

Analysis of a dataset using data mining techniques.

This notebook focuses on exploring and mining basic information from an anonymised retail transactions dataset given by the company Instacart.

Taking a first glance at the dataset

Source code for this section may be find in file `dist/first-glance.class.ts`

The dataset consists of information about 3.4 million grocery orders, distributed across 6 `.csv` files listed below:

```
In [1]: 1 import * as child_process from 'child_process';
2
3 /**
4  * Function executes a child_process listing the files in the instacart_basket_data folder
5  * and returns the listed files though a Promise.
6  */
7 function listFiles(): Promise<string[]> {
8   // Async behavior
9   return new Promise( (resolve,reject) => {
10    // Listing the files in the instacart_basket_data folder.
11    let ls: child_process.ChildProcess = child_process.exec('ls instacart_basket_data', (err: Error, stdout: string, stderr: string) => {
12      // Formatting ls output as an Array of strings (representing the file names)
13      resolve(stdout.match(/^[^\\r\\n]+/g));
14    });
15  });
16 }
17
18 function showFiles(): void {
19   // Async console for Jupyter
20   $$async$$ = true;
21
22   listFiles().then( (files: string[]) => {
23     console.log(files);
24     $$done$$();
25   }, (e) => console.log(e) );
26 }
27
28 showFiles();
29
30 [ 'aisles.csv',
31   'departments.csv',
32   'order_products__prior.csv',
33   'order_products__train.csv',
34   'orders.csv',
35   'products.csv' ]
```

Files composing the dataset are listed below:

```
[ 'aisles.csv',
  'departments.csv',
  'order_products__prior.csv',
  'order_products__train.csv',
  'orders.csv',
  'products.csv' ]
```

As a starting point, in order to have a first glance of the data we will actually be playing with throughout this entire report; let us display the first items composing each `.csv` file listed above.

We will use NodeJS's straight-forward File I/O `fs` to list, open and `'csv-parse'`'s `Parser` to parse `.csv` files.

```
In [3]: 1 import * as fs from 'fs';
2 import { Parser } from 'csv-parse';
3
4 /**
5  * Function reads a .csv file and returns it properly formatted.
6  */
7 function readFile<T>(filePath: string): Promise<{ data:Array<T>, file: string }> {
8     // Async console for Jupyter
9     $$async$$ = true;
10    // Async behavior
11    return new Promise( (resolve, reject) => {
12        let ret: Array<T> = [];
13
14        // 'csv-parse' Parser, columns options groups each row by column in an object
15        let parser: Parser = new Parser({
16            delimiter: ',',
17            columns: true
18        });
19
20        parser
21            .on('data', (data: T) => ret.push(data) )
22            .on('error', (error: any) => reject(error) )
23            .on('end', () => resolve({data: ret, file: filePath} ) );
24
25        // Reading the file and piping to parser
26        fs.createReadStream(filePath).pipe(parser);
27    });
28 }
29
30 listFiles().then( (files: string[]) => {
31     console.log('Parsing all the files, this may a while...');
32     // Creating promises array to run readFile() of files concurrently
33     let promises: Promise<any>[] = files.map( (file: string) => {
34         // Formatting to get absolute paths to files.
35         let filePath: string = `custom_data/small_${file}`;
36         // Reading file
37         return readFile<any>(filePath);
38     });
39     // Concurrent execution of promises
40     return Promise.all(promises);
41 }).then( (results: Array<any[]>) => {
42     // Reading first two elements of each file
43     results.forEach( (result: Array<any>) => {
44         console.log(`First two lines of file: ${result.file}`);
45         console.log(result.data.slice(0,2));
46     });
47     // Async console for Jupyter
48     $$done$$();
49 }).catch( (e) => console.log(e) );
```

Parsing all the files, this may a while...
First two lines of file: custom_data/small_aisles.csv
[{ aisle_id: '1', aisle: 'prepared soups salads' },
 { aisle_id: '2', aisle: 'specialty cheeses' }]
First two lines of file: custom_data/small_departments.csv
[{ department_id: '1', department: 'frozen' },
 { department_id: '2', department: 'other' }]
First two lines of file: custom_data/small_order_products_prior.csv
[{ order_id: '2',
 product_id: '33120',
 add_to_cart_order: '1',
 reordered: '1' },
 { order_id: '2',
 product_id: '28985',
 add_to_cart_order: '2',
 reordered: '1' }]
First two lines of file: custom_data/small_order_products_train.csv
[{ order_id: '1',
 product_id: '1',
 add_to_cart_order: '1',
 reordered: '1' },
 { order_id: '2',
 product_id: '101',
 add_to_cart_order: '2',
 reordered: '1' }]
First two lines of file: custom_data/small_orders.csv
[{ order_id: '2539329',
 user_id: '1',
 eval_set: 'prior',
 order_number: '1',
 order_dow: '2',
 order_hour_of_day: '08',
 days_since_prior_order: '' },
 { order_id: '2398795',
 user_id: '1',
 eval_set: 'prior',
 order_number: '2',
 order_dow: '3',
 order_hour_of_day: '07',
 days_since_prior_order: '15.0' }]
First two lines of file: custom_data/small_products.csv
[{ product_id: '1',
 product_name: 'Chocolate Sandwich Cookies',
 aisle_id: '61',
 department_id: '19' },
 { product_id: '2',
 product_name: 'All-Seasons Salt',
 aisle_id: '104',
 department_id: '13' }]

Out[3]: undefined

File aisles.csv is structured as such; and contains the two following rows:

aisle_id	aisle
1	prepared soups salads
2	specialty cheeses

File departments.csv is structured as such; and contains the two following rows:

department_id	department
1	frozen
2	other

File products.csv is structured as such; and contains the two following rows:

	product_id	product_name	aisle_id	department_id
	1	Chocolate Sandwich Cookies	61	19
	2	All-Seasons Salt	104	13

Files `aisles.csv`, `departments.csv`, `products.csv` and all (trivially) contain information about the aisles, products, and departments names respectively; which may not be of any interest to us other than for enhanced visualisation. We'll thus have to perform the proper `JOIN`s (a.k.a. `UNION`) between these tables and our future data / pattern collections when needed.

File `orders.csv` is structured as such; and contains the two following rows:

order_id	user_id	eval_set	order_number	order_dow	order_hour_of_day	days_since_prior_order
2539329	1	prior	1	2	08	
2398795	1	prior	2	3	07	15.0

File `orders.csv` is surely a bit more interesting, especially having in mind **sequential pattern mining**, as it lists all the orders, and contains information on **when** and **by who** it has been placed.

Finally, files `order_products__train.csv` and `order_products__prior.csv` contain the same, and most valuable information in regard to pattern mining, as they contain products ordered within each order. We will generate our transactions from this file, grouping `product_ids` by `order_id`.

These files are structured as below:

order_id	product_id	add_to_cart_order	reordered
1	49302	1	1
1	11109	2	1

Warming up: First statistics over the dataset

Main perk of using TypeScript could be the language's proximity to the DOM, for data visualisation. Upon this point, will be using `D3.js` library, as well as `zingchart` to visualise our data.

`D3.js` and `zingchart` are libraries allowing us to create an interactive charts with smooth transitions and interaction. If you're reading this of a `.pdf` version, please consider loading this notebook in a browser (application is dockerized), as many graphs only make sense when being interacted with (as Chord Diagrams for frequent itemset presentation for instance).

"Front-end" dependencies will be loaded below.

```
In [3]: 1 function loadDependencies(): void {
2         $.html(`
3             <script>
4                 // Zingchart Library
5                 $('head').append('<script src="./node_modules/zingchart/client/zingchart.min.js">');
6             </script>
7         `);
8     }
9
10    loadDependencies();
```

Out [3]:

We'll keep our first analysis of the data simple, and start by computing some simple statistics over the dataset, on the number of orders, products, products per order, etc. In addition to understanding more about the data, this will also allow us to find some appropriate and relevant criteria on which we could base (and reduce) our dataset for extended itemset analysis; as well as giving us a first glance at pattern redundancy.

Keeping in mind that each row of the `order_products__prior.csv` file is structured as stated above; an interesting approach would be to regroup these objects by both `order_id` and `product_id`, as it would allow us to have a glance over the product distribution through the orders on the first hand, as well as idea of each product's popularity on the other.

To do so, we need to transform the dataset consequently. Thus, let's start by defining `ProductOrder` as the structure of the data outputted by the `.csv` file parser. We will do the same for each data structure of each file composing our dataset:

```
1 /**
2  * Data structure we gather from CSV File.
3  */
4 interface ProductOrder {
5     order_id: string,
6     product_id: string,
7     add_to_cart_order: string,
8     reordered: string
9 }
```

From this point, as we focus on exploiting data from file `order_products__prior.csv`, which contains more than 1 million records; and for the sake of memory usage, we'll try to work on data streams as much as possible, rather than parsing 1-million-elements-cached `Arrays` when it comes to data transformation. We'll thus be using Reactive Programming library `RxJS` in that intent.

