Hashemite University
Prince Al-Hussein bin Abdullah II Faculty for
Information Technology
Department of Computer Information Systems

| For Instructor Use | |
|---|---|
| **Course Name** | Data Mining |
| **Course ID** | 151002351 |
| **Academic Year** | 2022/2023 |
| **Semester** | 2$^{nd}$ Semester |
| **Assignment** | 1 and 2 |
| **Due Date** | May-2023 |

| For Student Use | | | |
|---|---|---|---|
| Student Name | Student Id | Section Id | Seat Number |
| Rand Bassam Mohammad Wahdan | 2034877 | 2 | 19 |
| Ahmad Rawhi Mohammad Al-Qranawy | 2042497 | 1 | 59 |

# Instructor Name : Dr. Subhieh Elsalhi

**Part 1**

**QA.A1-**

| Number | Attribute Name | Attribute Type |
|--------|----------------|----------------|
| 1 | Age | Numeric (ratio) |
| 2 | Workclass | Nominal (symmetric) |
| 3 | Education | Nominal (ordinal) |
| 4 | Education-num | Numeric (ratio) |
| 5 | Marital_Status | Nominal |
| 6 | Occupation | Nominal (symmetric) |
| 7 | Relationship | Nominal (symmetric) |
| 8 | Race | Nominal (symmetric) |
| 9 | Gender | Nominal (symmetric) |
| 10 | Capital-gain | Numeric (ratio) |
| 11 | Capital-loss | Numeric (ratio) |
| 12 | Hours-per-week | Numeric (interval) |
| 13 | Native-country | Nominal |
| 14 | Fnlwgt | Nominal |

**A2-**

| Attribute Name | Location | Median | Mean | Variance | Mode |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Age | First | 37 | 38.58163 | 186.067 | 36 |
| Workclass | Second | - | - | - | Private |
| Education | Third | - | - | - | HS-grade |
| Education-num | Fourth | 10 | 10.08059 | 6.6188 | 36 |
| Marital_Status | Fifth | - | - | - | Married-civ-spouse |
| Occupation | Sixth | - | - | - | Prof-specialty |
| Relationship | Seventh | - | - | - | Husband |
| Race | Eighth | - | - | - | White |
| Gender | Nineth | - | - | - | Male |
| Capital-gain | Tenth | 0 | 1077.615 | 54544177.472 | 0 |
| Capital-loss | Eleventh | 0 | 87.30651 | 162381.6777 | 0 |
| Hours-per-week | Twelfth | 40 | 40.43747 | 152.4637 | 40 |
| Native-country | Thirteenth | - | - | - | United-States |
| Fnlwgt | Fourteenth | - | - | - | <=50k |

| Attribute Name | Max | Min | Range |
|:---:|:---:|:---:|:---|
| Age | 90 | 17 | 90-17=73 |
| Workclass | - | - | 8 |
| Education | - | - | 16 |
| Education-num | 16 | 1 | 16-1=15 |
| Marital_Status | - | - | 7 |
| Occupation | - | - | 14 |
| Relationship | - | - | 6 |
| Race | - | - | 5 |
| Gender | - | - | 2 |
| Capital-gain | 99999 | 0 | 99999 |
| Capital-loss | 4356 | 0 | 4356 |
| Hours-per-week | 99 | 1 | 99-1=98 |
| Native-country | - | - | 41 |
| Fnlwgt | - | - | 2 |

