# DATA WAREHOUSING AND DATA MINING

NAME                : MOTHISHWARAN C.
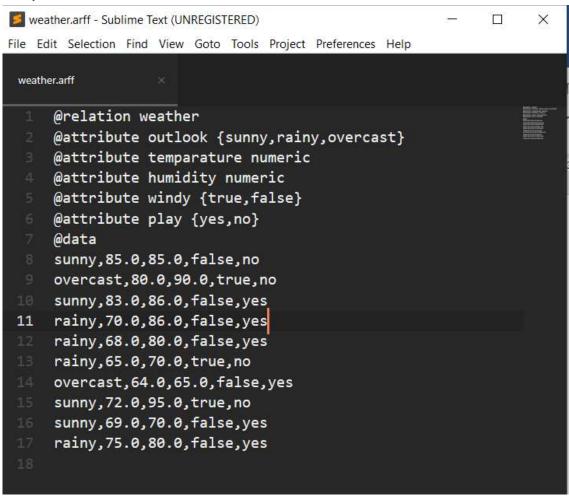
REG No             : 19MID0017

COURSE CODE     : CSI3010

SLOT                   : L39+L40

FACULTY           : DR. CHELLATHAMILAN T.

## DISCRETIZATION OF DATA USING WEKA

### 1) Weather dataset:



```
@relation weather
@attribute outlook {sunny,rainy,overcast}
@attribute temparature numeric
@attribute humidity numeric
@attribute windy {true,false}
@attribute play {yes,no}
@data
sunny,85.0,85.0,false,no
overcast,80.0,90.0,true,no
sunny,83.0,86.0,false,yes
rainy,70.0,86.0,false,yes
rainy,68.0,80.0,false,yes
rainy,65.0,70.0,true,no
overcast,64.0,65.0,false,yes
sunny,72.0,95.0,true,no
sunny,69.0,70.0,false,yes
rainy,75.0,80.0,false,yes
```

## Weka Explorer

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

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

**Filter**

Choose | None | Apply | Stop

**Current relation**

Relation: weather     Attributes: 5
Instances: 10     Sum of weights: 10

**Selected attribute**

Name: outlook     Type: Nominal
Missing: 0 (0%)     Distinct: 3     Unique: 0 (0%)

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | sunny | 4 | 4.0 |
| 2 | rainy | 4 | 4.0 |
| 3 | overcast | 2 | 2.0 |

**Attributes**

All | None | Invert | Pattern

| No. | Name |
|---|---|
| 1 | outlook |
| 2 | temparature |
| 3 | humidity |
| 4 | windy |
| 5 | play |

Remove

Class: play (Nom) | Visualize All

**Status**

OK     Log   x 0

---

## Viewer

Relation: weather

| No. | 1: outlook Nominal | 2: temparature Numeric | 3: humidity Numeric | 4: windy Nominal | 5: play Nominal |
|---|---|---|---|---|---|
| 1 | sunny | 85.0 | 85.0 | false | no |
| 2 | overcast | 80.0 | 90.0 | true | no |
| 3 | sunny | 83.0 | 86.0 | false | yes |
| 4 | rainy | 70.0 | 86.0 | false | yes |
| 5 | rainy | 68.0 | 80.0 | false | yes |
| 6 | rainy | 65.0 | 70.0 | true | no |
| 7 | overcast | 64.0 | 65.0 | false | yes |
| 8 | sunny | 72.0 | 95.0 | true | no |
| 9 | sunny | 69.0 | 70.0 | false | yes |
| 10 | rainy | 75.0 | 80.0 | false | yes |

Add instance | Undo | OK | Cancel

# Before Discretization (no bins added):

## After Discretization (4 bins added):

weka.gui.GenericObjectEditor                                    ✕

weka.filters.unsupervised.attribute.Discretize

**About**

An instance filter that discretizes a range of numeric attributes in the dataset into nominal attributes.

More

Capabilities

| | |
|---|---|
| attributeIndices | 2 |
| binRangePrecision | 3 |
| bins | 4 |
| debug | False ▼ |
| desiredWeightOfInstancesPerInterval | -1.0 |
| doNotCheckCapabilities | False ▼ |
| findNumBins | False ▼ |
| ignoreClass | False ▼ |
| invertSelection | False ▼ |
| makeBinary | False ▼ |
| spreadAttributeWeight | False ▼ |
| useBinNumbers | False ▼ |
| useEqualFrequency | False ▼ |

Open...      Save...      OK      Cancel

**Weka Explorer** (Preprocess tab)

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

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

**Filter**

Choose | Discretize -B 4 -M -1.0 -R 2 -precision 3 | Apply | Stop

**Current relation**

Relation: weather-weka.filters.unsupervised.attribute.Discretize... | Attributes: 5
Instances: 10 | Sum of weights: 10

**Selected attribute**

Name: temparature | Type: Nominal
Missing: 0 (0%) | Distinct: 4 | Unique: 1 (10%)

| No. | Label | Count | Weight |
|-----|-------|-------|--------|
| 1 | '(-inf-69.25]' | 4 | 4.0 |
| 2 | '(69.25-74.5]' | 2 | 2.0 |
| 3 | '(74.5-79.75]' | 1 | 1.0 |
| 4 | '(79.75-inf)' | 3 | 3.0 |

**Attributes**

All | None | Invert | Pattern

| No. | Name |
|-----|------|
| 1 | outlook |
| 2 | temparature |
| 3 | humidity |
| 4 | windy |
| 5 | play |

Remove

Class: play (Nom) | Visualize All

**Status**

OK | Log | x 0

---



**Weka Explorer** (Classify tab)

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 | % | 70

More options...