Reactive programming is nothing new as it only consists in programming with asynchronous data streams, which languages like JS are basically all about. `RxJS` yet provides us with an amazing and complete approach -as well as a great toolbox of functions- to combine, create and filter such streams easily.

Grouping by orders

Let us define a function allowing us to group `ProductOrder` objects by any key of the `ProductOrder` interface. We'll use in that intent `RxJS`'s `groupBy` method, which basically groups the items emitted by an `Observable` (a.k.a. stream; in our case, the `ReadStream` of the considered file) according to a specified criterion (in our case, either the `product_id` or the `order_id`), and emits these grouped items as `GroupedObservable`s, one `GroupedObservable` per group.

Let's start by defining a `Group<T>` as an object containing an `id` (the grouping criterion basically), as well as an `Array` of whatever item of type `T` we're grouping. This will be one of many products of our function:

```
1 interface Group<T> {
2     id: string,
3     items: Array<T>
4 }
```

We'll also forge ahead (keeping pattern mining in mind) by allowing this method to `.map()` the `ProductOrder` objects to whatever we want (its `id`, `product_name` ...) depending on our need.

Such a function is given below:

```

In [4]: 1 import { Observable } from 'rxjs/Observable';
2 import 'rxjs/add/operator/finally';
3 import 'rxjs/add/operator/groupBy';
4 import * as RxNode from 'rx-node';
5
6 /**
7  * Function returns an Observable of `ProductOrder` group by a defined criterion. You may map the parsed `ProductOrder` to whatever value you want.
8  */
9 function _readAndGroupBy<T>( key: keyof ProductOrder, map: (val: ProductOrder) => T ): Rx.Observable<Group<T>> {
10   /**
11    * 'csv-parse' Parser, columns options groups each row by column in an object.
12    */
13   let parser: Parser = new Parser({
14     delimiter: ',',
15     columns: true
16   });
17
18   // Turning native stream into Observable
19   return RxNode.fromStream( fs.createReadStream(`instacart_basket_data/order_products__train.csv`).pipe(parser) )
20     // Grouping objects by order
21     .groupBy( (data: ProductOrder) => data[key] )
22     // At this point, we basically have an Observable by group. Thus we need to flatten that.
23     .flatMap( (group: Rx.GroupedObservable<string, ProductOrder>) => {
24       return group
25         // Formatting the data
26         .map(map)
27         // And flattening the Observable array.
28         .reduce( (concat: Group<T>, current: T) => {
29           concat.items.push(current);
30           return concat;
31         }, {
32           id: group.key,
33           items: []
34         })
35     });
36 }

```

```

Out[4]: {}

```

Let us group the `ProductOrder` by their `order_id`. `ProductOrders` will be characterized by their `product_id` (We'll thus trivially have an list of Orders (or itemsets), as Arrays of `product_id`s).

Function above will return us with all the processed groups. We'll compute some basic statistics on these from there, such as:

- The number of orders (number of groups);
- The minimum number of product in an order (minimum of arrays length);
- The maximum number of product in an order (maximum of arrays length);
- The average product number per order (average array length);
- The number of records in the `order_products__prior.csv` file (sum of arrays length);

In [5]:

```
1 function statsOnOrders(): void {
2   $$async$$ = true;
3   console.log('Gathering data, this might take a while...');
4
5   /**
6    * All the groups.
7    */
8   let groups: Array<Group<string>> = [];
9
10  /**
11   * Product per order for Chart
12   */
13  let productsPerOrder: Array<number> = [];
14
15  let stats: any = {
16    max: 0,
17    min: Infinity,
18    sum: 0
19  }
20
21  /**
22   * Reads the file and groups `ProductOrders as intended`
23   */
24
25  _readAndGroupBy<string>('order_id', (productOrder: ProductOrder) => productOrder.product_id )
26    // Once all groups are loaded, displaying them.
27    .finally( () => {
28      console.log(`Maximum number of ProductOrders: ${stats.max}`);
29      console.log(`Minimum number of ProductOrders: ${stats.min}`);
30      console.log(`Average number of ProductOrders: ${stats.sum / groups.length}`);
31      console.log(`Total number of ProductOrders: ${stats.sum}`);
32      console.log(`Number of itemsets: ${groups.length}`);
33
34      let chartOptions: string = JSON.stringify({
35        id: 'productPerOrderChart',
36        data: {
37          type: 'bar',
38          title: {
39            text: 'Product per order sorted ascendingly'
40          },
41          scaleY: {
42            label: {
43              text: "Number of products"
44            },
45            item: {
46              fontSize: 10
47            }
48          },
49          scaleX: {
50            label: {
51              text: "Order id"
52            },
53            item: {
54              fontSize: 10
55            }
56          },
57          series: [{ values: productsPerOrder.sort( (a: number, b: number) => a - b) }]
58        }
59      });
60
61      $$html(`
62        <div id="productPerOrderChart"></div>
63        <script>
64          zingchart.render(${chartOptions});
65        </script>
66      `);
67
68      $$done$$();
69    })
70    // Note that this behaviour (induced by the flatMap of readAndGroupBy) makes everything pretty much blocking again.
71    .subscribe( (group: Group<string>) => {
72      // Computing some basic stats on the fly
73      stats.max = Math.max(group.items.length, stats.max);
74      stats.min = Math.min(group.items.length, stats.min);
75      stats.sum += group.items.length;
76
77      // Pushing group to groups.
78      groups.push(group);
79
80      // Chart display data
81      productsPerOrder.push(group.items.length);
82    });
83  }
84
85  statsOnOrders();
```

Gathering data, this might take a while...
Maximum number of ProductOrders: 80
Minimum number of ProductOrders: 1
Average number of ProductOrders: 10.552759338155157
Total number of ProductOrders: 1384617
Number of itemsets: 131209

Out[5]:

Grouping the records by their `order_id` enlights us of the following information:

Number of orders	Minimum product number per order	Maximum product number per order	Average product number per order	Total number of records
131,209	1	80	10.6	1,384,617

Eventhough the maximum number of product per order is 80, the graph above shows that approximatly 85% of all the orders only contains between 1 and 20 orders.

Grouping by product

Creating itemsets (*Groups*) of Orders, based on the `product_id` of `ProductOrders` may also be of interest to us, as it allow us "feel" a product "popularity" by counting the number of orders it appears in. This is a pretty good deal in regards to pattern mining, as:

- a "frequent" product will be more likely to appear in frequent itemsets;
- its number of appearance in the collection is, by definition, the maximum support over the dataset. Considering a product A being the most popular in a dataset such as ours, the itemset { A } will trivially be the absolute, most frequent itemset to be find in the entire dataset;
- if a product is too frequent, it may be of interest to ignore it, as the itemsets to be find may not be revealant enough.

Code is basically the same as before, thus won't be included in the notebook. Sources are however still available and results may be parsed from `custom_data/product_id__order_number.csv` .

Grouping the records by their `product_id` gives us the following results:

Number of products	Minimum order number per product	Maximum order number per product	Average order number per product	Total number of records
391,23	1	18,726	35.39	1,384,617

Let us determine the most popular products, by joining the retrieved data with table `products.csv` :

In [6]:

```
1 import * as CSVParser from './dist/class/csv-parser.class.js';
2 import * as join from './dist/function/join.function.js'
3
4 function joinResultsAndShowPopular(): void {
5     $$async$$ = true;
6     console.log('Gathering data, this might take a while...');
7
8     // Data to display.
9     let chartData: Array<any> = [];
10
11     // Loading the 'product.csv' table.
12     new CSVParser.CSVParser<Product>(`instacart_basket_data/products.csv`).loadAll()
13     .then( (products: Product[]) => {
14         console.log('Finished loading products');
15
16         // Now generating group of items based on the product_id
17         new CSVParser.CSVParser<ProductOrder>(`instacart_basket_data/order_products__train.csv`)
18         // Grouping items by product_id, and mapping every item composing these itemsets to their order_id.
19         .generateItemsets<string>('product_id', (productOrder: ProductOrder) => productOrder.order_id )
20         // Once execution is complete, creating chart.
21         .finally( () => {
22             console.log('Most popular products :');
23
24             chartData = chartData
25                 // Sorting desc.
26                 .sort( (a: any, b: any) => b.number - a.number )
27                 // Keeping 10 best products
28                 .splice(0,10);
29
30             let chartOptions: string = JSON.stringify({
31                 id: 'popularProductsChart',
32                 data: {
33                     type: 'bar',
34                     title: {
35                         text: 'The 10 most popular products (number of orders)'
36                     },
37                     scaleY: {
38                         label: {
39                             text: 'Times ordered'
40                         },
41                         item: {
42                             fontSize: 10
43                         }
44                     },
45                     scaleX:{
46                         values: chartData.map( (product: any) => product.name ),
47                         item: {
48                             'font-size':'6px'
49                         }
50                     },
51                     series: [{ values: chartData.map( (product: any) => product.number ) }]
52                 }
53             });
54
55             $$html(`
56                 <div id="popularProductsChart"></div>
57                 <script>
58                     zingchart.render({chartOptions});
59                 </script>
60             `);
61
62             $$done$$();
63
64         })
65         // Once a group is ready, pushing chart data to the proper array.
66         .subscribe( (group: Group<ProductOrder,string>) => {
67             chartData.push({
68                 name: join.join<Product>(products, 'product_id', group.id, 'product_name'),
69                 number: group.items.length
70             });
71         })
72     }, (e) => {});
73 }
74
75 joinResultsAndShowPopular();
```