# frequency of values

## (supervised→attribute→discretize)

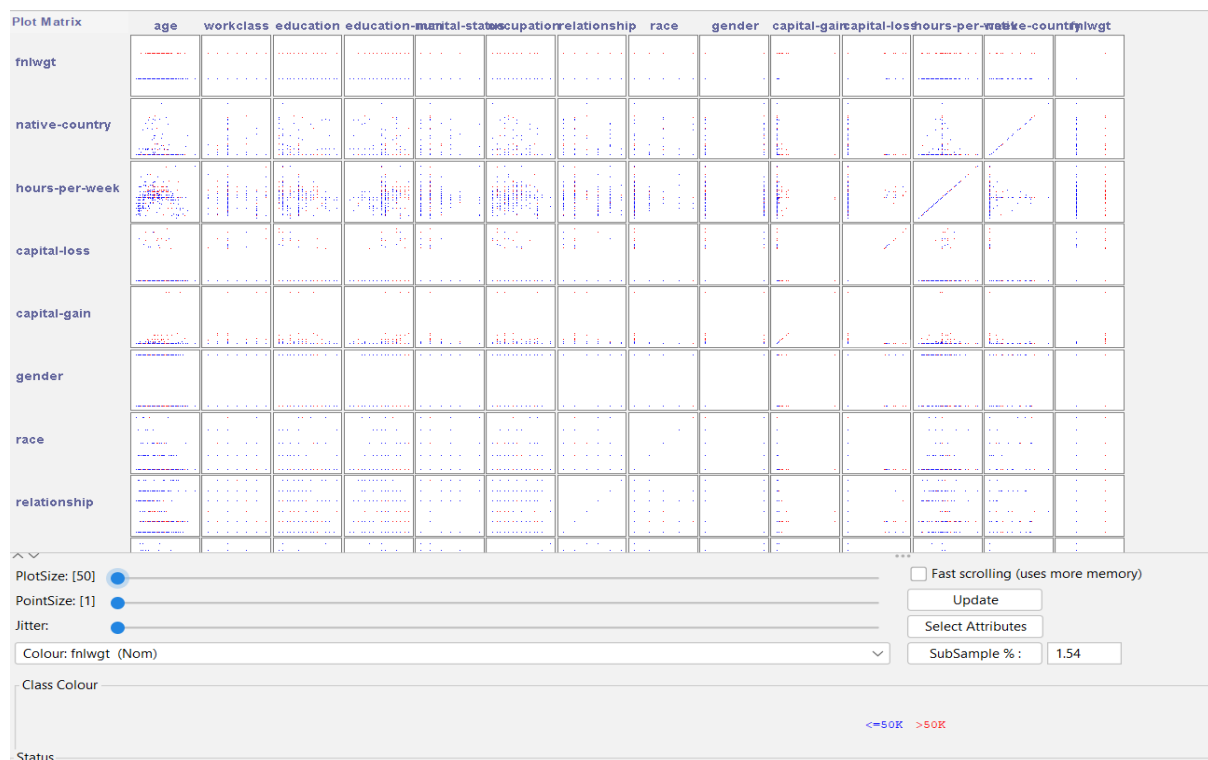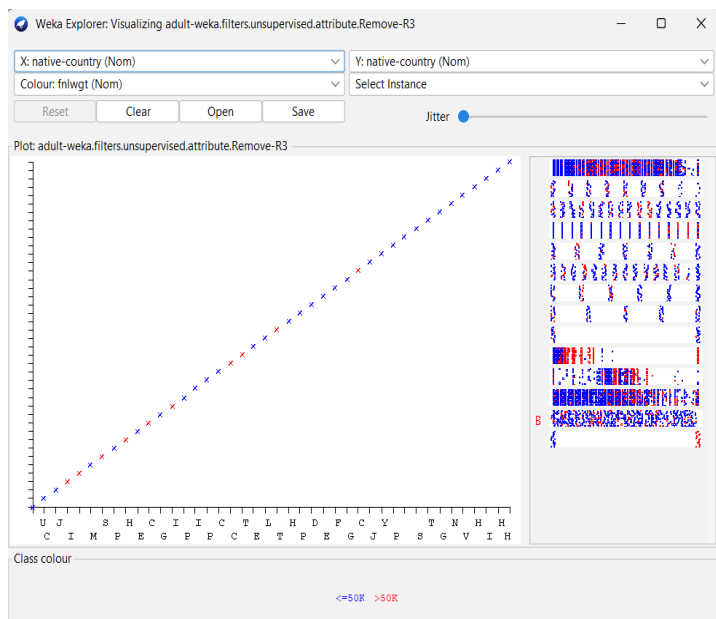| Attribute Name | frequency of values |
|---|---|
| Age | -∞-21.5:3130/21.5-23.5:1642/23.5-27.5:3259/27.5-29.5:1680/29.5-35.5:5214/35.5-43.5:6551/43.5-54.5:6577/54.5-61.5:2476/61.5-∞:**2032** |
| Workclass | Self-emp-not-inc:2541/Private:22696/State-gov:1297/federal-gov:960/Local-gov:2093/Self-emp-inc:1116/Without-pay:14/Never-worked:7 |
| Education | Bachelors:5354/HS:grade:10501/11th:1175/Masters:1723/9th:514/Some-colleage:7291/Assoc-acdm:1067/Assoc-voc:1382/7th-8th:646/Doctorate:413/Prof-school:576/5th-6th:333/10th:933/1st-4th:168/Preschool:51/12th:433 |
| Education-num | -∞-8.5:4253/8.5-9.5:10501/9.5-10.5:7291/10.5-12.5:2449/12.5-13.5:5355/13.5-14.5:1723/14.5-∞:**989** |
| Marital_Status | Maried-civ-supouse:14976/Divorced:4443/Married-spouse-absent:418/Never-married:10682/Separated:1025/Married-AF-spouse:23/Widowed:993 |
| Occupation | Exec-managerial:4066/Handlers-cleaners:1370/Prof-specialty:4140/Other-service:3295/Adm-clerical:3769/Sales:3650/Craft-repair:4099/Transport-moving:1597/Farming-fishing:994/Machine-op-inspct:2002/Tech-support:928/Protective-serv:649/Armed-Forces:9/Priv-house-serv:149 |
| Relationship | Husband:13193/Not-in-family:8304/wife:1568/own-child:5068/Unmarried:3446/Other-relative:981 |
