# Documentation Data Mining

by:

- عبدالرحمن رعدان
- مصعب الحبيشي
- هادي الجمعي
- مالك المصبحي
- عبدالعزيز عبدالغني
- ماجد النائب
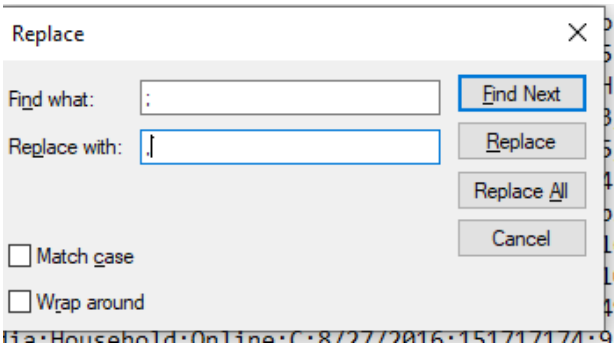
**اشراف :**

أ. ابراهيم الذارحي

# 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:** Clean the "Country" Column Data **Objective**: Ensure data quality by addressing issues with special characters. In our database, we found that the value "Cote d'Ivoire" in the "Country" column contains an apostrophe, which is causing an error in the Wiki.
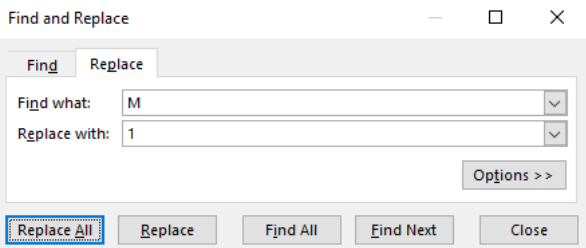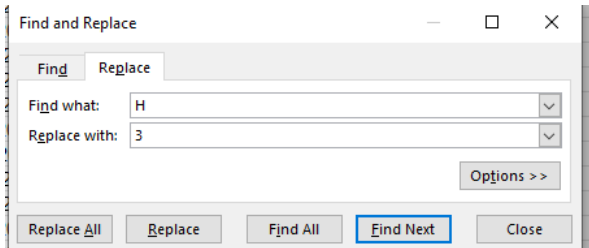
  **Problem**: The value "Cote d'Ivoire" includes an apostrophe that leads to errors during data analysis or query execution.

- **Step** 3: Re-encoding the "Order Priority" Column We performed re-encoding of the textual values in the "Order Priority" column into numerical values to facilitate processing by machine learning algorithms.

  The values were transformed as follows:

```
Critical → 3
High → 2
Medium → 1
Low → 0
```

  This step was necessary because some algorithms, such as Naïve Bayes and K-Means, require numerical data for analysis and processing.

# 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
Relation:     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: ٠.١ (1000 instances)
Minimum metric <confidence>: ٠.٥
Number of cycles performed: 18

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:

