matched words - bizarre
5 Polarity - negative
6
7 Tweet No. - 2
8 Tweet - Ts WCDsBNewsWire check out unplugged inclement weather materials on D2L Community of Practice wcdsbInnovates matt anderton mssuzan
9 Bag of words - {practice=1, ermalvisi=1, matt=1, anderton=1, mssuzannesmart=1, wcdsbnewsWire=1, check=1, community=1, out=1, unplugged=1,
10 Matched words - inclement
11 Polarity - negative
12
13 Tweet No. - 3
14 Tweet - sweater weather
15 Bag of words - {weather=1, sweater=1}
16 Matched words -
17 Polarity - neutral
18
19 Tweet No. - 4
20 Tweet - There is a flood watch in effect for our area today through Saturday Drivers are encouraged to practice caution The Riverwalk Trai
21 Bag of words - {through=1, saturday=1, practice=1, be=1, for=2, updates=1, our=1, trail=1, that=1, encouraged=1, are=1, today=1, contact=1
22 Matched words -
23 Polarity - neutral
24
25 Tweet No. - 5
26 Tweet - A storm system known as an Alberta clipper will move southeastward through North Dakota and into Minnesota on Saturday To the nort
27 Bag of words - {through=1, saturday=1, south=1, north=2, into=1, known=1, fall=1, and=2, of=2, where=1, clipper=2, mix=1, on=1, a=2, south
28 Matched words - fall
29 Polarity - negative
30
31 Tweet No. - 6
```


Task 3: Semantic Analysis

- Moncton was found in 0 articles. So, weather keyword is considered in place of Moncton.

Screenshot of sample output:

```
Run: Main x
/Users/prashitpatel/Library/Java/JavaVirtualMachines/azul-11.0.13/Contents/Home/bin/java -javaagent:/Applications/IntelliJ IDEA CE.app/Contents/lib/idea_rt.jar=54486:/App
Sample Output
Search Query | df | n/df | log10 n/df
canada | 60 | 33 | 1.5185139398778875
toronto | 9 | 222 | 2.346352974450639
weather | 14 | 142 | 2.1522883443830563

canada appeared in 60 documents
Article No. | Total Words | Frequency
1 | 814 | 1
2 | 234 | 2
3 | 153 | 1
... and so on.

Max Relative Frequency
Article No - 9
Article - 1t First Toronto Capital Corp said it completed an issue of a five mln dlrs convertible debenture to Arcalex B V a Dutch corporation that owns 54 3 pct of F
Key - toronto
Frequency - 0.05357142857142857

Semantic analysis performed on 2001 articles

Process finished with exit code 0
```

The complete output of the program can be found in the output file in program directory on GitLab.

```
output.txt x
1 | Search Query | df | n/df | log10 n/df
2 | canada | 60 | 33 | 1.5185139398778875
3 | toronto | 9 | 222 | 2.346352974450639
4 | weather | 14 | 142 | 2.1522883443830563
5
6 | canada appeared in 60 documents
7 | Article No. | Total Words | Frequency
8 | 1 | 814 | 1
9 | 2 | 234 | 2
10 | 3 | 153 | 1
11 | 4 | 413 | 6
12 | 5 | 66 | 1
13 | 6 | 292 | 1
14 | 7 | 195 | 1
15 | 8 | 195 | 1
16 | 9 | 292 | 1
17 | 10 | 110 | 1
18 | 11 | 183 | 3
19 | 12 | 106 | 2
20 | 13 | 99 | 1
21 | 14 | 82 | 1
22 | 15 | 339 | 3
23 | 16 | 203 | 1
24 | 17 | 307 | 7
25 | 18 | 187 | 1
26 | 19 | 331 | 3
27 | 20 | 175 | 1
28 | 21 | 74 | 1
29 | 22 | 879 | 2
30 | 23 | 491 | 1
31 | 24 | 191 | 1
-- | -- | -- | -- | --
Terminal Build Dependencies Event Log
```

References

- [1] “Positive words”, Ptrckprry [Online].
Available: <https://ptrckprry.com/course/ssd/data/positive-words.txt>
[Accessed on 24th November 2021]
- [2] “Negative words,” Github [Online].
Available: <https://gist.github.com/mkulakowski2/4289441> [Accessed on 24th November 2021]
- [3] “Print tables using formatter,” Javatpoint [Online].
Available: <https://www.javatpoint.com/how-to-print-table-in-java-using-formatter> [Accessed on 24th November 2021]
- [4] “Individual Company Sales Data,” Kaggle [Online]
Available: <https://www.kaggle.com/mickey1968/individual-company-sales-data> [Accessed on 25h November 2021]
- [5] “Cognos Analytics”, IBM [Online].
Available: <https://www.ibm.com/products/cognos-analytics> [Accessed on 25th November 2021]
- [6] Minqing Hu and Bing Liu. "Mining and Summarizing Customer Reviews." Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-2004), Aug 22-25, 2004, Seattle, Washington, USA [Accessed on 24th November 2021]