**Lab Number: 1**

**Title**
Preprocessing of Primary and Secondary Datasets Containing Dirty Data

**Objective**
To understand and practically perform data preprocessing on a dirty dataset using Weka Explorer.

**IDE/Tools Used**
Weka 3.8.6

**Theory**
**Dirty Data:** A dirty dataset refers to a collection of data that contains inaccuracies, inconsistencies, and errors, which can compromise its usefulness and reliability for analysis, reporting, or decision-making

**Types of Dirty Data**
- Missing Data
- Duplicate Records
- Inconsistent Values
- Noise
- Outliers

**Data Preprocessing:** Data preprocessing is the process of cleaning, transforming, and organizing raw data into a structured format that is ready for analysis or use in machine learning models.
- **Cleaning:** Involves handling missing values, removing duplicates, and correcting errors to make the data accurate and consistent.
- **Transformation:** Involves converting data into a suitable format. Examples include standardizing numerical features, normalizing data, or encoding categorical variables.
- **Integration:** Combines data from multiple sources into a single, unified dataset for analysis.
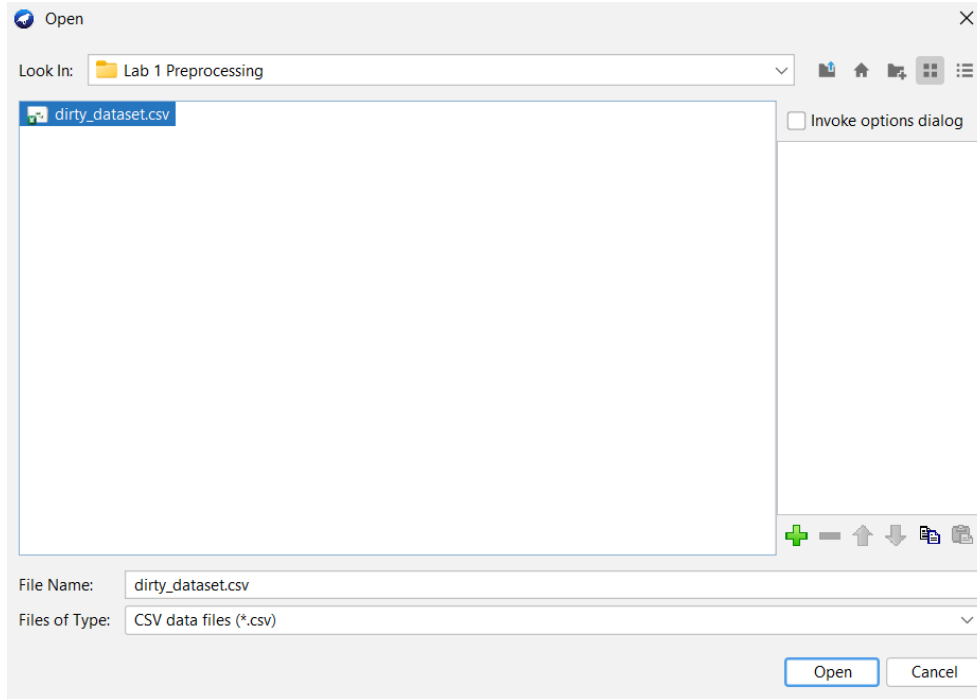
**Implementation**

**A. For Primary Dataset**

For primary dataset, a customer churn data was generated.

