

MSc in Data Science and Economics
Department of Economics, Management and Quantitative Methods
Università degli Studi di Milano

Market Basket Analysis-MeDal

Università degli Studi di Milano
Algorithms for massive data, cloud and distributed computing
JAMAL HIND – 989432



a.y. 2022/2023

I declare that this material, which I now submit for assessment, is entirely my own work and has not been taken from the work of others, save and to the extent that such work has been cited and acknowledged within the text of my work. I understand that plagiarism, collusion, and copying are grave and serious offences in the university and accept the penalties that would be imposed should I engage in plagiarism, collusion or copying. This assignment, or any part of it, has not been previously submitted by me or any other person for assessment on this or any other course of study.

Index

1 Abstract.....	4
2 Introduction.....	4
2.1 Data pre-processing.....	5
3 Market-Basket Analysis.....	7
3.1 A priori.....	7
3.2 FP-growth	8
4 Results.....	9
5 Conclusion.....	11
6 Bibliography.....	11

1 Abstract

The main aim of this project consists on the possibility to apply on a MeDal dataset some data mining techniques, in order to perform a Market-Basket analysis.

The algorithms were implemented using the Apache Spark computational framework, while the code was run on Google Colab.

In particular, considering our dataset, the column related to “texts” were identified as baskets while the respective “words” were items, and the purpose of this analysis was to output only the frequent itemset, so frequent pairs of words we had per text. For this aim we use the FP-Growth and the A-Priori algorithms.

After computing the frequent pairs and their related frequency (obviously obtaining the same result with both algorithms), the support of the frequent items were inspected, in order to output the confidence of the respective association rules. In fact, frequent itemset are useful to build meaningful association rules, from which one can infer that the presence of some items imply the presence of others.

In the design of the algorithms, the most typical Spark transformations and actions (e.g. 'map', 'flatMap', 'reduce', 'reduceByKey' etc.) were used. In the next section, a description of the two techniques used and their implementation will be presented.

2 Introduction

As we already said in the abstract, this project’s aim consists on describing and applying different techniques that are generally used to conduct market-basket analysis over huge datasets, in order to find frequent itemset.

In this specific case, the dataset is taken from the public repository of Kaggle and it contains more than 14 million rows. The folder is composed of different files representing a large dataset offering abbreviations and descriptions of all medical terms, and designed in a natural language understanding for pre-training in the medical domain. The MeDal (Medical Dataset for Abbreviation Disambiguation for Natural Language Understanding) dataset was published at the ClinicalNLP workshop at EMNLP. Actually, for the aim of this project, we are focusing on just one Excel file, which is the “full_data.csv” file, with 14 391 698 rows and 3 columns.

In the first column, nominated “TEXT”, we have the explanation of the medical terms and abbreviations; the second one is the “LOCATION”; and the last one is the “LABEL”, so one or two words used to identify a medical term. The information that we will use to develop our analysis are related to the “LABEL” and the “TEXT” column, focusing only on 2000 of the rows in our database.

Although in this setting only a sample of data is examined, the processing has been implemented to be suitable for any size of data, so for large-scale data.

To find frequent itemset we implemented different algorithms with the aid of the Map Reduce programming model.

Before applying the two algorithms we need to conduct a pre-processing phase, used to normalize the plain text of our dataset, providing a complete and cleaned input to the Market Basket Analysis algorithm.

2.1 Data pre-processing

The file of our interest can be retrieved via Kaggle API. To download the MeDal dataset we therefore need to use the relative command, which consists on putting as an input the exact name of the resource to download, which in our case is xhlulu/medal-emnlp.

```
!kaggle datasets download -d xhlulu/medal-emnlp
```

We will therefore unzip the file and before implementing and creating our algorithms, used to find frequent itemset, we need to convert our data and present them as a dataframe.

As a first step then we have to create the framework that enables us to perform parallel processing, and that is a SparkContext environment.

One of the most useful and interesting pieces of information that we can extract from our database consists on the frequency of the different labels describing each of our Medical definition. As we can see in figure 2.1, the most frequent ones appear with the words study, after, factors, development, cancer.

