**Aim:** To perform data preprocessing using WEKA

**Theory:**
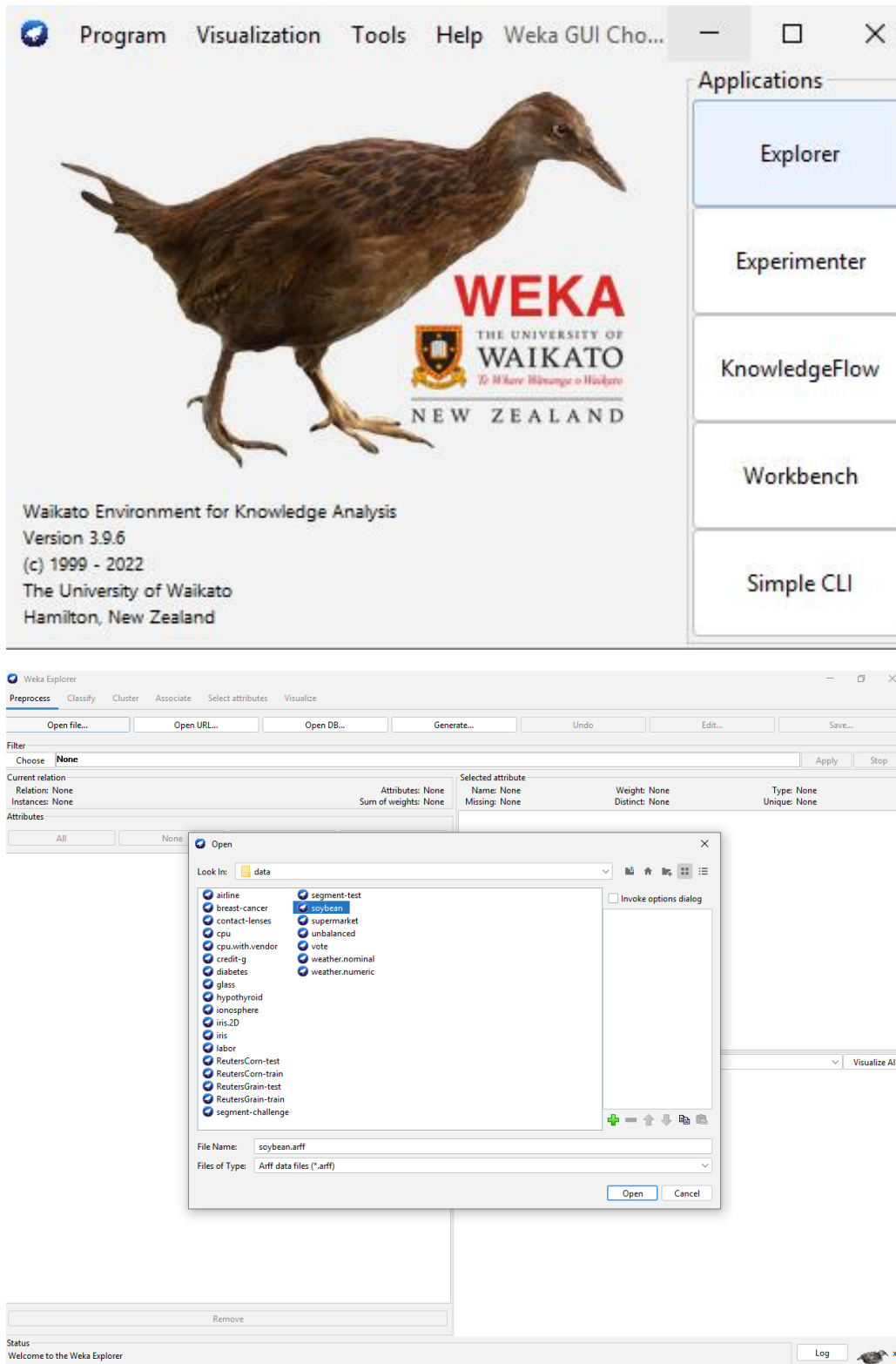
What is data preprocessing?

- Data preprocessing refers to the cleaning, transforming, and integrating of data in order to make it ready for analysis.
- The goal of data preprocessing is to improve the quality of the data and to make it more suitable for the specific data mining task.
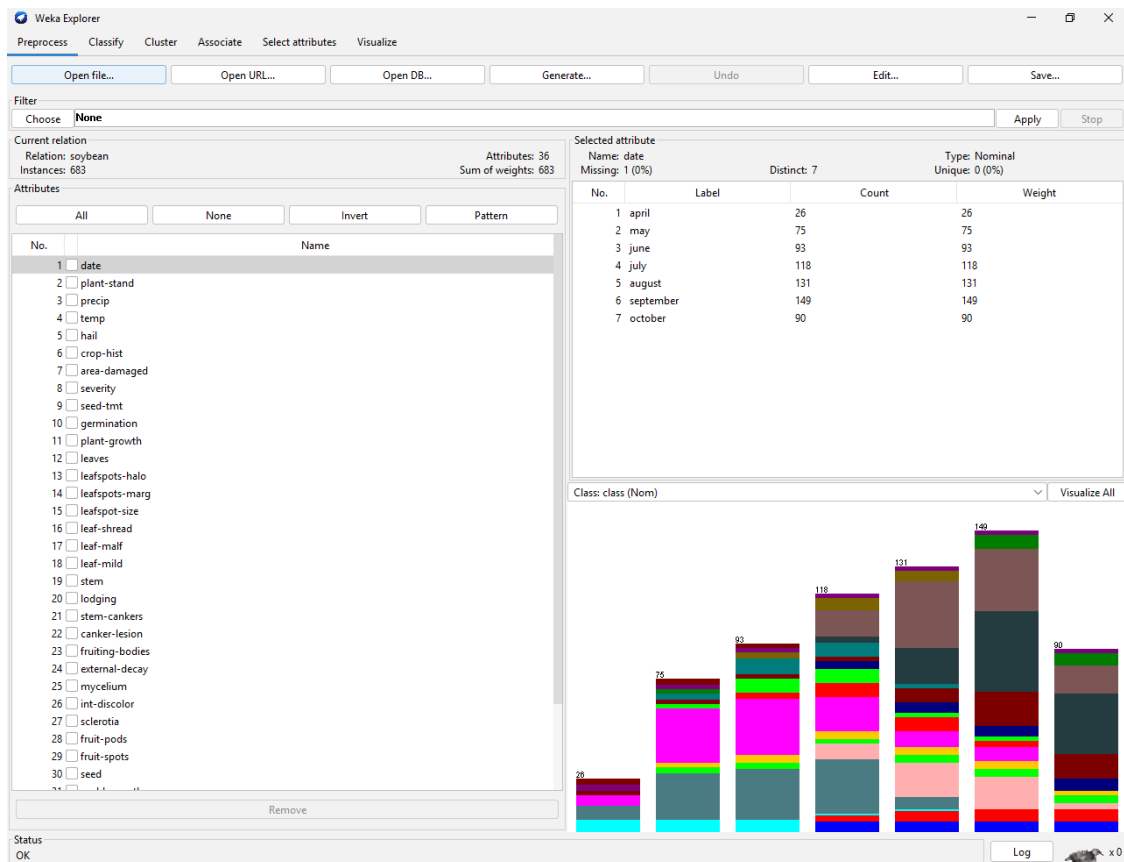
Common steps in data preprocessing include:

1. **Data cleaning:** this step involves identifying and removing missing, inconsistent, or irrelevant data. This can include removing duplicate records, filling in missing values, and handling outliers.
2. **Data integration:** this step involves combining data from multiple sources, such as databases, spreadsheets, and text files. The goal of integration is to create a single, consistent view of the data.
3. **Data transformation:** this step involves converting the data into a format that is more suitable for the data mining task. This can include normalizing numerical data, creating dummy variables, and encoding categorical data.
4. **Data reduction:** this step is used to select a subset of the data that is relevant to the data mining task. This can include feature selection (selecting a subset of the variables) or feature extraction (extracting new variables from the data).
