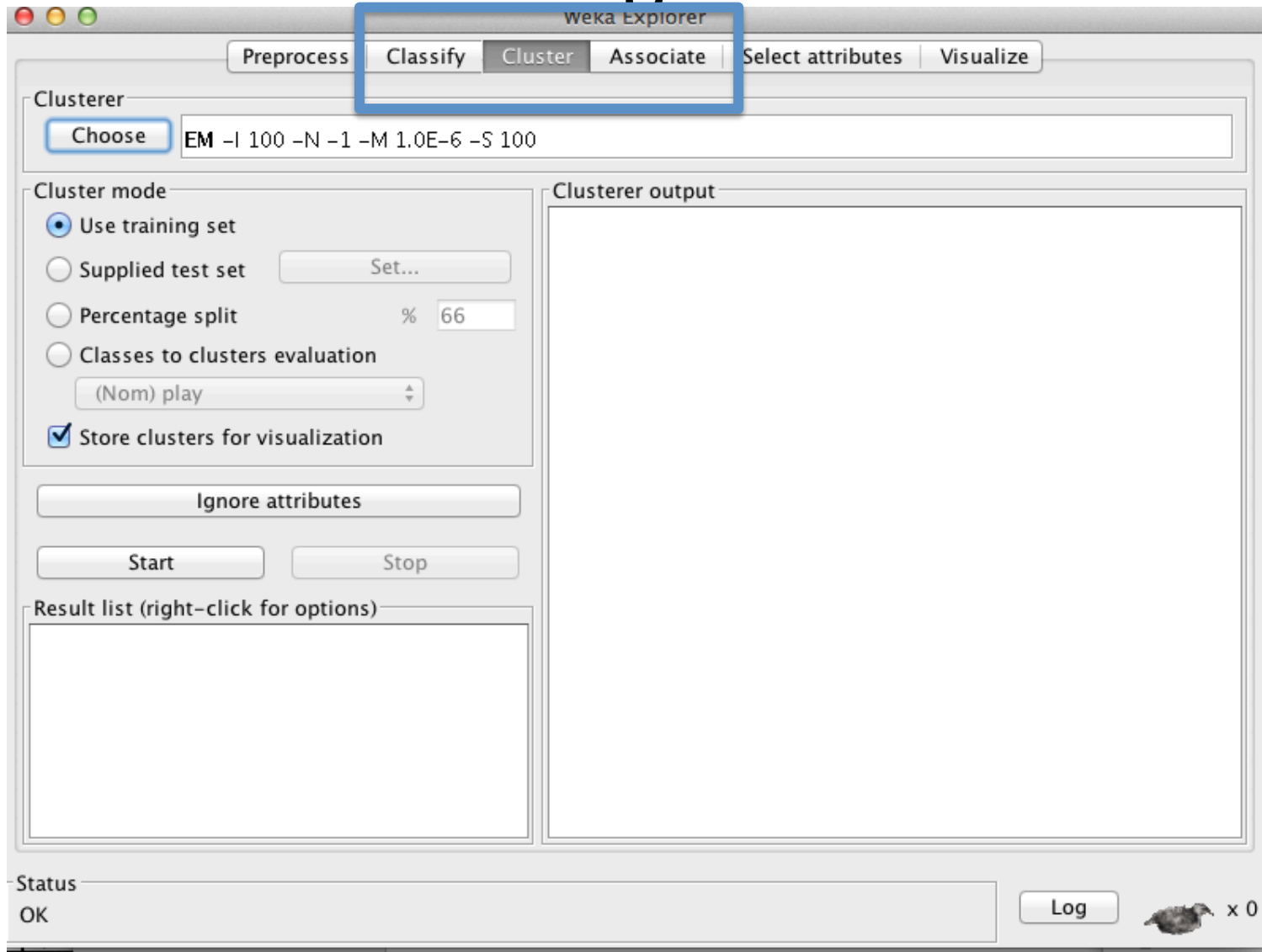# WEKA Explorer

Second part

# ML algorithms in weka belong to 3 categories

# Will see examples in each category (as we learn new algorithms)

1. **Classifiers** (given a set of categories, learn to assign each instance to a category. These are TRAINED methods): Decision Trees, decision tables, conjunctive rules..

2. **Clustering** (given a set of instances, group these instances in clusters according to some similarity function. These are UNTRAINED methods): Hierarchical clustering, DensityBased, etc)

3. **Association rules** (given a set of instances, find frequent patterns, e.g. rules that show dependencies among the data. These are UNTRAINED methods): Apriori, Filtered Associator, ..

4. Additional algorithms can be used, within the Experimenter (will see later)

# Weka Knowledge Explorer

Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize

## Classifier

Choose | **ZeroR**

## Test options

○ Use training set

○ Supplied test set    Set...

● Cross-validation    Folds | 10

○ Percentage split    % | 66

More options...

(Nom) class ▲▼

Start | Stop

## Result list (right-click for options)

## Classifier output

## Status

OK

Log | x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

- weka
  - ▼ classifiers
    - ▶ bayes
    - ▶ functions
    - ▶ lazy
    - ▶ meta
    - ▶ misc
    - ▼ trees
      - ▶ adtree
      - DecisionStump
      - Id3
      - ▼ j48
        - J48
      - ▶ lmt
      - ▶ m5
      - RandomForest
      - RandomTree
      - REPTree
      - UserClassifier
    - ▶ rules

## Status

OK

Log    x 0

# Weka Knowledge Explorer

| Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize |

## Classifier

[ Choose ]  **J48** -C 0.25 -M 2

## Test options

○ Use training set

○ Supplied test set    [ Set... ]

◉ Cross-validation   Folds  10

○ Percentage split    %   66

[ More options... ]

[ (Nom) class ▲▼ ]

[ Start ]    [ Stop ]

## Result list (right-click for options)

## Classifier output

## Status

OK

[ Log ]    x 0

# Weka Knowledge Explorer

**Preprocess** | **Classify** | **Cluster** | **Associate** | **Select attributes** | **Visualize**

## Classifier

| Choose | J48 -C 0.25 -M 2 |

## Test options

- ⚪ Use training set
- ⚪ Supplied test set — Set...
- ⚫ Cross-validation  Folds  10
- ⚪ Percentage split  %  66

More options...

(Nom) class

Start | Stop

## Result list (right-click for options)

## Classifier output

## Status

OK

Log   x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

[ Choose ] J48 -C 0.25 -M 2

**weka.gui.GenericObjectEditor**

weka.classifiers.trees.j48.J48

## Test options

- ( ) Use training set
- ( ) Supplied test set    [ Set... ]
- (•) Cross-validation   Folds   10
- ( ) Percentage split    %   66

[ More options... ]

(Nom) class

[ Start ]    [ Stop ]

## Result list (right-click for options)

| binarySplits | False |
| confidenceFactor | 0.25 |
| minNumObj | 2 |
| numFolds | 3 |
| reducedErrorPruning | False |
| saveInstanceData | False |
| subtreeRaising | True |
| unpruned | False |
| useLaplace | False |

[ Open... ]  [ Save... ]  [ OK ]  [ Cancel ]

## Status

OK

[ Log ]    x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |
|---|---|---|---|---|---|

**Classifier**

Choose  J48 -C 0.25 -M 2

**Test options**

- ○ Use training set
- ○ Supplied test set    Set...
- ● Cross-validation   Folds  10
- ○ Percentage split      %  66

More options...

