National Institute Of Technology – Calicut
# Data Mining
(Data pre-processing assignment)


**Que 02.**
**(AUTO MPG data set)**
**Data  Pre-processing Using Weka**

T.G. Deshan K. Sumanathilaka
B150413CS
Computer Science and Engineering
(B.Tech)

# Selected Data Set .............

# Step 1 :Loading Data Set to  Weka

# Step 2
## a)Cleaning up inconsistent spelling of terms

# b)Converting values that are text descriptions of numeric values to actual numeric values which are usable for analysis.

# c)Extracting and cleaning values for dates

# d)Fill in the missing values by the various options

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

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

Filter

Choose | Rep...

Current relatio...
Relation: aut...
Instances: 40(

**Viewer**
Relation: autoMPG-weka.filters.unsupervised.attribute.ReplaceMissing...

| No. | 1: mpg Numeric | 2: cylinders Numeric | 3: displacement Numeric | 4: horsepower Numeric | 5: Num |
|---|---|---|---|---|---|
| 1 | 18.0 | 8.0 | 307.0 | 130.0 | 350 |
| 2 | 15.0 | 8.0 | 350.0 | 165.0 | 369 |
| 3 | 18.0 | 8.0 | 318.0 | 150.0 | 343 |
| 4 | 16.0 | 8.0 | 304.0 | 150.0 | 343 |
| 5 | 17.0 | 8.0 | 302.0 | 140.0 | 344 |
| 6 | 15.0 | 8.0 | 429.0 | 198.0 | 434 |
| 7 | 14.0 | 8.0 | 454.0 | 220.0 | 435 |
| 8 | 14.0 | 8.0 | 440.0 | 215.0 | 431 |
| 9 | 14.0 | 8.0 | 455.0 | 225.0 | 442 |
| 10 | 15.0 | 8.0 | 390.0 | 190.0 | 385 |
| 11 | 23.51457 28... | 4.0 | 133.0 | 115.0 | 309 |
| 12 | 23.51457 28... | 8.0 | 350.0 | 165.0 | 414 |
| 13 | 23.51457 28... | 8.0 | 351.0 | 153.0 | 403 |
| 14 | 23.51457 28... | 8.0 | 383.0 | 175.0 | 416 |
| 15 | 23.51457 28... | 8.0 | 360.0 | 175.0 | 385 |
| 16 | 15.0 | 8.0 | 383.0 | 170.0 | 356 |
| 17 | 14.0 | 8.0 | 340.0 | 160.0 | 360 |
| 18 | 23.51457 28... | 8.0 | 302.0 | 140.0 | 335 |
| 19 | 15.0 | 8.0 | 400.0 | 150.0 | 376 |
| 20 | 14.0 | 8.0 | 455.0 | 225.0 | 308 |
| 21 | 24.0 | 4.0 | 113.0 | 95.0 | 237 |
| 22 | 22.0 | 6.0 | 198.0 | 95.0 | 283 |
| 23 | 18.0 | 6.0 | 199.0 | 97.0 | 277 |

Add instance | Undo | OK | Cancel

**Selected attribute**
Name: mpg            Type: Numeric
Missing: 0 (0%)    Distinct: 130    Unique: 73 (18%)

| Statistic | Value |
|---|---|
| Minimum | 9 |
| Maximum | 46.6 |
| Mean | 23.515 |
| StdDev | 7.738 |

Class: carname (Str)  ▼  | Visualize All

78  73  69  54  48  38  22  13  5  6
9        27.8        46.6

Remove

Status
OK

Log  × 0

---



**Weka Explorer** — Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter

Choose | **NumericTransform** -R 1,4 -C java.lang.Math -M floor    Apply

Current relation
Relation: autoMPG-we...                    9
Instances: 406                              406

Attributes

All

| No. | Name |
|---|---|
| 1 | mpg |
| 2 | cylinders |
| 3 | displaceme |
| 4 | horsepowe |
| 5 | weight |
| 6 | acceleratio |
| 7 | modelyear |
| 8 | origin |
| 9 | carname |