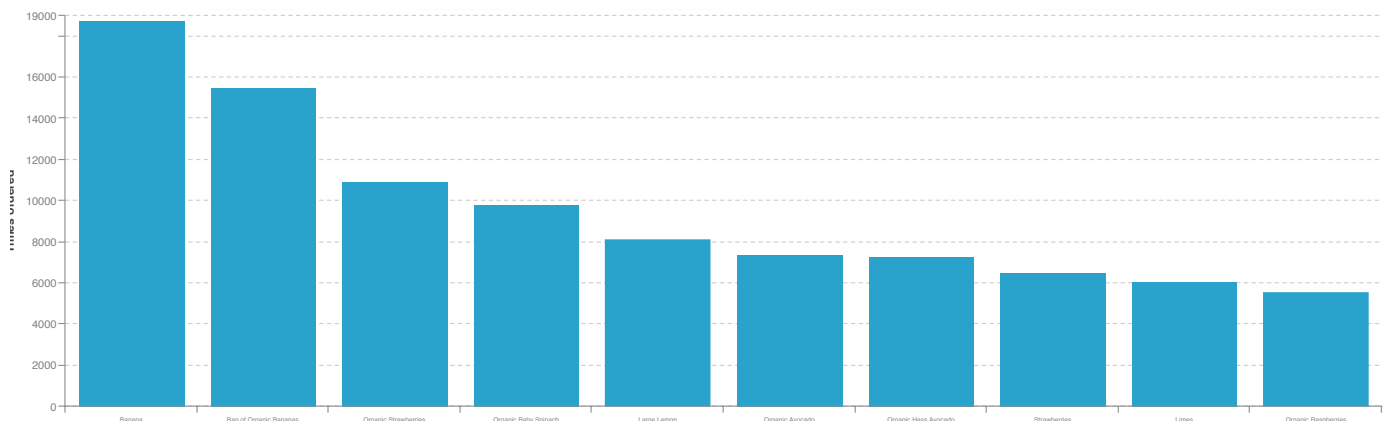
Gathering data, this might take a while...

Finished loading products

Most popular products :

Out[6]:

The 10 most popular products (number of orders)



Graph above show the 10 most popular products, based on the number of time they've been ordered; **Banana** being the most popular product, with 18,726 orders. From a pattern miner point of view, these products have more chance to be part of multiples itemsets, being that all singleton-itemsets of each of these products are the most frequents.

We can also see that many popular items fall into the same category of product (Banana & Bag of Organic Bananas / Organic Avocado & Organic Hass Avocado). If we had some time, regrouping these products would have been a good idea, as more interesting patterns may appear: Nobody would ever buy both Strawberries and Organic Strawberries in the same transaction, yet itemsets <{Banana , Strawberries}> and <{Bag of Organic Bananas , Rasberries}> **could** be considered as one (if you don't have interest in your consumer organic habits).

Grouping products per aisle (and exploring aisle-to-aisle itemsets) may be a less-time consuming alternative if the need of reducing the dataset becomes apparent: Results would appear as "Someone buying fruits have more chance buying vegetables"; though we obviously loose information on the products bought.

In the meantime, It is worth noticing that **46 products** including 100% Black Cherry & Concord Grape Juice , Breaded Popcorn Turkey Dogs or Lip Balm , are the least popular with only 1 order. Filtering less popular products may also be a solution to reducing the dataset; yet I guess this would not change a lot as these elements would still pruned pretty fast by data-mining algorithms when mining for patterns.

Frequent item sets

Knowing a little bit more about our data, we'll now move to mine and gather **frequent** itemsets from our training dataset (`order_products_train.csv`). In other words, we will focus on finding regularities in the shopping behavior of customers, based on what `product` `orders` (a.k.a. transactions) are composed of.

Vocabulary

Let us define some basic vocabulary we will be tremendously be using:

- Let I be a set of items. An itemset is a subset of I .
- Let D be a transaction database such that each transaction is an itemset.
- Frequency of an itemset is the number of transactions including the itemset.
- For a given support, an itemset is said to be frequent if its frequency is no less than the support.

SPMF Open-Source Library

Many (if not all) algorithms tested upon this point will be implementations from SPMF library (<http://www.philippe-fournier-viger.com/spmf/> (<http://www.philippe-fournier-viger.com/spmf/>)); an open-source data mining library written in Java, specialized in pattern mining.

A NodeJS JS/TypeScript wrapper, `class SPMF` for SPMF (spawning Java commands) has been written in order to execute these, while still exploiting the results with Typescript. Sources may be find in file `./src/class/spmf.ts` . Examples of its use are shown below.

Apriori

Keeping our goal in mind, we'll start by using the **Apriori Algorithm**. The object of this algorithm is to identify association between different sets of data, and to find out patterns in data using a "bottom up approach", as frequent subsets are extended one item at a time; and is based on a property called the anti-monotone property: It states that if an item set is not frequent, then none of its supersets can be frequent. As a result, the list of potential frequent itemset gets smaller as mining progresses [1].

Though the Apriori Algorithm is easy to implement and understand, it is far from being performant the reasons stated above.

Dataset formatting

Upon this point we will be using SPMF library's Java implementation of Apriori (through command lines), in order to mine frequent item sets from our dataset. This implementation needs the input data to be formatted as such:

- An item is represented by a positive integer.
- A transaction is a line in the file.
- In each line (transaction), items are separated by a single space.
- It is assumed that all items within a same transaction (line) are sorted according to a total order and that no item can appear twice within the same line.

For example, an input file is defined as follows:

```
1 2 5
2 7
```

Transforming this dataset into this format is pretty straightforward using the code we already wrote previously (same code, but encapsulated in a `CSVParser` class). To save us some time, already formatted data will be available to parse (file `custom_data/formatted-itemsets.csv`). Here is an example of `CSVParser` though, running on a custom 6-line dataset:

```
In [7]: 1 // Our custom parser class, looks like Jupyter/iTypescript doesn't like ES6 imports of custom Typescript classes.
2 import * as CSVParser from './dist/class/csv-parser.class.js';
3 import { ProductOrder } from './src/interface/product-order.interface';
4 import { Group } from './src/interface/group.interface';
5
6 function formatData(): void {
7     $$async$$ = true;
8     console.log('Generating transactions, this might take a while...');
9
10    let groups: Group<ProductOrder,number>[] = [];
11
12    new CSVParser.CSVParser<ProductOrder>(`custom_data/small_order_products__train.csv`)
13    // Grouping items by order_id, and mapping every item composing these itemsets to their product_id.
14    .generateItemsets<number>('order_id', (productOrder: ProductOrder) => Number.parseInt(productOrder.product_id))
15    // Once execution is complete, writing the formatted dataset into a proper file.
16    .finally( () => {
17        // Writing number of product per order_id in a new file : The array of already formatted rows is joined by a return carriage character.
18        console.log('Example of output:');
19        console.log(groups.join('\r\n'));
20        $$done$$();
21    })
22    // On group reception, formatting the items composing is as a ROW (joined by plain space character), and pushing it the output array.
23    .subscribe((group: Group<ProductOrder,number>) => groups.push(group.items.sort( (a: number, b: number) => a - b ).join(' ')))
24 }
25
26 formatData();
27
```

Generating transactions, this might take a while...

Example of output:

```
1
101 1002
10 100 2000
```

Parsing the results

This Apriori implementation enlights us with multiple information other than the frequent itemsets themselves, as it outputs multiple results through both `stdout` , including:

- The number of candidates;
- The maximum size of candidates the algorithm stopped at;
- The frequent itemsets count;
- Maximum memory usage;
- Total execution time.

Mined frequent itemsets may be find in the specified output file, formatted as below:

```
...A #SUP: B

-- EXAMPLES OF OUTPUT --
49628 #SUP: 186
49683 #SUP: 2413
13176 21137 #SUP: 164
13176 21903 #SUP: 175
```

Where `...A` represents the spread of `item_ids` composing the frequent itemset; and `B` the support of this itemset (in term of number of occurency).

Eventhough the output format is not rigorously `.csv` friendly, we can still easily parse the output file with our TypeScript application with no additionnal dependency or need of code; using a little hack, defining the string `#SUP:` as a `.csv` separator. This will result in us parsing an Array of 2-item rows:


```
[ ...item_ids, support ]
```

This is own our homemade SPMF wrapper inherently works.

EDIT: V2 of SPMF Wrapper uses Regular Expressions, see below

Last step will trivially be to parse those item_ids and joining these results with the products.csv table for enhanced visualisation.