**Steps used to clean the data:**

**1. Open the dataset in the pre-processor of WEKA**



**2. Visualize the data**

| No. | 1: CustomerID String | 2: Age Numeric | 3: Gender Nominal | 4: Income Numeric | 5: Region Nominal | 6: Spend Numeric | 7: SignupDate Nominal | 8: LastPurchase Nominal | 9: **Churn** Nominal |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 25.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 2 | 2 | | Female | | Europe | 850.0 | 2023-02-30 | | No |
| 3 | 3 | 45.0 | Male | 120000.0 | Asia | 5000.0 | 2023-03-10 | 2025-01-15 | Yes |
| 4 | 4 | 32.0 | F | 75000.0 | South A... | 3200.0 | 2023-04-05 | 2024-11-20 | No |
| 5 | 5 | 28.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 6 | 6 | 35.0 | | 62000.0 | Europe | | 2023-06-01 | 2024-10-10 | No |
| 7 | 7 | 999.0 | Male | 55000.0 | Africa | 300.0 | 2023-07-12 | 2023-07-12 | Yes |
| 8 | 8 | 41.0 | Male | 58000.0 | North A... | 1800.0 | 2023-08-20 | | No |
| 9 | 9 | 29.0 | Female | 48000.0 | Asia | 1100.0 | 2023-09-05 | 2024-09-05 | Yes |
| 10 | 10 | 33.0 | Male | | Oceania | 2200.0 | 2023-10-01 | 2024-12-10 | No |
| 11 | CUST11 | 31.0 | Female | 70000.0 | Europe | 1500.0 | 2023-11-11 | 2024-11-11 | Yes |
| 12 | 12 | 27.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 13 | 13 | 62.0 | Male | 85000.0 | Moon | 999999.0 | 2025-12-01 | 2025-12-01 | No |

## 3. Remove unwanted columns

In this case CustomerID was removed using the unsupervised.attribute.Remove filter



Weka Explorer

Preprocess  Classify  Cluster  Associate  Select attributes  Visualize

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

**Filter**
Choose | Remove -R 1 | Apply | Stop

**Current relation**
Relation: dirty_dataset    Attributes: 9
Instances: 13    Sum of weights: 13

**Selected attribute**
Name: CustomerID    Type: String
Missing: 0 (0%)    Distinct: 13    Unique: 13 (100%)

**Attributes**
All | None | Invert | Pattern

| No. | Name |
|---|---|
| 1 | CustomerID |
| 2 | Age |
| 3 | Gender |
| 4 | Income |
| 5 | Region |
| 6 | Spend |
| 7 | SignupDate |
| 8 | LastPurchase |
| 9 | Churn |

Class: Churn (Nom)    Visualize All

**Viewer**

Relation: dirty_dataset-weka.filters.unsupervised.attribute.Remove-R1

| No. | 1: Age Numeric | 2: Gender Nominal | 3: Income Numeric | 4: Region Nominal | 5: Spend Numeric | 6: SignupDate Nominal | 7: LastPurchase Nominal | 8: **Churn** Nominal |
|---|---|---|---|---|---|---|---|---|
| 1 | 25.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 2 | | Female | | Europe | 850.0 | 2023-02-30 | | No |
| 3 | 45.0 | Male | 120000.0 | Asia | 5000.0 | 2023-03-10 | 2025-01-15 | Yes |
| 4 | 32.0 | F | 75000.0 | South A... | 3200.0 | 2023-04-05 | 2024-11-20 | No |
| 5 | 28.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 6 | 35.0 | | 62000.0 | Europe | | 2023-06-01 | 2024-10-10 | No |
| 7 | 999.0 | Male | 55000.0 | Africa | 300.0 | 2023-07-12 | 2023-07-12 | Yes |
| 8 | 41.0 | Male | 58000.0 | North A... | 1800.0 | 2023-08-20 | | No |
| 9 | 29.0 | Female | 48000.0 | Asia | 1100.0 | 2023-09-05 | 2024-09-05 | Yes |
| 10 | 33.0 | Male | | Oceania | 2200.0 | 2023-10-01 | 2024-12-10 | No |
| 11 | 31.0 | Female | 70000.0 | Europe | 1500.0 | 2023-11-11 | 2024-11-11 | Yes |
| 12 | 27.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 13 | 62.0 | Male | 85000.0 | Moon | 999999.0 | 2025-12-01 | 2025-12-01 | No |

Right c

## 4. Remove any duplicate values

In this case this was done using the unsupervised.instance.RemoveDuplicates filter.



## 5. Replace any missing values

In this case it was done using the unsupervised.attribute.ReplaceMissingValues filter

## 6. Convert string into nominal values

This is done using the unsupervised.attribute.StringToNominal filter.

