



**ADDIS ABABA UNIVERSITY**  
**SCHOOL OF COMMERCE**  
**MASTERS OF BUSINESS INFORMATION**  
**SYSTEMS (BIS)**

**Mining Community Voices: Enhancing Humanitarian**  
**Accountability through Data Mining in Ethiopia**

A THESIS SUBMITTED TO THE SCHOOL OF COMMERCE ADDIS ABABA  
UNIVERSITY IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE  
DEGREE OF MASTER OF SCIENCE IN BUSINESS INFORMATION SYSTEMS

**BY**

**HAYELEGEGBREAL SEYOUM**

**JUNE 2025**

**ADDIS ABABA UNIVERSITY**  
**SCHOOL OF COMMERCE**  
**MASTERS OF BUSINESS INFORMATION**  
**SYSTEMS (BIS)**

**APPLYING DATA MINING TO ANALYZE COMMUNITY FEEDBACK**  
**FOR ENHANCED ACCOUNTABILITY IN HUMANITARIAN AID:**  
**INSIGHTS FROM THE NATIONAL ACCOUNTABILITY WORKING**  
**GROUP OF ETHIOPIA**

A THESIS SUBMITTED TO THE SCHOOL OF COMMERCE ADDIS ABABA  
UNIVERSITY IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE  
DEGREE OF MASTER OF SCIENCE IN BUSINESS INFORMATION SYSTEMS

**BY**

HAYELEGEGBREAL SEYOUM

Name and Signature of Members of Examining Board

NAME	SIGNATURE	DATE
Chairperson, Examining Board		
Advisor		
Examiner		
Examiner		

# CONTENTS

List of Figures, Tables.....	iv
List of abbreviations.....	1
Acknowledgement.....	2
Abstract.....	3
CHAPTER ONE: INTRODUCTION.....	5
1.1. Background of the Study.....	5
1.2. Statement of the Problem.....	7
1.3. Research Question.....	8
1.4. Research Objective.....	8
1.5. Significance of the Study.....	9
1.6. Scope of the Study.....	9
1.7. Definition of Terms.....	10
CHAPTER TWO: LITRATURE REVIEW.....	12
2.1.Data Mining Overview.....	12
2.2.Data Mining, OLAP, and Data Warehousing.....	12
2.3. Data Mining Process.....	13
2.3.1 Data Mining and Knowledge Discovery Process (KDP).....	13
2.3.2 The CRISP-DM Process.....	14
2.3.3 Hybrid Models.....	16
2.3.4 SEMMA.....	16
2.4 Data Mining Tasks.....	17
2.4.1 Classification.....	17
2.4.2 Clustering.....	18
2.4.3 Association Rule Mining.....	18

2.4.4 Sequential Pattern Mining .....	19
2.5 Data Mining Techniques for Humanitarian Contexts .....	19
2.5.1 Clustering Techniques in Humanitarian Feedback Analysis .....	19
2.5.2 Classification Techniques in Humanitarian Data Mining .....	20
2.5.3 Evaluating the Effectiveness of Segmentation and Classification Models .....	21
2.6 Application of Data Mining Techniques in the Humanitarian Industry .....	22
2.6.1 General Application of Data Mining in Humanitarian Operations.....	22
2.6.2 Data Mining in Feedback Management within Humanitarian Response .....	23
2.6.3 Research Gaps and Challenges in Feedback Management Studies .....	23
2.7 Related Works .....	24
CHAPTER THREE: METHODOLOGY .....	27
3.1. Description of the Study Area .....	27
3.2. Research Approach .....	27
3.3. Research Design .....	27
3.4. Research Framework.....	28
3.5. Population and Sample .....	29
3.6. Data Sources and Types.....	29
3.7. Data Collection Procedures .....	30
3.8. Ethical Considerations .....	30
3.9. Data Analysis.....	30
CHAPTER FOUR: BUSINESS UNDERSTANDING, DATA UNDERSTANDING AND DATA PREPROCESSING .....	32
4.1 Business understanding.....	32
4.1.1 Ethiopia Inter-Agency Accountability Working Group .....	32
4.1.2 Why community feedback.....	33
4.2 Data Understanding .....	34

4.2.1 Data Collection.....	34
4.2.2 Data Description.....	38
4.2.3 Data Mining Goal.....	42
4.3 Data pre-processing.....	43
4.3.1 Data Cleaning.....	43
4.3.2 Data Transformation.....	46
4.3.3 Data Reduction.....	47
4.3.4 Final Dataset Structure for Analysis.....	48
CHAPTER FIVE: MODEL BUILDING AND EVALUATION.....	50
5.1 Experimental Design.....	50
5.1.1 Format of The Dataset.....	51
5.2 Apriori Algorithm Model.....	52
5.3 K-Mean Clustering Model.....	57
5.4 Decision Tree Classification Model.....	60
5.5 Comparative Analysis of the Three Experiments.....	64
5.6 Evaluation by Domain Experts.....	69
CHAPTER SIX: SUMMARY, CONCLUSION AND RECOMMENDATIONS.....	71
6.1 Summary.....	71
6.2 Conclusion.....	72
6.3 Recommendations.....	73
References.....	75
Appendix 1 List of attributes selected.....	78
Appendix 2: Final data used.....	79
Appendix 3: Detailed Summary of Apriori Association Rule Mining Results.....	80
Appendix 4: Configuration Settings used in WEKA for each model.....	82

# LIST OF FIGURES, TABLES

## List of Figures

Figure 1:Steps of KDD. Adapted from Fayyad, Piatetsky-Shapiro & Smith (1996).....	14
Figure 2:CRISP-DM methodology Source: (Chapman, 2000).....	15
Figure 3: Screenshot of raw data before cleaning (sample for region column) .....	45
Figure 4: Screenshot of cleaned data after processing .....	45
Figure 5: Clustering and merging similar entries using OpenRefine “Cluster and Edit” .....	46
Figure 6: The transformed date field. ....	47
Figure 7: The community feedback data in ARFF format .....	52

## List of Tables

Table 1: Definition of terms.....	10
Table 2: List of feedback category .....	40
Table 3: List of feedback channels used by organisations .....	41
Table 4: Description of the final dataset variables used in the analysis .....	49
Table 5: Summary of centroids.....	60
Table 6: Class-Level Performance.....	62
Table 7: Confusion Matrix Summary .....	63

# LIST OF ABBREVIATIONS

## Abbreviation Full Form

AAP	Accountability to Affected Populations
ARFF	Attribute-Relation File Format
CHS	Core Humanitarian Standard
CRISP-DM	Cross-Industry Standard Process for Data Mining
CSV	Comma-Separated Values
DW	Data Warehouse
EDA	Exploratory Data Analysis
HSP	Humanitarian Partnership Standards
IAAWG	Inter-Agency Accountability Working Group
IDP	Internally Displaced Person
IOM	International Organization for Migration
KDP	Knowledge Discovery Process
K-Means	K-Means Clustering Algorithm
NGO	Non-Governmental Organization
OLAP	Online Analytical Processing
OLAM	Online Analytical Mining
ROC AUC	Receiver Operating Characteristic - Area Under Curve
SAS	Statistical Analysis System
SEMMA	Sample, Explore, Modify, Model, Assess
SOM	Self-Organizing Map
TP Rate	True Positive Rate
UN	United Nations
UNHCR	United Nations High Commissioner for Refugees
WEKA	Waikato Environment for Knowledge Analysis

## **ACKNOWLEDGEMENT**

I would want to express my deepest gratitude to the Almighty God for giving me the courage, tolerance, and perseverance to complete this research. I am extremely grateful to my advisor, Dr. Eyob N., for his invaluable expertise, relevant suggestions, and consistent motivation throughout the duration of this thesis. His expertise and guidance have been instrumental in shaping the direction and quality of this research. My heartfelt thanks are forwarded to the Inter-Agency Accountability Working Group (IAAWG) Ethiopia and its member agencies for providing access to the community feedback data and their commitment towards humanitarian accountability improvements. Their collaboration and openness facilitated this research. To family and friends, thank you for your unwavering encouragement, patience, and trust in me. You have been an ever-constant source of inspiration. Finally, I thank all the humanitarian practitioners and community members whose voices are given expression in this research. I hope that this contribution, however humble, is made to hearing their voices and improving the services they are provided.

## **ABSTRACT**

Listening to communities is not only humanitarian action's best practice, it's a necessity. But where things get tough is while aggregating the mountains of feedback collected from affected peoples, especially when analysis resources are limited. That is what this research sought to explore: how data mining tools can bridge that gap by turning raw community feedback into actionable insights. Based on information gathered by the Inter-Agency Accountability Working Group (IAAWG) in Ethiopia between 2021 and 2025, the research aimed to identify trends, group similar comments together, and predict the sentiment of the reaction. The ultimate purpose was to facilitate more accountable, transparent, and responsive humanitarian action by maximising the voices that are already being heard better. To achieve this, the study applied three basic data mining techniques, association rule mining, K-Means clustering, and decision tree classification, to a dataset of over 92,000 feedback entries. These methods were chosen to address specific research aims: identifying hidden patterns and trends, segmenting feedback along demographic and contextual parameters, and determining the viability of automated feedback classification. The study identified extremely strong associations among some areas, demographics, and feedback themes. For example, host communities within the Somali region typically worried about food security. It was possible for the clustering model to discern distinct community profiles, and the classification model was promising in the ability to determine if responses would be positive, negative, or neutral, potentially offering a tool for real-time triage of community grievances. Significantly, the findings were not only technically sound but also practically useful. Humanitarian practitioners debated the results and made sure the results agreed with their own practice knowledge. This ensured that the models were not just thought-provoking but also workable in actuality. The study also determined the strengths and weaknesses of each of the techniques, giving a balanced view of their suitability. While the classification model struggled somewhat with neutral feedback, overall model performance showed that data mining has the potential to be a valuable addition to humanitarian feedback systems, especially when coupled with human decision-making and contextual understanding. Lastly, in this research, it is illustrated that data mining is merely a technical answer, but rather a mechanism for amplifying and giving weight to the voices of the people. By transforming unstructured community feedback into structured, usable data and actionable

insights, humanitarian organizations can gain a deeper understanding of the needs and priorities of the people they serve, leading to more informed decisions, more targeted interventions, and stronger accountability. The study concludes by recommending the integration of data mining as a routine part of normal feedback analysis, the ethical handling of sensitive data, and continued exploration of these techniques through future research. By doing so, humanitarian actors can ensure that the voice of the community is not only heard but also translated into concrete action for the betterment of services and outcomes.

**Keywords:** data mining techniques, association rule mining, K-Means clustering, and decision tree classification

# CHAPTER ONE: INTRODUCTION

## 1.1. Background of the Study

The ongoing global humanitarian crisis is continued to be severe and it requires more cooperation, creative efforts, and stronger multilateral commitments to adequately address the problems of the affected populations. In 2025, the crisis is defined by intensifying needs that is impacting well over 300 million. They are brought about by armed conflict, climate variability, and economic decline, each adding to the suffering of millions worldwide. As conflicts sustain displacement and food insecurity stay elevated, while the impact of climate-related disasters worsens underlying vulnerabilities.(UNOCHA, 2024). In Ethiopia, an estimated 21.4 million individuals need humanitarian assistance due to ongoing conflict, drought, and economic adversity, with displacement and food insecurity being key concerns (UNOCHA, 2024)

Feedback mechanisms enables communities to voice their concerns, share their needs, and assess the quality of the services they receive, is also one of the primary tactics used by humanitarian organizations, which are essential in responding to crises by meeting the immediate as well as long-term needs of impacted populations. According to (Bonino, et al., 2014), feedback mechanisms are crucial for advancing accountability to affected populations (AAP), as they facilitate participation and trust between communities and humanitarian groups. They offer a channel of communication between humanitarian actors and affected communities, ensuring that they are heard and their interests are included in decision-making.

Ethiopia Inter-Agency Accountability Working Group (IAAWG-E) is a coalition of humanitarian agencies working in Ethiopia that aim to comply with the IASC Commitments to Accountability to Affected Populations (AAP), Core Humanitarian Standards (CHS), Sphere and the entire Humanitarian Partnership Standards (HSP). The working group was initiated in Addis Ababa in 2009 by organizations interested in advancing accountability in humanitarian and development work through the use of the HAP standards. The working group consists of 51 members (NGO and UN, Cluster Lead Agencies etc.), with IOM Ethiopia taking the lead Co-Chair and Plan International Ethiopia as supporting Co-Chair. The IAAWG offers a shared platform for gathering and acting on community feedback, thereby enabling a shared approach to

accountability. With a shared framework and fostering best practice, the IAAWG strives to make humanitarian action knowledge-informed by affected individuals.

Ethiopia is a country which is frequently impacted by conflict, drought, and displacement, so that, the delivery of essential services by humanitarian organizations is vital. Humanitarian organizations employ feedback mechanisms to assess and adjust their services in a bid to respond more effectively to the needs of vulnerable populations. Feedback mechanisms have emerged as a key tool for enhancing transparency and responsiveness within the humanitarian community. Input from the community is regularly gathered by organizations and sent to IAAWG for coordinated reporting and analysis. In order to better match community needs with humanitarian activities, different types and amounts of input are collected, depending on the organization's coverage area and feedback channels.

The huge amount of gathered data , data quality problems, and the absence of appropriate tools to process and analyze such data render effective feedback data analysis challenging despite the extensive adoption of feedback systems. This thus renders change implementation challenging. Humanitarian organizations often struggle and are unable to interpret this data to make informed decisions owing to limited resources and the sophistication of the humanitarian response context. They thus deny themselves the considerable knowledge that could more effectively provide service and advance their own accountability to affected populations. As such, while feedback is collected regularly, its potential to influence decision-making and improve humanitarian practices remains largely untapped (Bonino, et al., 2014).

Since 2021, the Inter-Agency Accountability Working Group (IAAWG) Ethiopia has been collecting and analysing feedback data monthly, resulting in regular reports that highlight the community's needs and concerns. The historical data that has been gathered over the last four years, however, has not been examined for any underlying trends or patterns offering a chance to use data mining techniques that allows more efficient and fact-based decision-making.

This research aimed to fill this gap by analyzing four years' worth of IAAWG Ethiopia feedback data to uncover patterns and trends that could enhance the efficiency and effectiveness of humanitarian responses in Ethiopia. Humanitarian organizations can gain a deeper understanding of recurring issues, address service gaps, and prioritize resource allocation based on actual, real-

time community needs by applying data mining methods to community feedback data. Data mining plays a crucial role in extracting hidden patterns within large datasets, revealing insights that might otherwise go unnoticed and enabling decision-makers to make data-driven, strategic interventions(Han, et al., 2012). The aim of the project was to increase the use of feedback data for future humanitarian crises through data mining tools and providing a more systematic approach to data analysis. The data from the past was used in the hope of discovering something that can help to make better decisions, be more accountable and ultimately increase effectiveness of humanitarian aid in Ethiopia.

