

When applying the Jaccard distance on the two sets, the entities will be paired based on the intersections of entities on the pages. In this case, when comparing the similarity of the sets, the numerator is then the number of entities intersecting over time in the different page versions within the set, and the denominator is the total number of entities of the different versions over time. The output shows the percentage of how similar each entity are between these sets. A is set one and B is set two where n stands for the page edits of the applicable set that are done over time. This will result in the following formula in Figure 3 (where X represents the entities of set A_n , while Y represents the entities of set B_n):

$$J(A_n, B_n) = 1 - \frac{|X \cap Y|}{|X \cup Y|}$$

Figure 3: Jaccard Distance for 2 Wikipedia pages

The Jaccard distance is chosen, because of the accessible appliance for MapReduce and Spark. MapReduce is a technique to handle large sets of data. Which in this case are the different page versions within the sets. Spark is a fast cluster computing system for processing both batch and real-time data processing. These systems and the applicable solutions are explained later on in the report (Section 3: Solution). For this project, it is chosen to apply the Spark solution to the generated dataset. But in the end, the solution should be accessible and reproducible for other programmers who are interested in calculating the Jaccard distance between different sets of Wikipedia pages.

The report is structured as follows. Section 2 frames how the dataset is generated. Section 3 describes the solution, where the MapReduce and Spark approaches are discussed. And finally, in Section 4, the result of the Spark solution is modelled in a similarity graph and described how the approach was done.

2 DATASETS

2.1 Generating Dataset

The dataset is generated by using the *Wikipedia-histories* library from Nathan Drezner. This library allows the user to scrape data from the applicable subject with all the previous page versions and the content made in the past.

For this report, it is chosen for six topics of the subcategory within the Business Intelligence subject of Wikipedia: Business intelligence, Competitive intelligence, Data mining, Market intelligence, Open-source intelligence and Product intelligence. Not all the subcategories could be generated due to the computation time of generating the pages and their previous versions. That is why it is decided to select six topics out of the eleven topics classified as a subcategory for the topic Business intelligence.

Business intelligence (BI) is defined by Negash and Gray [6] as "a data-driven decision support system that combines data storage, data gathering, and knowledge management with analysis to provide input to the decision process." As reported by Calof and Wright [2], Competitive intelligence (CI) is defined as the "system of environmental scanning which integrates the knowledge of everyone in the company."

Widom [8] describes a Data warehouse as that it "encompasses architectures, algorithms, and tools for bringing together selected data from multiple databases or other information sources into a single repository, which is suitable for direct querying or analysis." Marketing intelligence (MI) is defined by Jamil [3] as a "process designed to constantly produce knowledge for business sectors from dispersed data and information for strategic market positioning, as an organisational continuum that aims to answer typical decision problems faced by firms when competing in actual business environments."

Column name	Data type	Description
revid	Integer	Page revision ID
time	Datetime	Date and time of edits/creation
text	Varchar	Content of the page

Table 1: Data types of the dataframes

According to Steele [7], Open-source Intelligence (OI) is defined as "unclassified information that has been deliberately discovered, discriminated, distilled, and disseminated to a select audience to address a specific question." McFarlane et al. [5] describe Product Intelligence (PI) in an industrial context as "the linking of a physical order or product to information and rules governing the way it is intended to be made, stored or transported."

The information is scraped from these pages and filtered into three columns: revid (stands for revision ID), time, and text. In Table 1, a short description is provided and the applicable data types are described.

The dataframes were saved into several variables. Because Product intelligence was only edited 17 times, the last 17 edits were filtered from the topics BI, CI, Data warehouse, OI, and MI. When the data were filtered to these page versions and subsetting to these columns, the dataframes were merged into one master dataframe (stored as *master_df.csv*). This master dataframe contains 104 rows distributed over the three columns mentioned in Table 1.

3 SOLUTION

In this section, the two solutions (Spark dataframes and MapReduce) will be discussed in both ways, and for this report, it is chosen to make a temporal similarity graph out of a Spark dataframe.

3.1 MapReduce

For computing the set similarity when using MapReduce on the two sets of pages, it is required to count the number of links on each individual page A_n and then followed up by how similar it is with the links of the set of individual pages, B_n (links are represented in elements X and Y). The links from element Y that do not overlap with X should be subtracted. In the end, the number of entities will be divided by the amount elements which will result in the coefficient of how similar the sets are.

This is a step-by-step how the algorithm will run within a MapReduce manner:

Map: XY

The map method takes a key/value pair (X , Y) and returns these values unchanged in the algorithm. It will result in the following:

X , Y

Reduce: Y

The follow-up is the Reduce function, which takes the key and the set of all the values, and then computes the cardinality of values of set Z . This returns the input key with each input value and Z . This value will be used to compute the denominator of the Jaccard similarity.