(Nom) class

Start    Stop

**Result list (right-click for options)**

---

## weka.gui.GenericObjectEditor

weka.classifiers.trees.j48.J48

| | |
|---|---|
| binarySplits | False |
| confidenceFactor | 0.25 |
| minNumObj | 2 |
| numFolds | 3 |
| reducedErrorPruning | False |
| saveInstanceData | False |
| subtreeRaising | True |
| unpruned | False |
| useLaplace | False |

Open...  Save...  OK  Cancel

---

**Status**

OK

Log    x 0

# Weka Knowledge Explorer

| Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize |

## Classifier

( Choose ) **J48** -C 0.25 -M 2

## Test options

- ◯ Use training set
- ◯ Supplied test set    ( Set... )
- ⦿ Cross-validation   Folds  `10`
- ◯ Percentage split    %  `66`

( More options... )

[ (Nom) class ▼ ]

( Start )    ( Stop )

## Result list (right-click for options)

## Classifier output

## Status

OK

( Log )    x 0

# Weka Knowledge Explorer

Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize

## Classifier

**Choose** | J48 -C 0.25 -M 2

## Test options

- ◯ Use training set
- ◯ Supplied test set — Set...
- ◉ Cross-validation    Folds  10
- ◯ Percentage split    %  66

More options...

(Nom) class ⬍

Start | Stop

## Result list (right-click for options)

## Classifier output

## Status

OK

Log | 🐑 x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

Choose J48 -C 0.25 -M 2

## Test options

- Use training set
- Supplied test set    Set...
- Cross-validation    Folds 10
- Percentage split    % 66

More options...

(Nom) class

Start    Stop

## Result list (right-click for options)

## Classifier output

# Weka Knowledge Explorer

| Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize |

## Classifier

**Choose** | J48 -C 0.25 -M 2

## Test options

- ◯ Use training set
- ◯ Supplied test set    Set...
- ◯ Cross-validation    Folds    10
- ◉ Percentage split    %    66

More options...

(Nom) class

Start    Stop

## Result list (right-click for options)

## Classifier output

## Status

OK

Log    x 0

# Weka Knowledge Explorer

| Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize |

## Classifier

**Choose** | J48 -C 0.25 -M 2

## Test options

- ○ Use training set
- ○ Supplied test set    Set...
- ○ Cross-validation    Folds  10
- ● Percentage split    %  66

More options...

(Nom) class ▾

Start    Stop

## Result list (right-click for options)

## Classifier output

### ● ● ● Classifier evaluation opt

- ☑ Output model
- ☑ Output per-class stats
- ☐ Output entropy evaluation measures
- ☑ Output confusion matrix
- ☑ Store predictions for visualization
- ☐ Output text predictions on test set
- ☐ Cost-sensitive evaluation    Set...

Random seed for XVal / % Split    1

OK

## Status

OK

Log    x 0

# Weka Knowledge Explorer

Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize

## Classifier

Choose | J48 -C 0.25 -M 2

## Test options

- ○ Use training set
- ○ Supplied test set — Set...
- ○ Cross-validation — Folds 10
- ● Percentage split — % 66

More options...

(Nom) class

Start | Stop

## Result list (right-click for options)

## Classifier output

### Classifier evaluation options

- ☑ Output model
- ☑ Output per-class stats
- ☐ Output entropy evaluation measures
- ☑ Output confusion matrix
- ☑ Store predictions for visualization
- ☐ Output text predictions on test set
- ☐ Cost-sensitive evaluation — Set...

Random seed for XVal / % Split | 1

OK

## Status

OK

Log | x 0

## Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

### Classifier

( Choose )  **J48** -C 0.25 -M 2

### Test options

○ Use training set

○ Supplied test set        ( Set... )

○ Cross-validation  Folds  10

◉ Percentage split      %  66

( More options... )

(Nom) class  ▲▼

( Start )        ( Stop )

### Result list (right-click for options)

### Classifier output

### Status

OK

( Log )  🐑  x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

**Choose**  J48 -C 0.25 -M 2

## Test options

○ Use training set
○ Supplied test set    Set...
○ Cross-validation  Folds  10
● Percentage split    %  66

More options...

(Nom) class

Start    Stop

## Result list (right-click for options)

11:49:05 – trees.j48.J48

## Classifier output

```
=== Run information ===

Scheme:       weka.classifiers.trees.j48.J48 -C 0.25 -M 2
Relation:     iris
Instances:    150
Attributes:   5
              sepallength
              sepalwidth
              petallength
              petalwidth
              class
Test mode:    split 66% train, remainder test

=== Classifier model (full training set) ===

J48 pruned tree
------------------

petalwidth <= 0.6: Iris-setosa (50.0)
petalwidth > 0.6
|   petalwidth <= 1.7
|   |   petallength <= 4.9: Iris-versicolor (48.0/1.0)
|   |   petallength > 4.9
|   |   |   petalwidth <= 1.5: Iris-virginica (3.0)
|   |   |   petalwidth > 1.5: Iris-versicolor (3.0/1.0)
|   petalwidth > 1.7: Iris-virginica (46.0/1.0)

Number of Leaves  :     5
```
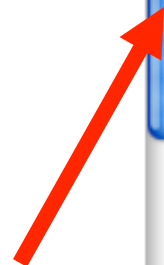
## Status

OK

Log    x 0

# Weka Knowledge Explorer

**Preprocess** | **Classify** | **Cluster** | **Associate** | **Select attributes** | **Visualize**

## Classifier

**Choose** | J48 -C 0.25 -M 2

## Test options

- ○ Use training set
- ○ Supplied test set — Set...
- ○ Cross-validation — Folds 10
- ● Percentage split — % 66

**More options...**

(Nom) class

**Start** | **Stop**

## Result list (right-click for options)

11:49:05 – trees.j48.J48

## Classifier output

```
=== Run information ===

Scheme:        weka.classifiers.trees.j48.J48 -C 0.25 -M 2
Relation:      iris
Instances:     150
Attributes:    5
               sepallength
               sepalwidth
               petallength
               petalwidth
               class
Test mode:     split 66% train, remainder test

=== Classifier model (full training set) ===

J48 pruned tree
------------------

petalwidth <= 0.6: Iris-setosa (50.0)
petalwidth > 0.6
|    petalwidth <= 1.7
|    |    petallength <= 4.9: Iris-versicolor (48.0/1.0)
|    |    petallength > 4.9
|    |    |    petalwidth <= 1.5: Iris-virginica (3.0)
|    |    |    petalwidth > 1.5: Iris-versicolor (3.0/1.0)
|    petalwidth > 1.7: Iris-virginica (46.0/1.0)

Number of Leaves  :      5
```

## Status

OK

**Log** | x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

**Choose** | J48 -C 0.25 -M 2

## Test options

- ◯ Use training set
- ◯ Supplied test set    Set...
- ◯ Cross-validation  Folds  10
- ● Percentage split    %  66

More options...