## **1.2. Statement of the Problem**

Humanitarian organizations rely on feedback that is coming from affected communities in order to improve their services and to be accountable. Over the past several years, there has been an increase in the size and type of feedback coming in, however, the technological, infrastructural, and operational issues have prevented this feedback from the communities being effectively analyzed, to the extent that it is not analyzed at all. As a result of this critical issues raised by affected populations may go unnoticed, leading to gaps in service provision and weakened trust between communities and humanitarian organizations. The challenge lies not in the availability of large volumes of feedback data but in the failure to leverage advanced technological methods to extract meaningful insights. While data mining has revolutionized sectors such as business intelligence and healthcare, its potential in humanitarian feedback analysis remains largely unexplored, resulting in missed opportunities for data-driven decision-making and enhanced accountability.(Kondraganti, et al., 2022)

Although the humanitarian sector acknowledges the importance of community feedback, much of the data collected is either underused or completely ignored (Bonino, et al., 2014). This gap between data collection and decision-making results in missed opportunities to improve response efforts. If feedback is not properly analyzed, key trends and issues may go unnoticed, leading to delays or ineffective responses to urgent community needs. Given the limited resources available in humanitarian operations, it is crucial to prioritize issues effectively, but without clear insights from feedback data, making informed decisions can be difficult(Bonino, et al., 2014).

The failure to properly analyze and respond to community feedback weakens accountability and ruin trust between humanitarian organizations and the people they serve. Affected communities may feel unheard when their feedback does not lead to a meaningful improvements in the services they receive. Repeatedly collecting feedback without a visible action can lead to frustration and reduce people's willingness to participate in future feedback efforts.

There is a lack of research on how data mining can be used to analyze community feedback in humanitarian operation. This study aimed to fill that gap by trying to explore how advanced data mining techniques can transform raw feedback in to useful insights, and also by ultimately improving service delivery and accountability in humanitarian responses.

### **1.3. Research Question**

The research question was formulated to address the primary problem identified in the study:

#### **Primary Research Question**

How can community feedback data be analyzed to enhance accountability and humanitarian response in Ethiopia using data mining techniques?

#### **Sub-Questions**

1. What patterns and trends can be identified in the community feedback data collected by the IAAWG Working Group?
2. How can these patterns inform decision-making to address the concerns of affected populations?
3. What are the challenges and limitations of applying data mining techniques in analyzing community feedback in humanitarian contexts?

### **1.4. Research Objective**

#### **General Objective**

The general objective the study wasto analyze community feedback data using data mining techniques to enhance accountability and improve humanitarian response in Ethiopia.

## **Specific Objectives**

To achieve the general objective, this study had the following specific objectives

1. To identify patterns, trends, and associations within community feedback entries collected by IAAWG.
2. To evaluate how the insights generated from data mining can be used to shape decision-making and improve service delivery for affected populations.
3. To assess the effectiveness and limitations of data mining techniques in processing large-scale community feedback datasets in the humanitarian sector.

### **1.5. Significance of the Study**

This study aimed to assess the applicability of data mining techniques in the humanitarian industry to identify patterns, trends, and associations of community feedback.

As this research was conducted for academic purpose the researcher gained an experience of conducting a research and also, because there aren't many studies in this field the results of this one might inspire other researchers to carry out more research. The result of the research also will help humanitarian organizations to improve community response and respond to feedback in a better way. The work of IAAWG can also be supported through advanced analysis that could help to setup scientific recommendations.

### **1.6. Scope of the Study**

This study explores how data mining techniques can be used to identify useful patterns and insights from the Inter-Agency Accountability Working Group (IAAWG) gathered community feedback data in Ethiopia. The dataset employed for this study spans four years, namely from January 2021 to February 2025, and provides a complete set of attributes varying from demographic details, feedback channels, feedback categories, date, and sector.

In order to guide the analysis of the research the study used the Cross-Industry Standard Process for Data Mining (CRISP-DM) and employed three main methods or algorithms namely K-Mean, J48 and Apriori. K-Means clustering was used in order to group similar feedback entries and

identify common themes between different sections of the population. J48 decision tree classification was used for type or category of feedback prediction from a given set of input features. Further, Apriori association rule mining was applied to find frequent patterns and correlations among various feedback attributes, including the way in which some types of concerns follow each other in a given area or amongst certain demographics.

## 1.7. Definition of Terms

<b>Data Mining:-</b>	A process of analyzing large datasets to find patterns, trends, or relationships that can be used to make decisions (Han, et al., 2012).
<b>Accountability to Affected Populations (AAP):-</b>	A framework to ensure humanitarian organizations listen to, engage with, and respond to the needs of the people they serve (CHS Alliance, 2020).
<b>Feedback Channel:-</b>	The methods or tools used to collect feedback from communities, such as surveys, hotlines, or community meetings.
<b>Feedback Category:-</b>	The methods to categories feedback from communities as positive and negative.
<b>Predictive Modelling:-</b>	A technique used to forecast outcomes based on historical data.
<b>Feedback Sector/Sector of Intervention:-</b>	The themes or topics under which feedback is classified, such as health, education, or protection.

Table 1: Definition of terms

## 1.9. Organization of the Study

Following this introductory chapter, the research will proceed with Chapter 2: Literature Review, which delves into previous studies on data mining within accountability frameworks and humanitarian contexts, providing a solid foundation for understanding the study's relevance and positioning it within the broader academic conversation; Chapter 3: Methodology will then outline the study area, research design, data collection techniques, and analytical procedures,

offering a comprehensive view of how the research was structured and conducted; Chapter 4: Business Understanding, Data Understanding, and Data Pre-processing will follow the CRISP-DM framework to define the study's context and detail the preparatory steps taken before modeling; Chapter 5: Model Building and Evaluation will focus on the development of the model and the evaluation process used to assess its performance and effectiveness; and finally, Chapter 6: Summary, Conclusion, and Recommendation will bring together the main findings, highlight their contributions to the field, and offer thoughtful suggestions for future research and practical application.

## **CHAPTER TWO: LITRATURE REVIEW**

### **2.1.Data Mining Overview**

Data mining is an evolving powerful process which integrates sets of techniques and tools to uncover rich patterns and connections between a very large sets of data and insights which can make informed and reliable predictions (Han, et al., 2012). It is a rapidly developing topic in business and academic communities and weaves together numerous fields and spans a wide variety of applications(Fayyad, et al., 1996). On a small leveldata mining facilitates improved business decision-making by improving the understanding of business environments and enabling organizations to gain a more competitive advantage (Turban, et al., 2018). On a bigger level, organizations leverage data mining algorithms and software to create and maintain a detailed data-based customer relationships to provide customers with more personalized products and services (Chen, et al., 2014).

Data mining is one the best technique that utilise statistical methods to examine big datasets in order to find patterns and trends that might otherwise go overlooked making it especially useful when dealing with large, complex datasets, where it can pull out insights that are both surprising and valuable (Han, et al., 2012). It's also a major part of what's known as the knowledge discovery process, which is all about turning raw data into something meaningful and useful (Fayyad, et al., 1996). To solve real-world problems, data mining uses different models, like grouping similar items, predicting future outcomes, or spotting trends over time, and these models help businesses make better decisions, prepare for what's coming, and fix issues before they grow into bigger problems (Han, et al., 2012)

### **2.2.Data Mining, OLAP, and Data Warehousing**

Different technologies are used to help analyze data and support better decision-making, and among them are tools like OLAP and data warehousing. OLAP, which stands for Online Analytical Processing, allows users to explore and analyze summarized data from databases or data warehouses, making it easier to spot trends, patterns, and insights that support strategic decisions.(Han, et al., 2012) explain that OLAP's capability to view data from multiple perspectives plays a crucial role in simplifying the data mining process. By providing

summarized data, OLAP facilitates the extraction of hidden insights, which can be used to inform decision-making. Before acting on patterns discovered through data mining, analysts typically use OLAP to understand their potential implications.

The integration of OLAP and data mining is known as OLAM (Online Analytical Mining), which is particularly important because data mining tools require clean and consistent data for accurate analysis (Han, et al., 2012). This integration streamlines data analysis and supports effective decision-making.

A data warehouse (DW) is a centralised data repository that stores integrated, historical data from multiple sources to aid in decision-making. (Han, et al., 2012) describe a data warehouse as a specialized database maintained separately from operational databases. Data warehousing is an evolving process that unfolds in stages, including the creation of a Virtual Data Warehouse, Data Mart, and Enterprise Data Warehouse. These stages allow for the aggregation of specialized or organization-wide data, enhancing decision-making at various levels.

## **2.3. Data Mining Process**

### **2.3.1 Data Mining and Knowledge Discovery Process (KDP)**

Since many scholars use the terms Data Mining (DM) and the Knowledge Discovery Process (KDP) interchangeably, there is sometimes confusion between the two. However, per (Cios, et al., 2007), data mining is only one phase of the larger KDP.

(Cios, et al., 2007) define the Knowledge Discovery Process (KDP) as the process of looking for fresh information in a particular field. In order to improve decision-making and obtain a competitive edge, KDP entails a number of operations intended to reveal hidden knowledge within databases. Finding legitimate, original, practical, and eventually intelligible patterns are the four fundamental criteria that constitute the KDP. Understanding data access and storage, applying effective algorithms for data analysis, making sure that interpretation and visualization are correct, choosing

In contrast, data mining is seen as one of the KDP's components. It's the stage where we dig into large sets of data using specific techniques to find useful information—kind of like uncovering

hidden gems based on certain rules or (Fayyad, et al., 1996). There are five main steps that usually make up the data mining process:

1. Selection:- Creating a target data set or focusing on specific variables or data samples to be analyzed.
2. Preprocessing:- Cleaning and preprocessing the target data to ensure consistency.
3. Transformation:- Applying dimensionality reduction or transformation methods to the data
4. Data Mining:- Searching for patterns in the data, typically aimed at prediction
5. Interpretation/Evaluation:-Analyzing and evaluating the mined patterns to extract useful insights.

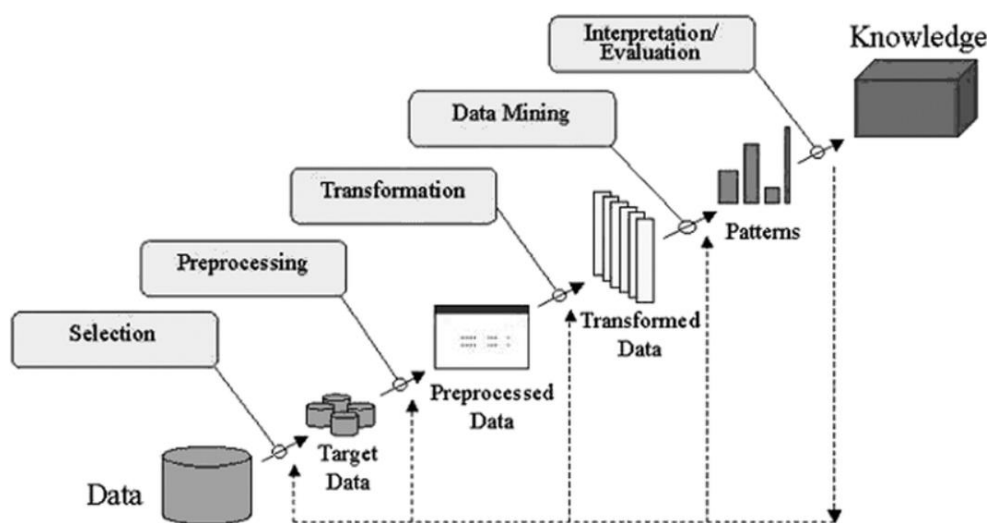


Figure 1: The steps of the Knowledge Discovery in Databases (KDD) process as adapted from Fayyad, Piatetsky-Shapiro, and Smith (1996)

### 2.3.2 The CRISP-DM Process

NCR, SPSS, and Daimler Chrysler formed a cooperation to develop a six steps cyclical process, which is designed to solve data mining problems called the CRISP-DM (CROSS-Industry Standard Process for Data Mining) process (Chapman, 2000). They are the following stages

1. **Business Understanding**:-This phase involves the comprehension of the objectives and requirements from an organizational point of view. It includes the transformation of the business issue into a data mining problem along with the determination of an initial approach to achieve the established goals.

2. **Data Understanding**:-During this stage, preliminary datasets are gathered, and the interrelations between the data's characteristics are analyzed to establish new knowledge. In addition, it is necessary to assess both the description and the quality of the data.
3. **Data Preparation**:-This process helps for getting the final dataset ready for use by the model in modeling software. It calls for the task of data purification to address the issues of noise, outliers, and missing values. The dataset is then formatted properly for the modeling phase without losing its original meaning.
4. **Modeling**:-Here, the appropriate modeling techniques, i.e., neural networks or decision trees, are selected.
5. **Evaluation**:-This process employ indicators which are specific and try to evaluate the performance and efficiency of the models.
6. **Deployment**:- In the last stage, the acquired knowledge is implemented in business activities, and documentation is done for future use.

These six steps guarantee a methodical approach to data mining, enhancing the procedure's efficacy and efficiency(Chapman, 2000).

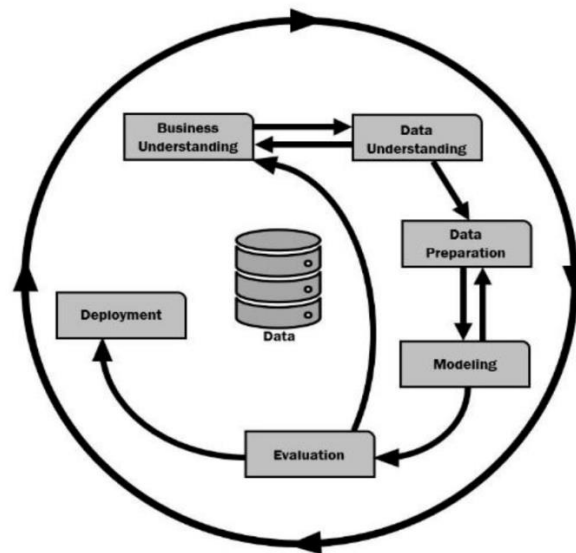


Figure 2:CRISP-DM methodology Source: (Chapman, 2000)

### 2.3.3 Hybrid Models

Hybrid models improve the process of knowledge discovery by integrating both industrial and academic perspectives within data mining endeavors. As offshoots of the CRISP-DM model, hybrid models are predominantly research-focused, with emphasis on the entire data mining process and not the modeling process alone. Feedback loops are incorporated in the six phases of hybrid models, and thus, results in one domain can be applied in other domains(Cios, et al., 2007). The six phases are:

- **Understanding the Problem Domain:** This step involves defining the problem, determining project goals, and consulting domain experts. The goals are then transformed into data mining (DM) objectives, and preliminary DM tools are selected.
- **Understanding the Data:** In this step, data is collected, and its quality is assessed. Tasks include checking for completeness, redundancy, and missing values, ensuring the data's relevance to the DM goals.
- **Preparation of the Data:** Data preparation involves sampling, checking correlations, cleaning, and handling missing or noisy data. Dimensionality reduction techniques like feature selection are used, and new attributes may be derived. The dataset is then prepared for modeling.
- **Data Mining:** DM methods are applied to the pre-processed data to discover patterns and insights.
- **Evaluation of the Discovered Knowledge:** The results of the DM models are evaluated to determine if the knowledge is novel and useful. Domain experts help interpret the findings, and the process may be revised to improve results.
- **Use of the Discovered Knowledge:** This final step involves planning how the discovered knowledge can be applied in other domains, implementing the findings, documenting the process, and deploying the model.