Maximum support

In the extend of running the Apriori Algorithm, we need to define a proper minimum support value, as defined above. This is a crucial step, as setting the support to an exagerated value would result no itemsets to be find; while setting it to a too low value will inherently result in a very long execution time.

As seen previously, Banana is the most frequent item in out dataset, being recorded in 18,726 of the 131,209 transactions i.e. 14.3% of all the orders. This number also represents the maximum support of any itemset in our dataset, as being the support of itemset composed of the Banana singleton.

First runs

Let us run Apriori a few times, with a different minimum support threshold, in order to compare the execution time as well as the number of return candidates. One small drawback of this implementation being that it returns single-item candidates (which are not of interest at all), we'll also specify the number on multiple-items candidates found.

```
In [8]: 1 import * as SPMF from './dist/class/spmf.class.js';
2 import { SPMFResults, ItemSet, Rule } from './src/class/spmf.class';
3
4 function aprioriRun(): void {
5     $$async$$ = true;
6     console.log('Mining patterns, this might take a while...');
7     // Our custom SPMF wrapper
8     new SPMF.SPMF('Apriori')
9         // Loading from file
10        .fromFile('custom_data/formatted_itemsets.txt`)
11        // Executes Apriori with 5% support
12        .exec<ItemSet>(5)
13        // Listening for results
14        .subscribe((results: SPMFResults<ItemSet>) => {
15            // Wrapper returns both the stats...
16            console.log('Stats:');
17            console.log(results.stats);
18            // ... and the frequent itemsets. Showing the first two itemsets:
19            console.log('First two frequent itemsets:');
20            console.log(results.output.slice(0,2));
21
22            $$done$$();
23        });
24    }
25
26 aprioriRun();
```

Mining patterns, this might take a while...

Stats:
{ candidates: 28,
 executionTime: 456,
 memory: 59.766456604003906 }
First two frequent itemsets:
[{ support: 15480, items: ['13176'] },
 { support: 10894, items: ['21137'] }]

Minimum support	Number of candidates	Frequent itemsets count	Multiple-items itemsets	Execution time
13%	1 (Banana, as expected)	1	0	422 ms
7%	10	4	0	423 ms
3%	153	17	0	748 ms
1%	5,467	120	14	8.2 s
0.5%	33,042	364	108	48.9 s
0.3%	210,892	1125	478	4 mn 58 s

As expected, the speed of Apriori gets in the way pretty fast when it comes to parsing a dataset of such size, and with a maximum support so low. In order to find more relevant itemsets, we may have to transform the dataset in order reduce the overall total number of candidates.

Our first frequent itemsets

Results of Apriori runs may be find in folder /output , labelled as such: output_approri_#SUPPORT.txt . Lets have a look to multiple-itemed frequent itemsets Apriori has gathered, for a minimum support of 0.3%. We want to display our itemsets through a Chord Chart (https://www.zingchart.com/docs/chart-types/chord-diagrams/) using zingchart .

With this intent, we need to format the data, zing-chart -way, as such:

```
{
  name: 'Banana',
  values '0, 1'
}, {
  name: 'Avocado',
  values '1, 0'
}
```

... where 1 is the support of the itemset <{ 'Banana' , 'Avocado' }> .

Such a transformation is quite tedious to explain, but I believe it's easy to understand line-by-line, though the code below (if not, I do apologize for any headache I may have cause).

Note that a class representing a Chord node has been created, in which you can easily manage relations of a serie of nodes. Code is available in the source codes ./src/class/chord-serie.class.ts .

```

In [9]: 1 import * as SPMF from './dist/class/spmf.class.js';
2 import * as CSVParser from './dist/class/csv-parser.class.js';
3 import * as ChordSerie from './dist/class/chord-serie.class.js';
4 import * as join from './dist/function/join.function.js'
5
6 function parseAprioriResults(): void {
7     $$async$$ = true;
8     console.log('Parsing Apriori results, this might take a while...');
9
10    // Loading the 'product.csv' table.
11    new CSVParser.CSVParser<Product>(`instacart_basket_data/products.csv`).loadAll()
12    .then( (products: Product[]) => {
13        console.log('Finished loading products');
14
15        // Our custom SPMF wrapper
16        new SPMF.SPMF('Apriori')
17        .loadResultsFromFile<ItemSet>('output/output_apriori_03.txt')
18        .subscribe((results: SPMFResults<ItemSet>) => {
19            results.output = results.output
20            // Only 2-item itemsets
21            .filter( (itemset: ItemSet) => itemset.items.length == 2 )
22            // Sorted, highest support first
23            .sort( (a: ItemSet, b: ItemSet) => b.support - a.support)
24            // Only the 20 most frequent itemsets
25            .slice(0,20);
26
27            // Gathering all items composing the itemsets.
28            let uniqueItems: string[] = [].concat.apply([], results.output.map( (itemset: ItemSet) => itemset.items ) )
29            // Filtering unique items.
30            .filter( (itemId: string, index: number, that: string[]) => that.indexOf(itemId) === index );
31
32            // Mapping unique ids to a ChordSerie node.
33            let series: ChordSerie.ChordSerie[] = uniqueItems.map( (uniqueItemId: string) => new ChordSerie.ChordSerie(uniqueItemId) );
34
35            // Now we need to populate each node, lets start by noticing all nodes of the size of the array, so they can create inner arrays of such size.
36            series.forEach( (serie: ChordSerie.ChordSerie) => serie.setSeriesNumber( series.length ) )
37
38            // Then for each itemset, we need to populate each node
39            results.output.forEach( (itemset: ItemSet) => {
40                // So for each items composing the itemset...
41                itemset.items.forEach( (itemId: string, index: number, that: string[]) => {
42                    // Copying array just to be sure, shouldn't be a problem though.
43                    let copy: string[] = that.slice(0);
44                    // Deleting the item id we're currently parsing
45                    copy.splice(index,1);
46                    // Creating an array of all the nodes related to this item id
47                    let indexesOfRelatedNodes: number[] = copy.map( (id: string) => uniqueItems.indexOf(id) )
48
49                    // Noticing the node of its relations
50                    series[uniqueItems.indexOf(itemId)].addRelation(itemset.support, indexesOfRelatedNodes);
51                })
52            });
53
54            let chartOptions: string = JSON.stringify({
55                id: 'aprioriItemSetChart',
56                data: {
57                    type: 'chord',
58                    title: {
59                        text: 'Top 20 most frequent 2-items itemsets'
60                    },
61                    legend: {
62                        overflow: 'scroll',
63                        layout: 'vertical',
64                        width: 200
65                    },
66                    options: {
67                        'angle-padding': 1,
68                        'band-width': 10,
69                        'band-space': 5,
70                        'radius': 190,
71                        style: {
72                            band: { },
73                            //chord: { 'border-width': 0, flat: true },
74                            tick: { visible: false },
75                            item: { visible: false },
76                            label: { visible: false }
77                        },
78                    },
79                    tooltip: {},
80                    // Populating with product names
81                    series: series.map( (serie: ChordSerie.ChordSerie[]) => serie.getSerie( id) => join.join<Product>(products, 'product_id', id, 'product
82                )
83            });
84
85            $$html(`
86            <div id="aprioriItemSetChart"></div>
87            <script>
88                zingchart.render({{chartOptions}});
89            </script>
90            `);
91
92            $$done$$();
93        });
94    });
95 }
96
97 parseAprioriResults();

```

Parsing Apriori results, this might take a while...
 Finished loading products

Out[9]:

Top 20 most frequent 2-items itemsets



Chart above has been populated with the 20 most popular 2-item itemsets mined by Apriori with a minimum support of 0.3%; the chords of the chart representing the itemsets, of which you can have the support by hovering the graph with your mouse.

It gives a visual idea of which items are often bought together (featuring the itemsets and their support through the chords), as well as the most popular items within these itemsets.

It is worth noticing that only 13 products share all top 20 most frequent 2-itemed itemset sets.

Fruits and vegetables appears to be over represented within the frequent itemsets, which makes sense, as fruits and vegetables are very often the products we buy altogether when going to the supermarket. Yet within the aisle, from a marketing point of view, it could be interesting to organise the products such as the most popular items are apart each other. One may be tempted to buy Apples on his way to find Strawberries after weighting his freshly bought Bananas; most frequent 2-item itemset being `<{"Banana", "Organic Strawberries"}>`.