$X, (Y, Z)$

Map: Switch (XY -> YX)

The next step is to take the key/value pairs, and then switch the strings between the key and the value. This will return the modified value and key of the pages. When using the Y element as the key, then this will result in combining the stages to collect all the X elements that correspond to each element of Y .

$Y, (X, Z)$

Reduce: Split

This reduce task continues with taking the key and the set of values, following up by splitting the set of values. This will output an empty key with each of the reduced lists of values (mapped as the value Q). The purpose is to load the balance of the collection of pairs in the last map method.

$Q, \{(X, Z)\}$

Map: Collect the pairs

This map takes in an empty key Q and the previously mapped list of values. Furthermore, it creates a series of X_i, X_j pairs with the first element in the list of values. This is also done in each of the following elements. The list of size n will produce $n - 1$ pairs. Except for the first one, one pair for each element in the list.

$(X_i, X_j), C_{ij} (C_{ij} = C_i + C_j)$

Reduce: Normalise the occurrences

Finally, this reduce will calculate the number of occurrences of X_i and X_j by counting the number of elements in the set of input values. This will be followed by dividing the numbers of occurrences by the total number of Y elements that occur with one or both of X_i and X_j .

$(X_i, X_j), |C_{ij}| / C_{ij} - |C_{ij}|$

Hereby the pseudocode of the solution (referring to algorithm 1).

After applying this algorithm on multiple pages and their previous pages (A_n, B_n), based on the output, a temporal similarity graph can be created based on the Jaccard distance (1 - Jaccard similarity) output of how similar each page within the set are with each other. The pseudocode for this solution is displayed in the "Jaccard distance - MapReduce" (referring to Algorithm 1).

3.2 Spark

For this project, it is chosen to make use of the Spark dataframes. For the Spark solution, the setup is slightly different. This solution has been approached by loading in the CSV file outputted from the generated dataset and convert it to a Spark dataframe. After converting it to a Spark dataframe, a function is created and used to calculate the Jaccard distance based on the text-similarity from the set of pages.

Algorithm 1: Jaccard distance - MapReduce

```
map(String key, String value) // key: set_pages1_links //
value: set_pages2_links
    for each set_pages1_link in key
    return: (set_pages1_links, set_pages2_links);

reduce(String key, Iterator values): // key: set_pages1_links
// values: set_pages2_links
    for each set_pages1_links in key:
        count of links += Count(set_pages2_links);
    return: set_pages1_links, (set_pages2_links, count of
links);

map_2(String key, String value) // key: set_pages1_links //
value: set_pages2_links, count of links
    for each link in key
    return: set_pages2_links, (set_pages1_links, count of
links);

reduce_2(String key, Iterator values): // key:
set_pages2_links // values: set_pages1_links, count of links
    for each set_pages2_links in values:
    return: empty set, {(set_pages1_links, count of links)};

map_3(String key, String value) // key: empty set // value:
set_pages1_links, count of links
    for each empty set in key
    return: (set_pages1_linksi, set_pages1_linksj), count
of linksij;

reduce_3(String key, Iterator values): // key:
(set_pages1_linksi, set_pages1_linksj) // values: count of
linksij
    for each link in key:
        result += ParseInt(count of linksij);
    return: set_pages1_linksi, set_pages1_linksj, results;
```

Algorithm 2: Jaccard distance - Spark

```
function pipeline (sparkdf);
Input : sparkdf
Output : sparkdf with distCol (Jaccard distance column)
start:
    create model =
        stage 1: HashingTF(sparkdf)
        stage 2: MinHash_LSH(sparkdf)

return: model

start:
    dfhashed = transform_model(sparkdf)
    dfmatches = calculate Jaccard distance

return: dfmatches (sparkdf with distCol);
```

Jaccard distance	Similarity level	Representing set
0.6 - 1	Similar	323
0.5	Somewhat similar	2448
0 - 0.4	Not similar	2585

Table 2: Pages similarity at different cutoffs

The pipeline function is using the *HashingTF* function from the PySpark library. The *HashingTF* maps a sequence of terms to the term frequencies using the hashing trick. Another function that has been applied for the pipeline function is by using the *MinHashLSH* function, which is also from the PySpark library. The *MinHashLSH* function inputs vectors, and then it calculates the Jaccard distance for the vectors. By combining the functions, it makes it possible to compute matches only once for each page, so that the cost of computation is growing linearly instead of exponentially.