### 2.3.4 SEMMA

The SEMMA (Sample, Explore, Modify, Model, Assess) model was created by SAS Institute to provide a data mining project methodology. The model is composed of five separate steps, as a

cyclical process of iterative experiment, and contains statistical techniques along with visualization tools to facilitate the data mining process (Olson & Delen, 2008). These steps are:

- **Sample:** Here a representative sample of data is extracted from a larger dataset which is large enough to include significant data but small enough for quick manipulation.
- **Explore:** Here in this phase, it involves exploring the data for trends and anomalies, helping better understand the dataset and identify potential insights.
- **Modify:** Here the necessary adjustments are made to the data, such as creating, selecting, and transforming variables for the modeling process.
- **Model:** Here the appropriate modeling techniques are applied to explain patterns in the data.
- **Assess:** Here on the last stage of the process the effectiveness and reliability of the models are evaluated to estimate their performance and usefulness.

## 2.4 Data Mining Tasks

We need to utilise specific data mining tasks to find useful patterns and insights of a data and the choice of which task to use really depends on kind of problem we're trying to solve and the type of patterns we expect to find. Due to no single method works for every situation, the technique has to be chosen carefully to match our goal. Some of the most common tasks in data mining include classification - which sorts data into categories; clustering - which groups similar items together; and association rule mining - which finds relationships between items like in market basket analysis.

### 2.4.1 Classification

Classification is a fundamental task in data mining and is one form of supervised learning, where the process is applied to predict both continuous numeric outputs and categorical (nominal) outputs based on input variables (Aggarwal, 2015). The model of classification delimits the input and output clearly, and they utilize inputs in predicting the result (output).

Classification constitutes two levels:

- **Classifier Construction:** a classification model is built using a training set where each sample belonging to a predefined class and the model is created by learning from the dataset and is represented through classification rules, decision trees, or mathematical formulas.

- **Usage of the Classifier:** Once the classifier is built, it is used to predict or classify unknown objects based on patterns observed in the training set.
  - In classification, data is mapped into predefined classes. The classes are defined before the analysis begins, which makes this a supervised learning method and the main goal is to use what's been learned from known examples to predict the category or class that a new instance is most likely to belong to, based on its observed characteristics or features.

### **2.4.2 Clustering**

Clustering works a bit differently than classification. It's used when we don't know how the data is organized and need to figure that out first. It helps by dividing a mixed group of data into a smaller group which are called clusters. The grouping is happen based on things they have in common. Unlike classification, clustering doesn't rely on any predefined labels or categories. As (Han, et al., 2012) noted, "the class labels are not present in the training data simply because they are not known to begin with."

Clustering aims to position objects in such a manner that every group is maximally similar and maximally dissimilar to each other. Usually used when we don't already know how the data is structured or organised and we need to discover it. The group, or clusters, formed through this process can be seen as categories of similar object from which useful rules and patterns can later be identified.

### **2.4.3 Association Rule Mining**

Association rule mining is an unsupervised data mining endeavour that aims to systematically discover hidden relationships or meaningful associations between various items or attributes that are present within a given database (Berry & Linoff, 2004). Unlike classification, association rule mining does not forecast a target output variable but rather detects patterns of co-occurrence.

Association rule mining is particularly useful for analyzing past events, acquisitions, or interactions to unearth hidden associations. A common application is market basket analysis, where relationships between products bought together are discovered. According to (Berry & Linoff, 2004), association discovery techniques help identify affinities between items based on their co-occurrence in records or transactions.

#### **2.4.4 Sequential Pattern Mining**

Sequential pattern mining is closely related to association rule mining but differs by considering the time sequence of events. Sequential patterns detect associations that occur in a specific order, taking into account the temporal element.

(Tsiptsis & Chorianopoulos, 2009)highlighted how sequential pattern mining takes into account the time order of events and accordingly recognizes patterns that result in a particular outcome. Sequential pattern mining is effective in scenarios where event timing is of extreme importance, like predicting disease development from a patient's past, analyzing user access patterns on websites, or inventory control in retail environments.

Sequential pattern mining is directed towards the identification of relationships among transactions in a sequence, whereas association rules primarily reveal relationships present within a single transaction.

### **2.5 Data Mining Techniques for Humanitarian Contexts**

Data mining techniques, often used in customer segmentation and classification in business contexts, can be equally valuable in analyzing community feedback and stakeholder behaviors within humanitarian settings. These methods allow organizations to track trends, classify answers, and establish experiential information that refines response and accountability mechanisms.

#### **2.5.1 Clustering Techniques in Humanitarian Feedback Analysis**

Clustering is a basic method of unsupervised learning that is applied in the grouping of similar data points without any pre-existing labels. Clustering in the case of humanitarian interventions may be utilized in the grouping of community feedback into specific categories depending on identifiable patterns or behaviour, such as various complaints, requests, or needs for services.

#### **K-Means Clustering**

This is one of the most widely used well-established clustering algorithms designed specifically to divide a dataset into a number of groups (K) which are predefined with the goal of keeping similar items grouped closely together within each cluster while ensuring that each group remains as distinct and separate as possible from the others(Berry & Linoff, 2004). In

humanitarian settings K-Means can be utilized to identify clusters of communities that share similar needs or reactions. This helps organizations tailor interventions based on the identified groups. As many studies have shown that K-Means clustering is effective in large-scale data processing, making it ideal for segmenting feedback from large affected populations (Berry & Linoff, 2004)

### **Hierarchical Clustering**

Unlike the K-Means clustering algorithm the hierarchical clustering algorithm does not require the user to specify the number of clusters in advance. This generates a tree-like structure (also known as a dendrogram) and which visually represents the nested groupings of data by illustrating how individual data points are progressively merged into clusters based on their similarity (Berry & Linoff, 2004). This can be useful when there is uncertainty about the exact number of segments or when feedback data needs to be grouped progressively. Hierarchical clustering has been successfully applied in humanitarian settings to identify subgroups within communities based on varying levels of vulnerability or need (Tsipsis & Chorianopoulos, 2009).

### **Self-Organizing Maps (SOM)**

The Self-Organizing Map (SOM) is an artificial neural network that projects high-dimensional data sets onto lower-dimensional grids. SOM use is most useful in representing intricate patterns in humanitarian data because it assists organizations in comprehending the relationship between various stakeholder groups or patterns of feedback along multiple dimensions (Kohonen, 2001). For instance, SOM can be employed to cluster feedback data based on sentiment or urgency, providing actionable insights for program improvement.

### **2.5.2 Classification Techniques in Humanitarian Data Mining**

Classification is the process of systematically sorting data into pre-defined categories based on their features and it is a key method commonly used in supervised learning within data mining (Berry & Linoff, 2004). In humanitarian work, classification models help predict or categorize new community feedback by recognizing patterns learned from past data.

- **Decision Tree Classifiers:** Decision trees are a popular method for classification tasks since they are interpretable and support both numerical and categorical variables. In humanitarian contexts, decision trees can classify feedback into predefined categories, e.g., the nature of

assistance needed or levels of satisfaction. The model works by recursively splitting the data based on attribute values that best separate the classes. Decision trees have also been utilized in humanitarian crises to forecast the probability of specific needs among a population, thereby facilitating optimal resource distribution (Han, et al., 2012).

- **Pruning Decision Trees:** In order to improve the generalizability of decision tree models, pruning is used to remove branches that may lead to overfitting. Pruning enhances the decision tree's accuracy in predicting outcomes for new, unseen data (Olson & Delen, 2008). In humanitarian data analysis, pruning can be employed to ensure that decision trees used for feedback classification remain effective even as new data is incorporated.
- **Random Forests and Ensemble Methods:** Random Forests, an ensemble method that combines numerous decision trees, can enhance classification accuracy and improve robustness by minimizing the possibility of overfitting and managing large datasets efficiently. Within the scope of humanitarian feedback analysis, Random Forests are used for classifying feedback along different dimensions like urgency, severity, and type of intervention needed. Ensemble methods' success lies in the fact that they aggregate the predictions of numerous models and hence generate more stable prediction (Breiman, 2001).

### 2.5.3 Evaluating the Effectiveness of Segmentation and Classification Models

The execution of clustering and classification models needs to be carefully assessed to ascertain their ability to yield valuable information. Assessment techniques comprise internal measures, for example, silhouette scores for clustering models, and external measures, for example, accuracy and F1 scores for classification models.

- **Clustering Evaluation:** In assessing clustering algorithms such as K-Means and hierarchical clustering, the assessment process usually involves quantifying how effective clusters are in distinguishing data points and how tight the clusters are. One of the methods applied for this is the silhouette score, which quantifies the level of separation between clusters. High silhouette scores show that the clustering process is effective at connecting similar data and neatly distinguishing dissimilar groups (Rousseeuw, 1987).
- **Classification Evaluation:** The effectiveness of models at classification can be assessed based on accuracy measures, precision, recall, and the F1 score. These measures give

practical information regarding the model's capacity to predict the correct class, which is highly relevant in humanitarian contexts where the speedy and correct classification of community needs is essential for effective intervention(Powers, 2011).

## **2.6 Application of Data Mining Techniques in the Humanitarian Industry**

Data mining techniques have received much attention as they can revolutionize many sectors, including humanitarian action where they are designed to improve informed decision-making, render responses more effective, and eventually lead to more impactful results for affected communities. While they have shown promising results in some domains, their utilization for community feedback management, a crucial domain for the improvement of emergency response and accountability practice, is yet not well understood.

### **2.6.1 General Application of Data Mining in Humanitarian Operations**

The humanitarian community is faced with various challenges in responding to crises due to the intricate working environments. Issues like logistical barriers, limited resources, and rapidly changing conditions require organizations to make timely, well-informed, and data-driven decisions. Increasingly, data mining techniques, in particular clustering and classification methodologies, are being employed to handle and interpret large sets of heterogeneous data such as demographic statistics, satellite images, and community needs assessments.

In the past few years, the humanitarian community has brought data mining into their processes to improve the distribution of resources, forecast imminent needs, and streamline supply chains. Predictive modeling among others is used to assess the chances of particular crises or forecast patterns in terms of displacement, disease spread, and health risks (Caldwell, et al., 2022). Furthermore, clustering algorithms have been used to classify affected populations into groups of vulnerability determinants in order to enable the prioritization of relief disbursement. Classification processes are also used to predict regions likely to achieve specific humanitarian consequences in order for intervention to be provided on time (Caldwell, et al., 2022).

The use of data mining techniques for feedback management in the humanitarian sector is still largely unexplored, despite the fact that similar applications have shown success in operational areas like resource distribution and supply chain management.

### **2.6.2 Data Mining in Feedback Management within Humanitarian Response**

Community feedback management is a critical part of the humanitarian response system, allowing organizations to gauge the needs and satisfaction of affected communities, thereby ensuring interventions are in accordance with their set priorities. But there are tremendous challenges in handling and analyzing community feedback due to the enormous volume, intricate nature, and heterogeneity of the data. Methods like surveying and interviewing, which can be termed as traditional, are time-consuming and can fail to address the broad spectrum of issues in the community in a timely manner.

Data mining techniques such as sentiment analysis, clustering, and classification are of tremendous value in this paradigm. Sentiment analysis, a text mining-specific field, is increasingly employed to analyze qualitative feedback data that can include complaints, suggestions, or queries. Sentiment analysis of community feedback assists organizations in identifying what new problems are arising, prioritizing interventions that are required, and enhancing their communications. Furthermore, clustering is employed in grouping feedback according to themes or topics, thereby assisting in identifying common issues within various segments of the affected population.

Despite the prospective usefulness of such techniques, research exclusively addressing the use of data mining in feedback management within the humanitarian field is still limited. The literature confirms that although various studies on the application of data mining in general humanitarian activities abound, few of them have dealt with its application in the management and analysis of feedback from affected communities. This constraint demonstrates a profound shortage in existing literature, as effective handling of feedback is instrumental to enhancing accountability and program design.

### **2.6.3 Research Gaps and Challenges in Feedback Management Studies**

The use of data mining techniques in the management of community feedback in humanitarian settings is confronted with numerous challenges. Among the key obstacles is heterogeneity and complexity of feedback data, which may consist of both structured and unstructured forms of data (e.g., open-ended questionnaire answers, social media, and text messages). Effective handling of such disparate data sources demands sophisticated text mining and natural language

processing methods, which have not been adequately explored in humanitarian feedback management (Aggarwal, 2015).

Moreover, Besides, feedback data is usually noisy and inconsistent, and it is hard to discover meaningful patterns from it. To circumvent this issue, preprocessing methods such as data cleansing and transformation are necessary. Nevertheless, limited studies have comprehensively investigated the efficacy of these preprocessing measures in the perspective of humanitarian feedback data (Delen, 2021)

Lastly, limited research exists in bringing together feedback data with other kinds of humanitarian data (e.g., demographic data and needs assessments) for the purpose of developing an aggregate picture of community needs. This approach of joining different data sources can offer a complete picture of a community's needs and help plan better and more effective aid efforts. Unfortunately, studies into how different data mining methods can be used with varied data sources in aid feedback management is lacking.

## **2.7 Related Works**

Data mining has demonstrated significant potential across a variety of industries offering valuable insights for decision-making, trend identification, and predictive modelling. The integration of this into humanitarian works specifically in the context of feedback management has not been fully explored. This section tries to provide a review of existing literature related to data mining techniques, their applications in related fields, and the gaps in humanitarian research.

As there is a limited direct research on applying data mining to accountability mechanisms in humanitarian sectors a broader review of data mining applications in industries such as customer relationship management, public service and assessment of insurance risk can provide a foundation for this research. For instance, (Denekew, 2003) applied clustering techniques to segment frequent flyer program members at Ethiopian Airlines, enhancing customer relationships by identifying high-value customers. This method of classifying large datasets using K-means clustering could also be used to categorize feedback from affected populations, enabling more targeted interventions. In addition, in the insurance sector, (Hintsay, 2002) used predictive modelling to assess risk in motor insurance. His application of decision trees and neural networks for risk assessment offers similarities for applying similar methods in

identifying at-risk populations based on feedback data in humanitarian settings. Similarly, (Fikre, 2005) demonstrated the use of predictive models at Ethiopian Insurance Corporation, where decision trees were used to classify customers according to risk, providing insights on how predictive analytics could significantly enhance the prioritization of feedback and response management efforts in emergencies.

These studies highlight the potential of using data mining techniques, such as clustering and predictive modeling, to gain insights from complex datasets—insights that can be adapted to analyze feedback from affected populations in humanitarian contexts.

Moreover, (Reganie, 2013) applied data mining techniques to customer segmentation in microfinance institutions using specifically K-means clustering. This work is important for understanding how data mining can enhance segmentation and enable targeted responses, which is directly applicable in identifying different needs within humanitarian populations and tailoring response strategies.

In spite of the widespread uses of data mining in numerous fields, there exists a significant shortage of direct usage of such techniques in humanitarian feedback management. Current research concentrates largely on public sector or commercial use, hence creating a evident gap for research that tailors these techniques to humanitarian aid and accountability contexts.