(Nom) class

**Start**    Stop

## Result list (right-click for options)

11:49:05 – trees.j48.J48

## Classifier output

```
Time taken to build model: 0.24 seconds

=== Evaluation on test split ===
=== Summary ===

Correctly Classified Instances          49              96.0784 %
Incorrectly Classified Instances         2               3.9216 %
Kappa statistic                          0.9408
Mean absolute error                      0.0396
Root mean squared error                  0.1579
Relative absolute error                  8.8979 %
Root relative squared error             33.4091 %
Total Number of Instances               51

=== Detailed Accuracy By Class ===

TP Rate    FP Rate   Precision   Recall   F-Measure   Class
  1          0          1          1          1        Iris-setosa
  1          0.063      0.905      1          0.95     Iris-versicolor
  0.882      0          1          0.882      0.938    Iris-virginica

=== Confusion Matrix ===

  a  b  c   <-- classified as
 15  0  0 |  a = Iris-setosa
  0 19  0 |  b = Iris-versicolor
  0  2 15 |  c = Iris-virginica
```

## Status

OK

Log    × 0

# Weka Knowledge Explorer

| Preprocess | **Classify** | Cluster | Associate | Select attributes | Visualize |

## Classifier

[ Choose ] J48 -C 0.25 -M 2

## Test options

○ Use training set
○ Supplied test set    [ Set... ]
○ Cross-validation    Folds [ 10 ]
● Percentage split    % [ 66 ]

[ More options... ]

(Nom) class ▾

[ Start ]    [ Stop ]

## Result list (right-click for options)

11:49:05 – trees.j48.J48

## Classifier output

```
Time taken to build model: 0.24 seconds


=== Evaluation on test split ===
=== Summary ===

Correctly Classified Instances        49              96.0784 %
Incorrectly Classified Instances       2               3.9216 %
Kappa statistic                        0.9408
Mean absolute error                    0.0396
Root mean squared error                0.1579
Relative absolute error                8.8979 %
Root relative squared error           33.4091 %
Total Number of Instances             51

=== Detailed Accuracy By Class ===
```

| View in main window | | Recall | F-Measure | Class |
| View in separate window | | 1 | 1 | Iris-setosa |
| Save result buffer | | 1 | 0.95 | Iris-versicolor |
| | | 0.882 | 0.938 | Iris-virginica |
| Load model | |
| Save model | |
| Re-evaluate model on current test set | |
| Visualize classifer errors | |
| **Visualize tree** | |
| Visualize margin curve | |
| Visualize threshold curve | ▶ |
| Visualize cost curve | ▶ |

## Status

OK

[ Log ]   x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

[ Choose ] J48 -C 0.25 -M 2

## Test options

- ○ Use training set
- ○ Supplied test set
- ○ Cross-validation
- ● Percentage split

[ More optic ]

(Nom) class

[ Start ]

### Result list (right-click for

11:49:05 – trees.j48.J

## Weka Classifier Tree Visualizer: 11:49:05 – trees.j48.J48 (iris)

### Tree View



```
petalwidth
   <= 0.6              > 0.6
Iris-setosa (50.0)     petalwidth
                   <= 1.7        > 1.7
              petallength     Iris-virginica (46.0/1.0)
          <= 4.9      > 4.9
 Iris-versicolor (48.0/1.0)   petalwidth
                         <= 1.5        > 1.5
              Iris-virginica (3.0)   Iris-versicolor (3.0/1.0)
```

96.0784 %
3.9216 %

ass
is-setosa
is-versicolor
is-virginica

```
 15  0  0 |  a = Iris-setosa
  0 19  0 |  b = Iris-versicolor
  0  2 15 |  c = Iris-virginica
```

## Status

OK

[ Log ]   x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Classifier

**Choose** `J48 -C 0.25 -M 2`

## Test options

- ○ Use training set
- ○ Supplied test set    Set...
- ○ Cross-validation    Folds    10
- ● Percentage split    %    66

**More options...**

(Nom) class

**Start**    **Stop**

## Result list (right-click for options)

11:49:05 – trees.j48.J48

## Classifier output

```
Time taken to build model: 0.24 seconds

=== Evaluation on test split ===
=== Summary ===

Correctly Classified Instances          49               96.0784 %
Incorrectly Classified Instances         2                3.9216 %
Kappa statistic                          0.9408
Mean absolute error                      0.0396
Root mean squared error                  0.1579
Relative absolute error                  8.8979 %
Root relative squared error             33.4091 %
Total Number of Instances               51

=== Detailed Accuracy By Class ===

TP Rate   FP Rate   Precision   Recall   F-Measure   Class
  1         0         1           1        1           Iris-setosa
  1         0.063     0.905       1        0.95        Iris-versicolor
  0.882     0         1           0.882    0.938       Iris-virginica

=== Confusion Matrix ===

  a  b  c   <-- classified as
 15  0  0 |  a = Iris-setosa
  0 19  0 |  b = Iris-versicolor
  0  2 15 |  c = Iris-virginica
```

## Status

OK                                                                          Log    🐑 x 0

# Explorer: clustering data

- WEKA contains "clusterers" for finding groups of similar instances in a dataset

- Implemented schemes are:

  - *k*-Means, EM, Cobweb, *X*-means, FarthestFirst

- Clusters can be visualized and compared to "true" clusters (if given)

- Evaluation based on loglikelihood if clustering scheme produces a probability distribution

# The K-Means Clustering Method

- Given *k*, the *k-means* algorithm is implemented in four steps:
  - Partition objects into *k* nonempty subsets
  - Compute seed points as the centroids of the clusters of the current partition (the centroid is the center, i.e., *mean point*, of the cluster)
  - Assign each object to the cluster with the nearest seed point
  - Go back to Step 2, stop when no more new assignment

Weka Clusterer Visualize: 16:12:19 - SimpleKMeans (iris-weka....