In addition, 3 itemsets have also be mined, and it can be of interest to take a glance at them:

```
In [10]: 1 import * as SPMF from './dist/class/spmf.class.js';
2 import * as CSVParser from './dist/class/csv-parser.class.js';
3 import * as join from './dist/function/join.function.js'
4
5 function threeItemSets(): void {
6     $$async$$ = true;
7     console.log('Parsing Apriori results, this might take a while...');
8
9     // Loading the 'product.csv' table.
10    new CSVParser.CSVParser<Product>(`instacart_basket_data/products.csv`).loadAll()
11    .then( (products: Product[]) => {
12        console.log('Finished loading products');
13
14        // Our custom SPMF wrapper
15        new SPMF.SPMF('Apriori')
16        .loadResultsFromFile<ItemSet>('output/output_apriori_03.txt')
17        .subscribe((results: SPMFResults<ItemSet>) => {
18            let rows: string = results.output
19            // Only 2-item itemsets
20            .filter( (itemset: ItemSet) => itemset.items.length > 2 )
21            // Sorted, highest support first
22            .sort( (a: ItemSet, b: ItemSet) => b.support - a.support)
23            // Creating HTML Rows
24            .map( (itemset: ItemSet) => {
25                let items: string[] = itemset.items.map( (id: string[]) => {
26                    return `<td>${join.join<Product>(products, 'product_id', id, 'product_name')}`
27                });
28                return `<tr>${items.join('')}<td>${itemset.support}</td></tr>`;
29            })
30            // Only the 10 most frequent itemsets
31            .slice(0,10);
32
33            $$$.html(`
34                <h3>Top 10 3-items frequent itemsets</h3>
35                <table style="width:100%">
36                    <tr>
37                        <th>Item 1</th>
38                        <th>Item 2</th>
39                        <th>Item 3</th>
40                        <th>Support</th>
41                    </tr>
42                    ${rows.join('')}
43                </table>
44            `);
45
46            $$done$$();
47        });
48    });
49 }
50
51 threeItemSets();
52
Parsing Apriori results, this might take a while...
Finished loading products
```

Out[10]: **Top 10 3-items frequent itemsets**

Item 1	Item 2	Item 3	Support
Bag of Organic Bananas	Organic Strawberries	Organic Hass Avocado	710
Bag of Organic Bananas	Organic Strawberries	Organic Raspberries	649
Bag of Organic Bananas	Organic Strawberries	Organic Baby Spinach	587
Bag of Organic Bananas	Organic Raspberries	Organic Hass Avocado	531
Bag of Organic Bananas	Organic Baby Spinach	Organic Hass Avocado	497
Organic Baby Spinach	Banana	Organic Avocado	484
Banana	Large Lemon	Organic Avocado	477
Banana	Limes	Large Lemon	452
Bag of Organic Bananas	Organic Strawberries	Organic Cucumber	424
Limes	Large Lemon	Organic Avocado	389

It's pretty fun to see that the top 5 (6 out of top 10) of most commonly bought-together products are all itemsets of **organic** products. This really represents the new trend and consumer habit of buying only organic products: One having in mind to buy organic products will only buy organic products. This behaviour result in creating interesting frequent datasets!

These frequent itemsets clearly points towards associations rules stating If one buys 2 organic products, he will buy 1 other organic products with a very high confidence. Let's keep this idea in mind for later.

LCM Algorithm

LCM is an algorithm known to be the fastest for mining frequent **closed** itemsets; knowing that an itemset X is closed in a dataset S if no proper super-itemset Y that has the same support count as X in S exists. Main reason of mining closed frequent datasets over all datasets is that such a set is usually much smaller than the set of frequent itemsets, even though no information is lost, as the entire set can be regenerated from the closed set[2].

We will be using SPMF library's JAVA implementation of LCM, in order to mine frequent closed itemsets from our dataset; An advantage being that this implementation inputs the same file format as the Apriori Algorithm of the same library, thus, there is no need to format our dataset in a different manner.

Running the algorithm

It would be interesting to compare the performances and results of the LCM and Apriori algorithms; thus, let's run LCM over the same dataset again, and with the same support. You can also try it by yourself by executing the code below:

```
In [11]: 1 import * as SPMF from './dist/class/spmf.class.js';
2
3 function lcmRun(): void {
4   $$async$$ = true;
5   console.log('Mining patterns, this might take a while...');
6   // Our custom SPMF wrapper
7   new SPMF.SPMF('LCM')
8     // Loading from file
9     .fromFile(`custom_data/formatted_itemsets.txt`)
10    // Executes LCM with 5% support
11    .exec<ItemSet>(5)
12    // Listening for results
13    .subscribe((results: SPMFResults<ItemSet>) => {
14      // Wrapper returns both the stats...
15      console.log('Stats:');
16      console.log(results.stats);
17      // ... and the frequent itemsets. Showing the first two itemsets:
18      console.log('First two frequent itemsets:');
19      console.log(results.output.slice(0,2));
20
21      $$done$$();
22    });
23 }
24
25 lcmRun();
```

Mining patterns, this might take a while...
Stats:
{ candidates: undefined, executionTime: 409, memory: undefined }
First two frequent itemsets:
[{ support: 15480, items: ['13176'] },
 { support: 10894, items: ['21137'] }]

Minimum support	Frequent closed itemsets count	Multiple-items itemsets	Execution time
13%	1	0	204 ms
7%	4	0	275 ms
3%	17	0	328 ms
1%	120	16	1.2 s
0.5%	364	108	2.3 s
0.3%	1125	339	4.9 s
0.1%	4444		24.4 s

Let's reduce the minimum support even more !

Enhanced performances of LCM over Apriori allows us to reduce the minimum even more while mining patterns on our dataset with decent executing times (See table above, with an execution for a minimum support of .1%). This allowed us to find some quite revealant 4-items itemsets, with still a decent support value:

```
In [12]: 1 import * as SPMF from './dist/class/spmf.class.js';
2 import * as CSVParser from './dist/class/csv-parser.class.js';
3 import * as join from './dist/function/join.function.js'
4
5 function lcmItemsets(): void {
6     $$async$$ = true;
7     console.log('Mining patterns, this might take a while...');
8
9     new CSVParser.CSVParser<Product>(`instacart_basket_data/products.csv`).loadAll()
10     .then( (products: Product[]) => {
11         console.log('Finished loading products');
12
13         // Our custom SPMF wrapper
14         new SPMF.SPMF('LCM')
15             // Loading from file
16             //.fromFile(`custom_data/formatted_itemsets.txt`)
17             // Executes LCM with .1% support
18             //.exec<ItemSet>(.1)
19             .loadResultsFromFile<ItemSet>('output/output_lcm_01.txt')
20             // Listening for results
21             .subscribe((results: SPMFResults<ItemSet>) => {
22                 let rows: string = results.output
23                 // Only 2-item itemsets
24                 .filter( (itemset: ItemSet) => itemset.items.length > 3 )
25                 // Sorted, highest support first
26                 .sort( (a: ItemSet, b: ItemSet) => b.support - a.support)
27                 // Creating HTML Rows
28                 .map( (itemset: ItemSet) => {
29                     let items: string[] = itemset.items.map( (id: string[]) => {
30                         return `<td>${join.join<Product>(products, 'product_id', id, 'product_name')}</td>`
31                     });
32
33                     return `<tr>${items.join('')}<td>${itemset.support}</td></tr>`;
34                 })
35                 // Only the 10 most frequent itemsets
36                 .slice(0,10);
37
38                 $$html(`
39                     <h3>4-items frequent itemsets</h3>
40                     <table style="width:100%">
41                         <tr>
42                             <th>Item 1</th>
43                             <th>Item 2</th>
44                             <th>Item 3</th>
45                             <th>Item 4</th>
46                             <th>Support</th>
47                         </tr>
48                         ${rows.join('')}
49                     </table>
50                 `);
51
52                 $$done$$();
53             });
54     });
55 }
56
57 lcmItemsets();
```

Mining patterns, this might take a while...
Finished loading products

Out[12]:

	Item 1	Item 2	Item 3	Item 4	Support
	Bag of Organic Bananas	Organic Strawberries	Organic Raspberries	Organic Hass Avocado	228
	Bag of Organic Bananas	Organic Strawberries	Organic Baby Spinach	Organic Hass Avocado	163
	Banana	Limes	Large Lemon	Organic Avocado	146
	Bag of Organic Bananas	Organic Strawberries	Organic Cucumber	Organic Hass Avocado	140
	Organic Baby Spinach	Banana	Large Lemon	Organic Avocado	134

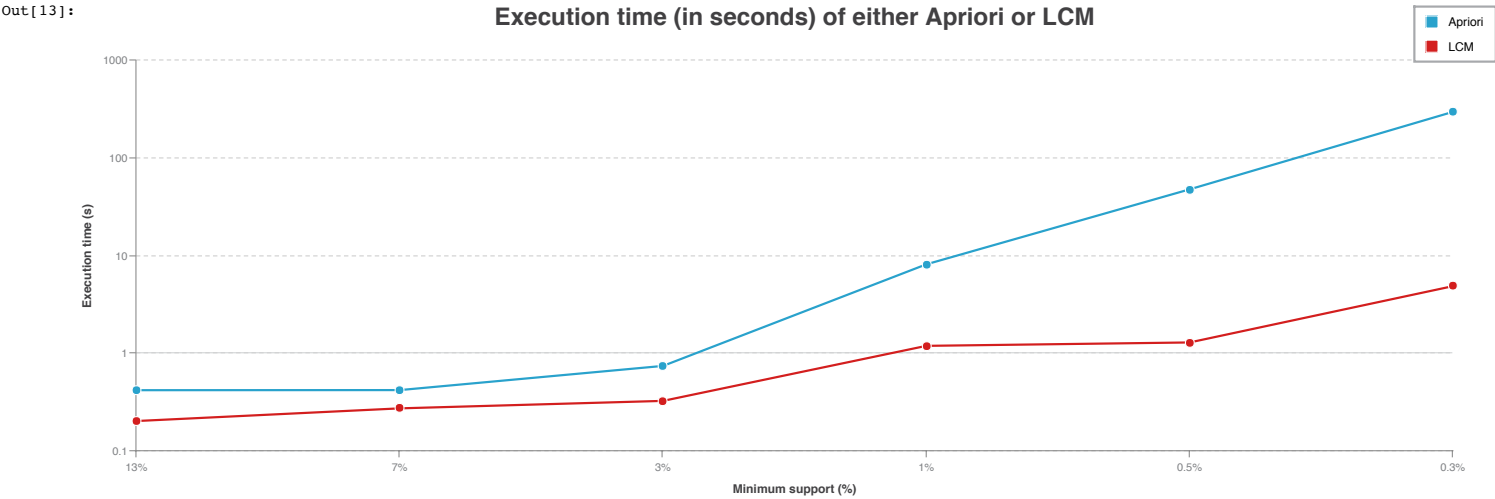
These results points in the same way than our previous observations: organic products tends to be bought altogether.