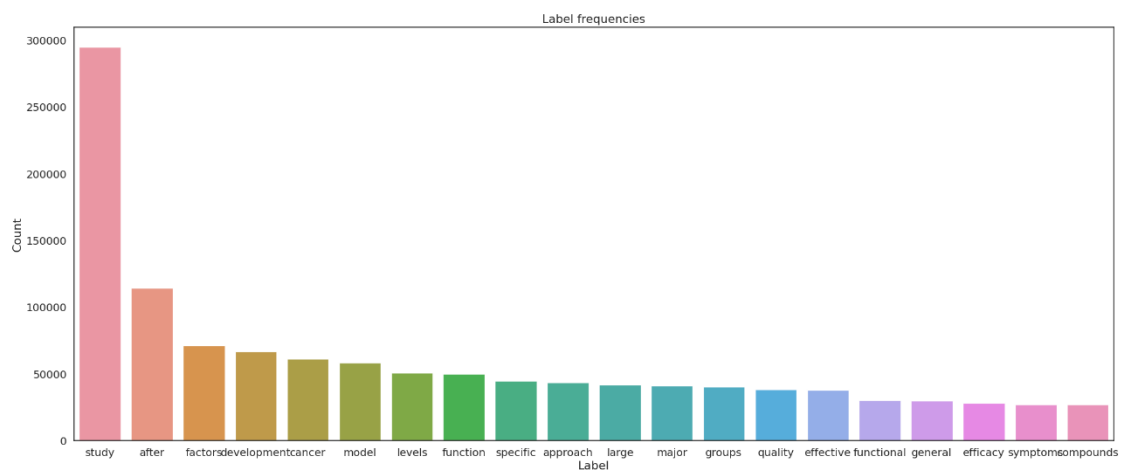


Figure 2.1. "LABEL" frequency-words.

As previously stated, our interest consists in finding the occurrences of n - items set in our file that satisfy a certain threshold, but before doing this we need to select the data of interest and execute some kind of functions to clean our data in the best possible way, so that the output is not affected by other additional inputs.

So, once we select the “TEXT” column, we are going to implement different map functions using SparkNLP, an open-source text processing library which integrates with both Python and Spark and works directly on Dataframes.

Thanks to this phase, we are able to create a set of key-value pairs containing words, on which different algorithms for market basket analysis could be applied. These pre-processing steps consist in tokenizer, required to split each text content into smaller shingles (or tokens), which corresponds to single words or punctuation marks; a map function converting each letter of each row in the “TEXT” column in a lower-case; a map function removing punctuation; and a map function removing stop-words, such as conjunctions or articles, which could negatively impact our results. We also need to apply an additional map function, required to remove any kind of words duplication per row.

What we obtain as a result is a column of cleaned words.

Figure 2.2 represents a table composed of three columns, in which we have the processed text, the cleaned-out words and a function counting the number of cleaned out words per row.

processed_text	words_clean	word_count
alphanisabolol ha...	[case, alteration...	26
a report is given...	[pronounced, deve...	123
the virostatic co...	[substance, prese...	33
rmi rmi and rmi a...	[hypotensive, lev...	105
a doubleblind stu...	[clear, obtained,...	93
stroma from eithe...	[lysis, pathway, ...	60
the effect of the...	[pathway, concent...	25
in one experiment...	[feeding, rumens,...	60
the presence of a...	[vmax, bindingdep...	96
the reaction of g...	[step, stable, ad...	81
choline acetyltra...	[using, ionexchan...	76
increasing concen...	[protein, globula...	65
the properties of...	[protein, builtin...	95
primary amines re...	[protein, react, ...	60
a purification pr...	[step, dhf, equil...	91
dihydrofolate red...	[chromatography, ...	115
ionization effect...	[indication, bett...	78
kinetic analyses ...	[dipeptidase, nac...	120
the nearultraviol...	[exhibit, fragmen...	93
the circular pola...	[appears, portion...	99

only showing top 20 rows

Figure 2.2. Complete table analysis.

Finally, as stated before, we need to select only part of the total rows, in particular given the large size of our dataset (14.393.619 rows) we will select only 2000 of the total rows, on which we will implement our algorithms.

3 Market-Basket analysis

The Market Basket Analysis is used to describe the relationship we have between two items. In particular, on one hand we have the items, while on the other one we have the baskets, and the Market Basket Analysis tries to understand which items can be considered to be frequent ones, by comparing their number of occurrences to a support threshold s and by making sure that the number of times such items appear in a basket is higher than s .