| Race | White:27815/Black:3124/Asian-pac-Islander:1039/Amer-Indian-Eskimo:311/Other:271 |
| Gender | Male:21789/Female:10771 |
| Capital-gain | -∞-57:29849/57-3048:472/3048-3120:97/3120-4243.5:309/4243.5-4401:70/4410-4668.5:65/4668.5-4826:26/48226-4973.5:18/4932.5-4973.5:7/4973.5-5119:70/5119-5316.5:97/5316.5-5505.5:11/5505.5-6618.5:37/6618.5-7073.5:34/7073.5-∞:1399 |
| Capital-loss | -∞-1551.5:31197/1551.5-15685:25/1568.5-1820.5:348/1820.5-1862:56/1862-1881.5:39/1881.5-1923:361/1923-1975.5:19/1975.5-1978.5:168/1978.5-2168.5:111/2168.5-2176.5:7/2176.5-2218.5:31/2218.5-2384.5:79/2384.5-2450.5:70/2450.5-3726.5:43/3726.5-∞:7 |
| Hours-per-week | -∞-34.5:5583/34.5-39.5:2180/39.5-41.5:15253/41.5-49.5:3083/49.5-65.5:5640/65.5-∞:**822** |
| Native-country | United_States:29169/Cuba:95/Jamaica:81/India:100/Mexico:643/South:80/Puerto-Rico:114/Honduras:13/England:90/Canada:121/Germany:137/Iran:43/Philippiens:198/Italy:43/Poland:60/Calumbia:59/Thailand:18/Ecuador:28/Laos:18/Taiwan:51/Haiti:44/Portugal:37/Dominican-Republic:70/El-Salvador:106/France:29/Guatemala:64/China:75/Japan:62/Yugoslavia:16/Peru:31/Outlying-US(Guam-USVI-etc):14/Scotland:12/Trinadad&Tobago:19/Greece:29/Nicaragua:34/Vietnam:67/Hong:20/Ireland:24/Hungary:13/Holand-Netherlands:1 |
| Fnlwgt | <=50k:24719/>=50k:7841 |