## 7. Removing Outliers

To remove Outliers, we perform the following steps:

### 7.1. Interquartile Range

Choose Interquartile Range filter from unsupervised.attribute.InterquartileRange and select the following settings. This will give outlier and extreme values in the dataset.

## 7.2. Remove rows with outliers

To remove the outlier, we use unsupervised.instance.RemoveWithValues filter and apply the following preferences and repeat for attribute indices Age_Outlier, Income_Outlier, and Spend_Outlier (i.e. 9, 11, 13). Here splitPoint is 0.5 because, "No" = 0 and "Yes" = 1, so anything beside "No" will be deleted.

## 7.3. Remove the columns created

Finally remove the columns from 9 to 14 by unsupervised.attribute.Remove filter.

## 8. Finalize

Data Cleaning Process is done. Visualize and save the clean data.

**Viewer**

Relation: dirty_dataset-weka.filters.unsupervised.attribute.Remove-R1-weka.filters.unsupervised.instance.Re

| No. | 1: Age<br>Numeric | 2: Gender<br>Nominal | 3: Income<br>Numeric | 4: Region<br>Nominal | 5: Spend<br>Numeric | 6: SignupDate<br>Nominal | 7: LastPurchase<br>Nominal | 8: **Churn**<br>Nominal |
|-----|------|--------|-----------|----------|--------|------------|------------|-------|
| 1 | 25.0 | Male | 45000.0 | North A... | 1200.0 | 2023-01-15 | 2024-12-01 | Yes |
| 2 | 133.2 | Female | 68666.666... | Europe | 850.0 | 2023-02-30 | 2024-12-01 | No |
| 3 | 45.0 | Male | 120000.0 | Asia | 5000.0 | 2023-03-10 | 2025-01-15 | Yes |
| 4 | 32.0 | F | 75000.0 | South A... | 3200.0 | 2023-04-05 | 2024-11-20 | No |
| 5 | 41.0 | Male | 58000.0 | North A... | 1800.0 | 2023-08-20 | 2024-12-01 | No |
| 6 | 29.0 | Female | 48000.0 | Asia | 1100.0 | 2023-09-05 | 2024-09-05 | Yes |
| 7 | 33.0 | Male | 68666.666... | Oceania | 2200.0 | 2023-10-01 | 2024-12-10 | No |
| 8 | 31.0 | Female | 70000.0 | Europe | 1500.0 | 2023-11-11 | 2024-11-11 | Yes |

**Save**

| Save In: | 📁 Lab 1 Preprocessing |

☐ Invoke options dialog

| File Name: | clean_dataset |
| Files of Type: | Arff data files (*.arff) |

Save   Cancel

## B. For Secondary Dataset

For secondary dataset, default data provided by the Weka, labor.arff was selected.

**Steps used to clean the data:**

1. **Open the dataset in the pre-processor of WEKA**



2. **Visualize the data**

## 3. Remove unwanted columns

This step was not necessary as all columns were needed.

## 4. Remove any duplicate values

In this case this was done using the unsupervised.instance.RemoveDuplicates filter.



## 5. Replace any missing values

In this case it was done using the unsupervised.attribute.ReplaceMissingValues filter

## 6. Convert string into nominal values

This is done using the unsupervised.attribute.StringToNominal filter.

## 7. Removing Outliers

To remove Outliers, we perform the following steps:

### 7.1. Interquartile Range

Choose Interquartile Range filter from unsupervised.attribute.InterquartileRange and select the following settings. This will give outlier and extreme values in the dataset.

## 7.2. Remove rows with outliers

To remove the rows with outliers, we use unsupervised.instance.RemoveWithValues filter and apply the following preferences and repeat for attribute indices 18 to 32. Here, splitPoint is 0.5 because, "No" = 0 and "Yes" = 1, so anything besides "No" will be deleted.

## 7.3. Remove the columns created

Finally remove the columns from 9 to 14 by unsupervised.attribute.Remove filter.

## 8. Finalize

Data Cleaning Process is done. Visualize and save the clean data.