

<a href="#">1. Aim of the project</a>	1
<a href="#">2. Description of the algorithm</a>	1
<a href="#">3. Datasets used for testing</a>	2
<a href="#">4. Results</a>	2
<a href="#">4.1. Custom dataset</a>	2
<a href="#">4.2. Sign dataset</a>	2
<a href="#">5. Conclusions</a>	3
<a href="#">5.1. Results analysis</a>	3
<a href="#">5.2. Potential uses</a>	3
<a href="#">6. Sources</a>	3

## 1. Aim of the project

Aim of the project is to prepare, implement and verify experimentally SPAM algorithm, regarding mining sequential patterns. The project should be implemented in Python language without using any of the datamining libraries available.

## 2. Description of the algorithm

The pseudo-code of the algorithm is presented below:

---

<b>SPAM</b> ( <i>SDB</i> , <i>minsup</i> )
1. Scan <i>SDB</i> to create $V(SDB)$ and identify $F_1$ , the list of frequent items.
2. <b>FOR</b> each item $s \in F_1$ ,
3. <b>SEARCH</b> ( $\langle\{s\}, F_1, \{e \in F_1 \mid e \succ_{\text{lex}} s\}, \text{minsup}\rangle$ ).

---

<b>SEARCH</b> ( <i>pat</i> , $S_n$ , $I_n$ , <i>minsup</i> )
1. Output pattern <i>pat</i> .
2. $S_{\text{temp}} := I_{\text{temp}} := \emptyset$
3. <b>FOR</b> each item $j \in S_n$ ,
4. <b>IF</b> the s-extension of <i>pat</i> is frequent <b>THEN</b> $S_{\text{temp}} := S_{\text{temp}} \cup \{i\}$ .
5. <b>FOR</b> each item $j \in S_{\text{temp}}$ ,
6. <b>SEARCH</b> (the s-extension of <i>pat</i> with $j$ , $S_{\text{temp}}$ , $\{e \in S_{\text{temp}} \mid e \succ_{\text{lex}} j\}$ , <i>minsup</i> ).
7. <b>FOR</b> each item $j \in I_n$ ,
8. <b>IF</b> the i-extension of <i>pat</i> is frequent <b>THEN</b> $I_{\text{temp}} := I_{\text{temp}} \cup \{i\}$ .
9. <b>FOR</b> each item $j \in I_{\text{temp}}$ ,
10. <b>SEARCH</b> (i-extension of <i>pat</i> with $j$ , $S_{\text{temp}}$ , $\{e \in I_{\text{temp}} \mid e \succ_{\text{lex}} j\}$ , <i>minsup</i> ).

---

### 3. Datasets used for testing

During the project two datasets have been used:

- small, generated by hand dataset with 4 sequences and 2 different items to verify proper implementation of the algorithm, as follows:
  - 1: {a}, {a,b}, {a}, {a,b}
  - 2: {a,b}, {b}
  - 3: {a,b}, {a}, {b}
  - 4. {a}, {b}
- Sign dataset, containing 731 sequences and 310 different items;

### 4. Results

Algorithm has been used to find sequential patterns in both datasets above:

#### 4.1. Custom dataset

Algorithm with *minSup* set to 3 returned following frequent patterns: {a}; {b}; {a,b}; {a}, {b}; {b}, {b}; {a,b}, {b};

Verifying results by hand yields the same results (Searching for frequent patterns with  $sup \geq 3$  implies that it has to be supported from at least one of transaction sets 2,4, which means that the only possible frequent patterns are those sequences which are subsequences of either of those two) .

#### 4.2. Sign dataset

The SPAM algorithm launched on our main dataset has found 100 frequent sequential patterns using a support threshold equal to 400 which corresponds to the relative support of around 0.547.

Instead, considering a higher threshold (i.e. 600), we can notice how the set of results is significantly reduced to just 7 frequent sequential patterns discovered. Example mined pattern is presented below:

```
6:
Sequence: {1},{253}
Absolute support: 616
Relative support: 0.842681
```

## 5. Conclusions

### 5.1. Results analysis

Tests show that although the algorithm works for the given dataset, this specific dataset contains only itemsets with one item each, therefore the i-extension part of the algorithm remains unused and one of the strengths of it is neutralized.

Despite this, SPAM performs well and is one of the fastest sequential pattern mining algorithm

### 5.2. Potential uses

**One of the fields that the algorithm really shines is finding a comprehensive list of all sequential patterns, which later can be processed via means of rule induction to generate sequential rules about the dataset.** One can obtain vital information from a dataset this way, assuming the dataset is appropriate for finding such rules. One of the examples would be the laboratory task where one had to find out what occurrences may cause hypoglycemic symptoms to appear.

## 6. Sources

- [http://www.philippe-fournier-viger.com/spmf/PAKDD2014\\_sequential\\_pattern\\_mining\\_CM-SPADE\\_CM-SPAM.pdf](http://www.philippe-fournier-viger.com/spmf/PAKDD2014_sequential_pattern_mining_CM-SPADE_CM-SPAM.pdf)
- <http://www.philippe-fournier-viger.com/spmf/index.php?link=datasets.php>
- <http://www.philippe-fournier-viger.com/spmf/datasets/SIGN.txt>