**Viewer**
Relation: autoMPG-weka.filters.unsupervised.attribute.ReplaceMissing...

| No. | 1: mpg Numeric | 2: cylinders Numeric | 3: displacement Numeric | 4: horsepow Numeric |
|---|---|---|---|---|
| 1 | 18.0 | 8.0 | 307.0 | 130 |
| 2 | 15.0 | 8.0 | 350.0 | 165 |
| 3 | 18.0 | 8.0 | 318.0 | 150 |
| 4 | 16.0 | 8.0 | 304.0 | 150 |
| 5 | 17.0 | 8.0 | 302.0 | 140 |
| 6 | 15.0 | 8.0 | 429.0 | 198 |
| 7 | 14.0 | 8.0 | 454.0 | 220 |
| 8 | 14.0 | 8.0 | 440.0 | 215 |
| 9 | 14.0 | 8.0 | 455.0 | 225 |
| 10 | 15.0 | 8.0 | 390.0 | 190 |
| 11 | 23.514572864321615 | 4.0 | 133.0 | 115 |
| 12 | 23.514572864321615 | 8.0 | 350.0 | 165 |
| 13 | 23.514572864321615 | 8.0 | 351.0 | 153 |
| 14 | 23.514572864321615 | 8.0 | 383.0 | 175 |
| 15 | 23.514572864321615 | 8.0 | 360.0 | 175 |
| 16 | 15.0 | 8.0 | 383.0 | 170 |
| 17 | 14.0 | 8.0 | 340.0 | 160 |
| 18 | 23.514572864321615 | 8.0 | 302.0 | 140 |
| 19 | 15.0 | 8.0 | 400.0 | 150 |
| 20 | 14.0 | 8.0 | 455.0 | 225 |
| 21 | 24.0 | 4.0 | 113.0 | 95 |
| 22 | 22.0 | 6.0 | 198.0 | 95 |
| 23 | 18.0 | 6.0 | 199.0 | 97 |

Add instance | Undo | OK | Cancel

**Selected attribute**
Name: mpg            Type: Numeric
Missing: 0 (0%)    Distinct: 130    Unique: 73 (18%)

| Statistic | Value |
|---|---|
| Minimum | 9 |
| Maximum | 46.6 |
| Mean | 23.515 |
| StdDev | 7.738 |

Class: carname (Str)  ▼  | Visualize All

78  73  69  54  48  38  22  13  5  6
9        27.8        46.6

Status
OK

Log  × 0

Viewer

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

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

**Filter**

Choose | **NumericTransform** -R 1,4 -C java.lang.Math -M floor | Apply

**Current relation**

Relation: autoMF...
Instances: 406

**Attributes**

All

| No. | Name |
|-----|------|
| 1 | mpg |
| 2 | cylinde |
| 3 | displac |
| 4 | horsep |
| 5 | weight |
| 6 | accele |
| 7 | model |
| 8 | origin |
| 9 | carnan |

Remove

**Selected attribute**

Name: mpg — Type: Numeric
Missing: 0 (0%) | Distinct: 36 | Unique: 3 (1%)

| Statistic | Value |
|-----------|-------|
| Minimum | 9 |
| Maximum | 46 |
| Mean | 23.342 |
| StdDev | 7.676 |

Class: carname (Str) ▾ | Visualize All

**Viewer**

Relation: autoMPG-weka.filters.unsupervised.attribute.ReplaceMissing...

No. 1: mpg 2: cylinders 3: displacement 4: horsepower 5: weight 6: a

| No. | Numeric | Numeric | Numeric | Numeric | Numeric |
|-----|---------|---------|---------|---------|---------|
| 1 | 18.0 | 8.0 | 307.0 | 130.0 | 3504.0 |
| 2 | 15.0 | 8.0 | 350.0 | 165.0 | 3693.0 |
| 3 | 18.0 | 8.0 | 318.0 | 150.0 | 3436.0 |
| 4 | 16.0 | 8.0 | 304.0 | 150.0 | 3433.0 |
| 5 | 17.0 | 8.0 | 302.0 | 140.0 | 3449.0 |
| 6 | 15.0 | 8.0 | 429.0 | 198.0 | 4341.0 |
| 7 | 14.0 | 8.0 | 454.0 | 220.0 | 4354.0 |
| 8 | 14.0 | 8.0 | 440.0 | 215.0 | 4312.0 |
| 9 | 14.0 | 8.0 | 455.0 | 225.0 | 4425.0 |
| 10 | 15.0 | 8.0 | 390.0 | 190.0 | 3850.0 |
| 11 | 23.0 | 4.0 | 133.0 | 115.0 | 3090.0 |
| 12 | 23.0 | 8.0 | 350.0 | 165.0 | 4142.0 |
| 13 | 23.0 | 8.0 | 351.0 | 153.0 | 4034.0 |
| 14 | 23.0 | 8.0 | 383.0 | 175.0 | 4166.0 |
| 15 | 23.0 | 8.0 | 360.0 | 175.0 | 3850.0 |
| 16 | 15.0 | 8.0 | 383.0 | 170.0 | 3563.0 |
| 17 | 14.0 | 8.0 | 340.0 | 160.0 | 3609.0 |
| 18 | 23.0 | 8.0 | 302.0 | 140.0 | 3353.0 |
| 19 | 15.0 | 8.0 | 400.0 | 150.0 | 3761.0 |
| 20 | 14.0 | 8.0 | 455.0 | 225.0 | 3086.0 |
| 21 | 24.0 | 4.0 | 113.0 | 95.0 | 2372.0 |
| 22 | 22.0 | 6.0 | 198.0 | 95.0 | 2833.0 |
| 23 | 18.0 | 6.0 | 199.0 | 97.0 | 2774.0 |

