

# Documentation

---

## Data Preprocessing Steps

### Handling Delimiters

- **Step 1:** Converting Commas to Semicolons in the Dataset  
**Objective:** Ensure the dataset is properly formatted by replacing all semicolons (👉) with commas (,).
- **Step 2:** Checking for Outliers, Missing Values, and Special Characters  
**Objective:** Ensure the data quality by identifying and handling the following issues:
  - **Outliers:** Perform statistical analysis (such as interquartile range or z-score) to detect extreme data points that may affect model performance.
  - **Missing Values:** Identify missing entries and handle them through appropriate methods, such as imputation (mean, median, or mode) or row removal when necessary.
  - **Special Characters:** Detect and clean special characters such as the apostrophe in entries like "Cote d'Ivoire" to avoid parsing errors. These can be replaced by suitable alternatives (e.g., removing the apostrophe or using escape sequences).

**Steps Taken:**

- Applied functions to identify and visualize outliers.
- Used conditional checks to replace or remove missing values.
- Implemented a find-and-replace mechanism for problematic characters.

## Algorithm Implementation

### Apriori Algorithm (Association Rule Mining)

- **Objective:** Discover frequent itemsets and generate association rules.
- **Selected Columns:** Item Type, Sales Channel.
- **Parameters:**
  - Support Threshold: A reasonable value based on dataset characteristics.
  - Confidence Threshold: Set a meaningful value to filter rules.
  - Lift: Evaluated to assess rule significance.
- **Steps:**
  - Preprocessed data by removing duplicates and inconsistencies.
  - Implemented Apriori using association rule mining tools in Weka.

```
=== Run information ===
Scheme:      weka.associations.Apriori -N 10 -T 0 -C 0.5 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
Relations:   10000 Sales Records-weka.filters.unsupervised.attribute.Remove-R1,6-14
Instances:   10000
Attributes:  4
             Country
             Item Type
             Sales Channel
             Order Priority
=== Associator model (full training set) ===

Apriori
=====
Minimum support: .\ (1000 instances)
Minimum metric <confidence>: .\
Number of cycles performed: 10

Generated sets of large itemsets:
Size of set of large itemsets L(1): 6
Size of set of large itemsets L(2): 8

Best rules found:
1. Order Priority=C 2555 ==> Sales Channel=Online 1326 <conf:(. . .)> lift:(. . .) lev:(.) [32] conv:(. . .)
1. Order Priority=C 2494 ==> Sales Channel=Online 1270 <conf:(. . .)> lift:(. . .) lev:(.) [7] conv:(. . .)
```

```
*. Order Priority=M 2503 ==> Sales Channel=Offline 1266 <conf:(-,-,-) lift:(-,-,-) lev:(-) [29] conv:(-,-,-)
i. Order Priority=M 2448 ==> Sales Channel=Online 1228 <conf:(-,-,-) lift:(-,-,-) lev:(-) [-10] conv:(-,-,-)
```

## Naïve Bayes (Classification)

- **Objective:** Build a probabilistic model to classify data into predefined classes.
- **Selected Columns:** Region, Country, Item Type, Sales Channel, Order Priority.
- **Steps:**
  - Split dataset into training (70%) and testing (30%) sets.
  - Assumed feature independence.
  - Evaluated model using accuracy, precision, recall, and F1-score.

Test options		Classifier output	
<input type="radio"/> Use training set		Run information ==	
<input type="radio"/> Supplied test set	Set	Scheme: weka.classifiers.bayes.NaiveBayes	
<input type="radio"/> Cross-validation	Folds: 10	Relation: 1000 Sales Records	
<input type="radio"/> Percentage split	% 70	Instances: 1000	
More options...		Attributes:	
		Region	
		Country	
		Item Type	
		Sales Channel	
		Order Priority	
		Order Date	
		Order ID	
		Ship Date	
		Units Sold	
		Unit Price	
		Unit Cost	
		Total Revenue	
		Total Cost	
		Total Profit	
		Test mode: split 70/30 train, remainder test	
		Classifier model (full training set) ==	
		Naive Bayes Classifier	
		Attribute	
		Class	
		H (-,-,-)	
		L (-,-,-)	
		C (-,-,-)	
		M (-,-,-)	
		Region	
		Middle East and North Africa	
		Europe	
		Sub-Saharan Africa	
		Central America and the Caribbean	
		Australia and Oceania	
		Asia	
		North America	
		[Total]	
		Country	
		Afghanistan	
		Albania	
		Algeria	
		Andorra	
		Angola	
		Antigua and Barbuda	
		Armenia	
		Australia	
		Austria	
		Azerbaijan	
		Bahrain	
		Bangladesh	
		Barbados	
		Belarus	
		Belgium	