X: Instance_number (Num)    Y: sepallength (Nom)    t attributes    Visualize

Colour: Cluster (Nom)    Select Instance

Reset    Clear    Open    Save    Jitter ○

Plot:iris-weka.filters.supervised.attribute.Discretize-Rfirst-last_clustered

Class colour

cluster0                                    cluster1

16:12:19 - SimpleKMeans

right click: visualize cluster assignement

on on test split ===

3
squared errors: 165.0
y replaced with mean/mode

                          Cluster#
      Full Data              0
        (99)               (35)          (
================================================
-inf-5.55]'    '(-inf-5.55]'    '(6.15-in
2.95-3.35]'    '(3.35-inf]'    '(-inf-2.9
(4.75-inf]'    '(-inf-2.45]'    '(4.75-in
(-inf-0.8]'    '(-inf-0.8]'    '(0.8-1.7
Iris-setosa    Iris-setosa Iris-virgin

Time taken to build model (percentage split) : 0.01 seconds

Clustered Instances

0      17 ( 33%)
1      34 ( 67%)

Status
OK                                              Log          x 0

# Explorer: finding associations

- WEKA contains an implementation of the Apriori algorithm for learning association rules
  - Works only with discrete data
- Can identify statistical dependencies between groups of attributes:
  - milk, butter $\Rightarrow$ bread, eggs (with confidence 0.9 and support 2000)
- Apriori can compute all rules that have a given minimum support and exceed a given confidence

# Basic Concepts: Frequent Patterns

| Tid | Items bought |
|-----|--------------|
| 10 | Beer, Nuts, Diaper |
| 20 | Beer, Coffee, Diaper |
| 30 | Beer, Diaper, Eggs |
| 40 | Nuts, Eggs, Milk |
| 50 | Nuts, Coffee, Diaper, Eggs, Milk |



**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

- itemset: A set of one or more items
- k-itemset $X = \{x_1, \ldots, x_k\}$
- *(absolute) support*, or, *support count* of X: Frequency or occurrence of an itemset X
- *(relative) support*, $s$, is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- An itemset X is *frequent* if X's support is no less than a *minsup* threshold

# Basic Concepts: Association Rules

| Tid | Items bought |
|-----|--------------|
| 10 | Beer, Nuts, Diaper |
| 20 | Beer, Coffee, Diaper |
| 30 | Beer, Diaper, Eggs |
| 40 | Nuts, Eggs, Milk |
| 50 | Nuts, Coffee, Diaper, Eggs, Milk |



**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

- Find all the rules $X \rightarrow Y$ with minimum support and confidence

  - support, $s$, probability that a transaction contains $X \cup Y$

  - confidence, $c$, conditional probability that a transaction having X also contains $Y$

*Let minsup = 50%, minconf = 50%*

*Freq. Pat.:* Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3

- Association rules: (many more!)
  - *Beer $\rightarrow$ Diaper* (60%, 100%)
  - *Diaper $\rightarrow$ Beer* (60%, 75%)

## Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Generate... | Undo | Edit... | Save...

**Filter**

Choose | Di

Apply

**Current relation**

Relation: supe...

Type: Nominal

Instances: 4627

ique: 0 (0%)

**Attributes**

All

---

### Open

📁 data

| Name | Date Modified |
|------|---------------|
| iris.arff | mercoledì 15 agosto 2012 0.12 |
| labor.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersCorn-test.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersCorn-train.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersGrain-test.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersGrain-train.arff | mercoledì 15 agosto 2012 0.12 |
| segment-challenge.arff | mercoledì 15 agosto 2012 0.12 |
| segment-test.arff | mercoledì 15 agosto 2012 0.12 |
| soybean.arff | mercoledì 15 agosto 2012 0.12 |
| supermarket.arff | mercoledì 15 agosto 2012 0.12 |
| vote.arff | mercoledì 15 agosto 2012 0.12 |
| weather.arff | mercoledì 15 agosto 2012 0.12 |
| weather.nominal.arff | mercoledì 15 agosto 2012 0.12 |

File Format: Arff data files (*.arff)

Cancel | Choose

Visualize All

| No. | Name |
|-----|------|
| 1 | depa |
| 2 | depa |
| 3 | depa |
| 4 | depa |
| 5 | depa |
| 6 | depa |
| 7 | depa |
| 8 | depa |
| 9 | depa |
| 10 | groce |
| 11 | depa |

Remove

**Status**

OK

Log | 🐦 x 0

# Weka Knowledge Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Undo | Save...

## Filter

Choose | None | Apply

## Current relation

Relation: vote
Instances: 435      Attributes: 17

## Selected attribute

Name: handicapped-infants      Type: Nominal
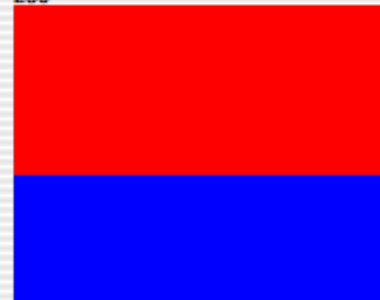Missing: 12 (3%)      Distinct: 2      Unique: 0 (0%)

| Label | Count |
| --- | --- |
| n | 236 |
| y | 187 |

## Attributes

| No. | Name |
| --- | --- |
| 1 | handicapped-infants |
| 2 | water-project-cost-sharing |
| 3 | adoption-of-the-budget-resolution |
| 4 | physician-fee-freeze |
| 5 | el-salvador-aid |
| 6 | religious-groups-in-schools |
| 7 | anti-satellite-test-ban |
| 8 | aid-to-nicaraguan-contras |
| 9 | mx-missile |
| 10 | immigration |
| 11 | synfuels-corporation-cutback |
| 12 | education-spending |
| 13 | superfund-right-to-sue |
| 14 | crime |
| 15 | duty-free-exports |
| 16 | export-administration-act-south-africa |
| 17 | Class |

Colour: Class (Nom) | Visualize All

236

187

## Status

OK

Log      x 0

# Weka Knowledge Explorer

**Preprocess** | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Undo | Save...

## Filter

Choose | None | Apply

### Current relation
Relation: vote
Instances: 435     Attributes: 17

### Selected attribute
Name: handicapped-infants     Type: Nominal
Missing: 12 (3%)     Distinct: 2     Unique: 0 (0%)

| Label | Count |
|-------|-------|
| n | 236 |
| y | 187 |

### Attributes

| No. | Name |
|-----|------|
| 1 | handicapped-infants |
| 2 | water-project-cost-sharing |
| 3 | adoption-of-the-budget-resolution |
| 4 | physician-fee-freeze |
| 5 | el-salvador-aid |
| 6 | religious-groups-in-schools |
| 7 | anti-satellite-test-ban |
| 8 | aid-to-nicaraguan-contras |
| 9 | mx-missile |
| 10 | immigration |
| 11 | synfuels-corporation-cutback |
| 12 | education-spending |
| 13 | superfund-right-to-sue |
| 14 | crime |
| 15 | duty-free-exports |
| 16 | export-administration-act-south-africa |
| 17 | Class |

Colour: Class (Nom) | Visualize All

236        187

## Status
OK

Log    x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | **Associate** | Select attributes | Visualize |

## Associator

[ Choose ] **Apriori** -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0

[ Start ] [ Stop ]