Quick comparison of LCM and Apriori

Let's have a quick look over the behaviours and execution of both tested algorithms:

Minimum support	Apriori execution time	LCM execution time
13%	422 ms	204 ms
7%	423 ms	275 ms
3%	748 ms	328 ms
1%	8.2 s	1.2 s
0.5%	48.9 s	1.3 s
0.3%	4 mn 58s	4.9 s

```
In [13]: 1 function executionTimeChart(): void {
2
3
4     let chartOptions: string = JSON.stringify({
5         id: 'executionTimeChart',
6         data: {
7             type: 'line',
8             scaleY: {
9                 progression: 'log',
10                logBase: 10,
11                label: {
12                    text: "Execution time (s)"
13                },
14                item: {
15                    fontSize: 10
16                },
17            },
18            scaleX: {
19                values: ['13%', '7%', '3%', '1%', '0.5%', '0.3%'],
20                label: {
21                    text: "Minimum support (%)"
22                },
23                item: {
24                    fontSize: 10
25                },
26            },
27            title: {
28                text: 'Execution time (in seconds) of either Apriori or LCM'
29            },
30            legend: {
31            },
32            series: [
33                { values: [.422,.423,.748,8.2,48,298], text: 'Apriori' },
34                { values: [.204,.275,.328,1.2,1.3,4.9], text: 'LCM' }
35            ]
36        }
37    });
38
39    $$html(`
40        <div id="executionTimeChart"></div>
41        <script>
42            zingchart.render(${chartOptions});
43        </script>
44    `);
45
46    };
47
48    executionTimeChart();
```



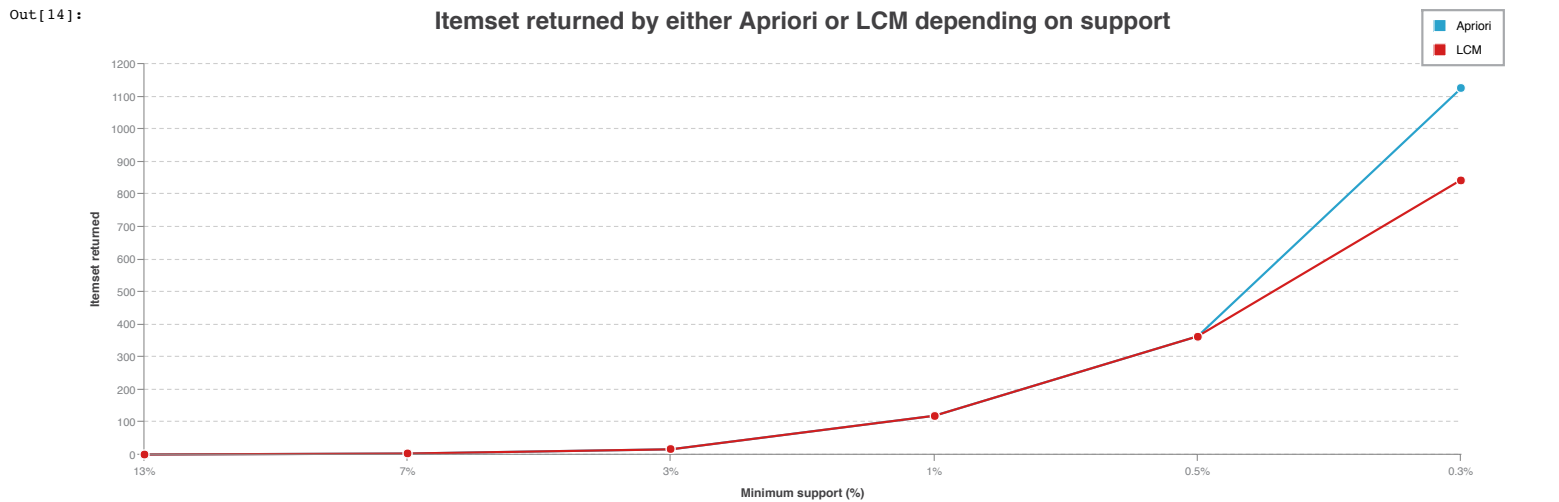
In terms of performance, LCM crushed Apriori when it comes to lower minimum support. This is anything but a surprise, as Apriori is nowadays outdated, as faster and more memory efficient algorithms for mining all datasets have been proposed since, like FPGrowth which we will use later. In the meantime, LCM is known to be one of the most efficient algorithm for mining closed frequent itemsets.

Minimum support	Apriori itemsets count	LCM itemsets count
13%	1	1
7%	4	4
3%	17	17
1%	120	120
0.5%	364	364
0.3%	1125	844

```

In [14]: 1 function itemSetChart(): void {
2
3     let chartOptions: string = JSON.stringify({
4         id: 'itemSetChart',
5         data: {
6             type: 'line',
7             scaleY: {
8                 label: {
9                     text: "Itemset returned"
10                },
11                item: {
12                    fontSize: 10
13                },
14            },
15            scaleX: {
16                values: ['13%', '7%', '3%', '1%', '0.5%', '0.3%'],
17                label: {
18                    text: "Minimum support (%)"
19                },
20                item: {
21                    fontSize: 10
22                }
23            },
24            title: {
25                text: 'Itemset returned by either Apriori or LCM depending on support'
26            },
27            legend: {
28            },
29            series: [
30                { values: [1,4,17,120,364,1125], text: 'Apriori' },
31                { values: [1,4,17,120,364,844], text: 'LCM' }
32            ]
33        }
34    });
35
36    $$html(`
37        <div id="itemSetChart"></div>
38        <script>
39            zingchart.render({chartOptions});
40        </script>
41    `);
42
43
44 };
45
46 itemSetChart();

```



In terms of returned itemsets, LCM globally returns less itemsets than Apriori. This was to be expected given that those Algorithms don't output the same set of itemsets; Apriori returning all itemsets while LCM only returns **closed itemsets**. However it may be of interest to notice that both algorithms return the **same** result upon a certain minimum support given the fact that itemsets returned are closed no matter what algorithm is used.

Frequent item sets association rules

Once patterns have been mined, we can also determine association rules, as each item may have several relations with others. These relations thus indirect relationships between the items composing the transactions; expressed with a certain **confidence**, as an example:

Diapers implies beer with 65% confidence

FPGrowth association mining

Association rules mining generally consist in two steps: the first step being to discover frequent itemsets (With Apriori, LCM or any other algorithm as we did previously); and the second, to generate rules by using these frequent itemsets.

We'll be using in that order yet another all-in-one Algorithm from the SPMF Library, the `FPGrowth_association_rules` algorithm, which basically mine item sets using a faster, more memory efficient algorithm than Apriori: `FPGrowth` ; before generating the association rules.

Algorithm takes 2 input values:

- The minimum support of mined itemsets, as before;
- The minimum confidence of resulting association rules;

Parsing the results

Parsing the results may defer a bit, as the output file now indicates rules instead of plain itemsets, the support of the association, as well as the confidence of the rule.

Example of output:

```

1 ==> 2 4 5 #SUP: 3 #CONF: 0.75
5 6 ==> 1 2 4 #SUP: 3 #CONF: 0.6
4 7 ==> 1 #SUP: 3 #CONF: 0.75

```

Such a file would result in being parsed the following way with our current code:

```

[ '1 ==> 2 4 5', '3 #CONF: 0.75' ]

```

Nothing to worry about though, as a really hacky to do it would be to `split()` both resulting strings using separators `==>` and `#CONF` . I'd really wish I would have use Regular expressions at this point though...