The gap is particularly significant given the complexity of managing large amounts of feedback data in crisis situations. Humanitarian agencies require practical ways of categorizing, analyzing, and acting on crisis-affected populations' feedback. The purpose of this study is to fill this gap by employing data mining techniques to humanitarian feedback mechanisms and create the basis for more accountable and responsive humanitarian action.

While there is no direct literature in the humanitarian sector, the methodologies used in related fields provide useful frameworks. For instance, decision trees, clustering techniques, and predictive modeling techniques that have been extensively studied in the fields of customer relationship management and insurance provide a theoretical basis to conduct the analysis of humanitarian feedback data. These techniques can be modified to address the unique challenges present in humanitarian situations, which comprise variability in data, missing records, and real-time information requirements.

This thesis built upon current methodologies, with an emphasis on applying them to the specific situation of humanitarian feedback management. The effectiveness of these approaches could be analyzed with respect to their capacity to increase accountability and refine response mechanisms during crises.

An exploratory investigation was conducted to provide preliminary information on the feasibility of using data mining in feedback mechanisms in humanitarian situations. This obtained valuable insights on the challenges and possibilities of using data mining in this environment, and the results could contribute towards further refining accountability processes.

In addition, this examination draws upon case studies within humanitarian agencies that have participated in the use of data mining methods. While these are not formal publications, the empirical knowledge of institutions like the UNHCR and ICRC is sure to provide important lessons and observations concerning the potential applications of data mining in the specific context of feedback management.

The absence of any directly relevant existing work on the application of data mining in feedback management in humanitarian situations makes this research even more innovative. By bridging this gap, this thesis also aims to be of use to both the humanitarians and the research community. The social contribution of the research can be realized in its usefulness to strengthen the accountability mechanisms, refine response processes, and thereby improve the effectiveness of humanitarian interventions.

## **CHAPTER THREE: METHODOLOGY**

In order to meet the goal of the study appropriate data mining methodology was selected based on the nature of the problem identified. Also, extensive effort was made to review related literatures from previous studies.

### **3.1. Description of the Study Area**

This study focused on the IAAWG feedback data collected in Ethiopia. The working group brings together various humanitarian organizations to ensure that the voices of affected populations are heard and that their feedback is used to improve service delivery. Feedback is collected from communities across multiple regions of Ethiopia, particularly in areas affected by ongoing crises, including displacement due to conflict and natural disasters as the country often face severe challenges, including limited access to basic services and infrastructure.

The study area specifically included feedback from the IAAWG common dataset which cover diverse sectors such as health, education, food security, and protection. The feedback is gathered through multiple channels including surveys, community meetings, mobile hotlines etc.

### **3.2. Research Approach**

The study relied on statistical techniques to interpret the feedback data, aiming to provide objective and measurable insights that can be applied in improving accountability to affected populations. It adopted a quantitative research approach, utilizing data mining techniques to analyze large datasets as this approach is suitable for uncovering patterns, trends, and associations in large-scale data that may not be easy through traditional analysis methods.

### **3.3. Research Design**

The research have both natures for exploratory and nature of descriptive aiming to provide a detailed understanding of the current situation while also uncovering new patterns, relationships, or insights that may not be immediately apparent. Descriptive research will allow for a comprehensive examination of the feedback data and provide an understanding of the overall trends within the dataset. Exploratory research will help uncover patterns and relationships that

may not be immediately obvious, providing deeper insights into the issues affecting affected populations.

This design is appropriate because as it aims to explore how data mining can be applied to humanitarian feedback data and identify patterns that can enhance accountability mechanisms.

### **3.4. Research Framework**

For the purpose of directing the analysis in this research CRISP-DM (Cross-Industry Standard Process for Data Mining) approach was chosen as the research framework. CRISP-DM is considered to be one of the most widely used and well-respected data mining techniques available today (Chapman, 2000). It's practical, versatile, and easy to implement in many different industries and fields of study—like humanitarian interventions.

One of the primary reasons to choose CRISP-DM is that it is structured but iterative. It divides the data mining process into six stages: Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment. The steps are not linear in nature; rather, they provide for forward-and-back movement, which is extremely useful when dealing with real data that might mean working back through previous steps (Chapman, 2000).

Below is a brief explanation of how each step was carried out in this research

#### **Business Understanding**

This step was to know the objectives of the Inter-Agency Accountability Working Group (IAAWG) and how communities' input may be utilized more effectively to improve humanitarian response and accountability.

#### **Data Understanding**

During this stage, IAAWG feedback data collected were analyzed to assess their quality, format, and functionality. This identified issues or gaps early on.

#### **Data Preparation**

Raw data were cleaned, converted, and reorganized to get the data ready for examination. This involved replacing missing values, standardizing formats, and selecting the most informative variables.

## **Modelling**

A variety of data mining techniques were applied, including K-Means clustering to group similar feedback, decision trees to classify feedback types, and association rule mining to identify patterns and relationships.

## **Evaluation**

Models were evaluated for accuracy and usefulness, as well as for whether or not they were beneficial in the humanitarian setting, based on domain expert feedback.

## **Deployment**

Finally, the results of the analysis were interpreted and translated into practical recommendations that could guide more responsive and accountable humanitarian response.

CRISP-DM also exists at different levels of granularity—from high-level phases to specific tasks and actual implementations (Chapman, 2000). This allowed us to tailor the process to this research's particular needs without deviating from established methodology.

In brief, CRISP-DM provided us with a clear, flexible, and organized method of examining community feedback data. It allowed us to guarantee that the research was not only methodologically sound, but also grounded in the realities of humanitarian action in Ethiopia.

### **3.5. Population and Sample**

The population for this study consists of the feedback entries collected by IAAWG, during 2021-2025, approximately 93,000 feedback, across various programs.

### **3.6. Data Sources and Types**

The primary data source for this study was the feedback data collected by IAAWG Ethiopia. The data is first organised in a tabular format. Community feedback by nature is unstructured however organisations use different tools or templates to make a structured tabular format. The IAAWG nationally provides a simple template that can be utilised by organisations that doesn't have their own templates. The categorisation of information depends on the person entering the feedback. The feedback accepted populated in to the actual feedback section of the template.

However, some organisations also categories actual feedbacks for privacy or simplification purposes. The actual feedback then get analysed and the remaining attributes of the template is filled. This study focuses on the tabular data, that is already populated on the templates, to provide valuable insights into the concerns of affected populations.

### **3.7. Data Collection Procedures**

The community feedback data is collected using several feedback channels such as suggestion boxes, hotlines, helpdesks, face to face during focus group discussions etc. The collected data will be collected and populated by different organisation based on their area of response and geographic coverages. At the end of each month, they're expected to send this information to the IAAWG. This research is based on the final monthly reports that these organizations submit.

### **3.8. Ethical Considerations**

Ethical considerations take top precedence to this research, given the fact that the data under discussion directly involve vulnerable groups. Foremost among the ethical issues are:

- **Confidentiality:**-Prior to conducting the analysis, all personally identifiable information were stripped or anonymized to ensure individuals' privacy. Furthermore, details such as the organization's name and specific feedback were removed to ensure confidentiality.
- **Data Security:**-Data erased after anonymisation.
- **Non-Discrimination:**-The research ensured the perspectives of all the groups impacted, particularly from marginalized communities, are represented within the evaluation.

### **3.9. Data Analysis**

The examination was conducted using data mining techniques for pattern, trend, and relationship discovery in the feedback data. The process outlined below employed:

- **Data Preprocessing:** The first task was cleaning and preprocess the data using OpenRefine 3.9.0. This included removing duplicate records, filling or managing missing values, and normalizing text so that the data was prepared and uniform for analysis.

- Exploratory Data Analysis (EDA): As an initial step, an exploratory examination was carried out with the assistance of WEKA 3.9.6. This showed easy patterns—such as how often various feedback types occur, the most common categories, and the spread of feedback between and across groups and time.
- Pattern Discovery: For further investigation, data mining algorithms such as clustering, classification, and association rule mining were applied with the help of WEKA. Clustering enabled the collection of similar feedback, and overall themes emerged. Classification (through the J48 decision tree algorithm) enabled the prediction of categories of feedback from certain attributes. Association rule mining showed relationships between feedback attributes.

## **CHAPTER FOUR: BUSINESS UNDERSTANDING, DATA UNDERSTANDING AND DATA PREPROCESSING**

### **4.1 Business understanding**

The first step in applying data mining techniques to humanitarian feedback management is to have a solid understanding of the business environment, such as knowing the roles of the key stakeholders, the purpose of the feedback system, and the decision-making implications of the results.

#### **4.1.1 Ethiopia Inter-Agency Accountability Working Group**

The Ethiopia Inter-Agency Accountability Working Group (IAAWG) is a group of humanitarian agencies that come together to improve the way in which aid agencies address the needs of people affected by crises. From 2009, the working group has focused on ensuring the voice of the community in the humanitarian response. It has over 50 members from NGOs, UN organizations, and cluster leads, and IOM Ethiopia and Plan International Ethiopia serve as co-chairs.

One of the group's main roles is to collect and coordinate feedback from communities across the country. This feedback comes from hotlines, helpdesks, community meetings, and other channels. Member organizations share this information every month. The aim is to make sure that people's concerns, suggestions, and grievances are heard and considered during decision-making.

While IAAWG has made progress in gathering feedback, there's still a gap when it comes to analyzing the data in depth. Many of the data are put together into reports, but very little is put to the use of the newest tools to find patterns or trends that can be leveraged to improve services. That's where data mining can make a difference.

In this research, understanding how IAAWG operates is important because it unearths the potential and the limitations of the current system. With the use of data mining methodologies, the research aims to help IAAWG maximize the feedback that it receives—turning raw data into information that can generate more responsive and accountable humanitarian action.

#### **4.1.2 Why community feedback**

Community feedback is a central element of accountability in humanitarian action since it is the vehicle whereby communities voice their needs, concerns, and satisfaction regarding interventions that are underway. In crisis situations, where crisis and fast-paced change tend to generate emergent needs and frequently changing demands, feedback from affected populations provides essential information regarding the effectiveness and appropriateness of response activities.

Incorporating community feedback ensures, in addition to assuring aid agencies are acting with urgency, increasing the long-term impact of their interventions. Feedback can be utilized to identify where services aren't being provided and an early signal of impending issues, as well as effective strategies that could be scaled or replicated. Therefore, the integration of community feedback in humanitarian programming is not just a method for maximizing program outcomes but also for building trust among affected populations.

The Core Humanitarian Standard on Quality and Accountability (CHS) is an international standard that establishes nine commitments to make humanitarian responses effective, accountable, and grounded in the needs and rights of affected populations, places community feedback at the center of its strategy, with a requirement for organizations to establish mechanisms through which affected people can make their contributions, voice concerns, and leave their signature throughout the project cycle (CHS Alliance, 2020). By incorporating feedback mechanisms into the CHS, humanitarian aid agencies are able to gain transparency, genuine participation, and respect for the agency and the dignity of the population they work with. The CHS thus improves program quality and accountability, as well as affirming the key role of feedback in informing relevant and responsive humanitarian action.

In the context of this thesis, community feedback analysis is ranked among the key aspects that would immensely be contributed to by data mining methods. By analyzing large volumes of feedback data, relief organizations will be able to gain more insights into patterns, recognize areas of need, and make more informed decisions that are also representative of the community's concerns.

## 4.2 Data Understanding

A good understanding of the data on hand is essential before doing any data mining activity. This will help an effective utilisation of techniques and in return provides proper results. In this section we will look at how the working group gathers the data from the organisations, how the data is structured and what kind of gaps might exist.

### 4.2.1 Data Collection

The data used in this research was collected through the Inter-Agency Accountability Working Group (IAAWG). This group includes several organizations working within the humanitarian sector. These organizations gather feedback from affected communities using a range of methods. The goal is to capture diverse perspectives that can help shape humanitarian responses and strengthen accountability.

#### Feedback Channels Used by Organizations

All organizations that are part of the working group use different feedback collection methods to listen and connect with the disaster affected communities. The tools are designed to reach a different type of audience as they collect input from men, women, children, and vulnerable groups. After the gathering of the feedback it is submitted to the IAAWG platform to be compiled together and analysed.

The main feedback channels used by the IAAWG member organizations are the following categories:

#### Face-to-Face Channels

<b>Community Meetings and Focus Groups</b>	These are organized at distribution points, community centers, or during field visits. Affected populations participate in discussions where they express concerns, provide feedback, or ask questions regarding the ongoing humanitarian response.
<b>Community Events</b>	These are public gatherings or information sessions where humanitarian workers engage with the community. Feedback is often gathered informally during these events through discussions and surveys.
<b>Place of Worship</b>	Feedback is also gathered in places of worship where community members

	may feel more comfortable discussing their concerns in a familiar environment
--	---

### Digital Channels

<b>SMS and Telegram</b>	These platforms allow the community to send and receive messages related to the humanitarian response. SMS and Telegram are often used due to their accessibility, particularly in areas with limited internet connectivity. This method helps gather real-time feedback.
<b>WhatsApp</b>	This widely used messaging platform facilitates direct communication with communities, enabling immediate responses to questions or issues raised by affected individuals

### Phone-Based Channels

<b>Hotlines and Phone Calls</b>	Hotlines provide a direct communication line for community members to reach out to humanitarian organizations with their concerns, inquiries, or feedback regarding the ongoing response.
---------------------------------	---

### Written and Anonymous Channels

<b>Suggestion and Complaints Boxes</b>	These boxes are placed in community centers or distribution sites, allowing individuals to submit feedback anonymously encouraging individuals to share sensitive concerns without fear of retributions
<b>Logbooks</b>	Feedback is recorded in logbooks by staff or volunteers at community centers or during distributions. This method serves as a written record of verbal feedback received on-site.

### Community-Based Mechanisms

<b>Community Volunteers and Information Boards</b>	Feedback is collected through the efforts of community volunteers who engage with individuals at the grassroots level. Information boards are also placed in strategic locations to allow individuals to post their feedback or concerns
<b>Child-Friendly Spaces</b>	These spaces are designed to give children a safe environment to provide feedback, ensuring that their voices are heard, especially in emergencies where children are among the most vulnerable

### Data Storage and Management

After the feedback has been collected from the various channels, the data of every organization is being compiled and reported to the IAAWG on a monthly basis. The feedback is analysed and made available for follow-up action. To ensure an organizational outlook, the feedback is entered into a standard Excel spreadsheet. This spreadsheet structures the data in the same format so that analysis will be easier in the future.

Each row in the spreadsheet must be of significance to have one feedback record, while columns hold different information—such as the date on which the feedback was received, the type of concern reported, the medium used to provide the feedback, and demographic information of the person who provided the feedback. This systematic organization makes it easier to sort, filter, and analyze by organizations and geographies.