Naive Bayes Classifier		Class			
Attribute		H (-,-,-)	L (-,-,-)	C (-,-,-)	M (-,-,-)
Region					
Middle East and North Africa		321.0	315.0	323.0	309.0
Europe		666.0	614.0	705.0	652.0
Sub-Saharan Africa		652.0	667.0	640.0	648.0
Central America and the Caribbean		257.0	271.0	248.0	247.0
Australia and Oceania		185.0	200.0	207.0	209.0
Asia		369.0	378.0	382.0	344.0
North America		60.0	56.0	57.0	46.0
[Total]		2510.0	2501.0	2562.0	2455.0
Country					
Afghanistan		15.0	13.0	21.0	14.0
Albania		17.0	15.0	14.0	13.0
Algeria		14.0	15.0	13.0	10.0
Andorra		15.0	16.0	10.0	17.0
Angola		15.0	16.0	12.0	15.0
Antigua and Barbuda		13.0	15.0	16.0	9.0
Armenia		14.0	9.0	12.0	15.0
Australia		15.0	12.0	10.0	18.0
Austria		16.0	13.0	16.0	17.0
Azerbaijan		16.0	11.0	16.0	17.0
Bahrain		18.0	17.0	17.0	16.0
Bangladesh		14.0	16.0	22.0	15.0
Barbados		15.0	11.0	17.0	17.0
Belarus		12.0	15.0	17.0	12.0
Belgium		9.0	9.0	19.0	10.0

Classifier output		Class			
Units Sold		H (-,-,-)	L (-,-,-)	C (-,-,-)	M (-,-,-)
mean		1487.7111	8111.1111	8111.1111	8111.1111
std. dev.		1487.7111	1487.7111	1487.7111	1487.7111
weight sum		1487.7111	1487.7111	1487.7111	1487.7111
precision		1.0000	1.0000	1.0000	1.0000
Unit Price					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000
Unit Cost					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000
Total Revenue					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000
Total Cost					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000
Total Profit					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000

Total Cost					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000
Total Profit					
mean		111.1111	111.1111	111.1111	111.1111
std. dev.		111.1111	111.1111	111.1111	111.1111
weight sum		111.1111	111.1111	111.1111	111.1111
precision		1.0000	1.0000	1.0000	1.0000
Time taken to build model: 0.11 seconds					
Time taken to test model on test split: 0.11 seconds					
Summary ==					
Correctly Classified Instances		11.1111	11.1111	11.1111	11.1111
Incorrectly Classified Instances		11.1111	11.1111	11.1111	11.1111
Kappa statistic		1.0000	1.0000	1.0000	1.0000
Mean absolute error		1.0000	1.0000	1.0000	1.0000
Root mean squared error		1.0000	1.0000	1.0000	1.0000
Relative absolute error		1.0000	1.0000	1.0000	1.0000
Root relative squared error		1.0000	1.0000	1.0000	1.0000
Total Number of Instances		11.1111	11.1111	11.1111	11.1111
Detailed Accuracy By Class ==					

## ID3 Algorithm (Decision Trees)

- **Objective:** Create decision trees based on information gain.
- **Selected Columns:** Region, Item Type, Sales Channel, Order Priority, Order Date.
- **Steps:**
  - Used entropy and information gain to construct the tree.
  - Visualized decision tree structure.

- Evaluated accuracy using cross-validation.

## K-Means Algorithm (Clustering)

- **Objective:** Partition the data into clusters based on similarity.
- **Selected Columns:** Units Sold, Unit Price, Total Revenue, Total Profit.
- **Parameters:**
  - Number of Clusters (K): Determined using the elbow method.
  - Initialization: Used k-means++ to enhance convergence.
- **Steps:**
  - Standardized the data.
  - Applied K-Means clustering algorithm.
  - Visualized clusters and centroids.

```

=== Run information ===

Scheme: weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -R 2 -A
Relation: 10000 Sales Records-weka.filters.unsupervised.attribute.Remove-RI-3,11
Instances: 10000
Attributes: 5
    Unit Sold
    Unit Price
    Total Revenue
    Total Cost
    Total Profit
Test mode: evaluate on training data

=== Clustering model (full training set) ===

kMeans
=====

Number of iterations: 17
Within cluster sum of squared errors: 1784.5004961635325

Initial starting points (random):

Cluster 0: 1000,1000,1000,1000,1000
Cluster 1: 1000,1000,1000,1000,1000
Cluster 2: 1000,1000,1000,1000,1000
Cluster 3: 1000,1000,1000,1000,1000
Cluster 4: 1000,1000,1000,1000,1000

Missing values globally replaced with mean/mode

```

```
Initial starting points (random):  
Cluster 0: 7879, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78  
Cluster 1: 1147, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78, 136, 78  
  
Missing values globally replaced with mean/mode  
  
Final cluster centroids:  
  
Attribute          Full Data      Clusters  
                  (10000.0)    (2470.0)           0           1  
-----  
Units Sold         18797.4444     1135.1111     6943.3333     6943.3333  
Unit Price         191.8444      226.4111     136.1111     136.1111  
Total Revenue       36099.7778    254911.1111    94555.5556    94555.5556  
Total Cost          18797.4444     1135.1111     6943.3333     6943.3333  
Total Profit        17302.3333    141359.9999    87612.2223    87612.2223  
  
Time taken to build model (full training date) : ~1.5 seconds  
  
=== Model and evaluation on training set ===  
  
Clustered Instances  
  
(N=1) 115.  
(N=1) 78.  
(N=1) 78.
```

## Evaluation Metrics

- Used metrics such as accuracy, precision, recall, F1-score, Silhouette Score, and Inertia.

## Results and Insights

- Visualizations provided insights into sales trends, classification outcomes, and cluster formations.
- Performance scores highlighted the effectiveness of different algorithms for the dataset.