Note that this kind of analysis is linked to the association rule $I \rightarrow j$, where I is a set of items and j is an element, and according to this rule we have that whenever in a basket all the elements in the subset I are contained, also the element j is contained in that basket, so, in other terms, it consists on the probability of j being in the basket with I . The association rule is evaluated related to two measures:

- 📊 Confidence. It consists on a fraction between the support of the itemset I and the element j , and the support of the itemset I . Confidence consists therefore on the probability that an element j appears in some baskets, given that subset I is in them.

$$Conf(I \rightarrow j) = \frac{supp(I \cup \{j\})}{supp(I)}$$

- 📊 Interest. It consists on the difference between the confidence and the fraction of baskets that contain j . If I has no influence on j , then the interest is equal to 0 and the fraction of baskets including I that contain j is the same as the fraction of all the baskets that contain j ; if I has a high interest, then the presence of I in a basket encourages the presence of j in that same basket; and if I has a low interest the presence of I in a basket discourages the presence of j in that same basket.

$$Int(I \rightarrow j) = Conf(I \rightarrow j) - \frac{supp(\{j\})}{total\ baskets}$$

In order to find frequent itemset, two algorithms have been considered in this project, that obviously give us back the same results.

In particular, we will use Frequent Pattern Growth and A-priori algorithms.

3.1 A-priori

To implement the A-priori algorithm we defined a function considering as an input the RDD of baskets and based on the Map-Reduce programming model.

In particular, we will set each key to 1 and the values will contain a list of words of each text. The basic idea of this algorithms consists on the possibility to avoid computing the frequency of all possible itemsets generated from our dataset, by

creating a set of candidates whose frequency will be checked to avoid false positives.

This approach is in particular facilitated thanks to the monotonicity property, according to which if an itemset is frequent, all its subsets must be frequent as well. Therefore, it cannot exist a frequent itemset generated by items which are not frequent themselves, and thanks to such property we will avoid getting also false negatives.

The Apriori algorithm presents two phases.

First Phase: each item is set with a counter equal to one, so that we are forming some key-value pair (item,1). Then, the algorithm is going to sum the occurrences of each item collected in the first step and keep only the items whose number of occurrences exceeds the threshold.

Second Phase: those words whose count are greater than the support threshold are combined to create new itemset, and this is the set of candidate pairs. For each new pair created the algorithm is going to check if such pairs are included in the basket and in an affirmative case the respective counter will be increased by one. If the resulting counter of a given candidate pair is higher than the support, they are appended to the frequent item list.

3.2 FP-Growth

FP-Growth is a highly developed algorithm that allow us to fit a model on a larger number of observations of the dataset requiring a smaller amount of time if compared to the previous approach.

The FP-Growth algorithm can overcome the two shortcomings we have identified with the Apriori algorithm, which are the extremely large size of the candidate itemsets, since this set is actually no longer required for the FP-Growth algorithm; and the high costs given by the fact that with A-priori algorithm we need to scan and do some computations over the same item sets.

In the first step this algorithm computes the frequency of the items and compare it with a minimum threshold; in the second step the algorithm creates a data structure called frequent pattern tree, where only the frequent singletons and their count is stored.

4 Results

As previously stated, the two different methods adopted to perform the Market-Basket Analysis give us back the same result.

The output from the FP Growth Algorithm is represented in figure 4.1 and 4.2. These tables show the frequent item sets and their corresponding frequencies, which are computed using a minimum support threshold of 0.1.

items	freq		items	freq
[ph]	777		[activity, ph]	271
[activity]	587		[enzyme, activity]	251
[effect]	439		[acid, ph]	215
[two]	426		[enzyme, ph]	209
[found]	426		[c, ph]	202
[also]	421			
[acid]	403			
[enzyme]	392			
[results]	385			
[studied]	346			
[observed]	341			
[increased]	325			
[c]	322			
[concentration]	320			
[presence]	317			
[one]	301			
[may]	288			
[increase]	276			
[activity, ph]	271			
[similar]	269			

Figure 4.1. Frequent singletons and pairs of items.

The table below allow us to gain information about the relationships between different items in the same dataset.

In particular, such representation shows us the association rule generated by the FP Growth algorithm, including the antecedent (items that occur before) and the relative consequent (items that occur after).