<b>Title</b>	<b>Description</b>	<b>Data Type</b>
<b>Date Feedback Was Received</b>	The exact date on which the feedback was received.	Date (DD/MM/YYYY)
<b>Feedback Status (New or Pending)</b>	Indicates whether the feedback is newly received this month or carried over from a previous month.	Categorical ( <i>New/Pending</i> )
<b>Organisation</b>	Name of the humanitarian organization receiving and responding to the feedback.	Text
<b>Region</b>	Regional location where the feedback	Text

	was collected.	
<b>Zone</b>	Administrative zone within the region.	Text
<b>Woreda</b>	The district (woreda) where the feedback originated.	Text
<b>Gender</b>	Gender of the individual providing feedback.	Categorical ( <i>Male/Female/Other</i> )
<b>Age</b>	Age of the individual providing feedback.	Categorical ( <i>9yrs and under, 10-14yrs, 15-17yrs, 18-24yrs, 25-49yrs, 50-59yrs, 60-69yrs and above 70yrs</i> )
<b>Community Type</b>	Classification of the community	Categorical ( <i>IDP, Host, Refugee, Returnee</i> ).
<b>Vulnerability</b>	Perceived or assessed vulnerability level of the individual.	Categorical ( <i>People with Disability, Child Headed, Chronically Ill, None</i> )
<b>Language</b>	Language used in delivering the feedback.	Text
<b>Actual Feedback</b>	Verbatim text or summary of the feedback received from the community member.	Text (Long)
<b>Feedback Channel</b>	Mode through which the feedback was collected	Categorical ( <i>e.g., Helpdesk, Face-to-Face, Hotline etc..see section 4.2.1</i> ).
<b>Feedback Category</b>	Categorization and thematics of feedback. Further detail seen in table below	Categorical ( <i>e.g. Appreciation, Request for Assistance, Suggestion, Concerns etc... see section 4.2.2</i> )
<b>Feedback Concern by Sector/Subsector</b>	More specific breakdown of the feedback issue by sector and subsector.	Categorical ( <i>e.g., WASH, Protection, Education</i> )

<b>Feedback Status</b>	Current status of the community feedback	Categorical ( <i>Referred, Closed, Pending, or Unanswered</i> ).
<b>Actions Taken</b>	Description of the steps taken in response to the feedback.	Text
<b>Responsibility for Follow-Up</b>	Assigned individual or team responsible for follow-up and resolution.	Text
<b>Expected Feedback Closure Date</b>	Target date by which the feedback should be resolved.	Date (DD/MM/YYYY)
<b>Reason for Missed Closure Date</b>	Explanation if the feedback was not resolved by the expected date.	Text

Table 2: The attributes IAAWG feedback collection template

**4.2.2 Data Description**

The community feedback data set used in the study is mostly time-based data and categorical data. Gender, region, category of feedback, and community type fall under categorical variables, with the day on which feedback is given offering a temporal factor to the data set. Although age exists as a variable, it's also categorized according to pre-established age groups, making it ordinal in nature compared to a continuous numerical value. Because of such an architecture, analysis involves investigation of patterns and correlations across such categories, rather than statistical analysis requiring sequential numerical data.

**Date Feedback Was Received *DD/MM/YYYY* (Time Stamp)**

Every entry has a timestamp on when feedback was received. This is helpful for trend analysis over time so that changing community needs or emerging issues with critical need would be able to be highlighted.

### **Region, Zone and Woreda (Geographical Information)**

There is also location-based data, e.g., region, zone, or woreda/district where feedback was received. This is helpful for geographic pattern identification—whether certain problems are universal or localized.

### **Gender, Age, Community type, Language and Vulnerability (Demographic Information)**

Each feedback entry may have demographic data such as age, gender, location, and even socio-economic status sometimes. This is to know who is providing feedback and to make sure that the analysis gets the opinions of diverse groups—especially those often marginalized groups like women, children, people with disabilities, and marginalized communities.

### **Actual Feedback and Feedback Category (Feedback Types)**

"Actual Feedback" column notes the first comment or message provided by the individual—a complaint, suggestion, Appreciation, or other form of feedback. It is the respondent's voice in their own words and can contain rich context or emotion that fixed fields would not capture.

To allow for analysis of the data, each piece of feedback is also assigned a "Feedback Category" selected from a pre-established list and are typically selected by the person recording the feedback or by the organisation receiving the feedback. Capturing the categories allow for the structuring of the information and support more intense analysis of recurring patterns or concerns expressed by affected communities.

<b>Feedback Category</b>
Appreciation
Complaints about project/program service
Concerns
General Feedback
Questions incl Requests for Information
Request for Assistance

Suggestions
-------------

Table 3: List of feedback category

## Feedback Channels

This column indicates the channel by which each feedback was obtained. Feedback can come from many sources, including hotlines, surveys, public meetings, and suggestion boxes. Each channel provides a different perspective. Face-to-face communication, for example, will receive more textured and detailed responses, while electronic mechanisms may capture more structured or categorized feedback. Knowing what channel is utilized is significant since it has the potential to impact both the level and nature of information obtained.

<b>Feedback Channel</b>
Complaints/ Suggestions Box
Face to face (at distribution/service sites)
Face to Face (interviews, focus group discussions, community meetings, etc)
Face to face (Community events)
Face to face (Place of worship)
Community volunteers
Community information boards
Helpdesk
Youth Centre
Child friendly spaces
Hotline - open
Phone calls to communities
SMS
E-mail

Facebook
Twitter
WhatsApp
Other social media (specify one)
Other internet chat platforms (Please specify)
IEC materials
Radio
Other

Table 4: List of feedback channels used by organisations

### **Feedback Status**

This column shows the present status of each feedback entry, how far it has progressed in the response process. The statuses are as follows:

- Closed – Feedback has been solved and addressed.
- Pending – The comment has been received but is still under review or awaiting action.
- Referred – The complaint has been referred to some other unit/organization for follow-up.
- Unanswered – There is nothing done yet, and the feedback still isn't answered.

### **Action Taken**

For the majority of cases, the dataset includes a record for what response or action has been taken by the organization in response to the feedback. This allows measurement of how well the

feedback loop is being closed and whether community concerns are being well addressed in a timely manner.

### **Data Size and Volume:**

The dataset is expected to be large, especially during peak periods of humanitarian response when the affected population may submit a higher volume of feedback. As feedback is collected over time, the data size may grow significantly, making it important to implement efficient data mining and analysis techniques to manage large datasets.

### **4.2.3 Data Mining Goal**

The primary goal of applying data mining techniques to this dataset is to enhance the accountability and responsiveness of humanitarian organizations by deriving actionable insights from community feedback. To be more specific this study sought to achieve the following objectives:

1. **Categorization and Classification:** Use clustering and classification techniques to categorize the feedback based on its content (e.g., complaints, suggestions, requests for assistance). This will allow humanitarian organizations to prioritize issues and identify patterns in feedback.
2. **Identifying Trends:** Detect recurring issues or emerging trends in the feedback data helping organizations to adjust their response strategies in real time and address the most pressing concerns raised by the affected population.
3. **Predictive Modeling:** Develop predictive models to forecast the likelihood of certain issues arising in the future. e.g, predictive models could be used to identify which types of feedback are likely to be associated with higher urgency or risk, allowing organizations to take preemptive actions.
4. **Segmentation of Populations:** Apply segmentation techniques to group the feedback by demographic factors (e.g., age, gender, location) to understand if certain groups are facing unique challenges or have specific concerns that need targeted interventions.

## **4.3 Data pre-processing**

### **4.3.1 Data Cleaning**

#### **Missing Data Handling**

One of the first and most important steps in preparing the data for analysis as part of preparing the dataset was missing data handling. Because of the nature of humanitarian feedback gathering—occasionally conducted in challenging environments and in diverse regions—it is not surprising that some of the entries would be missing.

For the sake of providing the quality and consistency of the analysis, rows with high missing values were examined carefully and removed. These were the instances where some significant key fields were missing, making them unqualified for meaningful interpretation or modeling. This was done with care not to miss valuable data but also considering the fact that missing data can introduce noise and lower the effectiveness of data mining techniques.

After this filtering exercise, the dataset was then reduced from the original 94,013 feedback records to 92,895. While this reduced it by a small extent, it greatly improved the overall consistency and usability of the data. The resulting dataset continues to represent a rich and diverse foundation upon which to identify patterns and insights that can support more responsible and responsive humanitarian action.

#### **Duplicates**

On raw feedback processing into categorized forms, it was revealed that some of the entries were found to be duplicates during processing. For example, several feedback messages such as "we need additional water support" and "we need sanitation support" both fell into the category of "Request for Assistance" under the WASH sector. Since these entries also shared similar metadata (such as gender and submission date), they became duplicates in the processed dataset. While these are not in the literal sense duplicates, they are alternative feedback by alternative individuals. As the aim of this research is to analyze patterns and trends within categorized feedback, these were retained. To keep them on record guarantees that frequency and distribution of feedback types are accurately represented within the analysis.

## **Further cleaning process**

To tidy the dataset for analysis, I used OpenRefine, which is a powerful and specially designed tool for reshaping and cleaning data. Given the size and nature of the feedback dataset OpenRefine made the more interactive and hands-on process of identifying and correcting inconsistencies easier.

The first task was to validate the dataset for inconsistencies within categorical fields such as gender, region, and categories of feedback. Using OpenRefine's faceting feature, I was able to quickly spot spelling, case, and capitalization variations (e.g., "female" vs. "Female" or "wash" vs. "WASH"). Those were normalized to ensure consistency across the dataset.

Next, I addressed missing values. Faceting columns to identify blank entries allowed me to ascertain which fields had gaps and if gaps should be filled, flagged, or the affected rows deleted entirely. This made the data more complete as a whole without compromising integrity.

I also made use of OpenRefine's clustering and editing feature to merge near-duplicate values—useful for cleaning out entries that were virtually identical but worded slightly differently. This was a critical process in order not to fragment categories like sectors and feedback channels by variable naming.

Finally, I normalized the date fields so that they took on a uniform format, which would later allow for time-based analysis.

OpenRefine Compiled Data csv [Permalink](#)

Facet / Filter Undo / Redo 10 / 10 **92894 rows**

Refresh Reset all Remove all Show as: rows records Show: 5 10 25 50 100 500 1000 rows

Region	Date Feedback Was Received (DD/MM/YYYY)	Is this a new feedback from this month or still pending/answered from last month?	Region	Gender
Addis Ababa	2024-07-01	New	Tigray	Male
Afan Oromo	2024-07-01	New	Tigray	Male
AFAR	2024-07-02	New	Tigray	Male
Afar	2024-07-02	New	Tigray	Male
Amahara	2024-07-03	New	Gambela	Female
Amhara	2024-07-03	New	Somali	Male
Anyuak	2024-07-10	New	Oromia	Male
Arabic	2024-07-10	New	Gambela	Male
Benishangul Gumuz	2024-07-10	New	Somali	Male
Benishangul-Gumuz	2024-07-12	New	Amhara	Male

Figure 3: Screenshot of raw data before cleaning (sample for region column)

OpenRefine Compiled Data csv [Permalink](#)

Facet / Filter Undo / Redo 39 / 39 **92894 rows**

Refresh Reset all Remove all Show as: rows records Show: 5 10 25 50 100 500 1000 rows

Region	Date Feedback Was Received (DD/MM/YYYY)	Is this a new feedback from this month or still pending/answered from last month?	Region	Gender
Addis Ababa	2024-07-01	New	Tigray	Male
Afan Oromo	2024-07-01	New	Tigray	Male
AFAR	2024-07-02	New	Tigray	Male
Afar	2024-07-02	New	Tigray	Male
Amahara	2024-07-03	New	Gambela	Female
Amhara	2024-07-03	New	Somali	Male
Anyuak	2024-07-10	New	Oromia	Male
Arabic	2024-07-10	New	Gambela	Male
Benishangul Gumuz	2024-07-10	New	Somali	Male
Benishangul-Gumuz	2024-07-12	New	Amhara	Male

Figure 4: Screenshot of cleaned data after processing

## Cluster and edit column "Feedback Category"

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. [Find out more...](#)

Method Key collision Keying function Fingerprint 15 clusters found

Cluster size	Row count	Values in cluster	Merge?	New cell value
8	22175	<ul style="list-style-type: none"><li>Request for Assistance (19677 rows)</li><li>Request for assistance (2430 rows)</li><li>Request for assistance (43 rows)</li><li>request for Assistance (13 rows)</li><li>request for assistance (7 rows)</li><li>Request For Assistance (2 rows)</li><li>Request for assistance (2 rows)</li><li>Request for assistance</li></ul>	<input checked="" type="checkbox"/>	Request for Assistan
5	28648	<ul style="list-style-type: none"><li>Concerns (28006 rows)</li><li>concerns (482 rows)</li><li>Concerns (156 rows)</li><li>Concerns (2 rows)</li><li>concerns (2 rows)</li></ul>	<input type="checkbox"/>	Concerns
3	259	<ul style="list-style-type: none"><li>Request for Information (194 rows)</li><li>Request for information (43 rows)</li><li>Request for information (22 rows)</li></ul>	<input type="checkbox"/>	Request for Informati
3	1620	<ul style="list-style-type: none"><li>Suggestions (1582 rows)</li><li>Suggestions (30 rows)</li><li>suggestions (8 rows)</li></ul>	<input type="checkbox"/>	Suggestions

Select all Deselect all Export clusters Merge selected & re-cluster Merge selected & Close Close

Figure 5: Clustering and merging similar entries using OpenRefine "Cluster and Edit"

### 4.3.2 Data Transformation

To better understand when the feedback was submitted, I split the initial date column (Date Feedback Was Received) into more readable and analyzable fields: day of the week. This was done with basic date formatting functions that split the entire date into its weekday name (e.g., "Friday").

This conversion facilitated simpler access for investigation of temporal patterns—i.e., if there are weekly trends, e.g., whether additional feedback is incurred at the start or end of the week. Such knowledge can help humanitarian agencies better interpret community engagement behavior and potentially streamline when and how they receive feedback.

	A	B	C	D	E
	Date Feedback Was Received (DD/MM/YYYY)	Day	Is this a new feedback from this month or still pending/unansw	Region	Gender
1					
2	1-Jul-24	Monday	New	Tigray	Male
3	1-Jul-24	Monday	New	Tigray	Male
4	2-Jul-24	Tuesday	New	Tigray	Male
5	2-Jul-24	Tuesday	New	Tigray	Male
6	3-Jul-24	Wednesday	New	Gambela	Female
7	3-Jul-24	Wednesday	New	Somali	Male
8	10-Jul-24	Wednesday	New	Oromia	Male
9	10-Jul-24	Wednesday	New	Gambela	Male
10	10-Jul-24	Wednesday	New	Somali	Male
11	12-Jul-24	Friday	New	Amhara	Male
12	12-Jul-24	Friday	New	Tigray	Male
13	15-Jul-24	Monday	New	Amhara	Male
14	15-Jul-24	Monday	New	Gambela	Male
15	15-Jul-24	Monday	New	Amhara	Male
16	16-Jul-24	Tuesday	New	Amhara	Male
17	17-Jul-24	Wednesday	New	Amhara	Male
18	17-Jul-24	Wednesday	New	Gambela	Male
19	17-Jul-24	Wednesday	New	Oromia	Male
20	18-Jul-24	Thursday	New	Amhara	Male
21	18-Jul-24	Thursday	New	Gambela	Male
22	19-Jul-24	Friday	New	Amhara	Male
23	19-Jul-24	Friday	New	Gambela	Male
24	19-Jul-24	Friday	New	Tigray	Male
25	22-Jul-24	Monday	New	Amhara	Male
26	23-Jul-24	Tuesday	New	Afar	Male
27	23-Jul-24	Tuesday	New	Amhara	Male
28	24-Jul-24	Wednesday	New	Amhara	Male
29	25-Jul-24	Thursday	New	Afar	Male
30	26-Jul-24	Friday	New	Tigray	Male

Figure 6: The transformed date field.