Late update: Regular expressions

That was I ended up doing eventually, keeping in mind that we will need to parse output files from even more algorithms (sequential pattern mining, etc...). I'm thus using the following Regular expressions to match results as they are gathered from output files; and using even more regular expressions within the matches of expressions to parse for numeric values.

```
Matching itemsets:
/^(\\d* )+#SUP: (\\d*( )*)+\\t*n*r*/g

Matching association rules:
/^(\\d* )+==> (\\d* )+#SUP: (\\d* )+#CONF: (\\d*+\\..*\\d*\\t*n*r*/g
```

Running the algorithm

You can try the algorithm by yourself, by editing and executing the code below:

```
In [15]: 1 import * as SPMF from './dist/class/spmf.class.js';
2
3 function associationRulesMining(): void {
4     $$async$$ = true;
5     console.log('Mining association rules, this might take a while...');
6     // Our custom SPMF wrapper
7     new SPMF.SPMF('FPGrowth_association_rules')
8         // Loading from file
9         .fromFile('custom_data/formatted_itemsets.txt`)
10        .exec<Rule>(0.1,25)
11        // Listening for results
12        .subscribe((results: SPMFResults<Rule>) => {
13            // Wrapper returns both the stats...
14            console.log('Stats:');
15            console.log(results.stats);
16            // ... and the frequent rules. Showing the first two rules:
17            console.log('First two rules mined:');
18            console.log(results.output.slice(0,2));
19
20            $$done$$();
21        });
22 }
23
24 associationRulesMining();
```

Mining association rules, this might take a while...

Stats:

{ candidates: undefined, executionTime: 54, memory: undefined }

First two rules mined:

```
[ { support: 198,
  confidence: 0.3168,
  items: [ '45' ],
  results: [ '24852' ] },
  { support: 136,
  confidence: 0.3919308357348703,
  items: [ '46149' ],
  results: [ '196' ] } ]
```

Association rules & Aisles relations

Itemset association rules are great candidates to Chord Diagram representation. We'll use these mined associations to determine how aisles interact with each other. The idea we have in mind is simple: We will map every item and result (which are `product_ids`) of each rule to their `aisle_id` , in order to express associations rules as such:

If customer buys a product from aisle A (and B) then he will buy a product from aisle C (and D)

From a marketer point of view, we want aisle A (and B) as far as possible from aisle C (and D)... Thus when parsing our newly parsed rules, keeping our example in mind: Each time we encounter id of `Aisle A` , we will increment the value linked to the band (a.k.a. edge) between it and `Aisle C` .

The resulting graph will show the weighted relations between each aisle.

Note that we will not take `confidence` into account. Given the fact that minimal confidence as been set to 25%, we believe confidence value should not be a discriminative factor for mapping aisles (probability of each is high).


```

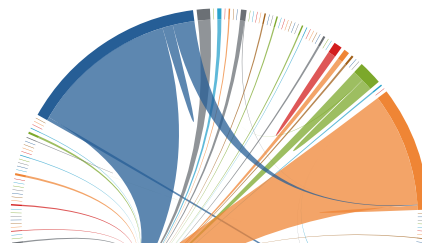
In [16]: 1 import * as SPMF from './dist/class/spmf.class.js';
2 import * as CSVParser from './dist/class/csv-parser.class.js';
3 import * as ChordSerie from './dist/class/chord-serie.class.js';
4 import * as join from './dist/function/join.function.js'
5
6 function parseAssociationRules( ignoreIds: string[] = [] ): void {
7     $$async$$ = true;
8     console.log('Parsing association rules, this might take a while...');
9
10    let products: Product[] = [];
11    let aisles: Aisle[] = []
12
13    // Loading the 'product.csv' table.
14    new CSVParser.CSVParser<Product>(`instacart_basket_data/products.csv`).loadAll()
15    .then( (p: Product[]) => {
16        products = p;
17        console.log('Finished loading products');
18
19        // Loading the 'aisles.csv' table.
20        return new CSVParser.CSVParser<Aisle>(`instacart_basket_data/aisles.csv`).loadAll()
21    })
22    .then( (a: Aisle[]) => {
23        aisles = a;
24        console.log('Finished loading aisles');
25
26        // Our custom SPMF wrapper
27        new SPMF.SPMF('FPGrowth_association_rules')
28        .loadResultsFromFile<Rule>('output/output_fpgrowth_assoc_01_25.txt')
29        .subscribe((results: SPMFResults<Rule>) => {
30            // Gather all the aisle_ids
31            let aislesIds: number[] = aisles.map( (aisle: Aisle) => aisle.aisle_id );
32
33            // Series will be aisles, as planned.
34            let series: ChordSerie.ChordSerie[] = aisles.map( (aisle: Aisle) => new ChordSerie.ChordSerie(aisle.aisle_id).setSeriesNumber(aisles.length) );
35
36            // For each rule...
37            results.output.forEach( (rule: Rule) => {
38                // Map all aisle_ids the rule is pointing towards.
39                let indexesOfResults: number[] = rule.results.map( (id: string) => aislesIds.indexOf( join.join<Product>(products, 'product_id', id, 'aisl
40
41                // For each item composing the rule
42                rule.items.forEach( (id: string) => {
43                    // Get the aisle_id from the product_id
44                    let aisleId: number = join.join<Product>(products, 'product_id', id, 'aisle_id');
45                    // If it's ignored, then meh.
46                    if(ignoreIds.includes(aisleId)) return;
47
48                    // Else notice the corresponding aisle it has a relation with others.
49                    let index: string = aislesIds.indexOf( aisleId );
50                    series[index].addRelation(1,indexesOfResults);
51                });
52            })
53
54            let chartOptions: string = JSON.stringify({
55                id: 'associationRulesChart' + ignoreIds.join('_'),
56                data: {
57                    type: 'chord',
58                    title: {
59                        text: 'Implied relations between aisles',
60                    },
61                    subtitle: {
62                        text: ignoreIds.length ? `Aisles ignored : ${ignoreIds.join(' ')}` : ''
63                    },
64                    legend: {
65                        overflow: 'scroll',
66                        layout: 'vertical',
67                        width: 200
68                    },
69                    options: {
70                        'angle-padding': 1,
71                        'band-width': 10,
72                        'band-space': 0,
73                        'radius': 190,
74                        "style": {
75                            "chord": {
76                                "border-width": 0,
77                                "flat": false
78                            },
79                            "band": { "border-width": 0 },
80                            "tick": { "visible": false },
81                            "item": { "visible": false },
82                            "label": { "visible": false }
83                        }
84                    },
85                    tooltip: {},
86                    // Populating with product names
87                    series: series.map( (serie: ChordSerie.ChordSerie[]) => serie.getSerie( (id) => join.join<Aisle>(aisles, 'aisle_id', id, 'aisle') ) ),
88                }
89            });
90
91            $$html(`
92                <div id="associationRulesChart${ignoreIds.join('_')}"></div>
93                <script>
94                    zingchart.render({${chartOptions}});
95                </script>
96            `);
97
98            $$done$$();
99        });
100    });
101 }
102
103 parseAssociationRules();

```

Parsing association rules, this might take a while...
 Finished loading products
 Finished loading aisles

Out[16]:

Implied relations between aisles



- specialty cheeses
- bulk grains rice dried goods
- energy granola bars
- instant foods
- marinades meat preparation
- other
- packaged meat
- bakery desserts
- pasta sauce
- kitchen supplies
- cold flu allergy

Graph above enlightens us with many interesting information:

- You want to keep fresh vegetables, packed vegetables fruits and fresh fruits aisles as far as possible from each other (which is definitely a difficult task I concede);
- Many rules associate fresh fruits products to products of itself (as we saw previously): You will have to deal inner aisle organisation !

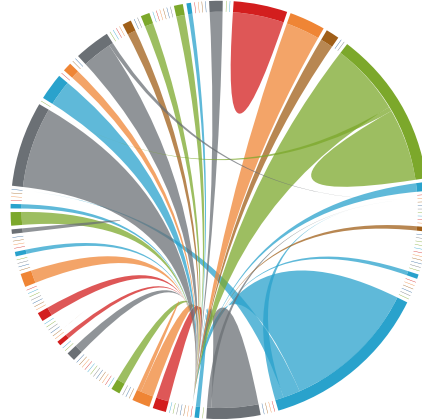
Let's ignore aisles fresh vegetables, packed vegetables fruits and fresh fruits relations and recompile the graph:

```
In [21]: 1 parseAssociationRules(['123','24','83']);
Parsing association rules, this might take a while...
Finished loading products
Finished loading aisles
```

Out[21]:

Implied relations between aisles

Aisles ignored : 123 24 83



- bulk grains rice dried goods
- prepared soups salads
- energy granola bars
- instant foods
- marinades meat preparation
- other
- packaged meat
- bakery desserts
- pasta sauce
- kitchen supplies
- cold flu allergy
- fresh pasta
- prepared meals
- tofu meat alternatives
- packaged seafood
- fresh herbs
- baking ingredients
- bulk dried fruits vegetables
- oils vinegars
- oral hygiene
- packaged cheese

Many outer relations of associations from other aisles still point toward fresh fruit aisle products. *omnibus viis Romam pervenitur*.

