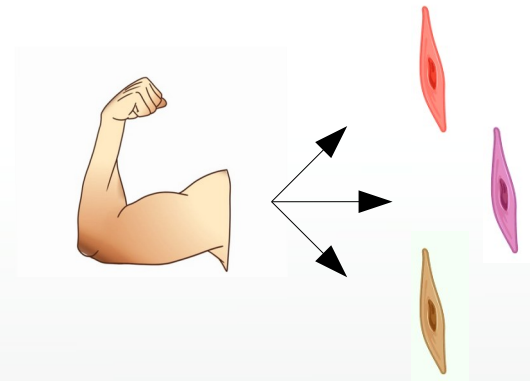
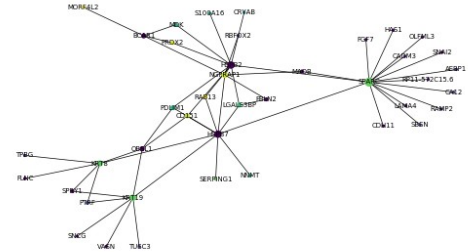


How to find genes network into single cell RNA data

By Romuald MARIN

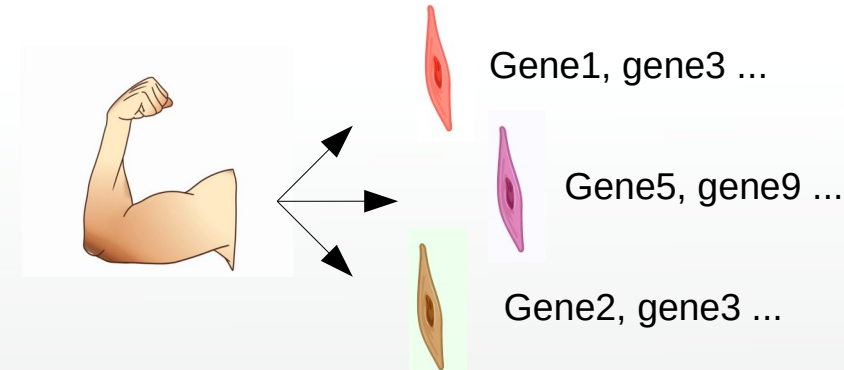
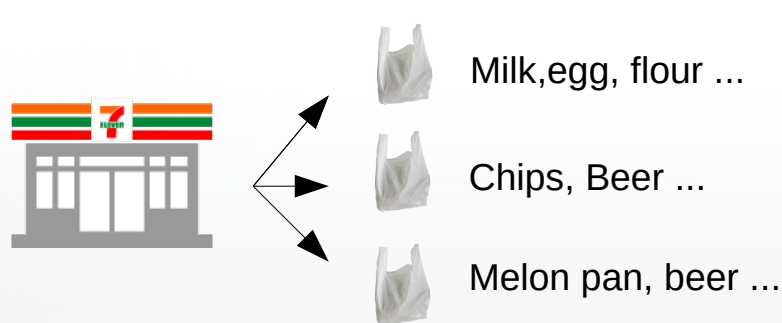
Why do we want to find genes network?

- Genes network :
 - Groupe of genes which interact with each other
 - Transcription factor gene
- Single cell RNA :
 - Cellular level study
 - Big variation
 - Non-expressed gene

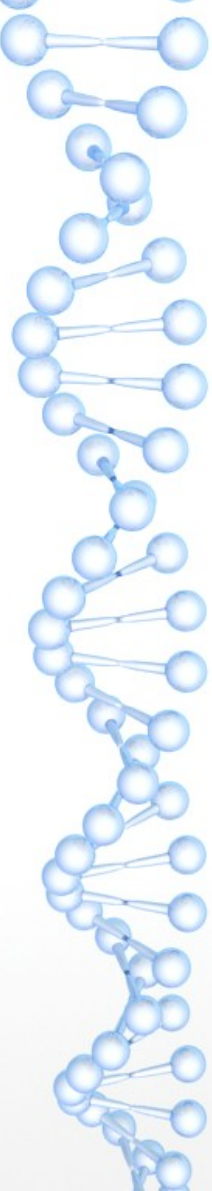


My approach?

- Association rule learning :
 - Discovering interesting relations between variables in large databases
 - Discovering patterns in large datasets
- Apriori Algorithm
 - Find objects bought together
 - Find relations between different items



How does it work ? (1)



	A1BG	ACTB	A1CF	...	FAT1
Cell 1	0	364	1014	...	40
Cell 2	0	909	0	...	0
Cell 3	0	501	590	...	2046
Cell 4	0	107	0	...	0

Count Matrix with
expressed genes

Normalize and bool :

- Remove no expressed gene
- Remove too small count

	ACTB	A1CF	...	L3HYPDH
Cell 1	True	True	...	False
Cell 2	True	False	...	False
Cell 3	True	True	...	True
Cell 4	True	False	...	False

Boolean Matrix with
expressed genes

Support : Fraction of cells who expressed genes
(or a group of genes)