### Result list (right-click for options)

### Associator output

## Status

OK

[ Log ] x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Associator

Choose | Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0

Start | Stop

### Result list (right-click for optic

16:29:37 – Apriori

### Associator output

```
Minimum metric <confidence>: 0.9
Number of cycles performed: 11

Generated sets of large itemsets:

Size of set of large itemsets L(1): 20

Size of set of large itemsets L(2): 17

Size of set of large itemsets L(3): 6

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

Best rules found:

 1. adoption-of-the-budget-resolution=y physician-fee-freeze=n 219 ==> Class=democrat
 2. adoption-of-the-budget-resolution=y physician-fee-freeze=n aid-to-nicaraguan-con
 3. physician-fee-freeze=n aid-to-nicaraguan-contras=y 211 ==> Class=democrat 210
 4. physician-fee-freeze=n education-spending=n 202 ==> Class=democrat 201    conf:(1
 5. physician-fee-freeze=n 247 ==> Class=democrat 245    conf:(0.99)
 6. el-salvador-aid=n Class=democrat 200 ==> aid-to-nicaraguan-contras=y 197    conf
 7. el-salvador-aid=n 208 ==> aid-to-nicaraguan-contras=y 204    conf:(0.98)
 8. adoption-of-the-budget-resolution=y aid-to-nicaraguan-contras=y Class=democrat 20
 9. el-salvador-aid=n aid-to-nicaraguan-contras=y 204 ==> Class=democrat 197    conf
10. aid-to-nicaraguan-contras=y Class=democrat 218 ==> physician-fee-freeze=n 210
```

## Status

OK

Log | x 0

# Additional features of Explorer: Attribute Selection and Visualization

# Explorer: attribute selection

- Panel that can be used to investigate which (subsets of) attributes are the most predictive ones

- Attribute selection methods contain two parts:
  - A search method: best-first, forward selection, random, exhaustive, genetic algorithm, ranking
  - An evaluation method: correlation-based, wrapper, information gain, chi-squared, …

- Very flexible: WEKA allows (almost) arbitrary combinations of these two

- Will see in more detail in dedicated labs

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | **Select attributes** | Visualize |

## Attribute Evaluator

[ Choose ] **CfsSubsetEval**

## Search Method

[ Choose ] **BestFirst** -D 1 -N 5

## Attribute Selection Mode

- ⦿ Use full training set
- ◯ Cross-validation    Folds   10
               Seed   1

(Nom) Class

[ Start ]   [ Stop ]

## Result list (right-click for options)

## Attribute selection output

## Status

OK

[ Log ]   x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | **Select attributes** | Visualize |

## Attribute Evaluator

[ Choose ]  **CfsSubsetEval**

## Search Method

[ Choose ]  **BestFirst** -D 1 -N 5

## Attribute Selection Mode

- ⦿ Use full training set
- ◯ Cross-validation    Folds [ 10 ]
-                        Seed  [ 1 ]

[ (Nom) Class ▾ ]

[ Start ]    [ Stop ]

## Result list (right-click for options)

16:39:40 – BestFirst + CfsSubsetEval

## Attribute selection output

```
                duty-free-exports
                export-administration-act-south-africa
                Class
Evaluation mode:    evaluate on all training data



=== Attribute Selection on all input data ===

Search Method:
        Best first.
        Start set: no attributes
        Search direction: forward
        Stale search after 5 node expansions
        Total number of subsets evaluated: 83
        Merit of best subset found:    0.729

Attribute Subset Evaluator (supervised, Class (nominal): 17 Class):
        CFS Subset Evaluator

Selected attributes: 4 : 1
                    physician-fee-freeze
```

## Status

OK

[ Log ]    🐑 x 0

## Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | **Select attributes** | Visualize |

### Attribute Evaluator

[ Choose ] **CfsSubsetEval**

### Search Method

[ Choose ] **BestFirst -D 1 -N 5**

### Attribute Selection Mode

○ Use full training set
○ Cross-validation  Folds  10
  Seed  1

(Nom) Class

[ Start ]  [ Stop ]

### Result list (right-click for options)

16:39:40 – BestFirst + CfsSubsetEval

### Attribute selection output

```
                duty-free-exports
                export-administration-act-south-africa
                Class
Evaluation mode:    evaluate on all training data



=== Attribute Selection on all input data ===

Search Method:
        Best first.
        Start set: no attributes
        Search direction: forward
        Stale search after 5 node expansions
        Total number of subsets evaluated: 83
        Merit of best subset found:    0.729

Attribute Subset Evaluator (supervised, Class (nominal): 17 Class):
        CFS Subset Evaluator

Selected attributes: 4 : 1
                physician-fee-freeze
```

### Status

OK

[ Log ]  x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

## Attribute Evaluator

- weka
  - ▼ attributeSelection
    - CfsSubsetEval
    - ClassifierSubsetEval
    - WrapperSubsetEval
    - ConsistencySubsetEval
    - ReliefFAttributeEval
    - InfoGainAttributeEval
    - GainRatioAttributeEval
    - SymmetricalUncertAttributeEval
    - OneRAttributeEval
    - ChiSquaredAttributeEval
    - PrincipalComponents
    - SVMAttributeEval

```
                     duty-free-exports
                     export-administration-act-south-africa
                     Class
uation mode:    evaluate on all training data




Attribute Selection on all input data ===

ch Method:
      Best first.
      Start set: no attributes
      Search direction: forward
      Stale search after 5 node expansions
      Total number of subsets evaluated: 83
      Merit of best subset found:    0.729

ibute Subset Evaluator (supervised, Class (nominal): 17 Class):
      CFS Subset Evaluator


Selected attributes: 4 : 1
                     physician-fee-freeze
```

## Status

OK

Log    x 0

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

**Attribute Evaluator**

[ Choose ] **InfoGainAttributeEval**

**Search Method**

- 📁 weka
  - ▼ 📁 attributeSelection
    - 📄 BestFirst
    - 📄 ForwardSelection
    - 📄 RaceSearch
    - 📄 GeneticSearch
    - 📄 RandomSearch
    - 📄 ExhaustiveSearch
    - 📄 Ranker
    - 📄 RankSearch

E308 -N -1

te selection output

```
            duty-free-exports
            export-administration-act-south-africa
            Class
uation mode:    evaluate on all training data



Attribute Selection on all input data ===

ch Method:
    Best first.
    Start set: no attributes
    Search direction: forward
    Stale search after 5 node expansions
    Total number of subsets evaluated: 83
    Merit of best subset found:    0.729

ibute Subset Evaluator (supervised, Class (nominal): 17 Class):
    CFS Subset Evaluator


cted attributes: 4 : 1
            physician-fee-freeze
```

**Status**

OK

[ Log ]    🐑 x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | **Select attributes** | Visualize |

## Attribute Evaluator

[ Choose ]  **InfoGainAttributeEval**

## Search Method

[ Choose ]  **Ranker** -T -1.7976931348623157E308 -N -1

## Attribute Selection Mode

( ● ) Use full training set
( ○ ) Cross-validation    Folds [ 10 ]
                          Seed  [ 1 ]

[ (Nom) Class                        ▲▼ ]

[ Start ]   [ Stop ]