5. **Data discretization:** this step is used to convert continuous numerical data into categorical data, which can be used for decision tree and other categorical data mining techniques.

## Steps to perform preprocessing in WEKA:

Open Weka tool

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

Filter
Choose  None                                                                                    Apply   Stop

Current relation
Relation: soybean                                  Attributes: 36
Instances: 683                                     Sum of weights: 683

Selected attribute
Name: date                                         Type: Nominal
Missing: 1 (0%)        Distinct: 7                  Unique: 0 (0%)

Attributes

All   None   Invert   Pattern

| No. | Name |
|---|---|
| 1 | date |
| 2 | plant-stand |
| 3 | precip |
| 4 | temp |
| 5 | hail |
| 6 | crop-hist |
| 7 | area-damaged |
| 8 | severity |
| 9 | seed-tmt |
| 10 | germination |
| 11 | plant-growth |
| 12 | leaves |
| 13 | leafspots-halo |
| 14 | leafspots-marg |
| 15 | leafspot-size |
| 16 | leaf-shread |
| 17 | leaf-malf |
| 18 | leaf-mild |
| 19 | stem |
| 20 | lodging |
| 21 | stem-cankers |
| 22 | canker-lesion |
| 23 | fruiting-bodies |
| 24 | external-decay |
| 25 | mycelium |
| 26 | int-discolor |
| 27 | sclerotia |
| 28 | fruit-pods |
| 29 | fruit-spots |
| 30 | seed |