Though we can still glance some interesting things happening in yoghurt and sparkling water aisles, as it seems products from these aisles are often bought together ! If it actually definitely sounds silly that people want to buy different kind of sparkling water at the same time.

As a general conclusion to itemset pattern mining... Never have short storage of bananas, and just put the fresh fruits aisle all across the whole store :)

Sequential pattern mining

Sequential pattern mining consists of discovering interesting subsequences in a set of sequences (where, has a comparison, mining frequent itemsets consisted in discovering itemsets in a set of transaction). In the context of our example, sequential pattern mining can be used to find the sequences of items frequently bought by customers, as we acknowledged that transaction data could have been ordered by time, with the proper table joins (orders.csv). This can be useful to understand the behavior of customers to take marketing decisions.

Such a transformation is however pretty heavy, and an already transformed dataset, output/transactions_seq.txt is already provided with this notebook; a 4-columns csv-like file (separators are tabulation characters) as such:

- Column 1: user id;
- Column 2: order number (not an order id, but a number in the sequence of orders);
- Column 3: size of order;
- Column 4: list of products in text format (without spaces), separated by a comma;

In order to mine sequential pattern from our data, we will use the CloSpan Algorithm, from the SPMF library. This algorithm is used for discovering **closed** sequential patterns in sequence databases. As a closed frequent itemset is to a frequent itemset; a closed sequential pattern is a sequential pattern such that it is not strictly included in another pattern having the same support.

Data formatting

The input file format for SPMF CloSpan is defined as follows:

- Each line represents a sequence from a sequence database, and ends with the value -2;
- Each item from a sequence is a positive integer, and items from the same itemset within a sequence are order ascendingly, and separated by single space;
- Each itemset is separated by the value -1;

Example of input from the SPMF Documentation :

```
1 -1 1 2 3 -1 1 3 -1 4 -1 3 6 -1 -2
1 4 -1 3 -1 2 3 -1 1 5 -1 -2
5 6 -1 1 2 -1 4 6 -1 3 -1 2 -1 -2
5 -1 7 -1 1 6 -1 3 -1 2 -1 3 -1 -2
```

Back to our example, the file transactions_sec.txt has yet be found quite unusable, as the 'comma' has been used as a separator to separate items within the itemsets, making it impossible to map the product_names to their actual id.

```
96980 15 10
Asparagus,GreekExtraVirginOliveOil,Jelly,Blackberry,OrganicButterhead(Boston,Butter,Bibb)Lettuce,OrganicIcebergLettuce,RiceWhip,RipeLargePittedOlives
```

In the exemple above, Jelly,Blackberry and OrganicButterhead(Boston,Butter,Bibb)Lettuce are two separate products. Thus total count of products is 8, not 10.

Only way to compute this file has been to map all unique strings to a new integer (in a certain way, to compute our own products.csv (re-joining both orders and order_product is not an easy process). It's a bit sad though as we won't be able to compute more stats on aisles or departments from our results.

Here will be the transformations to apply :

- Group each order per client_id: This will be our sequences;
- Map the product to an integer;
- Order and format the sequence as specified.

Algorithm below executes this exact process, for you to tinker with. Two files have been generated and available for you though.

- `products_transactions_seq.csv` , a map of products to their 'new' unique id;
- `formatted_transactions_seq.txt` , the formatted `transactions_sec.txt` as intended.

Few interfaces had to be declaed too:

```
1 export interface TransactionSeq {
2   user_id: string,
3   order_number: string,
4   order_size: string,
5   products: string
6 }
7
8 export interface ItemSetSeq {
9   productIds: number[],
10  order: number
11 }
```

```
In [22]: 1 import * as CSVParser from './dist/class/csv-parser.class.js';
2 import { TransactionSeq } from './src/interface/transaction-seq.interface';
3 import { ItemSetSeq } from './src/interface/item-set-seq.interface';
4
5 function formatTransactionsSeq(): void {
6   console.log('Gathering data, this might take a while...');
7   $$sync$$ = true;
8
9   let uniqueProducts: string[] = [];
10  let sequences: string[] = [];
11
12  // Open the transaction_sec file.
13  new CSVParser.CSVParser<TransactionSeq>('custom_data/small_transactions_seq.txt', { delimiter: '\t', columns: true })
14    // Group by user_id and map TransactionSeq to ItemSetSeq.
15    .generateItemsets<ItemSetSeq>('user_id', (transactionSeq: TransactionSeq) => {
16      // Parsing products for product_ids.
17      let productIds: number[] = transactionSeq.products.split(',').map( (productName: string) => {
18        // Looking for the productName.
19        let index: number = uniqueProducts.indexOf(productName.replace('\r',''));
20        // If the as already been incountered, returning its index.
21        if(index > -1) return index;
22
23        // Else adding it to our 'uniqueProducts' array.
24        uniqueProducts.push(productName)
25        return uniqueProducts.length - 1;
26      });
27
28      return {
29        // Keep in mind that items must be sorted asc. in the itemsets.
30        productIds: productIds.sort( (a: number, b: number) => a - b ),
31        order: Number.parseInt(transactionSeq.order_number)
32      };
33    })
34    // Once all the group are gathers and sequences formatted.
35    .finally( () => {
36      // Example of output
37      console.log('Example of output:');
38      console.log(sequences.join('\r\n'));
39
40      setTimeout( () => $$done$$(), 300);
41    })
42    // Once one group is available.
43    .subscribe( (group: Group<TransactionSeq,ItemSetSeq>) => {
44      // Let's format groups to a sequence as specified for input
45      let sequence: string = group.items
46        // Within the groups, we need to sort ItemSetSeq by order asc.
47        .sort( (a: ItemSetSeq, b: ItemSetSeq) => a.order - b.order )
48        // Keep in mind that items are separated by a plain space character within the itemsets.
49        .map( (itemSetSeq: ItemSetSeq) => itemSetSeq.productIds.join(' ') )
50        // Each itemset is separated by a ' -1 ' value
51        .join(' -1 ');
52
53      // A sequence ends with ' -2'
54      sequences.push(sequence + ' -2')
55    });
56 }
57
58 formatTransactionsSeq();
```

Gathering data, this might take a while...

Out[22]: undefined

Example of output:

0 1 2 3 4 -1 0 2 3 5 6 7 -1 2 3 7 8 9 -1 2 3 4 7 9 -1 2 3 5 7 9 10 11 12 -1 2 3 7 9 -1 2 3 6 7 9 -1 2 3 7 9 13 14 -1 2 3 7 9 13 14 -1 2 3 6 7 9 14 15 16 17 -2

Parsing output

The output file introduces yet another format as such:

- Each line being a frequent sequential pattern;
- Each itemset being separated by the value -1;
- Support of the sequence is indicated.

Nothing major though, it's only a matter of yet another Regular Expression to match the results with our current existing code:

Where looking for at least 1 ensemble of a spread of any length of positive integers `\d+` followed by '-1' (`(\d+)+-1)+` , and then followed by '#SUP:' and a positive integer `#SUP: \d+ .`

Regular expression:
`((\d+)+-1)+#SUP: \d+`

Running the algorithm

```
In [23]: 1 import * as SPMF from './dist/class/spmf.class.js';
2
3 function sequencesMining(): void {
4     $$async$$ = true;
5     console.log('Mining sequences, this might take a while...');
6     // Our custom SPMF wrapper
7
8     new SPMF.SPMF('CloSpan')
9     .fromFile('custom_data/formatted_transactions_seq.txt')
10    .exec<Sequence>(20)
11    .subscribe((results: SPMFResults<Sequence>) => {
12        // Wrapper returns both the stats...
13        console.log('Stats:');
14        console.log(results.stats);
15        // ... and the frequent sequences. Showing the first two sequences:
16        console.log('First two sequences mined:');
17        results.output.slice(0,2)
18        console.log(results.output.slice(0,2));
19
20        $$done$$();
21    });
22 }
23
24 sequencesMining();
```

```
Mining sequences, this might take a while...
Stats:
{ candidates: undefined, executionTime: 588, memory: undefined }
First two sequences mined:
[ { itemsets: [ [Array] ], support: 6084 },
  { itemsets: [ [Array], [Array] ], support: 4045 } ]
```

Sequential patterns

File output/output_clospan_2.txt provided contains mined sequential patterns with cloSpan and a minimum support of 2%. Let's see what we can find about them...

No time to finish visualisation :(

References

[1]: "Kumar, Dr K Ramesh & D, Usha. (2013). Frequent Pattern Mining Algorithm for Crime Dataset: An Analysis. INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH.". Avalaible on https://www.researchgate.net/publication/266498986_Frequent_Pattern_Mining_Algorithm_for_Crime_Dataset_An_Analysis
https://www.researchgate.net/publication/266498986_Frequent_Pattern_Mining_Algorithm_for_Crime_Dataset_An_Analysis

[2]: "Philippe Fournier Viger. Mining Frequent Closed Itemsets using the LCM Algorithm. SPMF documentation.". Avalaible on <http://www.philippe-fournier-viger.com/spmf/LCM.php> (<http://www.philippe-fournier-viger.com/spmf/LCM.php>)