

Mining Massive Data Sets Midterm Report

1st 522H0036 - Luong Canh Phong
Faculty of Information Technology
Ton Duc Thang University
Ho Chi Minh City, Vietnam
522H0036@student.tdtu.edu.com

2nd 522H0092 - Cao Nguyen Thai Thuan
Faculty of Information Technology
Ton Duc Thang University
Ho Chi Minh City, Vietnam
522H0092@student.tdtu.edu.com

3rd 522H0075 - Tang Minh Thien An
Faculty of Information Technology
Ton Duc Thang University
Ho Chi Minh City, Vietnam
522H0075@student.tdtu.edu.com

4th 522H0167 - Truong Tri Phong
Faculty of Information Technology
Ton Duc Thang University
Ho Chi Minh City, Vietnam
522H0167@student.tdtu.edu.com

5th Instructor: Nguyen Thanh An
Faculty of Information Technology
Ton Duc Thang University
Ho Chi Minh City, Vietnam
nguyenthanhan@tdtu.edu.com

Abstract—This document is a model and instructions for L^AT_EX. This and the IEEEtran.cls file define the components of your paper [title, text, heads, etc.]. *CRITICAL: Do Not Use Symbols, Special Characters, Footnotes, or Math in Paper Title or Abstract.

Index Terms—component, formatting, style, styling, insert

I. INTRODUCTION

The project consists of three distinct tasks that cover different aspects in mining large datasets. The first task involves the practice of running a MapReduce program in Java and implementation of the A-Priori algorithm with the data given by the lecturer. The second task entails the implementation of the PCY algorithm OOP form to identify frequent pairs and generate relevant information. The last task involves the student conducting implementation of the MinHashLSH algorithm and a manual method of calculating the Jaccard distance to identify similar pairs of dates defined in the requirement, as well as make comparison where needed.

II. FIRST TASK: A-PRIORI ALGORITHM FOR FREQUENT CUSTOMERS

A. Overview of MapReduce

1) *What is MapReduce*: MapReduce is a yarn-based system commonly used for processing massive dataset:

- Performs concurrent processing by dividing the dataset into multiple chunks on the Hadoop commodity servers.
- Instead of sending the data to the machine with the logic to execute, we send the logic to the data to execute, specifically, the server.

2) *How MapReduce works*: MapReduce is executed in the following order:

- Split: Divides the dataset into multiple data batches.
- Map: Maps every element within each data batch to a $\langle key, value \rangle$ pair.
- Await Completion: Wait for all data batches to finish mapping the pairs.

- Combine: Generates $\langle key, value \rangle$ pairs in the form of a list (e.g., $[[A, 1], [A, 1]]$)
- Partition: Determines which reducer should handle each key. It uses a hash function (e.g., $hash(key) \% num_reducers$) to distribute keys evenly.
- Reduce: Processes every data assigned to it and return the output.

B. First subtask

In the first subtask, we are assigned to store the data on Hadoop Distributed File System (HDFS). After which we will implement a Hadoop MapReduce program in Java to discover groups of customers going shopping at the same date.

1) The Mapper Class: CustomerGroupByDateMapper

The Mapper class is responsible for reading the input data and emitting key-value pairs. Key aspects of its implementation include:

- Input Processing: The input is a CSV file where each line contains multiple fields, including Member_number (customer ID) and Date (transaction date).
- Filtering Headers: The Mapper skips header lines by checking if the first token equals Member_number.
- Emitting Key-Value Pairs: The key is the transaction date, and the value is the customer ID. This allows transactions to be grouped by date in the shuffle and sort phase.

Example Output from Mapper:

```
(01/01/2014, 11111)
(01/01/2014, 22222)
(02/01/2014, 11111)
```

2) The Reducer Class: CustomerGroupByDateReducer

The Reducer aggregates the values emitted by the Mapper for each unique date. Key implementation features include:

- Collecting Unique Customers: Customer IDs are added to a HashSet to remove duplicates.
- Joining Values: The set of unique customer IDs is converted to a comma-separated string.

- **Emitting Results:** The final output consists of the transaction date as the key, and a list of unique customer IDs as the value.

Example Output from Reducer:

```
(01/01/2014, 11111,22222)
(02/01/2014, 11111)
```

3) *Driver Program (Main Method):* The driver program configures and runs the MapReduce job. It performs the following actions:

- **Job Setup:** A Hadoop job is created with the name `Customer Date Groups`, using `GroupMapReduce` as the main class.
- **Mapper and Reducer Assignment:** The appropriate Mapper and Reducer classes are set. The Reducer is also used as a Combiner.
- **I/O Paths:** The input and output paths are taken from command-line arguments.
- **Job Execution:** The job is submitted to Hadoop, and the program exits based on success or failure.

C. Second subtask

In the second subtask, we implement the A-Priori algorithm to identify frequent customer pairs. This is achieved using two MapReduce passes.

