# CSA1668-DATA WAREHOUSING AND DATA MINING

T, MANOHAR
192110219

## DATA PRE PROCESSING

# DISCRETIZE

# CLUSTERING

# ASSOCIATE



Weka Explorer — Associate

```
=== Associator model (full training set) ===


Apriori
=======

Minimum support: 0.75 (214 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 5

Generated sets of large itemsets:

Size of set of large itemsets L(1): 4

Size of set of large itemsets L(2): 5

Size of set of large itemsets L(3): 2

Best rules found:

 1. ExtremeValue=no 206 ==> Outlier=no 206    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 2. Outlier=no 206 ==> ExtremeValue=no 206    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 3. node-caps=no 230 ==> Outlier=no 230    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 4. node-caps=no 230 ==> ExtremeValue=no 230    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 5. node-caps=no ExtremeValue=no 230 ==> Outlier=no 230    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 6. node-caps=no Outlier=no 230 ==> ExtremeValue=no 230    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 7. node-caps=no 230 ==> Outlier=no ExtremeValue=no 230    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 8. irradiat=no 218 ==> Outlier=no 218    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 9. irradiat=no 218 ==> ExtremeValue=no 218    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
10. irradiat=no ExtremeValue=no 218 ==> Outlier=no 218    <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
```
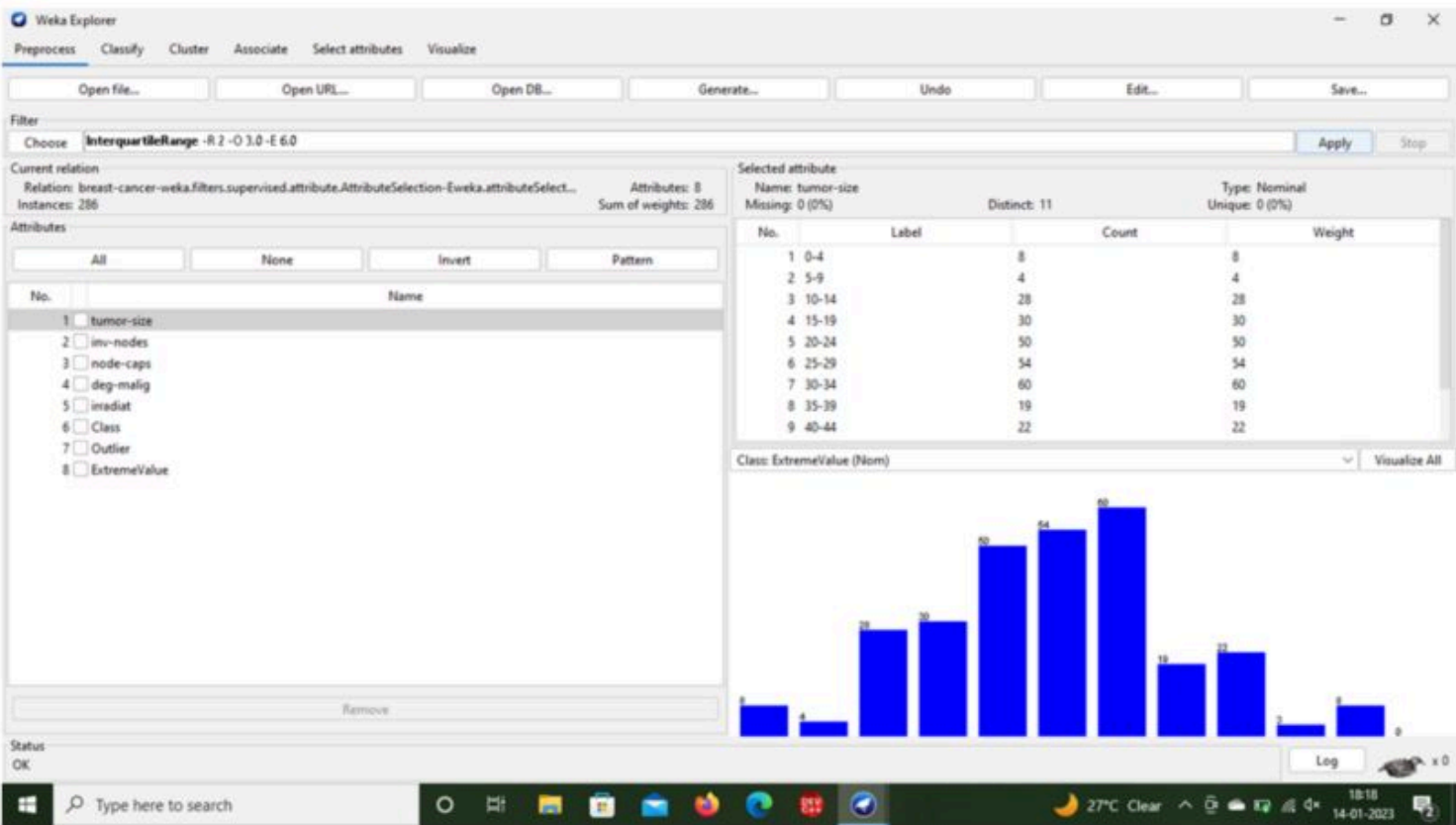
# INTERQUARTILE RANGE

# REPLACE MISSING VALUES

# CLASSIFICATION

# EM

Preprocess   Classify   Cluster   Associate   Select attributes   Visualize

**Clusterer**

Choose   EM -I 100 -N -1 -X 10 -max -1 -ll-cv 1.0E-6 -ll-iter 1.0E-6 -M 1.0E-6 -K 10 -num-slots 1 -S 100

**Cluster mode**
- Use training set
- Supplied test set          Set...
- Percentage split        %  66
- Classes to clusters evaluation
  (Nom) ExtremeValue
- ☑ Store clusters for visualization

Ignore attributes

Start          Stop

**Result list (right-click for options)**
18:21:22 - SimpleKMeans
18:21:36 - EM

**Status**
OK

**Clusterer output**

```
    [total]                212.8894   79.1106
irradiat
  yes                       27.0213   42.9787
  no                       184.8681   35.1319
    [total]                211.8894   78.1106
Class
  no-recurrence-events     169.7021   33.2979
  recurrence-events         42.1873   44.8127
    [total]                211.8894   78.1106
Outlier
  no                       210.8894   77.1106
  yes                            1          1
    [total]                211.8894   78.1106
ExtremeValue
  no                       210.8894   77.1106
  yes                            1          1
    [total]                211.8894   78.1106


Time taken to build model (full training data) : 0.93 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      211 ( 74%)
1       75 ( 26%)


Log likelihood: -5.35071
```

# NORMALIZE