## Result list (right-click for options)

16:39:40 – BestFirst + CfsSubsetEval
16:43:05 – Ranker + InfoGainAttributeEval

## Attribute selection output

```
        Information Gain Ranking Filter

Ranked attributes:
    0.7078541    4 physician-fee-freeze
    0.4185726    3 adoption-of-the-budget-resolution
    0.4028397    5 el-salvador-aid
    0.34036     12 education-spending
    0.3123121   14 crime
    0.3095576    8 aid-to-nicaraguan-contras
    0.2856444    9 mx-missile
    0.2121705   13 superfund-right-to-sue
    0.2013666   15 duty-free-exports
    0.1902427    7 anti-satellite-test-ban
    0.1404643    6 religious-groups-in-schools
    0.1211834    1 handicapped-infants
    0.1007458   11 synfuels-corporation-cutback
    0.0529956   16 export-administration-act-south-africa
    0.0049097   10 immigration
    0.0000117    2 water-project-cost-sharing

Selected attributes: 4,3,5,12,14,8,9,13,15,7,6,1,11,16,10,2 : 16
```

## Status

OK

[ Log ]   🐑 x 0

# Explorer: data visualization

- Visualization very useful in practice: e.g. helps to determine difficulty of the learning problem

- WEKA can visualize single attributes (1-d) and pairs of attributes (2-d)

- Color-coded class values

- "Jitter" option to deal with nominal attributes (and to detect "hidden" data points). (Jittering occurs when you have too many instances placed on the same point, see http://blogs.sas.com/content/iml/2011/07/05/jittering-to-prevent-overplotting-in-statistical-graphics.html)

- "Zoom-in" function

# Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Generate... | Undo | Edit... | Save...

## Filter

Choose | Di

## Open

data

| Name | Date Modified |
|------|---------------|
| contact-lenses.arff | mercoledì 15 agosto 2012 0.12 |
| cpu.arff | mercoledì 15 agosto 2012 0.12 |
| cpu.with.vendor.arff | mercoledì 15 agosto 2012 0.12 |
| diabetes.arff | mercoledì 15 agosto 2012 0.12 |
| glass.arff | mercoledì 15 agosto 2012 0.12 |
| ionosphere.arff | mercoledì 15 agosto 2012 0.12 |
| iris.arff | mercoledì 15 agosto 2012 0.12 |
| labor.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersCorn-test.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersCorn-train.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersGrain-test.arff | mercoledì 15 agosto 2012 0.12 |
| ReutersGrain-train.arff | mercoledì 15 agosto 2012 0.12 |
| segment-challenge.arff | mercoledì 15 agosto 2012 0.12 |

File Format: Arff data files (*.arff)

Cancel | Choose

Apply

## Current relation

Relation: Glas
Instances: 214

ype: Numeric
que: 145 (68%)

## Attributes

All

| No. | Name |
|-----|------|
| 1 | RI |
| 2 | Na |
| 3 | Mg |
| 4 | Al |
| 5 | Si |
| 6 | K |
| 7 | Ca |
| 8 | Ba |
| 9 | Fe |
| 10 | Type |

Visualize All

Remove

3   4                        4   3   3      0   1   1

1.51                         1.52                    1.53

## Status

OK

Log      x 0

# Weka Knowledge Explorer

| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

Open file...  Open URL...  Open DB...  Undo  Save...

## Filter

Choose  **None**  Apply

## Current relation

Relation: Glass
Instances: 214          Attributes: 10

## Attributes

| No. | Name |
|-----|------|
| 1 | RI |
| 2 | Na |
| 3 | Mg |
| 4 | Al |
| 5 | Si |
| 6 | K |
| 7 | Ca |
| 8 | Ba |
| 9 | Fe |
| 10 | Type |

## Selected attribute

Name: RI                          Type: Numeric
Missing: 0 (0%)     Distinct: 178     Unique: 145 (68%)

| Statistic | Value |
|-----------|-------|
| Minimum | 1.511 |
| Maximum | 1.534 |
| Mean | 1.518 |
| StdDev | 0.003 |

Colour: Type (Nom)          Visualize All

45  48  46
18
15
3  3  1  5        7    9   5      1  1    4    1  0  0  1  0  1

1.51              1.52              1.53

## Status

OK          Log          x 0

Weka Knowledge Explorer: Visualizing Glass

X: Al (Num)    Y: Ca (Num)

Colour: Type (Nom)    Select Instance

Reset    Clear    Save    Jitter

Plot: Glass

16.19

10.81

5.43

0.29    1.9    3.5

Class colour

build wind float    build wind non-float    vehic wind float

vehic wind non-float    containers    tableware    headlamps

X: Al (Num)

