

Mining High-Utility Patterns in Utility Databases

What is High-Utility pattern mining?

High utility pattern mining aims to discover all the patterns with utility of pattern is no less than user-specified **minimum utility** threshold **minutil**. **minUtil** controls the minimum utility of patterns should have.

Reference: Hong Yao and Howard J. Hamilton. 2006. Mining itemset utilities from transaction databases. Data Knowl. Eng. 59, 3 (December 2006), 603–626.

<https://doi.org/10.1016/j.datak.2005.10.004>

What is a utility database?

A utility database is a collection of transaction, where each transaction contains a set of items and a positive integer called **internal utility** respectively. And each unique item in database is also associated with another positive number called **external utility**.

A hypothetical utility database with items **a, b, c, d, e, f and g** and its **internal utility** is shown below at right side and items with its **external utility** is presented at left side.

Transactions	Item	Profit
(a,2) (b,3) (c,1) (g,1)	a	4
(b,3) (c,2) (d,3) (e,2)	b	3
(a,2) (b,1) (c,3) (d,4)	c	6
(a,3) (c,2) (d,1) (f,2)	d	2
(a,3) (b,1) (c,2) (d,1) (g,2)	e	5
(c,2) (d,2) (e,3) (f,1)	f	2
(a,2) (b,1) (c,1) (d,2)	g	3
(a,1) (e,2) (f,2)		
(a,2) (b,2) (c,4) (d,2)		
(b,3) (c,2) (d,2) (e,2)		

Note: Duplicate items must not exist in a transaction.

What is acceptable format of a utility databases in PAMI?

Each row in a utility database must contain only items, total sum of utilities and utility values.

```
a b c g:7:2 3 1 1
b c d e:10:3 2 3 2
a b c d:10:2 1 3 4
a c d f:7:3 2 1 2
a b c d g:9:3 1 2 1 2
c d e f:8:2 2 3 1
a b c d:6:2 1 1 2
a e f:5:1 2 2
a b c d:10:2 2 4 2
b c d e:9:3 2 2 2
```

Understanding the statistics of database

To understand about the database. The below code will give the detail about the transactional database.

- Total number of transactions (Database size)
- Total number of unique items in database
- Minimum length of transaction that existed in database
- Average length of all transactions that exists in database
- Maximum length of transaction that existed in database
- Minimum utility value exists in database
- Average utility exists in database
- Maximum utility exists in database
- Standard deviation of transaction length
- Variance in transaction length
- Sparsity of database

The below sample code prints the statistical details of a database.

```
In [5]: import PAMI.extras.dbStats.utilityDatabaseStats as stats
obj = stats.utilityDatabaseStats('sample_input.txt', ' ')
obj.run()
obj.printStats()
```

```
Database size : 10
Number of items : 7
Minimum Transaction Size : 3
Average Transaction Size : 4.0
Maximum Transaction Size : 5
Minimum utility : 3
Average utility : 11.714285714285714
Maximum utility : 19
Standard Deviation Transaction Size : 0.4472135954999579
Variance : 0.2222222222222222
Sparsity : 0.42857142857142855
```

What are the input parameters?

The input parameters to a frequent pattern mining algorithm are:

- **Utility database**

Acceptable formats:

- String : E.g., 'utilityDatabase.txt'
- URL : E.g., https://u-aizu.ac.jp/~udayrage/datasets/transactionalDatabases/transactional_T10
- DataFrame with the header titled 'Transactions', 'Utility' and 'TransactionUtility'

- **minUtil**

specified in

- **count**

- **seperator**

default seperator is '\t' (tab space)

How to store the output of a high utility pattern mining algorithm?

The patterns discovered by a high utility pattern mining algorithm can be saved into a file or a data frame.

How to run the high utility pattern mining algorithms in a terminal?

- Download the PAMI source code from github.
- Unzip the PAMI source code folder and enter into high utility pattern folder.
- You will find folder like **basic**
- Enter into the basic folder and execute the following command on terminal.

syntax: python3 algorithmName.py <path to the input file> <path to the output file> <minUtil> <seperator>

Example: python3 EFIM.py inputFile.txt outputFile.txt \$20\$ ' '

How to execute a High utility pattern mining algorithm in a Jupyter Notebook?

- Install the PAMI package from the PYPI repository by executing the following command: **pip3 install PAMI**
- Run the below sample code by making necessary changes

```
In [ ]: import PAMI.highUtilityPatterns.basic.EFIM as alg

iFile = 'sample_Input.txt' #specify the input utility database <
minUtil = 20               #specify the minSupvalue
seperator = ' '           #specify the seperator. Default seperator i
oFile = 'utilityPatterns.txt' #specify the output file name

obj = alg.EFIM(iFile, minUtil, seperator) #initialize the algorithm
obj.startMine()                          #start the mining process
obj.savePatterns(oFile)                   #store the patterns in file
df = obj.getPatternsAsDataFrame()         #Get the patterns discovered into a
obj.printStats()                          #Print the statistics of mining pro
```

The utilityPatterns.txt file contains the following patterns (format: pattern:utility):!cat utilityPatterns.txt

```
In [3]: !cat utilityPatterns.txt
```

```
e d c : 20
a b d : 23
a b d c : 33
a b c : 30
a d : 22
a d c : 34
a c : 27
b d : 25
b d c : 39
b c : 29
d c : 35
```

The dataframe containing the patterns is shown below:

In [4]:

```
df
```

Out[4]:

	Patterns	Utility
0	e d c	20
1	a b d	23
2	a b d c	33
3	a b c	30
4	a d	22
5	a d c	34
6	a c	27
7	b d	25
8	b d c	39
9	b c	29
10	d c	35