# A3- Distributions of each Attributes:

**Scatter plot**

- Provides a first look at bivariate data to see clusters of points, outliers .
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane.

## By choose visualize

positively correlated.
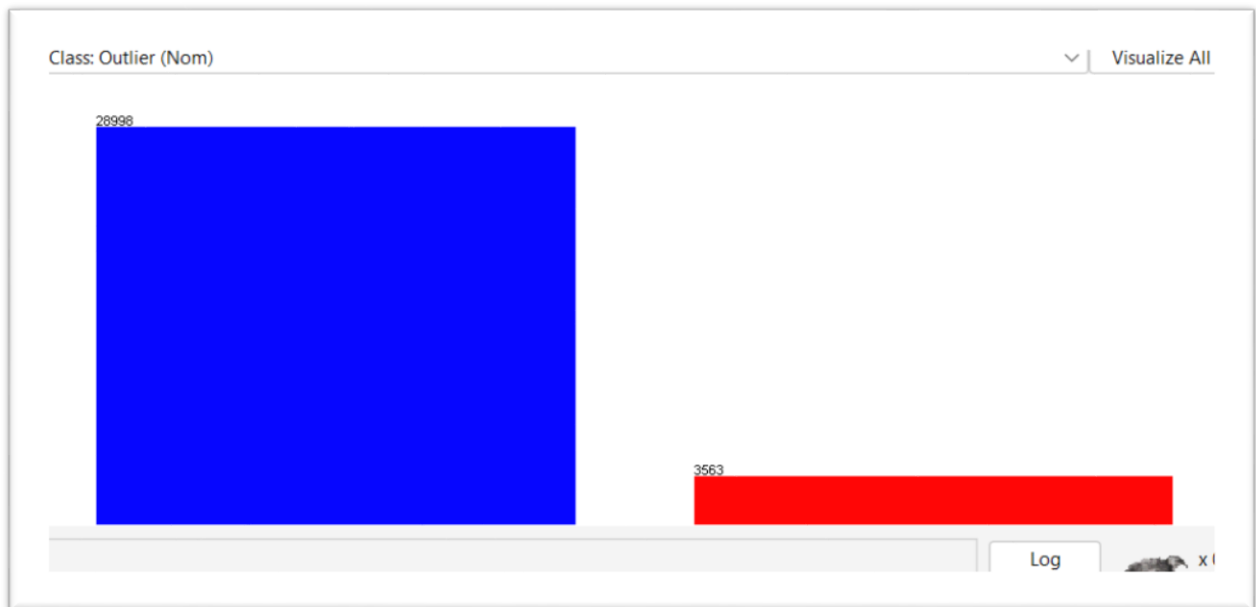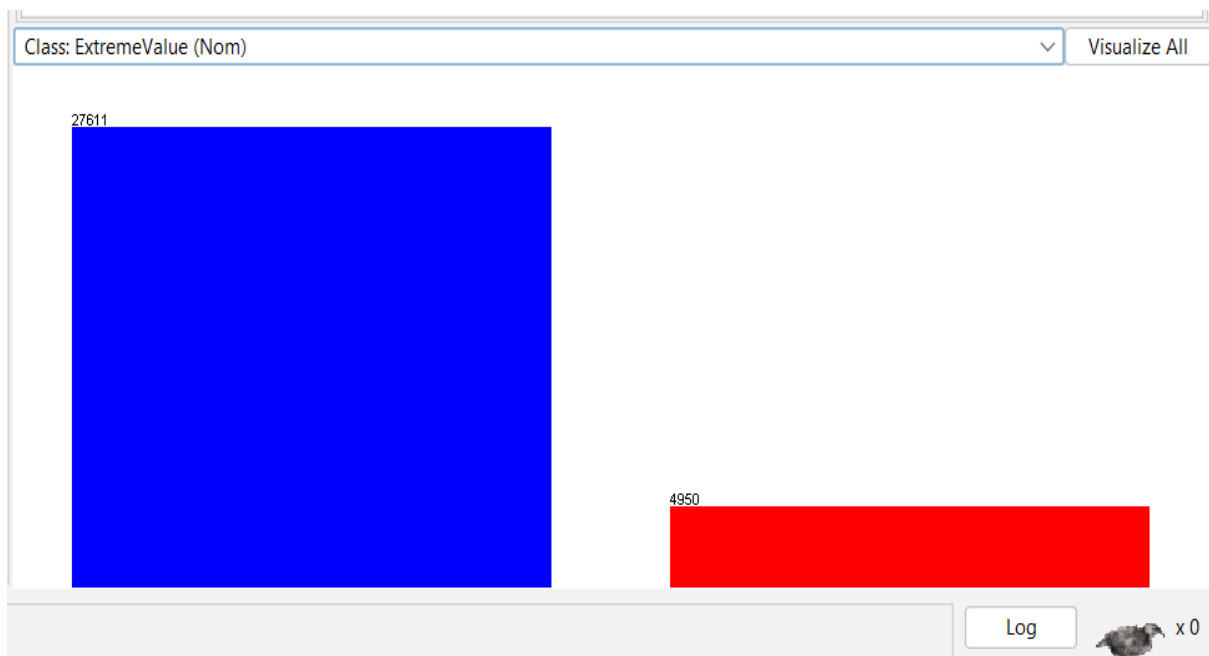