Y: Ca (Num)
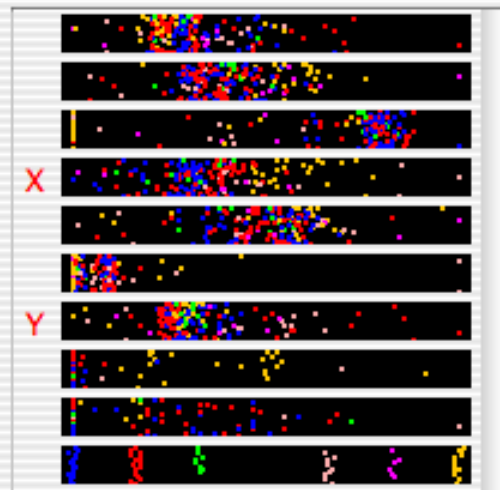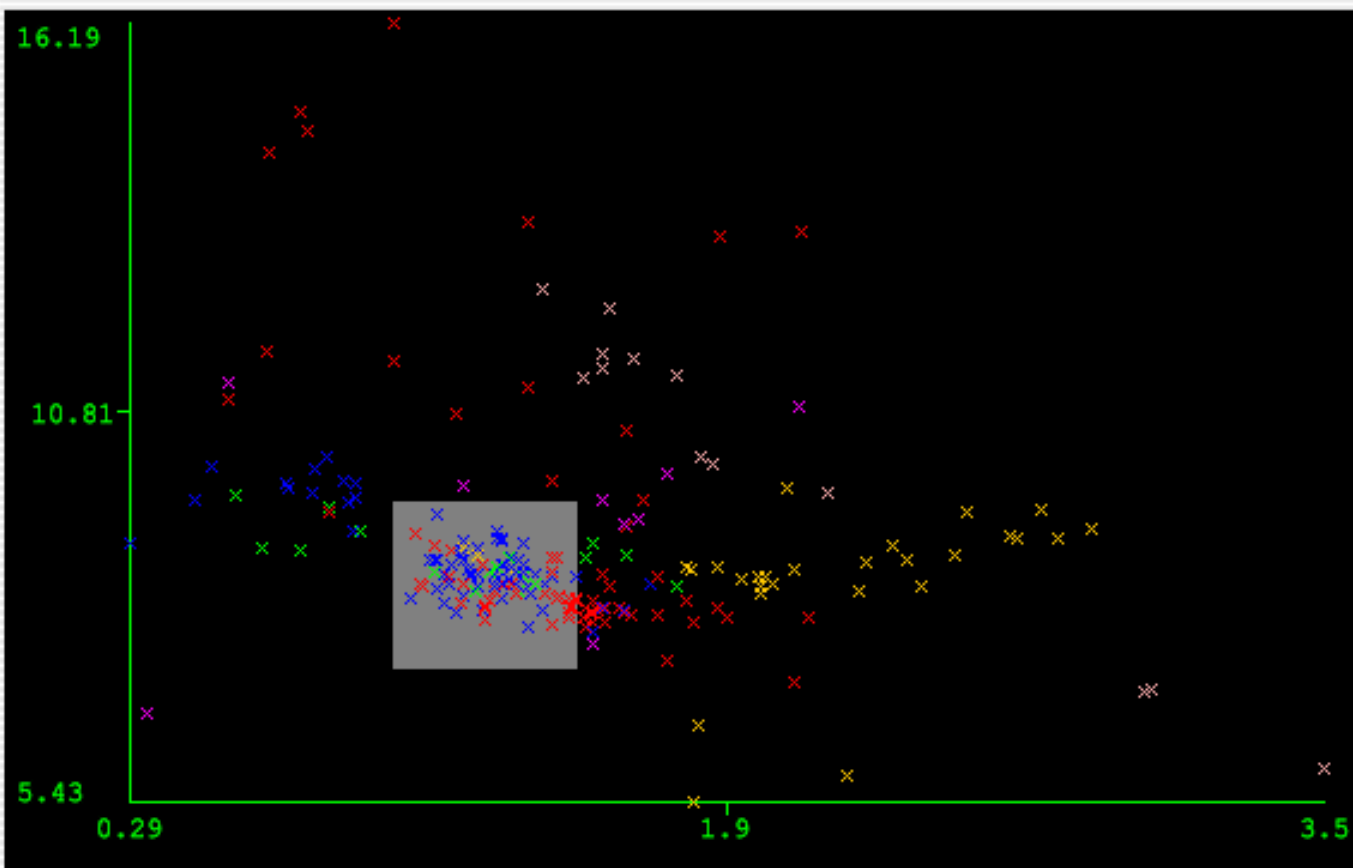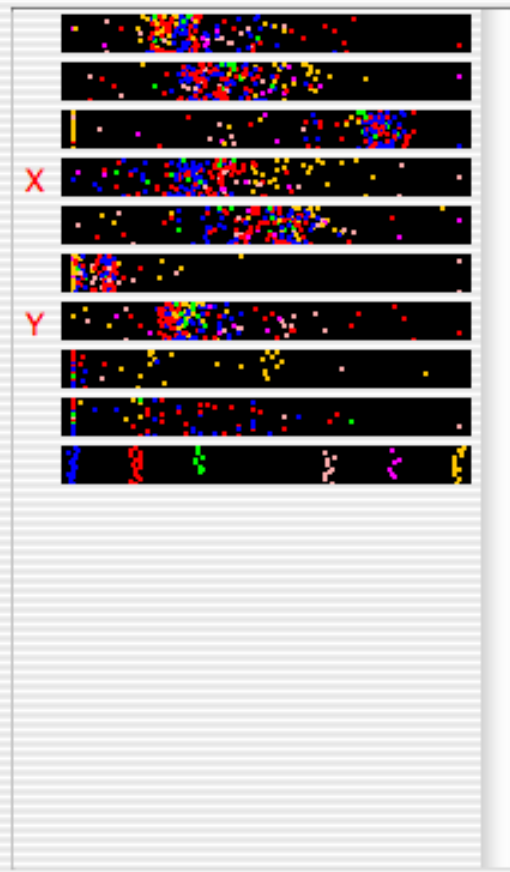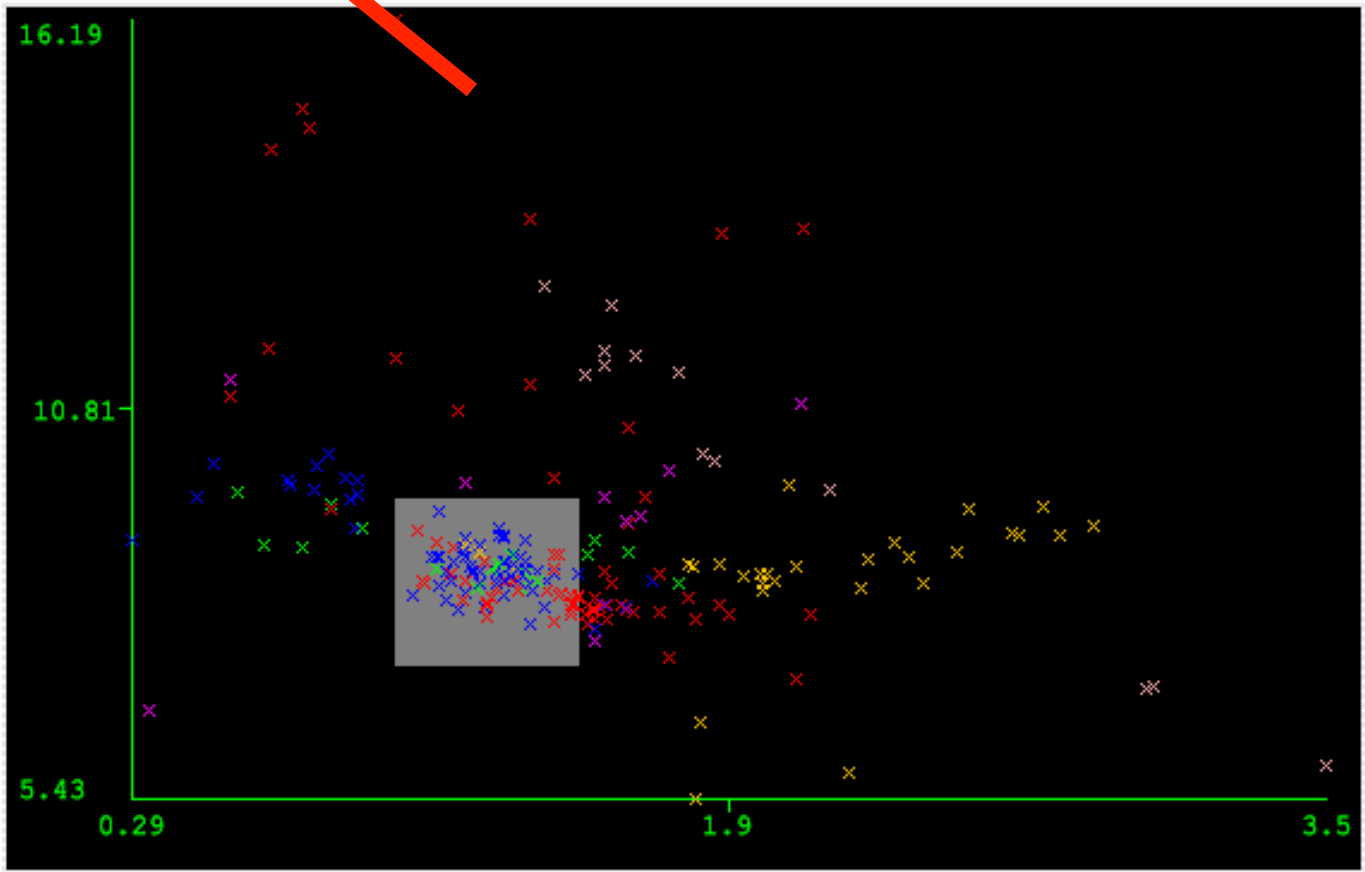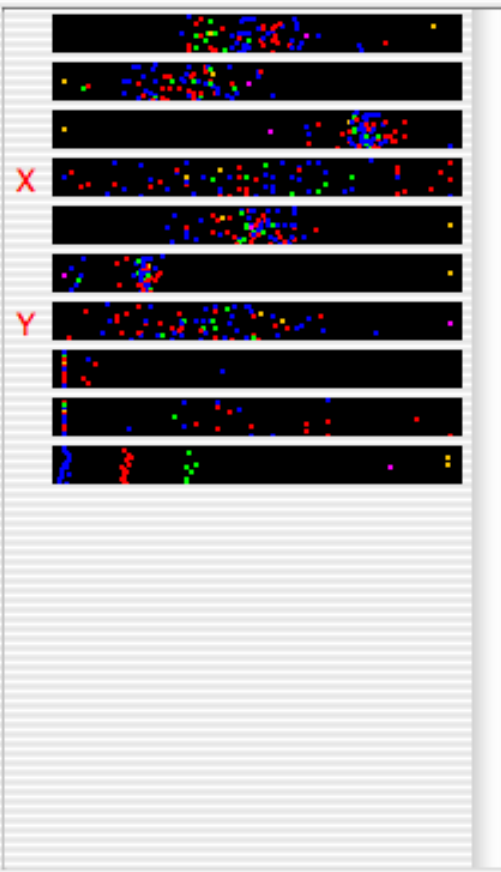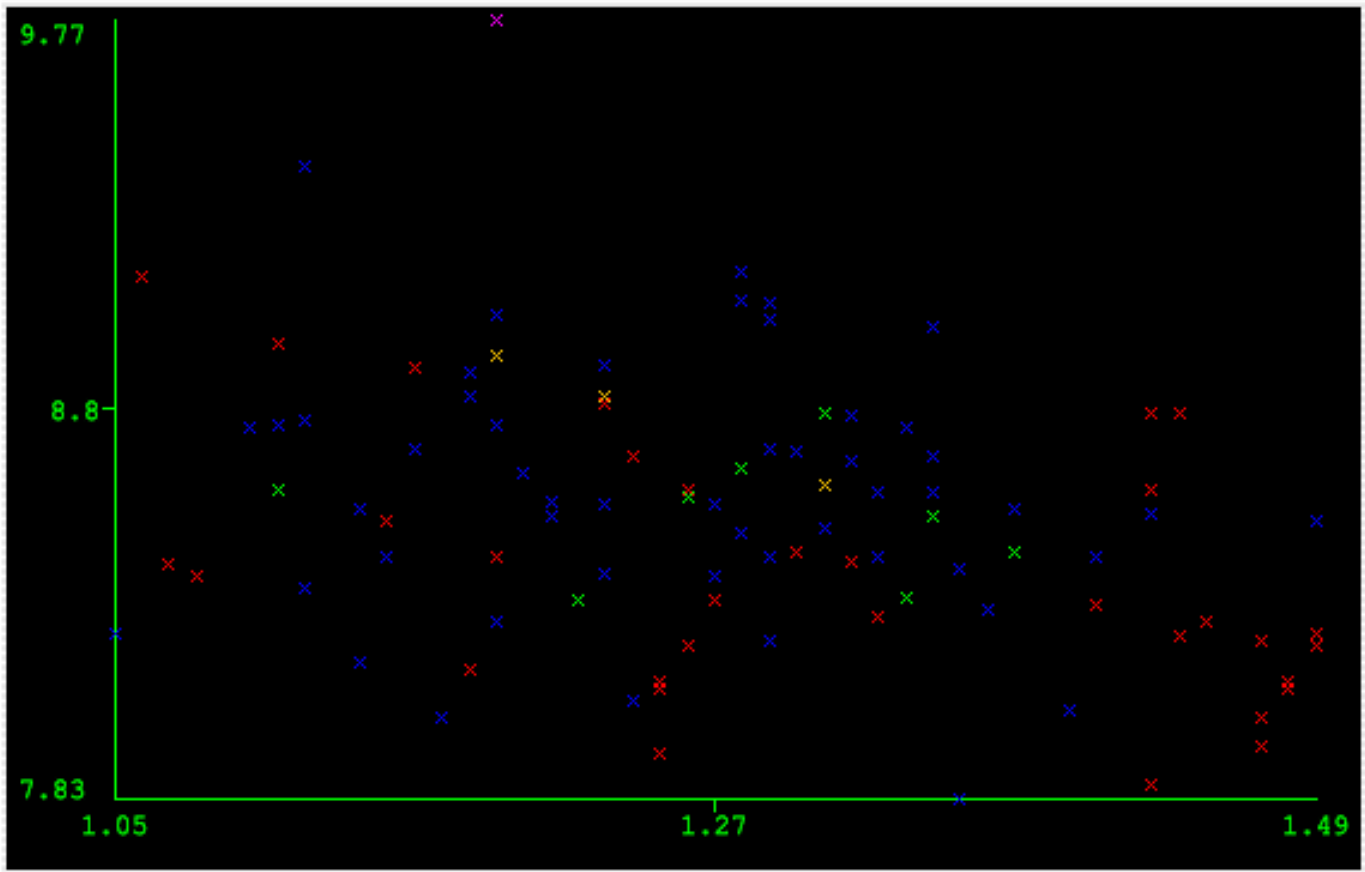
Colour: Type (Nom)

Rectangle

Submit    Clear    Save

Jitter

Plot: Glass



16.19

10.81

5.43

0.29                1.9                3.5

X

Y

Class colour

build wind float build wind non-float vehic wind float vehic wind non-float containers tableware headlamps

# References and Resources

- References:
  - WEKA website: http://www.cs.waikato.ac.nz/~ml/weka/index.html
  - WEKA Tutorial:
    - Machine Learning with WEKA: A presentation demonstrating all graphical user interfaces (GUI) in Weka.
    - A presentation which explains how to use Weka for exploratory data mining.
  - WEKA Data Mining Book:
    - Ian H. Witten and Eibe Frank, Data Mining: Practical Machine Learning Tools and Techniques (Second Edition)
  - WEKA Wiki: http://weka.sourceforge.net/wiki/index.php/Main_Page
  - Others:
    - Jiawei Han and Micheline Kamber, Data Mining: Concepts and Techniques, 2nd ed.