Note: The output of the first subtask (Customer Grouping by Date) will serve as the input for both passes of the second subtask (A-Priori Algorithm):

- **First Subtask (Customer Grouping by Date):**
 - `args[0]` – Input path to the raw transaction CSV file.
 - `args[1]` – Output path where grouped customer data by date will be written.
- **Second Subtask (A-Priori Algorithm):**
 - `args[0]` – Input path to the grouped customer data (output from the first subtask).
 - `args[1]` – Output path for the first pass (frequent individual customers).
 - `args[2]` – Output path for the second pass (frequent customer pairs).

This approach allows the workflow to seamlessly transition from the first subtask (grouping by date) to the second subtask (identifying frequent customers and pairs).

1) *The First Pass: Identifying Frequent Individual Customers*

- **Mapper Class: *AprioriFirstPassMapper***
 - **Function:** Reads grouped customer data (output of first subtask), splits the customer list, and emits each customer ID with a value of 1.
 - **Input Format:** Each line is a tab-separated pair where the key is a date and the value is a comma-separated list of customers.
 - **Filtering:** Skips malformed lines where the customer list is missing.
 - **Example Input:**

```
01/01/2014 12345,67890
```

– **Example Output from Mapper:**

```
(12345, 1)
(67890, 1)
```

- **Reducer Class: *AprioriFirstPassReducer***

- **Function:** Aggregates the occurrences of each customer ID.
- **Filtering:** Only customers meeting the support threshold (minimum occurrences) are retained.
- **Example Output from Reducer:**

```
(12345, 2)
(67890, 1)
```

2) *The Second Pass: Identifying Frequent Customer Pairs*

- **Mapper Class: *AprioriSecondPassMapper***

- **Setup:** Loads the list of frequent customers from the first pass output using Hadoop's distributed cache.
- **Processing:** For each transaction line, splits the list of customer IDs, filters only frequent customers, and generates all valid customer pairs.
- **Emitting:** Outputs each pair of frequent customers with a count of 1.
- **Example Output from Mapper:**

```
(12345,67890, 1)
(12345,54321, 1)
```

- **Reducer Class: *AprioriReducer***

- **Function:** Aggregates the counts for each customer pair.
- **Filtering:** Emits only those pairs whose frequency is greater than or equal to the support threshold.
- **Example Output from Reducer:**

```
(12345,67890, 3)
```

3) *Driver Program (Main Method):*

- **First Pass Execution:**

- The first MapReduce job is run with `AprioriFirstPassMapper` to compute individual customer frequencies.
- The output is saved and later loaded into memory for the second pass.

- **Second Pass Execution:**

- The second job uses `AprioriSecondPassMapper`, which loads the frequent customers using Hadoop's cache mechanism.
- It then computes the frequency of customer pairs and applies the support threshold.

4) *Helper Method: createJob*

A reusable helper method named `createJob` is implemented to reduce code repetition when setting up MapReduce jobs.

- **Parameters:**

- `jobName:` A string representing the name of the job.
- `mapperClass:` The class to be used as the Mapper.

- inputPath, outputPath: Paths for input and output directories.
- Functionality:
 - Configures the job with the specified name and sets the `AprioriReducer` as both the Combiner and Reducer.
 - Assigns key and value output types and adds file paths.
 - Returns a configured `Job` instance ready for execution.

III. SECOND TASK: PARK-CHEN-YU (PCY) ALGORITHM FOR FREQUENT ITEMS

A. Overview of PCY

Frequent itemset mining is essential for discovering item associations in transactional data, such as market basket analysis. The PCY algorithm improves efficiency by using hash buckets to reduce the computational cost of finding frequent item pairs. This project applies the PCY algorithm to mine frequent itemsets and generate association rules based on support and confidence thresholds, using PySpark for scalable data processing.

The PCY algorithm is based on two key passes through the data. In the first pass, frequent individual items are identified and counted. In the second pass, frequent item pairs are counted, and hash buckets are used to prune less frequent pairs. The hash function maps item pairs to buckets, and only pairs that have a sufficient bucket count are considered frequent. This approach significantly reduces the number of pair comparisons and improves algorithm efficiency.

- Step 1: Count individual items using the support threshold.
- Step 2: Count pairs of frequent items and hash them into buckets.
- Step 3: Prune item pairs that are not frequent based on the bucket counts.
- Step 4: Generate association rules using the confidence threshold.

B. Implementation Details

1) *Data Loading and Preprocessing*: Data is loaded using PySpark's `read.csv` function. Each transaction is represented as a basket, and the data is grouped by customer and date. The `collect_set` function is used to create a list of items bought together in each transaction.

2) *First Pass: Counting Frequent Items*: In the first pass, each item's frequency is counted, and only those items that meet the support threshold are considered frequent. The item counts are stored in a dictionary, sorted in descending order.

3) *Second Pass: Counting Frequent Pairs*: During the second pass, the algorithm generates pairs from frequent items and counts their occurrences. Hashing is applied to map pairs into buckets, and the bucket counts are used to prune pairs that do not meet the minimum support threshold.

4) *Association Rule Generation*: For each frequent pair, confidence is calculated as the ratio of the pair's count to the individual item count. Association rules are generated if the confidence meets the given threshold. The rules are sorted by confidence.