### 4.3.3 Data Reduction

A meticulous inspection of the dataset was undertaken before initiating the data cleaning task to guarantee that only clean, ethical, and useful analytically data would remain available for analysis. There were initially **94,013** feedback comments in the dataset, each with a set of **20** attributes. All these columns were not open to analysis for practical and ethical purposes, though.

To begin with, several columns were removed to protect the confidentiality of groups and individuals. Included in these were the **names of organizations** that collected comments, the **actual text** of the comments themselves, **actions taken** as follow-up, and the names or titles of the individual who conducted **follow-up**. These columns, although potentially useful information, contained sensitive data that would be breaching confidentiality. Since as per the data sharing guidelines of the participating organizations, these columns were omitted in order to uphold ethical standards as well as the anonymity of all the parties involved.

In addition, certain columns were also deleted because they were not specifically relevant to the research objectives. For example, columns like the **expected close date** for feedback and the **reason for delay** in closing it were operational in nature and did not provide much value to the analytical intent of this research.

The **Language** column was also removed due to a high rate of missing data. In most cases, the language that the feedback was provided in was not recorded, and therefore the field was of no use in analysis.

Lastly, the **Zone** and **Woreda** columns, which were very fine-grained geographic data, were removed in an effort to preserve a coarser regional focus. Since the analysis needed to be run at a regional level, retaining these very fine-grained geographic identifiers was not required and could have added unnecessary complexity without creating value.

These choices were undertaken to clean up the dataset, minimize noise, and guarantee that the analysis remained concise, ethical, and consistent with the goals of the study.

#### 4.3.4 Final Dataset Structure for Analysis

After data cleaning and transformation, the list of variables to be used in analysis was finalized. The columns were chosen according to their relevance to the research objectives and their potential contribution to meaningful insights. The Date Feedback Was Received column was excluded from the final dataset, as its most significant components—month and day—were already split and stored for temporal analysis.

The final dataset includes the following columns:

<b>Column</b>	<b>Description</b>
Day	– to identify patterns by weekday.
New or unresolved	– to identify new and outstanding feedback.
Region	– to analyze geographic distribution of feedback
Gender	– to search for gender variation in feedback
Age	– to observe patterns by age group.
Community type	– to divide IDPs, host communities, refugees, etc.
Feedback Channel	– to analyze how feedback is being collected.
Feedback Type	– to classify the type of feedback (e.g., concern, suggestion).

Feedback Concern by Sector/Subsector	– to identify thematic areas of concern.
Feedback Status	– to track whether feedback has been resolved or is still pending.

Table 5: Description of the final dataset variables used in the analysis

This dataset is used as the foundation for the subsequent analysis using data mining techniques in WEKA.

## **CHAPTER FIVE: MODEL BUILDING AND EVALUATION**

In the last chapter, we discussed the data understanding and preparation stages extensively. These processes were crucial in defining a clean and organized dataset for analysis. The researcher's understanding of the dataset and the study objectives informed the process of preparing data. The final result of that stage was a finalized dataset, specifically designed for model-building purposes that would be capable of discovering meaningful information.

This chapter then proceeds to the next critical stage: selecting and applying the appropriate data mining techniques and measuring how effectively they operate. Tools and algorithms were not chosen in a random manner—it was guided by the previously defined objectives of the business understanding phase, augmented by the nature of the data itself.

To perform all the steps of modeling and analysis, the researcher used Weka 3.9.6, an extremely easy-to-use and powerful data mining software. The decision to use Weka was based on its ease of use, flexibility, and excellent support for clustering and classification that fits the research. The research flows from data preprocessing through pattern discovery and model building that can aid humanitarian organizations in heightening their understanding and response to community feedback using Weka.

### **5.1 Experimental Design**

It is not always a simple choice to select the right model for a data mining project. There is not one "best" algorithm that does well at everything; that is why this is an applied-research undertaking, applying different models and assessing their performances to find out which work well with an individual dataset.

The key objective in this case is to uncover underlying correlations and patterns in the community feedback data, information that may not necessarily be obvious at first glance but could be incredibly valuable to humanitarian response and accountability. In order to do this, the research introduces two main learning approaches: supervised and unsupervised.

Supervised learning is used when we already have some idea of what we are looking for—for example, predicting the kind of feedback from familiar categories. To that end, the J48 decision tree algorithm was chosen. It's very commonly used because it's easy to understand and interpret, which is very helpful when dealing with stakeholders who don't have technical background.

Alternatively, unsupervised learning is focused on discovering natural clusters in the data without labels. Toward this end, the K-Means clustering algorithm was selected. It's a simple but powerful way to bunch similar feedback entries together to identify shared concerns or themes across several communities.

In addition to these two approaches, the study also incorporates association rule mining using the Apriori algorithm. This method is particularly useful for discovering frequent patterns and co-occurrences within the data, such as identifying which types of feedback tend to appear together in specific regions or among certain demographic groups.

The following sections describe why and how these two algorithms were used in the research.

### **5.1.1 Format of The Dataset**

Before going to the actual construction of models, the data had to be formatted in a way that is consumable by the Weka software. The data were initially saved in a file with Comma-Separated Values (CSV) which is a common and convenient method of table-based data storing. However, Weka manipulates data most effectively in a format called ARFF (Attribute-Relation File Format), so the CSV file was converted to this format.

This was translated with a plain text editor. Essentially, the rows and columns of the CSV file (CVD\_IAAWG\_Data.csv) were reorganized into two main parts that form an ARFF file (CVD\_IAAWG\_Data.arff):

**Header Section:** This section declares the dataset format and starts with a @RELATION tag that defines the name of the dataset, followed by a series of attribute using '@ATTRIBUTE' declarations that describe each column (or feature) of the dataset and its data type (i.e., numeric, nominal).

**Data Section:** This is where the actual records are listed. The values in each row correspond to a single feedback entry, which are aligned according to the attributes defined in the header.

This ARFF format made loading the dataset in Weka easier and starting the modelling process.

```

C:\WAAG_Data\aff - Notepad
File Edit Format View Help
Relation CVD_IWAAG_Data

@attribute Month {January, March, July, August, September, October, November, December, February, April, May, June}
@attribute Day {Thursday, Monday, Tuesday, Wednesday, Saturday, Friday, Sunday}
@attribute 'New or unresolved' {New, unresolved}
@attribute Region {Gambela, Somali, Oromia, Amhara, Afar, Tigray, 'Benishangul Gumz', 'South Ethiopia', 'South West Ethiopia', 'Central Ethiopia', 'Addis Ababa', Sidama, 'Dire Dawa', Harari}
@attribute Gender { 'Prefer not to say', Male, Female}
@attribute Age {10-14yrs, 15-17yrs, 18-24yrs, 25-49yrs, 50-59yrs, 60-69yrs, '9yrs and under', 'Above 70yrs', 'Preferred not to say'}
@attribute 'Community type' {Refugee, Host, IDP, Returnee, Other}
@attribute 'Feedback Channel' {'Child friendly spaces', Hotline, 'Face to Face', Helpdesk, 'Community Volunteers', Phone, 'Community information boards', 'Youth Centre', 'Social Media', 'Suggestion Box', SMS, E-mail}
@attribute 'Feedback Category' {Concerns, 'Request for Assistance', 'Complaints about project/program service', Appreciation, 'Questions incl Requests for Information', 'General Feedback', Suggestions}
@attribute 'Feedback Concern by Sector, subsector' {'Child Protection', 'Food Security', WASH, Livelihoods, Nutrition, Education, 'Multipurpose Cash', Health, Operations, 'Non Food Items', Protection, 'Gender Based Violence', CCM, Shelter}
@attribute 'Feedback Status' {Closed, Open, Pending, Referred}

@data
January, Thursday, New, Gambela, 'Prefer not to say', 10-14yrs, Refugee, 'Child friendly spaces', Concerns, 'Child Protection', Closed
January, Monday, New, Somali, Male, 10-14yrs, Host, Hotline, Concerns, 'Food Security', Open
March, Tuesday, New, Oromia, Male, 10-14yrs, Host, Hotline, 'Request for Assistance', 'Food Security', Closed
March, Tuesday, New, Somali, Male, 10-14yrs, Host, Hotline, 'Complaints about project/program service', 'Food Security', Closed
March, Wednesday, New, Gambela, Male, 10-14yrs, Refugee, 'Child friendly spaces', 'Request for Assistance', 'Child Protection', Closed
March, Wednesday, New, Gambela, Male, 10-14yrs, Refugee, 'Child friendly spaces', 'Request for Assistance', 'Child Protection', Closed
July, Tuesday, New, Oromia, Male, 10-14yrs, IDP, 'Face to Face', 'Request for Assistance', WASH, Closed
July, Tuesday, New, Oromia, Male, 10-14yrs, IDP, 'Face to Face', 'Request for Assistance', WASH, Closed
August, Saturday, New, Somali, Female, 10-14yrs, Host, 'Face to Face', Appreciation, Livelihoods, Closed
September, Thursday, New, Somali, Female, 10-14yrs, Host, Helpdesk, 'Request for Assistance', 'Food Security', Pending
September, Friday, New, Somali, Female, 10-14yrs, IDP, Helpdesk, Appreciation, 'Food Security', Closed
September, Friday, New, Somali, Female, 10-14yrs, Host, Helpdesk, 'Request for Assistance', 'Food Security', Closed

```

Figure 7: The community feedback data in ARFF format

## 5.2 Apriori Algorithm Model

### Experiment 1: High-Confidence Association Rule Mining on a 92,894-Instance Community Feedback Dataset Using the Apriori Algorithm

```

Associator output

=== Run information ===

Scheme: weka.associations.Apriori -N 10 -I 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -cc -1
Relation: CVD_IWAAG_Data
Instances: 92894
Attributes: 8
  Day
  Region
  Gender
  Age
  Community type
  Feedback Channel
  Feedback Type
  Sector

=== Associator model (full training set) ===

Apriori
=====
Minimum support: 0.1 (9289 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 16

Generated sets of large itemsets:

Size of set of large itemsets L(1): 22
Size of set of large itemsets L(2): 79
Size of set of large itemsets L(3): 84
Size of set of large itemsets L(4): 31
Size of set of large itemsets L(5): 1

Best rules found:

1. Region=Somali Feedback_Channel=Face to Face Sector=Food Security 9547 ==> Community_type=Host 9625 <conf:(0.97)> lift:(1.25) lev:(0.02) [1909] conv:(6.91)
2. Region=Somali Gender=Female Feedback_Type=Positive 15529 ==> Community_type=Host 14907 <conf:(0.96)> lift:(1.24) lev:(0.03) [2862] conv:(5.6)
3. Region=Somali Age=25-49yrs Feedback_Channel=Face to Face ==> Community_type=Host 9783 <conf:(0.96)> lift:(1.24) lev:(0.02) [1869] conv:(5.44)
4. Region=Somali Feedback_Channel=Face to Face Feedback_Type=Positive 11582 ==> Community_type=Host 11088 <conf:(0.96)> lift:(1.23) lev:(0.02) [2104] conv:(5.25)
5. Region=Somali Feedback_Type=Positive Sector=Food Security 12472 ==> Community_type=Host 11875 <conf:(0.95)> lift:(1.23) lev:(0.02) [2201] conv:(4.68)
6. Region=Somali Gender=Female Feedback_Channel=Face to Face 14532 ==> Community_type=Host 13811 <conf:(0.95)> lift:(1.23) lev:(0.03) [2539] conv:(4.52)
7. Region=Somali Feedback_Channel=Face to Face 15271 ==> Community_type=Host 17356 <conf:(0.95)> lift:(1.22) lev:(0.03) [3264] conv:(4.45)
8. Region=Somali Age=25-49yrs Feedback_Type=Positive 11429 ==> Community_type=Host 10851 <conf:(0.95)> lift:(1.22) lev:(0.02) [1986] conv:(4.43)
9. Region=Somali Feedback_Type=Positive 19517 ==> Community_type=Host 18899 <conf:(0.95)> lift:(1.22) lev:(0.04) [3450] conv:(4.39)
10. Community_type=Host Feedback_Type=Positive Sector=Food Security 12524 ==> Region=Somali 11875 <conf:(0.95)> lift:(2.09) lev:(0.07) [6202] conv:(10.54)

```

## Overview of the Model Configuration

This was one of the experiments done to study and understand trends in community feedback data by employing the association rule mining methods. For the first experiment, this study used the Apriori algorithm by Weka tool weka.associations.Apriori for frequent pattern and

association discovery from the dataset. The model was configured with a minimum support of 10 percent, meaning any rule would need to cover at least 9,289 instances in order to be considered. A confidence level of 90 percent was also established so that only highly reliable rules would be retained.

The data used in this analysis contained 92,894 cases and consisted of eight principal attributes: Day, Region, Gender, Age, Community Type, Feedback Channel, Feedback Type, and Sector. These attributes were considered to include a broad swath of demographic and contextual variables applicable to the subject of community feedback.

The main goal here was to find strong, general rules that highlight frequent patterns in the data. In other words, with a high confidence level, the model prefers associations that are not only statistically valid but also practically relevant for decision-making.

### **Summary of Discovered Itemsets**

The Apriori algorithm found 217 frequent itemsets of all sizes. These include 22 itemsets of size one, 79 of size two, 84 of size three, 31 of size four, and just one itemset of size five. This monotonic decrease in the number of itemsets as the size increases is a typical output of the Apriori process. It reflects the increasing difficulty of meeting the support constraint as the number of attributes increases..

### **Key Association Rules and Interpretations**

All the top ten association rules of the Apriori algorithm predicted the community type to be "Host." Most of these rules were conditioned on Somali being the region, which followed a strong geographic pattern. Here are some of the rules with their explanation :

#### **Rule 1**

*If Region = Somali, Feedback Channel = Face to Face, and Sector = Food Security, then Community Type = Host*

- Confidence: 97%
- Interpretation: Individuals in the Somali region who provide face-to-face feedback on food security are highly likely to be from the host community. This rule reflects a strong and reliable association.

**Rule** **2**

*If Region = Somali, Gender = Female, and Feedback Type = Positive, then Community Type = Host*

- Confidence: 96%
- Interpretation: Positive feedback from women in the Somali region is predominantly associated with the host community. This may suggest program effectiveness or satisfaction among this demographic.

**Rule** **3**

*If Region = Somali, Feedback Type = Positive, and Sector = Food Security, then Community Type = Host*

- Confidence: 95%
- Interpretation: Positive feedback on food security in the Somali region is strongly linked to host community members, reinforcing the importance of this sector in that context.

**Rule** **4**

*If Region = Somali, Feedback Channel = Face to Face, and Feedback Type = Positive, then Community Type = Host*

- Confidence: 95%
- Interpretation: This rule highlights the effectiveness of face-to-face communication in eliciting positive feedback from host communities in the Somali region.

**Rule** **5**

*If Region = Somali, Gender = Female, and Sector = Food Security, then Community Type = Host*

- Confidence: 94%
- Interpretation: Women in the Somali region who engage with food security programs are most likely from the host community, suggesting targeted engagement.