uncorrelated.

**Step 1:(filters→unsupervised→ordinal to numeric→apply)**

**Step2:( filters→unsupervised→intterquartilerang→apply)**

## Outlier



## ExtremeValue

- Cluster(From "Weka" ):We have get cluster from weka by:

(a) From cluster tab .

(b) Choose simple EM (Expectation Maximization ) class.

(c) Finally Hit "Apply".



Clusterer

| Choose | **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 |

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)
20:49:41 - SimpleKMeans

Clusterer output

```
Attribute          Full Data          0              1
                   (32561.0) (17360.0) (15201.0)
=================================================
age                 38.5816    40.6457    36.2244
workclass              2.31     2.5814          2
education            3.4245     3.6289      3.191
education-num      10.0807    10.1817     9.9653
marital-status     1.0838      1.159     0.9979
occupation          4.6664     4.9976     4.2881
relationship       1.5424     1.5636     1.5182
race               0.2217     0.4158          0
gender             0.3308     0.3238     0.3388
capital-gain     1077.6488 2021.2744          0
capital-loss       87.3038     163.75          0
hours-per-week     40.4375    40.7862    40.0392
native-country     1.2903     2.4202          0
fnlwgt               <=50K      <=50K      <=50K
Outlier                 no         no         no
ExtremeValue           yes        yes         no




Time taken to build model (full training data) : 1.25 seconds

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

Clustered Instances

0      17360 ( 53%)
1      15201 ( 47%)
```

Status
OK

Log        🐦 x 0

**QB - B1.**

# • Equi-width binning (3 bins).

(a) From "Pre-process" Tab.

(b) Click Filter–>Unsupervised–>Attribute–>Discretize.

(c) Open Discretize editor.

(d) Change number of bins to (3).

(e) Change attribute indices to (3).

(f) Make sure that useEqualFrequency must be false.

(g) Hit Apply.

# • Equi-depth binning (3 bins).

(a) From "Pre-process" Tab.

(b) Click Filter–>Unsupervised–>Attribute–>Discretize.

(c) Open Discretize editor.

(d) Change number of bins to (3).

(e) Change attribute indices to (3).

(f) Make sure that useEqualFrequency must be true.

(g) Hit Apply.

**B2**

• **min-max normalization to transform the values onto the range [0.0,1.0].(Normalization)**

(a) From "Preprocess" Tab.

(b) Click Filter–>Unsupervised–>Attribute–>Normalize.

(c) Open Normalize editor.