Add instance | Undo | OK | Cancel

Histogram values: 13, 79, 79, 51, 66, 50, 32, 24, 6, 6
(9 — 27.5 — 46)

**Status**

OK

Log | × 0

# e) Removing duplicate rows

# f)Using a scatter plot to visualize relationships between values in different columns

(using same attribute as x-axis and y-axis)

# g)Exporting cleaned data to Excel

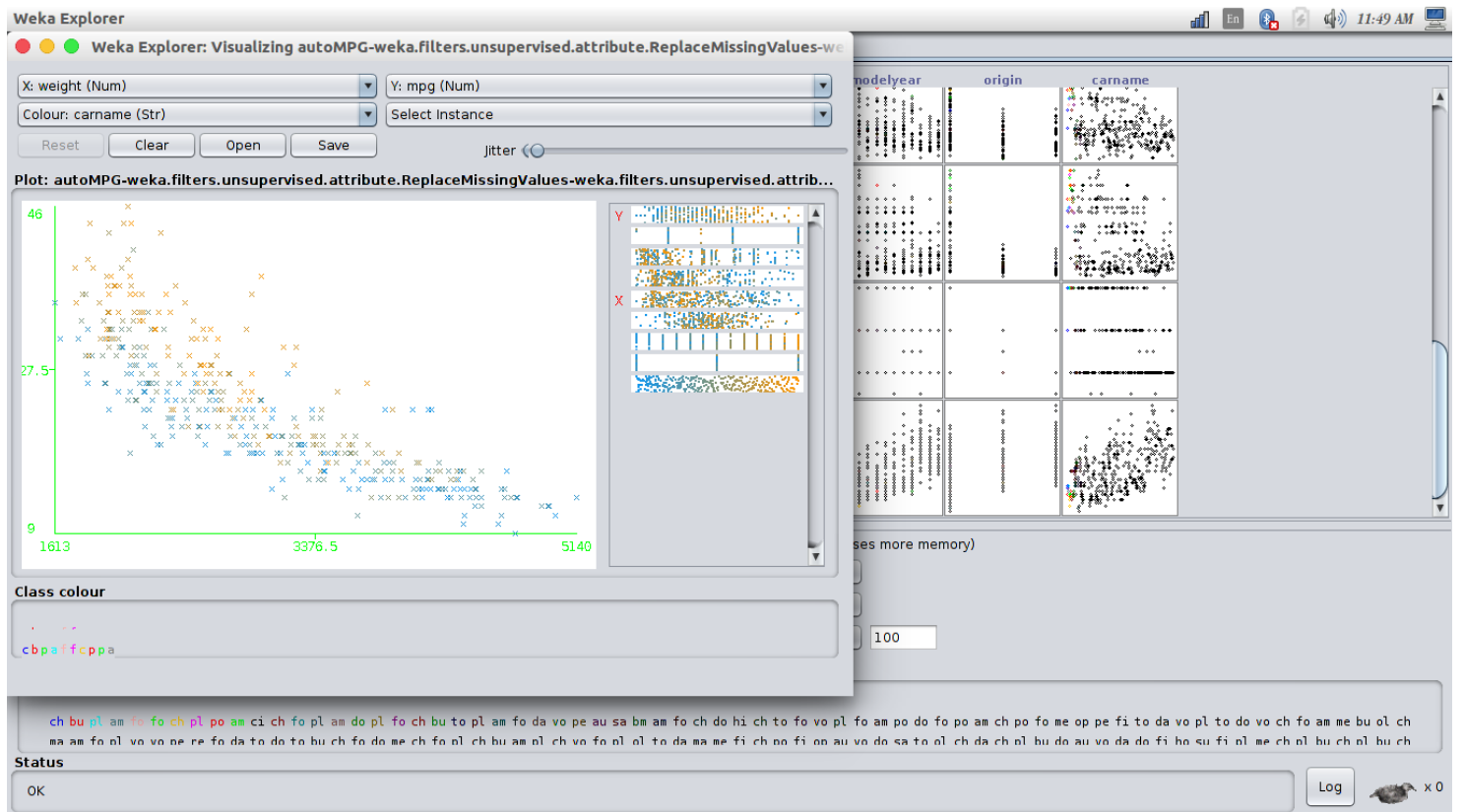| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 13 | | | | | | | | | |
| 14 | 15 | 8 | 307 | 130 | 3504 | 12 | 70 | 1 | 'chevrolet chevelle malibu' |
| 15 | 15 | 8 | 350 | 165 | 3693 | 11.5 | 70 | 1 | 'buick skylark 320' |
| 16 | 18 | 8 | 318 | 150 | 3436 | 11 | 70 | 1 | 'plymouth satellite' |
| 17 | 16 | 8 | 304 | 150 | 3433 | 12 | 70 | 1 | 'amc rebel sst' |
| 18 | 17 | 8 | 302 | 140 | 3449 | 10.5 | 70 | 1 | 'ford torino' |
| 19 | 15 | 8 | 429 | 198 | 4341 | 10 | 70 | 1 | 'ford galaxie 500' |
| 20 | 14 | 8 | 454 | 220 | 4354 | 9 | 70 | 1 | 'chevrolet impala' |
| 21 | 14 | 8 | 440 | 215 | 4312 | 8.5 | 70 | 1 | 'plymouth fury iii' |
| 22 | 14 | 8 | 455 | 225 | 4425 | 10 | 70 | 1 | 'pontiac catali ' |
| 23 | 15 | 8 | 390 | 190 | 3850 | 8.5 | 70 | 1 | 'amc ambassador dpl' |
| 24 | 23 | 4 | 133 | 115 | 3090 | 17.5 | 70 | 2 | 'citroen ds-21 pallas' |
| 25 | 23 | 8 | 350 | 165 | 4142 | 11.5 | 70 | 1 | 'chevrolet chevelle concours (sw)' |
| 26 | 23 | 8 | 351 | 153 | 4034 | 11 | 70 | 1 | 'ford torino (sw)' |
| 27 | 23 | 8 | 383 | 175 | 4166 | 10.5 | 70 | 1 | 'plymouth satellite (sw)' |
| 28 | 23 | 8 | 360 | 175 | 3850 | 11 | 70 | 1 | 'amc rebel sst (sw)' |
| 29 | 15 | 8 | 383 | 170 | 3563 | 10 | 70 | 1 | 'dodge challenger se' |
| 30 | 14 | 8 | 340 | 160 | 3609 | 8 | 70 | 1 | 'plymouth \'cuda 340' |
| 31 | 23 | 8 | 302 | 140 | 3353 | 8 | 70 | 1 | 'ford mustang boss 302' |