١. Order Priority=C 2555 ==> Sales Channel=Online 1326    <conf:(٠.٥٢)> lift:(١.٠٣) lev:(٠) [32] conv:(١.٠٣)
٢. Order Priority=L 2494 ==> Sales Channel=Online 1270    <conf:(٠.٥١)> lift:(١.٠١) lev:(٠) [7] conv:(١.٠١)
٣. Order Priority=H 2503 ==> Sales Channel=Offline 1266   <conf:(٠.٥١)> lift:(١.٠٢) lev:(٠) [29] conv:(١.٠٢)
٤. Order Priority=M 2448 ==> Sales Channel=Online 1228    <conf:(٠.٥)> lift:(٠.٩٩) lev:(٠) [-10] conv:(٠.٩٩)
```
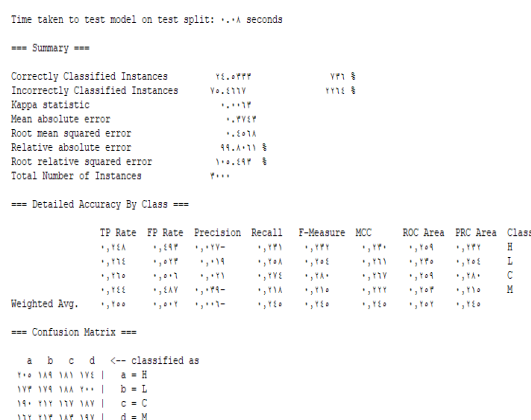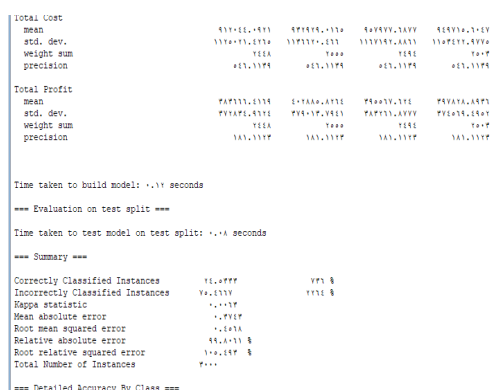
## **Algorithm and Reasoning:**

The Apriori algorithm is used to identify frequent itemsets in a dataset and generate association rules. It helps in discovering interesting relationships between variables in large databases.

**Parameters:**

- **Support Threshold:** A minimum percentage of records in the dataset that contain the itemset. In our case, with 10,000 records, a support threshold of 0.05 means any itemset must appear in at least 500 transactions to be considered significant.
- **Confidence Threshold:** The likelihood that a rule is true for the dataset. A threshold of 0.5 means we are interested in rules where the likelihood is at least 50%.
- **Lift:** Measures the importance of a rule. Lift > 1 indicates a strong association.

## 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.

### Algorithm and Reasoning:

Naïve Bayes is a probabilistic classification algorithm based on Bayes' Theorem. It assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

**Relevant Columns:**

- **Region:** Useful for understanding geographic influence on purchasing behavior.
- **Country:** To identify differences between countries.
- **Item Type:** To know the types of products that are bought together.
- **Sales Channel:** To understand the differences between various sales channels (online and offline).
- **Order Priority:** To understand how order priority affects product relationships.

---

## 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 -N 2 -A '
Relation:     10000 Sales Records-weka.filters.unsupervised.attribute.Remove-R1-8,11
Instances:    10000
Attributes:   5
              Units 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: 1784.5004861635325

Initial starting points (random):

Cluster 0: ۷۴۱۴,۱۱۸.۲۷,۱۹۰۷۱۰۲.۱۱,۴٦۹۰۱۰۹.۲۲,۱۲۱۱٦۵۵.۴٩
Cluster 1: ٤٤٩١,٤٢١.٨٩,۱۸۹۱۸۱۷.٤٤,۱۱۴٩٦٤٦.۲٤,٥٥۷۱۷۱.۲

Missing values globally replaced with mean/mode
```

```
Initial starting points (random):

Cluster 0: ۷۴۱۴,۱۱۸.۲۷,۱۹۰۷۱۰۲.۱۱,۴٦۹۰۱۰۹.۲۲,۱۲۱۱٦۵۵.۴٩
Cluster 1: ٤٤٩١,٤٢١.٨٩,۱۸۹۱۸۱۷.٤٤,۱۱۴٩٦٤٦.۲٤,٥٥۷۱۷۱.۲

Missing values globally replaced with mean/mode

Final cluster centroids:
                              Cluster#
Attribute        Full Data         0            1
                 (10000.0)     (2470.0)     (7530.0)
====================================================
Units Sold        ٤٥۴۳.٩۰۰۲    ۱٤۴.٤۰۴٦     ۸۰۰۲.٨۸٥۹
Unit Price        ۱۷٤.۳۸۸۷     ٥٥۴.۴۱۱٥     ۲٦۸.۱٤۴۱
Total Revenue     ۱۳۹٤۷۰.۲۴۰۹  ۴٤٤۸۱۸٥.٤۱۱٤  ۱۴۴۴٥٥۵.۷۴۱٤
Total Cost        ۴٩۱۴۷۳.۸۸۰۲  ۴۱۰۰۱٦۴.٤۱٤۸  ۴۴۸۱٥۰.۷۸۱٩
Total Profit      ۲٤٤٩٦.٤٥۵۷   ٤۲۷۷۸٥.٥٥۱۱   ۴٥۰۰۸٩.۴٤۲۲


Time taken to build model (full training data) : ۰.۲۸ seconds

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

Clustered Instances

0      ۲٤۷۰     ( ٢٥٪)
1      ۷٥۳۰     ( ٧٥٪)
```

# Evaluation Metrics

**Feature Importance**



**Cluster Visualization**



# Results and Insights

**Cluster Distribution**



**Decision Rules**