**Rule** **6**

*If Region = Somali, Age = 25–49, and Feedback Type = Positive, then Community Type = Host*

- Confidence: 94%
- Interpretation: Positive feedback from individuals aged 25 to 49 in the Somali region is strongly associated with the host community, indicating a key demographic for engagement.

**Rule** **7**

*If Region = Somali, Sector = Food Security, and Gender = Female, then Community Type = Host*

- Confidence: 93%
- Interpretation: This rule reinforces the intersection of gender and sectoral focus in identifying host community members.

**Rule** **8**

*If Region = Somali, Feedback Channel = Face to Face, and Gender = Female, then Community Type = Host*

- Confidence: 93%
- Interpretation: Face-to-face feedback from women in the Somali region is a strong indicator of host community affiliation.

**Rule** **9**

*If Region = Somali, Feedback Type = Positive, and Gender = Female, then Community Type = Host*

- Confidence: 92%
- Interpretation: This rule further confirms the trend of positive feedback from women in the Somali region being linked to host communities.

**Rule** **10**

*If Community Type = Host, Feedback Type = Positive, and Sector = Food Security, then Region = Somali*

- Confidence: 95%

- Interpretation: This reverse rule indicates that when host community members provide positive feedback on food security, it is most likely to originate from the Somali region. This insight is particularly useful for geographically targeted programming.

### **Emerging Patterns and Insights**

Several important patterns came out of the analysis. One of the key finding was that the community type "Host" consistently appeared as the outcome in the strongest rules. This suggests that host communities are most active in providing feedback, particularly in the Somali region.

Second, face-to-face emerged as the dominant feedback channel. This may reflect a desire for direct interaction in regions with fewer digital infrastructures. It also highlights the importance of leaving in-person feedback channels open to ensure inclusivity and accountability.

Third, food security was the most frequently mentioned sector in the regulations. This underscores the prevalence of food concerns in community feedback and suggests that food security interventions need to remain a top priority for monitoring and evaluation.

### **Practical Implications for Accountability Programming**

The findings of this analysis have several practical implications. To begin with, face-to-face communication in the Somali region needs to be prioritized in community engagement initiatives, particularly among host community members and women in the 25-49 years age group. This is possibly the most valid and actionable feedback.

The Somali region's high levels of positive feedback offer the possibility of that portray elements of program success that could be replicated in other regions. Understanding what is working well can inform broader program design and implementation.

The association rules can be used to enhance monitoring and evaluation systems. By identifying patterns of region, feedback type, and sector that are predictive of particular community responses, organizations can develop early warning systems and more accurately target their interventions.

### **Limitations and Recommendations**

While the results are intriguing, there are some limitations to consider. The abundance of rules related to the Somali region may either indicate a data bias, in which responses from this region are over-represented, or it may reflect a genuine heightened responsiveness. Inter-regional comparison is recommended to determine the cause.

In addition, the "Positive" feedback category is very broad. Dividing this category into more specific categories—like appreciation, satisfaction, or constructive suggestions—may provide more insight into community sentiment and increase the detail level for future analyses.

To the future, integrating these association rules in a real-time feedback dashboard can significantly enhance program responsiveness. Data visualization tools such as Power BI can be leveraged to visualize emerging patterns and support data-driven decision-making.

## **Conclusion**

The Apriori analysis yielded a subset of high-confidence, statistically significant association rules with a strong emphasis on the Somali region and host communities. The findings underscore the requirement for face-to-face feedback mechanisms and recognize food security as one of the primary concerns. The findings can immediately be used for the improvement of accountability programming, to help design community engagement strategies, and to more focused and effective interventions.

## **5.3 K-Mean Clustering Model**

### **Experiment 2: Clustering Community Feedback Using K-Means on a 92,894-Instance Dataset**

```

Clusterer output
--- 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 "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10
Relation: CVD_ILMMQ_Data
Instances: 92894
Attributes: 8
  Day
  Region
  Gender
  Age
  Community_type
  Feedback_Channel
  Feedback_Type
  Sector
Test mode: evaluate on training data

--- Clustering model (full training set) ---

KMeans
-----
Number of iterations: 3
Within cluster sum of squared errors: 309775.0
Initial starting points (random):
Cluster 0: Thursday,Tigray,Female,25-49yrs,IDE,Hotline,Negative,'Food Security'
Cluster 1: Tuesday,Somali,Female,18-24yrs,Host,'Face to Face',Neutral,Nutrition
Missing values globally replaced with mean/mode

Final cluster centroids:
Attribute          Full Data          Cluster#
(92894.0)          (45356.0)          (47538.0)
-----
Day                Monday             Thursday            Tuesday
Region             Somali              Tigray              Somali
Gender              Female              Male                 Female
Age                 25-49yrs           25-49yrs            25-49yrs
Community_type      Host                Host                 Host
Feedback_Channel    Helpdesk            Hotline              Face to Face
Feedback_Type       Neutral              Negative              Neutral
Sector              Food Security       Food Security         Food Security

```

In the second experiment, the objective was to reveal natural groupings in the community feedback dataset with unsupervised learning methods. Here, in particular, the K-Means cluster algorithm was used in the same dataset from Experiment 1 that comprised 92,894 records and eight most important attributes.

The aim was to identify particular clusters or trends in the feedback—sets of answers that share similar characteristics. These sets can reveal unique community profiles or emerging concerns, offering useful insights to make better-informed, more efficient, and more effective programming and engagement strategies.

## Model Configuration

The clustering was performed using Weka’s SimpleKMeans implementation. The model was configured with the following parameters:

- Number of clusters: 2
- Distance function: Euclidean
- Maximum iterations: 500
- Random seed: 10
- Missing values were globally replaced with mean or mode

The algorithm converged after just three iterations, indicating a relatively stable clustering structure. The within-cluster sum of squared errors was 309,775, suggesting a moderate level of compactness within the clusters.

### Initial Cluster Seeds

The algorithm began with the following randomly selected initial centroids:

- **Cluster 0:** Thursday, Tigray, Female, 25–49 years, IDP, Hotline, Negative feedback, Food Security
- **Cluster 1:** Tuesday, Somali, Female, 18–24 years, Host, Face to Face, Neutral feedback, Nutrition

These initial seeds reflect diverse demographic and contextual profiles, which helped guide the clustering process.

The **Within-Cluster Sum of Squared Errors (SSE)** was calculated at 309,775, which indicates a moderate level of compactness within the identified clusters.

### Final Cluster Characteristics

After training, the model produced two clusters of nearly equal size:

- **Cluster 0:** 45,356 instances (49%)
- **Cluster 1:** 47,538 instances (51%)

The final centroids for each cluster are summarized below:

Attribute	Cluster 0	Cluster 1
Day	Thursday	Tuesday
Region	Tigray	Somali
Gender	Male	Female
Age	25–49 years	25–49 years
Community Type	Host	Host
Feedback Channel	Hotline	Face to Face
Feedback Type	Negative	Neutral
Sector	Food Security	Food Security

Table 6: Summary of centroids

## **Interpretation of Clusters**

The below cluster results can be utilized to create tailored communication strategies and service delivery models that better mirror the needs and aspirations of each cluster.

The cluster analysis revealed two community feedback profiles that differed in terms of the following:

**Cluster 0:** This cluster has a higher likelihood of using hotline services and providing negative feedback. It is most strongly associated with the Tigray region and has a higher proportion of male respondents.

**Cluster 1:** Feedback in this cluster is predominantly collected via face-to-face and the majority of the respondents are female from the Somali region. The tone of feedback is predominantly neutral, and the age group is evenly spread like in Cluster 0.

Although both clusters are interested in topics of concern to the food security community, differences in communication channels, gender balance, and geographic distribution point to strong segmentation within the community.

## **5.4 Decision Tree Classification Model**

### **Experiment 3: Decision Tree Classification of Community Feedback Using the J48 Algorithm**

```

Classifier output
| Region = Oromia: Positive (0.0)
| Region = Benishangul Gumz: Positive (10.0/1.0)
| Region = Somali: Positive (0.0)
| Region = South Ethiopia: Nutral (1.0)
| Region = South West Ethiopia: Positive (0.0)
| Region = Central Ethiopia: Positive (0.0)
| Region = Addis Ababa: Positive (0.0)
| Region = Sidama: Positive (0.0)
| Region = Harari: Positive (0.0)
| Region = Dire Dawa: Positive (0.0)
Feedback_Channel = Suggestion Box: Nutral (26.0/8.0)
Feedback_Channel = SMS: Negative (151.0/5.0)

Number of Leaves : 1898
Size of the tree : 2274

Time taken to build model: 0.15 seconds

=== Evaluation on training set ===

Time taken to test model on training data: 0.05 seconds

=== Summary ===

Correctly Classified Instances 62345 67.1141 %
Incorrectly Classified Instances 30549 32.8859 %
Kappa statistic 0.5105
Mean absolute error 0.2911
Root mean squared error 0.3815
Relative absolute error 66.2538 %
Root relative squared error 81.3964 %
Total Number of Instances 92894

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC  ROC Area  PRC Area  Class
          0.830  0.200  0.602  0.830  0.698  0.578  0.886  0.684  Positive
          0.731  0.163  0.704  0.731  0.717  0.563  0.863  0.756  Negative
          0.508  0.124  0.722  0.508  0.596  0.420  0.771  0.723  Nutral
Weighted Avg.  0.671  0.158  0.683  0.671  0.665  0.512  0.833  0.724

=== Confusion Matrix ===

      a   b   c  <-- classified as
20561 1445 2764 | a = Positive
 4366 23514 4284 | b = Negative
 9232  8458 18270 | c = Nutral

```

The third experiment comprised the categorization of community feedback using a supervised approach. The J48 decision tree classifier, which is Weka's implementation of the C4.5 classifier, was used on the same data used in the previous experiments. The data contained 92,894 instances and eight attributes. The overall objective was to construct an interpretable model that would predict the category of feedback—positive, negative, or neutral—on the basis of demographic and contextual information.

### Model Configuration

To set up the model, the J48 classifier was set up with a confidence factor of 0.25 to prune, and a minimum of 2 instances per leaf. Missing values in the data set were addressed before training by being replaced with the mean in the case of numerical attributes and the mode in the case of categorical attributes, thereby having a complete data set to model.

### Model Performance

The decision tree generated from the process had 1,898 leaves and 2,274 nodes and took approximately 0.15 seconds to construct. The entire dataset was utilized for training and testing the model. On testing on the train set, it gave the following performance statistics:

- **Correctly classified instances:** 62,345 (67.11%)
- **Incorrectly classified instances:** 30,549 (32.89%)
- **Kappa statistic:** 0.5105
- **Mean absolute error:** 0.2911
- **Root mean squared error:** 0.3815

These indicate a moderate level of predictive accuracy, with the model performing best on positive and negative feedback classes, and less effectively on neutral responses.

### Class-Level Accuracy

The detailed accuracy by class is summarized below:

Class	TP Rate	FP Rate	Precision	Recall	F1-Score	ROC AUC
<b>Positive</b>	0.830	0.200	0.602	0.830	0.698	0.886
<b>Negative</b>	0.731	0.163	0.704	0.731	0.717	0.863
<b>Neutral</b>	0.508	0.124	0.722	0.508	0.596	0.771
<b>Weighted Avg.</b>	0.671	0.158	0.683	0.671	0.665	0.833

Table 7: Class-Level Performance

### Interpretation

- The model performs **best on Positive and Negative feedback**, with high recall (ability to find all relevant instances).
- **Neutral feedback** is harder to classify, with lower recall (50.8%), suggesting ambiguity or overlap with other classes.

- **Precision for Neutral is high (72.2%)**, meaning when the model predicts Neutral, it's usually correct—but it misses many actual Neutral cases.
- **ROC AUC scores** are strong across all classes, indicating good class separability.

#### Confusion Matrix

Actual \ Predicted	Positive	Negative	Neutral
Positive	20,561	1,445	2,764
Negative	4,366	23,514	4,284
Neutral	9,232	8,458	18,270

Table 8: Confusion Matrix Summary

The model confuses Neutral with both Positive and Negative, which explains the lower recall for Neutral. Positive and Negative classes are more clearly separated, with relatively fewer misclassifications.

#### Key Decision Patterns

The decision tree revealed several notable patterns:

- **Feedback Channel as a Primary Split:** The root of the tree was based on the feedback channel, with "Face to Face" interactions leading to more granular splits involving region, sector, age, and gender.
- **Regional and Sectoral Influence:** Regions such as Somali and Amhara, and sectors like Food Security and Education, frequently appeared in the decision paths, indicating their strong influence on feedback classification.
- **Demographic Differentiation:** Age and gender were significant in determining the sentiment of feedback. For instance, younger females in the Somali region providing face-to-face feedback on food security were more likely to give positive responses.

#### Interpretation and Implications

Experiment 3's decision tree model is a straightforwardly interpretable and transparent model for learning about the conditions under which various demographic and contextual features induce feedback sentiment. The model's 67.11% overall accuracy is poor, yet it does provide informative insight into the conditions under which various types of feedback—positive, negative, or neutral—are most likely to occur. For instance, the model identified the feedback channel as the single largest driver, and then region, gender, and sector as key determinants.

Such findings can be extremely useful for program managers who need to tailor communication interventions and strategies. For instance, being aware that face-to-face feedback in certain regions is more likely to be positive can inform the development of better engagement strategies and accountability systems.

### 5.5 Comparative Analysis of the Three Experiments

Three distinct data mining techniques—association rule mining, clustering, and classification—were applied in this study to analyze the community views collected by the Inter-Agency Accountability Working Group (IAAWG) in Ethiopia. Each of them was chosen relative to fulfilling a particular analysis goal and with the aim to produce complementary findings on the patterns, trends, and predictive associations within the data.

#### Purpose and Methodological Differences

<p><b>Experiment 1:</b> Association Rule Mining</p>	<p>the Apriori algorithm was used to mine association rules with the purpose of discovering frequent and high-confidence patterns in the data. This was to determine attribute combinations—region, gender, and type of feedback—that best predicted some outcome, primarily type of community. The most reliable patterns were those related to the Somali region and host community, whose rules were greater than 90% confidence. These trends consistently correlated face-to-face feedback, food security challenges, and female respondents to the host community, offering actionable results for geographically and demographically targeted programming.</p>
<p><b>Experiment 2:</b> K-Means Clustering</p>	<p>Experiment 2 employed K-Means clustering to partition the dataset into naturally occurring segments regardless of the application of pre-defined labels. This unsupervised learning technique detected two clusters: one of negative remarks by male interviewees in the Tigray region via hotline</p>

	channels, and another of neutral remarks by female interviewees in the Somali region via face-to-face contact. Although both clusters pertained to food security, the differences in sentiment, place, and communication mode revealed the diversity of community experience and likes.
<b>Experiment 3:</b> Decision Tree Classification	In Experiment 3, supervised learning was used to create a predictive model with the J48 decision tree algorithm. The model achieved a moderate accuracy of 67.11% and performed best with positive and negative feedback categories. The tree root node was established on the feedback channel, emphasizing its high predictive power and the model emphasised the importance of region, gender, and sector in the identification of feedback sentiment.