Confidence : Fraction of gene expressed if an
other gene (or a group of genes) is expressed

Lift : Measure of correlation, if a gene or a groupe
of genes are expressed and an other gene

How does it work ? (2)

1st step : Define threshold of Support and Confidence

2nd step : Find gene who pass different threshold

3rd step : Create group of gene of length L+1 and check values

4th step : Repeat until length threshold

Threshold

- Min Support = 0,1
- Max Support = 0,7
- Min Confidence = 0.9
- Length = 3

	ACTB	A1CF	...	FAT1
Cell 1	True	True	...	False
Cell 2	True	False	...	False
Cell 3	True	True	...	True
Cell 4	True	False	...	False

Gene	Support
AACS	0.05
ABHD	0.13
COQ4	0.62
LCP2	0.20
SAT1	0.96

New group :

ABHD , COQ4
ABHD, LCP2
LCP2, COQ4

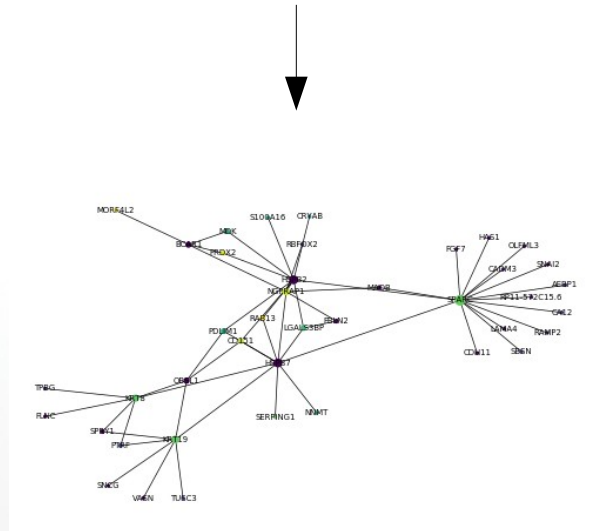
Calculate the value of support and confidence

Check value of support and confidence

Create new groupe of gene length +1

- results json format
- List in txt format
- Log file _____
- Network representation

- results json format
- List in txt format
- Log file _____
- Network representation



```

    "CNN3": {
      "support": 0.29493183473826623,
      "LY6E": {
        "support": 0.2819231970028099,
        "confidence": 0.955892731122089,
        "lift": 1.3944395405119405
      },
      "EEF2": {
        "LY6E": {
          "support": 0.2680820064522843,
          "confidence": 0.9509043927648578,
          "PDIA6": {
            "lift": 0.6450125367716212,
            "support": 0.254761161411177,
            "confidence": 0.9503105590062111
          },
          "OSTC": {
            "lift": 0.6416995356966715,
            "support": 0.25548964512436256,
            "confidence": 0.953027950310559
          }
        }
      }
    }
  },
}

```

```

1 workspace(do="datainling", input=matrixGSE146026/GSE146026_Izar_HGSC0C_ascites_10x_log.tsv",
2   max_length=4, max_support=0.7, min_confidence=0.95, min_support=0.2, normalize=False,
3   output=resultat_GSE146026MAXSUPPORT0.7L4", processor=8, rowremove="", transpose=True)
4 Loading data from file matrixGSE146026/GSE146026_Izar_HGSC0C_ascites_10x_log.tsv
5 Transpose matrix
6
7 Cell_ID AL627309.1 LINC00115 SAMD11 ... AL354822.1 PNRC2 SRSF101
8 | 10x_1 False False False ... False False False
9 | 10x_2 False False False ... False False False
10 | 10x_3 False False False ... False False False
11 | 10x_4 False False False ... False False False
12 | 10x_5 False False False ... False False False
13 |
14 | [5 rows x 11548 columns]
15 |
16 | Maximun lenght of tnsset is :
17 | 54
18 |
19 | Launch new apriori
20 |
21 | Generate C1
22 |
23 | number of itemsets find :
24 |
25 | 2405
26 |
27 | new number of itemsets find :
28 |
29 | 2405
30 |
31 | Generate C2
32 |
33 | Remove clone. Old number :
34 |
35 | 62
36 |
37 | New number :
38 |
39 | 62
40 |
41 | number of itemsets find :
42 |
43 | 62
44 |
45 | Generate C3
46 |
47 | Remove clone. Old number :
48 |
49 | 310
50 |
51 | New number :
52 |
53 | 278
54 |
55 | number of itemsets find :
56 |
57 | 278

```



Importance of different thresholds

Support threshold:

- Maximum support allow to remove genes who expressed in all cell (default=0.8)
- Minimum support allow to discovered low frequent network (default=0.3)

Confidence and lift

- allow to keep only gene network interesting (default confidence=0.9 & lift > 1)

Lenght

- Reduce compute time (default=4)

-

- Threshold :

- Results :

- Example of pattern :

-



Conclusion

ADVANTAGE

- Easy to use
- Very fast compare to other method
- Print graph and different results files
- Allow to study different level with threshold

DISADVANTAGE

- Only up regulated genes network are find
- More accurate methods are available