The applied solution is constructed in the pseudocode "Jaccard distance - Spark" (referring to Algorithm 2). It is describing how the pipeline function works for calculating the Jaccard distance in the context of Spark.

Based on the newly created dataframe, it is possible to calculate how similar the pages are with each other. A 1 is representing a similarity between the sets of pages, and a 0 represents the dissimilarity of the sets of pages. When applying different cutoffs in the Spark solution, it will result that 323 pages of the different sets are similar (values between 0.6 - 1) to each other, 2448 pages of the different sets are somewhat similar (0.5), and that 2585 pages of the different sets are not similar at all (values between 0 - 0.4), see Table 2 and Figure 4.

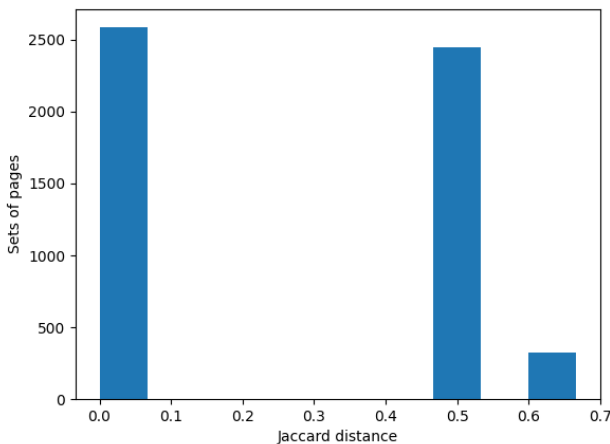


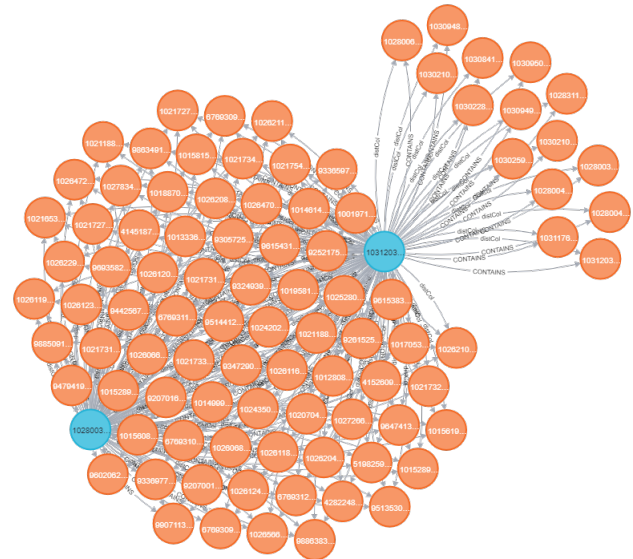
Figure 4: Jaccard distance column distribution

This concludes the Spark solution that has been applied for this assignment. For the source code, it can be found in the "SRC" folder.

4 STORAGE

4.1 Temporal Similarity Graph

Based on the Spark dataframe solution, the following graph has been modelled within Neo4J, see Figure 5. The nodes represent the page IDs for each edit, while the edges have the values of the distance column stored in them.



REFERENCES

- [1] Bela Bollobas, Gautam Das, Dimitrios Gunopulos, and Heikki Mannila. 1997. Time-series similarity problems and well-separated geometric sets. In *Proceedings of the thirteenth annual symposium on Computational geometry*. 454–456.
- [2] Jonathan L Calof and Sheila Wright. 2008. Competitive intelligence. *European Journal of marketing* (2008).
- [3] George Leal Jamil. 2013. Approaching Market Intelligence concept through a case analysis: Continuous knowledge for marketing strategic management and its complementarity to competitive intelligence. *Procedia Technology* 9 (2013), 463–472.
- [4] Włodzimierz Lewoniewski, Krzysztof Węcel, and Witold Abramowicz. 2020. Modeling Popularity and Reliability of Sources in Multilingual Wikipedia. *Information* 11, 5 (2020), 263.
- [5] Duncan McFarlane, Vaggelis Giannikas, Alex CY Wong, and Mark Harrison. 2013. Product intelligence in industrial control: Theory and practice. *Annual Reviews in Control* 37, 1 (2013), 69–88.
- [6] Solomon Negash and Paul Gray. 2008. Business intelligence. In *Handbook on decision support systems 2*. Springer, 175–193.
- [7] Robert David Steele. 2007. Open source intelligence. *Handbook of intelligence studies* 42, 5 (2007), 129–147.
- [8] Jennifer Widom. 1995. Research problems in data warehousing. In *Proceedings of the fourth international conference on Information and knowledge management*. 25–30.