Remove

| No. | Label | Count | Weight |
|---|---|---|---|
| 1 | april | 26 | 26 |
| 2 | may | 75 | 75 |
| 3 | june | 93 | 93 |
| 4 | july | 118 | 118 |
| 5 | august | 131 | 131 |
| 6 | september | 149 | 149 |
| 7 | october | 90 | 90 |

Class: class (Nom)                                 Visualize All

Status
OK

Log   x 0

# Check for missing values

Viewer

Relation: soybean

| No. | 1: date Nominal | 2: plant-stand Nominal | 3: precip Nominal | 4: temp Nominal | 5: hail Nominal | 6: crop-hist Nominal | 7: area-damaged Nominal | 8: severity Nominal | 9: seed-tmt Nominal | 10: germination Nominal | 11: plant-growth Nominal | 12: leaves Nominal | 13: leafspots-halo Nominal |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 40 | june | lt-normal | gt-norm | lt-norm | yes | same-lst-yr | low-areas | severe | none | 80-89 | abnorm | abnorm | absent |
| 41 | june | lt-normal | gt-norm | norm | | same-lst-yr | low-areas | | | | abnorm | abnorm | |
| 42 | may | lt-normal | gt-norm | norm | | same-lst-yr | low-areas | | | | abnorm | abnorm | |
| 43 | april | lt-normal | gt-norm | norm | yes | same-lst-s... | low-areas | pot-severe | none | 90-100 | abnorm | abnorm | absent |
| 44 | april | lt-normal | norm | norm | no | same-lst-t... | low-areas | severe | fungicide | 90-100 | abnorm | abnorm | absent |
| 45 | july | lt-normal | gt-norm | lt-norm | yes | same-lst-yr | low-areas | severe | fungicide | 90-100 | abnorm | abnorm | absent |
| 46 | june | lt-normal | gt-norm | gt-norm | | same-lst-s... | low-areas | | | | abnorm | abnorm | |
| 47 | april | lt-normal | gt-norm | norm | yes | same-lst-t... | low-areas | pot-severe | none | 80-89 | abnorm | abnorm | absent |
| 48 | june | lt-normal | norm | gt-norm | | same-lst-t... | low-areas | | | | abnorm | abnorm | absent |
| 49 | june | lt-normal | gt-norm | norm | no | same-lst-yr | low-areas | severe | none | lt-80 | abnorm | abnorm | absent |
| 50 | april | lt-normal | gt-norm | norm | yes | same-lst-s... | low-areas | pot-severe | none | lt-80 | abnorm | abnorm | absent |
| 51 | may | lt-normal | gt-norm | norm | yes | diff-lst-year | low-areas | severe | fungicide | 80-89 | abnorm | abnorm | absent |
| 52 | may | lt-normal | gt-norm | norm | | diff-lst-year | low-areas | | | | abnorm | abnorm | |
| 53 | july | lt-normal | gt-norm | norm | | same-lst-yr | low-areas | | | | abnorm | abnorm | |
| 54 | june | lt-normal | gt-norm | norm | | same-lst-yr | low-areas | | | | abnorm | abnorm | absent |
| 55 | july | lt-normal | gt-norm | gt-norm | | same-lst-t... | low-areas | | | | abnorm | abnorm | |
| 56 | may | lt-normal | gt-norm | norm | no | same-lst-s... | low-areas | sev... | none | 80-89 | abnorm | abnorm | absent |
| 57 | july | lt-normal | norm | norm | | same-lst-s... | low-areas | | | | abnorm | abnorm | absent |
| 58 | june | lt-normal | gt-norm | gt-norm | | same-lst-yr | low-areas | | | | abnorm | abnorm | |
| 59 | july | lt-normal | norm | gt-norm | | same-lst-t... | low-areas | | | | abnorm | abnorm | absent |
| 60 | may | lt-normal | gt-norm | gt-norm | | same-lst-yr | low-areas | | | | abnorm | abnorm | |
| 61 | june | lt-normal | gt-norm | gt-norm | | same-lst-s... | low-areas | | | | abnorm | abnorm | absent |
| 62 | july | lt-normal | norm | norm | | diff-lst-year | low-areas | | | | abnorm | abnorm | absent |

Add instance   Undo   OK   Cancel

## Create missing values by selecting any one and deleting data



## Choose a filter to apply:

### 1) Using replace missing values

Click Apply to apply the filter to the data



All missing values are now replaced

## 2) Using Add Cluster filter



## Click Apply to apply the filter to the data

Click edit to view the last column created when clusters were created



**Conclusion:** Successfully performed data preprocessing using WEKA.