| 32 | 15 | 8 | 400 | 150 | 3761 | 9.5 | 70 | 1 | 'chevrolet monte carlo' |
| 33 | 14 | 8 | 455 | 225 | 3086 | 10 | 70 | 1 | 'buick estate wagon (sw)' |
| 34 | 24 | 4 | 113 | 95 | 2372 | 15 | 70 | 3 | 'toyota coro  mark ii' |
| 35 | 22 | 6 | 198 | 95 | 2833 | 15.5 | 70 | 1 | 'plymouth duster' |
| 36 | 18 | 6 | 199 | 97 | 2774 | 15.5 | 70 | 1 | 'amc hornet' |
| 37 | 21 | 6 | 200 | 85 | 2587 | 16 | 70 | 1 | 'ford maverick' |
| 38 | 27 | 4 | 97 | 88 | 2130 | 14.5 | 70 | 3 | 'datsun pl510' |
| 39 | 26 | 4 | 97 | 46 | 1835 | 20.5 | 70 | 2 | 'volkswagen 1131 deluxe sedan' |
| 40 | 25 | 4 | 110 | 87 | 2672 | 17.5 | 70 | 2 | 'peugeot 504' |
| 41 | 24 | 4 | 107 | 90 | 2430 | 14.5 | 70 | 2 | 'audi 100 ls' |
| 42 | 25 | 4 | 104 | 95 | 2375 | 17.5 | 70 | 2 | 'saab 99e' |
| 43 | 26 | 4 | 121 | 113 | 2234 | 12.5 | 70 | 2 | 'bmw 2002' |
| 44 | 21 | 6 | 199 | 90 | 2648 | 15 | 70 | 1 | 'amc gremlin' |

# Differences of Weka And Open Refine

1) visualization of data in weka is much more clear compare to refine and easy. The plot matrix of weka will give a good visualization and easy to understand the properties of data set.

2)weka supports histogram visualized method.but not in refine.

3)Operations like Removing Duplicates,  removing inconsistent spelling of terms can be easily perform in weka.

4)Tuple format data consideration can perform in weka while column format data classification can perform in open refine.

5)Clustering and merging can perform easily in open refine compared to Weka.

6) Replacing  missing Values can be perform only in weka while removing data tuples with   missing value can be perform in both open refine and weka.

7) No inbuild, direct function to remove duplicate rows in open refine.but in weka .