(d) Make sure that translation = 0.0, which is the minimum value in the range.

(e) Make sure that scale = 1.0, which is the maximum value in the range.

(f) Hit Apply.

- **z-score normalization to transform the values. (Standardize)**
  **(unsupervised→attribute→Standardize)**

Selected attribute

| | |
|---|---|
| Name: age | Type: Numeric |
| Missing: 0 (0%) | Distinct: 73 | Unique: 2 (0%) |

| Statistic | Value |
|---|---|
| Minimum | -1.582 |
| Maximum | 3.77 |
| Mean | -0 |
| StdDev | 1 |

B3- Discretize the Age attribute into the following categories:

- Teenager = 1-16 ▯ ifelse(A>=1, ifelse(A<=16,1,2),2)

weka.gui.GenericObjectEditor ✕

weka.filters.unsupervised.attribute.MathExpression

**About**

Modify numeric attributes according to a given mathematical expression.

[More]
[Capabilities]

| | |
|---|---|
| debug | False ▾ |
| doNotCheckCapabilities | False ▾ |
| expression | ifelse(A>=1,ifelse(A<=16,1,2),2) |
| ignoreClass | False ▾ |
| ignoreRange | |
| invertSelection | False ▾ |

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

## After discretizing the values will be like this:

Selected attribute

Name: age — Type: Nominal
Missing: 0 (0%) — Distinct: 1 — Unique: 0 (0%)

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | 'All' | 32561 | 32561 |

## And it means there is no age between→[1,16].

- Young = 17-35 ▯ ifelse(A>=17,ifelse(A<=35,1,2),2)

Selected attribute

Name: age — Type: Nominal
Missing: 0 (0%) — Distinct: 2 — Unique: 0 (0%)

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | '(-inf-1.5]' | 14925 | 14925 |
| 2 | '(1.5-inf)' | 17636 | 17636 |

## 14925 of them is between [17,35]

• Mid_Age = 36-55 🡪 ifelse(A>=36, ifelse(A<=55,1,2),2)

```
Selected attribute
  Name: age                                                Type: Nominal
  Missing: 0 (0%)                  Distinct: 2              Unique: 0 (0%)
```

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | '(-inf-1.5]' | 13547 | 13547 |
| 2 | '(1.5-inf)' | 19014 | 19014 |

13546 of them is between [36,55]

• Mature = 56-70 🡪 ifelse(A>=56, ifelse(A<=70,1,2),2)

```
Selected attribute
  Name: age                                                Type: Nominal
  Missing: 0 (0%)                  Distinct: 2              Unique: 0 (0%)
```

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | '(-inf-1.5]' | 3549 | 3549 |
| 2 | '(1.5-inf)' | 29012 | 29012 |

3549 of them is between [56,70]

• Old = 71+ 🡪 ifelse(A>=71,1,2)

```
Selected attribute
  Name: age                                                Type: Nominal
  Missing: 0 (0%)                  Distinct: 2              Unique: 0 (0%)
```

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | '(-inf-1.5]' | 540 | 540 |
| 2 | '(1.5-inf)' | 32021 | 32021 |

540 of them is between [71, ∞]

## B4 -
## Convert the "Gender" variable into binary variables with values ["0" or "1"]. (NominaltoBinary )
→RenameNominalvalue →selectedAttribute→9→replacments(Male:0,Female:1)



→unsupervised→attribute→NominalToBinary→AttributeIndices→10→Apply.