C. Experimental Results

The PCY algorithm was executed on a transactional dataset of retail transactions, where each transaction (basket) represented a set of items purchased by a customer. The algorithm was applied with the following parameters:

- Support threshold: 2 (minimum count for items to be considered frequent).
- Confidence threshold: 0.5 (minimum confidence for association rules).

1) *Frequent Items*: The first pass of the algorithm counted the occurrence of each item in the transactions. The following are the top 30 frequent items identified based on the support threshold:

TABLE I
MOST FREQUENT ITEMS WITH THEIR COUNTS

Item	Count
Whole milk	2363
Other vegetables	1827
Rolls/buns	1646
Soda	1453
Yogurt	1285
Root vegetables	1041
Tropical fruit	1014
Bottled water	908
Sausage	903
Citrus fruit	795
Pastry	774
Pip fruit	734

2) *Frequent Item Pairs*: In the second pass, the algorithm counted pairs of frequent items across all transactions. The top 30 frequent item pairs, sorted by frequency, are as follows:

TABLE II
FREQUENT ITEM PAIRS WITH THEIR COUNTS

Pair	Count
('Whole milk', 'Other vegetables')	222
('Whole milk', 'Rolls/buns')	209
('Whole milk', 'Soda')	174
('Whole milk', 'Yogurt')	167
('Rolls/buns', 'Other vegetables')	158
('Soda', 'Other vegetables')	145
('Whole milk', 'Sausage')	134
('Whole milk', 'Tropical fruit')	123
('Yogurt', 'Other vegetables')	121
('Rolls/buns', 'Soda')	121
('Yogurt', 'Rolls/buns')	117
('Whole milk', 'Root vegetables')	113

3) *Association Rules*: Once frequent item pairs were identified, association rules were generated based on the confidence threshold of 0.5. The confidence for each rule was computed by dividing the pair count by the count of the antecedent item. The top 30 association rules, sorted by confidence, are presented below:

TABLE III
VALIDATED ASSOCIATION RULES

Rule	Confidence
Preservation products → Soups	1.00
Kitchen utensil → Pasta	1.00
Kitchen utensil → Bottled water	1.00
Kitchen utensil → Rolls/buns	1.00
Bags → Yogurt	0.50

4) *Performance Evaluation*: The PCY algorithm effectively identified frequent items, pairs, and association rules with reasonable execution time and memory usage. Given the dataset size, the distributed nature of Spark ensured that the computation was scalable.

- **Time Complexity**: The use of hashing significantly reduces the complexity of item pair generation, making the PCY algorithm faster than traditional algorithms such as the Apriori algorithm.
- **Memory Usage**: Memory usage was managed well by leveraging the distributed processing capabilities of Spark.

IV. THIRD TASK: MINHASHLSH FOR SIMILAR DATES

ACKNOWLEDGMENT

The preferred spelling of the word “acknowledgment” in America is without an “e” after the “g”. Avoid the stilted expression “one of us (R. B. G.) thanks ...”. Instead, try “R. B. G. thanks...”. Put sponsor acknowledgments in the unnumbered footnote on the first page.

REFERENCES

Please number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in [3]—do not use “Ref. [3]” or “reference [3]” except at the beginning of a sentence: “Reference [3] was the first ...”

Number footnotes separately in superscripts. Place the actual footnote at the bottom of the column in which it was cited. Do not put footnotes in the abstract or reference list. Use letters for table footnotes.

Unless there are six authors or more give all authors’ names; do not use “et al.”. Papers that have not been published, even if they have been submitted for publication, should be cited as “unpublished” [4]. Papers that have been accepted for publication should be cited as “in press” [5]. Capitalize only the first word in a paper title, except for proper nouns and element symbols.

For papers published in translation journals, please give the English citation first, followed by the original foreign-language citation [6].

REFERENCES

- [1] G. Eason, B. Noble, and I. N. Sneddon, “On certain integrals of Lipschitz-Hankel type involving products of Bessel functions,” *Phil. Trans. Roy. Soc. London*, vol. A247, pp. 529–551, April 1955.
- [2] J. Clerk Maxwell, *A Treatise on Electricity and Magnetism*, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [3] I. S. Jacobs and C. P. Bean, “Fine particles, thin films and exchange anisotropy,” in *Magnetism*, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.
- [4] K. Elissa, “Title of paper if known,” unpublished.
- [5] R. Nicole, “Title of paper with only first word capitalized,” *J. Name Stand. Abbrev.*, in press.
- [6] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, “Electron spectroscopy studies on magneto-optical media and plastic substrate interface,” *IEEE Transl. J. Magn. Japan*, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetism Japan, p. 301, 1982].
- [7] M. Young, *The Technical Writer’s Handbook*. Mill Valley, CA: University Science, 1989.

IEEE conference templates contain guidance text for composing and formatting conference papers. Please ensure that all template text is removed from your conference paper prior to submission to the conference. Failure to remove the template text from your paper may result in your paper not being published.