Each item is also described by specific measures, like the confidence, which is the number of baskets that contain the antecedent item set; the lift, which measures the strength of the association between the antecedent and the subsequent item, compared to what would be expected; and the support, which measures the frequency of occurrences of the item sets in baskets.

antecedent	consequent	confidence	lift	support
[ph]	[enzyme]	0.26898326898326896	1.3723636172615763	0.1045
[ph]	[activity]	0.3487773487773488	1.1883384966860266	0.1355
[ph]	[c]	0.25997425997426	1.6147469563618633	0.101
[ph]	[acid]	0.2767052767052767	1.3732271796787927	0.1075
[enzyme]	[activity]	0.6403061224489796	2.1816222299482	0.1255
[enzyme]	[ph]	0.5331632653061225	1.3723636172615765	0.1045
[activity]	[enzyme]	0.42759795570698467	2.181622229948194	0.1255
[activity]	[ph]	0.4616695059625213	1.1883384966860264	0.1355
[c]	[ph]	0.6273291925465838	1.614746956361863	0.101
[acid]	[ph]	0.533498759305211	1.3732271796787927	0.1075

Figure 4.2. Association rule.

We then applied the Apriori Algorithm. In the following table we are showing some of the items with the highest frequency.

```
[('effect', 439),
 ('activity', 587),
 ('obtained', 261),
 ('t3', 217),
 ('determined', 204),
 ('blood', 234),
 ('one', 301),
 ('well', 207),
 ('conditions', 257),
 ('different', 263)]
```

```
[(('acid', 'ph'), 215),
 (('enzyme', 'ph'), 209),
 (('activity', 'enzyme'), 251),
 (('activity', 'ph'), 271),
 (('c', 'ph'), 202)]
```

Figure 4.3. Frequent Item sets.

As we expected, frequent pairs are composed by frequent singletons.

In both algorithms, there are no triple items with high frequency, and the number of pairs is limited to 5. This is because of the value of the threshold, which is set to a higher value with respect to the frequencies that we have, and therefore the algorithms cannot return back as frequent triples or pairs some item sets with a number of occurrences lower than the threshold value s (equal to 200); and because of the size of the sample that we are considering, which is extremely small with respect to the total amount of rows that we have in our table.

Since our A-priori algorithm is requiring a huge amount of time to perform such operations in a larger sample, we will use the FP-Growth algorithm on a sample of 35985 rows. As we can see from figure 4.4, the frequent singletons are, approximately (with a different order, since ph is now ranked 7th in terms of frequency), still the same ones, with some additional frequent pairs and triples.

items freq	items freq	
[activity] 7843	[cell, cells] 2642	
[found] 7466	[enzyme, activity] 2487	
[results] 7343	[one, two] 2211	
[two] 7304	[ph, activity] 2096	
[patients] 7206	[found, activity] 1955	
[also] 6994	[also, activity] 1868	
[ph] 6643	[two, found] 1801	
[cells] 6589	[also, found] 1777	
[may] 6033	[acid, ph] 1738	
[effect] 5951	[enzyme, ph] 1706	
[one] 5727	[two, activity] 1706	
[studied] 5577	[degrees, c] 1700	
[observed] 5147	[cells, found] 1684	
[increased] 4845	[effect, activity] 1677	
[used] 4783	[results, activity] 1672	
[acid] 4676	[treatment, patie...] 1670	
[normal] 4583	[cells, activity] 1669	
[different] 4555	[patients, two] 1659	
[showed] 4506	[also, results] 1632	
[treatment] 4474	[two, results] 1626	
		items freq
		[enzyme, ph, acti...] 1098
		[degrees, c, ph] 999
		[purified, enzyme...] 737

Figure 4.4. Frequent elements with a larger sample.

5 Conclusion

In this paper we just presented a few of the different approaches that could be used to deal with frequent item sets, but it is important to keep in mind that we still have many other algorithms and functions extremely important for our aim.

Comparing the results of the two algorithms, most of the item sets are in both outputs. We can therefore conclude that, despite the fact that the two algorithms computing frequent item sets are using different methods, since the outputs are similar the results of the two algorithms can be considered reliable.

However, because of the higher amount of computations required by the A-priori algorithm, it is basically impossible to increase the size of the sample.

FP-Growth algorithm can give us a deep overview of the different measures that could be used to understand and perform some other analysis on our data, by computing some operations on a larger sample.

6 Bibliography

<https://towardsdatascience.com/fp-growth-frequent-pattern-generation-in-data-mining-with-python-implementation-244e561ab1c3>

Mustakim, D. Herianda, and A. Ilham. Market basket analysis using apriori and fp-growth for analysis consumer expenditure patterns at berkah mart in pekanbaru riau. *Journal of Physics*, pages 1–10, 2018.

<https://towardsdatascience.com/market-basket-analysis-using-pysparks-fpgrowth-55c37ebd95c0>