# 2 Part 2

1. Data Set Information: 10 Attributes as shown in (Figure 1)



**In Pre-processing →**
**We convert numeric data to nominal by discretize attribute with 3-bins**

**After that we rename all nominal value by (unsupervised→attribute→RenameNominalValues).**

Selected attribute

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

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | a1 | 19926 | 19926 |
| 2 | a2 | 11477 | 11477 |
| 3 | a3 | 1158 | 1158 |

**Then, we replace missing value by(unsupervised→attribute→ReplaceMissingValues) and merge missing value with most frequent value**

weka.gui.GenericObjectEditor ✕

weka.filters.unsupervised.attribute.MergeManyValues

About

Merges many values of a nominal attribute into one value.    More    Capabilities

| attributeIndex | 2 |
| debug | False |
| doNotCheckCapabilities | False |
| ignoreClass | False |
| label | private |
| mergeValueRange | 3,6 |

Open...    Save...    OK    Cancel

weka.gui.GenericObjectEditor ✕

weka.filters.unsupervised.attribute.MergeManyValues

About

Merges many values of a nominal attribute into one value.    More    Capabilities

| attributeIndex | 7 |
| debug | False |
| doNotCheckCapabilities | False |
| ignoreClass | False |
| label | Prof-speciaty |
| mergeValueRange | 4,12 |

Open...    Save...    OK    Cancel

weka.gui.GenericObjectEditor ✕

weka.filters.unsupervised.attribute.MergeManyValues

About

Merges many values of a nominal attribute into one value.    More    Capabilities

| attributeIndex | 14 |
| debug | False |
| doNotCheckCapabilities | False |
| ignoreClass | False |
| label | United-States |
| mergeValueRange | 1,5 |

Open...    Save...    OK    Cancel

**By these steps we cope-with missing value (No missing values in our set).**



**Then we sortLabels by(unsupervised→attribute→ sortLabels)**

Nearest Neighbour Learning and Decision Trees

KNN classification algorithm 'IBk'

1. KNN=1

(a) From "classify" Tab.

(b) loaded the dataset and ran the classifier with default options.

(c) Click choose–>lazy–>IBK

(d) keep all options at their default values.

(e) run the classifier to obtain the initial results.

2. KNN=10



3. KNN=20

4.Our summery

| K= | Time taken to build model | Accuracy | Error rate |
|---|---|---|---|
| 1 | 0.01 seconds | 82.3961% | 17.6039% |
| 10 | 0.04 seconds | 83.3727% | 16.6273% |
| 20 | 0.02 seconds | 83.3543% | 16.6457% |

Best K we will use is K=10 because the accuracy =**83.3727%**

4. J48

(a) From "classify" Tab.

(b) loaded the dataset and ran the classifier with default options.

(c) Click choose–>trees–>J48

(d) keep all options at their default values.

(e) hit start



1. **Identify the decision tree model that was obtained in terms of the number of nodes, branches, and levels in the tree:**

| Number of nodes | Number of branches | Number of levels |
|---|---|---|
| 533 | 533-1=532 | 533-452=81 |

2.  **Visualize the decision tree.**

*3. Determine the accuracy of the model in terms of contingency matrix:*

Correctly classified instances 27136

Total number of instances  32561

Accuracy=( classified instances/total  instances)*100

Accuracy=(27136/32561)*100=83.3389

5. As shown in (figure 4) the confusion matrix for Decision tree is:

=== Confusion Matrix ===

   a   b   <-- classified as

22970  1750 |   a = <=50K

 3675  4166 |   b = >50K

For more explanation:

- 22970 instances correctly classified as <=50k
- 1750 instances incorrectly classified as >=50k
- 3675 instances incorrectly classified as <=50k
- 4166 instances correctly classified as >50k
→That's why Error rate =16.661%
→accuracy explained in previous question

**Confusion Matrix:**

| A\P | C | ¬C | |
|-----|-----|-----|-----|
| C | TP | FN | P |
| ¬C | FP | TN | N |
| | P' | N' | All |

→ Final summery :

| chosen technique | Accuracy | Error rate |
|------------------|----------|------------|
| Lazy classifier – IBK with k=1 | 82.3961% | 17.6039% |
| Lazy classifier – IBK with k=10 | 83.3727% | 16.6273% |
| Lazy classifier – IBK with k=20 | 83.3543% | 16.6457% |
| Decision tree classifier | 83.339% | 16.661% |

Best technique is Lazy classifier – IBK with k=10
because the accuracy =83.3727%