(Nom) play

Start | Stop

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

21:09:05 - trees.J48
21:16:27 - trees.J48

**Classifier output**

```
=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances           1               33.3333 %
Incorrectly Classified Instances         2               66.6667 %
Kappa statistic                          0
Mean absolute error                      0.5714
Root mean squared error                  0.6061
Relative absolute error                102.8571 %
Root relative squared error            104.9781 %
Total Number of Instances                3

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Cl
                 1.000    1.000    0.333      1.000   0.500      ?      0.500     0.333     ye
                 0.000    0.000    ?          0.000   ?          ?      0.500     0.667     no
Weighted Avg.    0.333    0.333    ?          0.333   ?          ?      0.500     0.556

=== Confusion Matrix ===

 a b   <-- classified as
 1 0 | a = yes
 2 0 | b = no
```

**Status**

OK | Log | x 0

Accuracy varies when we change the split percentages.

| Split % | Before Discretization | After Discretization |
|---------|----------------------|----------------------|
| 70% | 33.33% | 33.33% |
| 50% | 40% | 40% |

# 1) Pima Diabetes Dataset:

# Before Discretization (no bins added):

**Accuracy: 76.52%**

## After Discretization (5 bins added):

weka.gui.GenericObjectEditor ✕

weka.filters.unsupervised.attribute.Discretize

**About**

An instance filter that discretizes a range of numeric attributes in the dataset into nominal attributes.

[ More ]

[ Capabilities ]

| | |
|---|---|
| attributeIndices | 2,5| |
| binRangePrecision | 3 |
| bins | 5 |
| debug | False ▼ |
| desiredWeightOfInstancesPerInterval | -1.0 |
| doNotCheckCapabilities | False ▼ |
| findNumBins | False ▼ |
| ignoreClass | False ▼ |
| invertSelection | False ▼ |
| makeBinary | False ▼ |
| spreadAttributeWeight | False ▼ |
| useBinNumbers | False ▼ |
| useEqualFrequency | False ▼ |

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

Weka Explorer

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

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

**Filter**

Choose | **Discretize** -B 5 -M -1.0 -R 2,5 -precision 3 | Apply | Stop

**Current relation**

Relation: pima_diabetes-weka.filters.unsupervised.attribute.Di...  Attributes: 9
Instances: 768  Sum of weights: 768

**Selected attribute**

Name: plas  Type: Nominal
Missing: 0 (0%)  Distinct: 5  Unique: 0 (0%)

| No. | Label | Count | Weight |
|-----|-------|-------|--------|
| 1 | '(-inf-39.8]' | 5 | 5.0 |
| 2 | '(39.8-79.6]' | 36 | 36.0 |
| 3 | '(79.6-119.4]' | 367 | 367.0 |
| 4 | '(119.4-159.2]' | 258 | 258.0 |
| 5 | '(159.2-inf)' | 102 | 102.0 |

**Attributes**

All | None | Invert | Pattern

| No. | Name |
|-----|------|
| 1 | preg |
| 2 | plas |
| 3 | pres |
| 4 | skin |
| 5 | insu |
| 6 | mass |
| 7 | pedi |
| 8 | age |
| 9 | class |

Remove

Class: class (Nom) | Visualize All

**Status**

OK | Log | x 0

---

Weka Explorer

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

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

**Filter**

Choose | **Discretize** -B 5 -M -1.0 -R 2,5 -precision 3 | Apply | Stop

**Current relation**

Relation: pima_diabetes-weka.filters.unsupervised.attribute.Di...  Attributes: 9
Instances: 768  Sum of weights: 768

**Selected attribute**

Name: insu  Type: Nominal
Missing: 0 (0%)  Distinct: 5  Unique: 0 (0%)

| No. | Label | Count | Weight |
|-----|-------|-------|--------|
| 1 | '(-inf-169.2]' | 642 | 642.0 |
| 2 | '(169.2-338.4]' | 100 | 100.0 |
| 3 | '(338.4-507.6]' | 17 | 17.0 |
| 4 | '(507.6-676.8]' | 6 | 6.0 |
| 5 | '(676.8-inf)' | 3 | 3.0 |

**Attributes**

All | None | Invert | Pattern

| No. | Name |
|-----|------|
| 1 | preg |
| 2 | plas |
| 3 | pres |
| 4 | skin |
| 5 | insu |
| 6 | mass |
| 7 | pedi |
| 8 | age |
| 9 | class |

Remove

Class: class (Nom) | Visualize All

**Status**

OK | Log | x 0

**Accuracy: 75.65%**

Accuracy varies when we change the split percentages.

| Split % | Before Discretization | After Discretization |
|---------|----------------------|----------------------|
| 70% | 76.52% | 75.65% |
| 50% | 74.21% | 71.35% |