**Key Findings**

<b>Association Rule Mining</b>	revealed that the most frequent and reliable patterns were associated with the Somali region and the host community type. Rules with confidence levels above 90% consistently linked face-to-face feedback, food security concerns, and female respondents to the host community. These patterns offer actionable insights for geographically and demographically targeted programming.
<b>K-Means Clustering</b>	produced two distinct clusters. One cluster was characterized by negative feedback from male respondents in the Tigray region using hotline channels, while the other featured neutral feedback from female respondents in the Somali region using face-to-face communication. Despite both clusters focusing on food security, the differences in sentiment, region, and communication style highlighted the diversity of community experiences.
<b>Decision Tree Classification</b>	achieved a moderate accuracy of 67.11%, with the model performing best on positive and negative feedback categories. The root node of the tree was based on the feedback channel, indicating its strong predictive power and also the model emphasised the importance of region, gender, and sector in determining feedback sentiment.

**Comparative Insights**

<b>Aspect</b>	<b>Experiment 1: Association Rules</b>	<b>Experiment 2: Clustering</b>	<b>Experiment 3: Classification</b>
<b>Learning Type</b>	Unsupervised (Descriptive)	Unsupervised (Descriptive)	Supervised (Predictive)
<b>Primary Goal</b>	Discover frequent, high-confidence patterns	Identify natural groupings in feedback	Predict feedback sentiment
<b>Interpretability</b>	High (human-readable rules)	Moderate (requires centroid analysis)	High (decision tree structure)
<b>Strengths</b>	Reveals strong, actionable associations	Segments population for targeted response	Enables automated classification
<b>Limitations</b>	Focuses only on frequent patterns	May oversimplify with few clusters	Moderate accuracy, struggles with neutral class
<b>Best Use Case</b>	Strategic planning and policy design	Community segmentation and profiling	Real-time feedback triage and prioritization

## **Synthesis and Implications**

Combined, the three experiments—association rule mining, clustering, and classification—provide an integrated analytical system for studying and reacting to community sentiment in humanitarian situations. Association rule mining is particularly well-suited to detect robust, consistent patterns in the data that can influence strategic choice and inform resource allocation. Clustering, on the other hand, enables humanitarian practitioners to value the diversity of community life, allowing for the possibility of creating more targeted and context-specific interventions that address the concerns of specific segments of society. Meanwhile, classification systems open up the possibility of automating sorting out feedback, possibly greatly accelerating and making responses to decision-making much faster.

By integrating these techniques, humanitarian organizations can move beyond basic reporting and toward a more data-driven, accountable, and community-centered approach to program design and evaluation.

## **Analysis in Light of the Research Objectives**

This study set out to explore how data mining techniques can be used to better understand and respond to community feedback in humanitarian settings with guided analysis by three specific objectives, each of which is discussed below in relation to the findings from the three experiments.

### **Objective 1: Identifying Patterns, Trends, and Associations in Community Feedback**

The initial aim was to uncover insightful data hidden in the feedback data. All three experiments independently made a contribution towards this goal:

Associations were most clearly revealed by the Apriori algorithm in Experiment 1. This technique picked out high-confidence rules that had high degrees of correlation between demographic features, modes of feedback, and thematic concerns. For instance, face-to-face feedback from women between 25–49 years old from the Somali group was always accompanied by positive feedback on food security. These links were not only statistically significant but also aligned with the operational realities observed by humanitarian practitioners in the field.

Patterns actively appeared in Experiment 2, which applied K-Means clustering to divide the dataset into groups. Two clusters were observed to predominate the analysis: one of male respondents in Tigray providing negative hotline service feedback, and the other with female respondents in Somali preferring face-to-face communication and giving neutral feedback. The patterns of repeated behavior and communications modes in community segments were observed in these clusters, offering improved insights into differences in feedback based on region, gender, and channel.

Trends were recorded across all three experiments, particularly through the temporal and demographic features of the data set. Transformation of date fields into month and weekday variables enabled the identification of weekly patterns and seasonality in feedback frequency and type. In addition, Experiment 3's decision tree model indicated how sentiment in feedback varied by region, age group, and sector over time, indicating shifting community concerns and styles of interaction.

Taken together, these findings demonstrate how data mining techniques can effectively reveal latent knowledge—trends, patterns, and relationships—that could otherwise remain unobserved in huge, complex data collections. These findings provide the foundation for more informed, adaptive, and ethical humanitarian interventions.

### **Objective 2: Using Insights to Improve Decision-Making and Service Delivery**

The second objective focused on how these insights could be used to improve humanitarian response. The insights from these experiments are not just theoretical, they can directly inform how humanitarian programs are designed, delivered, and improved. The association rules from **Experiment 1** offer clear guidance for program design. For instance, knowing that women in Somali are more likely to give positive feedback on food security through face-to-face channels suggests that these channels should be prioritized in that region.

The clusters identified in **Experiment 2** can help organizations tailor their communication and engagement strategies. If one group prefers hotlines and another prefers in-person discussions, then services can be adapted accordingly to ensure inclusivity and effectiveness.

**Experiment 3**, which used a decision tree model to classify feedback sentiment, showed that it's possible to automate the categorization of feedback with reasonable accuracy which could help humanitarian organizations respond more quickly to urgent concerns, especially in high-volume situations.

### **Objective 3: Evaluating the Effectiveness and Limitations of Data Mining Techniques**

The study evaluated the effectiveness of these data mining techniques within a real-world humanitarian context and the overall findings were encouraging with each method contributing distinct and valuable insights to the analysis of community feedback. Association rule mining proved to be particularly effective in uncovering clear and interpretable patterns that can directly inform strategic planning. Clustering revealed natural groupings within the data that were not predefined, offering a deeper understanding of the underlying structure and diversity of community responses. Meanwhile, classification emerged as a practical tool for predicting the sentiment of feedback, thereby supporting real-time decision-making and enabling more agile and responsive humanitarian interventions.

That said, there were also some limitations. The decision tree model, for example, had difficulty accurately classifying neutral feedback—likely because such responses are often vague or context-dependent. The clustering model was limited to just two groups, which may not fully capture the diversity of community experiences. And because the Somali region was heavily represented in the dataset, there’s a risk that the findings may not generalize to other regions without further validation.

Despite these challenges, the techniques proved to be effective overall. They handled a large dataset of nearly 93,000 entries, produced actionable insights, and were validated by domain experts who confirmed that the results aligned with what they were seeing on the ground.

## **5.6 Evaluation by Domain Experts**

To make sure the models developed in this study were not only technically sound but also practically useful, feedback was gathered from professionals working in the humanitarian field—people who regularly engage with community feedback and accountability systems. Their insights helped assess how well the models aligned with real-world needs and whether the results could actually support better decision-making on the ground.

Their evaluation focused on two types of models: **descriptive models**, which help explain what’s happening in the data, and **predictive models**, which aim to forecast or classify future feedback.

### **Descriptive Models: Making Sense of the Data**

The descriptive models used in this study—namely the **K-Means clustering** and **Apriori association rule mining**—were designed to uncover patterns and groupings in the feedback data.

Experts found these models especially helpful for understanding the bigger picture. For example, the association rules clearly showed that feedback from the Somali region, particularly from women in host communities, was strongly linked to food security concerns matching what many field workers had already observed, which gave them confidence in the model’s accuracy.

The clustering results were also well received as experts appreciated how the model grouped feedback into two distinct clusters, each with its own characteristics. One cluster, for instance, was dominated by hotline-based negative feedback from Tigray, while the other reflected more

neutral, face-to-face feedback from Somali. These insights may help organizations think about how different groups prefer to communicate and what issues matter most to them.

That said, some experts felt that limiting the clustering to just two groups might oversimplify the diversity of community voices. They suggested that future models could explore more nuanced groupings to capture a wider range of experiences.

### **Predictive Models: Supporting Real-Time Decisions**

The decision tree classification model was evaluated as a predictive model to identify the sentiment of feedback, positive, negative, or neutral. Experts saw a lot of potential in this model, especially for organizations that receive large volumes of feedback and need to respond quickly. The ability to automatically classify feedback could help teams prioritize urgent issues and streamline their response processes.

They also appreciated how easy the decision tree was to understand. Unlike more complex machine learning models, the tree structure made it clear how decisions were being made, which helped build trust in the results.

However, there were some concerns about the model's performance on neutral feedback, which was harder to classify accurately. Experts noted that neutral answers are typically imprecise or contextual, and that dividing this category into more specific subtypes could perform better in the future.

The domain experts were positive about the models developed in this study as they found the results relevant, easy to interpret, and aligned with what they were seeing in the field. Most importantly, they saw clear opportunities to use these insights to improve how feedback is collected, analyzed, and acted upon.

Their feedback also highlighted the importance of keeping models simple and transparent, especially when they're being used by non-technical staff. They encouraged continued collaboration between data analysts and field teams to ensure that future models remain grounded in the realities of humanitarian work.

# CHAPTER SIX: SUMMARY, CONCLUSION AND RECOMMENDATIONS

## 6.1 Summary

This research aimed to consider how data mining techniques can be utilized to examine feedback data within humanitarian settings, focusing on the work of the Inter-Agency Accountability Working Group (IAAWG) in Ethiopia. The motivation for this research came from a realization that while large quantities of feedback are collected systematically from affected populations, much of this is not being utilized because there is limited analytical capacity and tools.

In order to bridge this gap, the study utilized three of the core data mining techniques, association rule generation, clustering, and classification, with the aim of uncovering hidden insights, grouping similar feedback records, and predicting feedback sentiment based on demographic and contextual information.

The research utilized a systematic strategy, beginning with a thorough understanding of the humanitarian setting and the nature of the feedback data. The database of over 92,000 feedback records collected from 2021 to 2025 was cleaned and prepared using OpenRefine and mined using the WEKA software.

The following are some of the primary findings of the analysis :

- Association rule mining corroborated evidence of meaningful and interpretable relationships, with host communities and the Somali region as most frequently the providers of feedback—chiefly on food security.
- Clustering allowed for the identification of the natural groupings within the data, revealing how different regions and demographic segments engage with feedback mechanisms in different ways.
- Classification, using a decision tree model, was found to have potential for the automation of comments' categorization into positive, negative, or neutral sentiments with reasonable accuracy.

Each of the models was tested using conventional performance measures and tested again by validation by domain experts. Their comments made sure that the results were not just technically right but were also relevant to practice and in line with observations at the field level.

## **6.2 Conclusion**

To address the challenge, the research employed three key data mining techniques—association rule mining, clustering, and classification. These techniques not only stood the test of being feasible but were surprisingly successful in extracting enlightening information from complex feedback data sets. Each of the techniques contributed something unique: association rule mining uncovered significant and interpretable associations between variables; clustering pointed to natural groupings in the data; and classification demonstrated potential in the potential for an automatic sentiment classification in feedback.

One of the most compelling findings was the consistent prioritization of the Somali region and host communities as being at the center as sources of input, particularly for food security. This finding, which is based on high-confidence association rules, demonstrates the importance of region-specific and demographically sensitive programming. Similarly, the clustering analysis highlighted how different groups—such as women in Somali and men in Tigray—engage with feedback mechanisms differently, showing in greater detail community behavior and tendencies.

The classification model, while not ideal, was promising in supporting real-time decision making. Through predicting the sentiment of feedback entries, it offers a helpful tool for prioritizing responses—especially in high-volume or crisis scenarios where swift action is needed.

Apart from the technical innovations, the study also underlined the significance of interactions between field practitioners and data analysts. The involvement of domain experts in evaluating the models ensured that the results were not only statistically significant but also operationally appropriate. Their observations supported that the models could be suitably integrated into existing accountability mechanisms.

More broadly, this research draws upon the growing awareness that data-driven interventions can be a strong driver of humanitarian response. With its capacity to convert raw feedback into

actionable knowledge, organizations can be more transparent, responsive, and accountable to the populations they are seeking to serve. This is particularly important where humanitarian needs are vast, as in Ethiopia, and available resources are scarce.

Lastly, this study demonstrates that data mining is not merely a technical exercise—it is empowerment. Properly wielded ethically and responsibly, it has the capacity to re-sound the voices of affected groups, facilitate better decisions, and strengthen the relationship between humanitarian actors and the people they aim to assist.

### **6.3 Recommendations**

With the observations of this study, a variety of recommendations is proposed in order to assist humanitarian practitioners as well as researchers in using community feedback more effectively using data mining approaches. These suggestions are directed at ensuring accountability, promoting responsiveness, and supporting stronger evidence-based decision-making in humanitarian interventions.

#### **1. Make Data Mining a Part of Standard Feedback Work**

Humanitarian agencies are encouraged to embed data mining techniques—such as clustering, classification, and association rule mining—into their routine feedback mechanisms. These methods can help identify meaningful trends, classify feedback on chosen population variables, and even partially automate feedback processing. This would allow agencies to better understand communities and respond to issues with greater speed and relevance.

#### **2. Make Feedback Analysis More Timely**

To be responsive, one has to go beyond intermittent reviews and towards real-time or near-real-time feedback analysis. Having systems with feedback processing and analysis capabilities to operate continuously would allow organizations to pick up on pressing issues early on, react faster, and change interventions quickly—generally improving the overall quality of humanitarian aid.

#### **3. Increase and Diversify Feedback Collection**

While this study focused on formal feedback, care should be taken to assure that input is gathered from diverse channels and community groups. Humanitarian players have to continue investing in both face-to-face and distant methods, making sure excluded voices—e.g., women, youth, and persons with disabilities—are heard and included in the data. Inclusion guarantees stronger relevance and equity of response efforts.

#### 4. Strengthen interaction among analysts and field teams

While data mining can be very insightful, value is a function of the extent to which they are interpreted and acted upon. Frontline staff and data experts thus need to collaborate as field teams offer rich contextual data that can be used to anchor analytical findings in real-world wisdom to make interventions more suitable and effective.

#### 5. Explore More Advanced Machine Learning Methods

This study used classification via decision trees, which worked fine. Future research could attempt something more sophisticated such as Random Forests, Support Vector Machines, or Neural Networks that could potentially offer more accuracy, especially when dealing with larger datasets or more complex types of feedback, such as open-ended comments or multimedia.

#### 6. Handle Feedback Data Ethically and Responsibly

Community responses carry sensitive or personal data. It is imperative that organizations handle such information to the highest ethical standards. This includes anonymizing personal identifiers, securely storing systems, limiting access, and being transparent to communities about using their feedback. Ethical practices ensure data trust and safeguard rights of affected individuals.

#### 7. Test and Adapt the Models in Other Settings

This exercise used feedback data from Ethiopia by pilot-testing these models in other nations and crises can make it feasible to envision how widely they could be applied.

## REFERENCES

- Aggarwal, C. C.**, 2015. *Data mining: The textbook*. 1st ed. New York: Springer..
- Berry, M. J. & Linoff, G. S.**, 2004. *Data mining techniques: for marketing, sales, and customer relationship management*. 2nd ed. Indiana: Wiley Publishing.
- Bonino, F., Jean, I. & Clarke, P. K.**, 2014. *Humanitarian feedback mechanisms: Research, evidence and guidance*. London: ALNAP/ODI (ALNAP is a well-known organization focused on humanitarian learning.).
- Breiman, L.**, 2001. Random Forests. *Machine Learning*, 10, 45(1), pp. 5-32.
- Caldwell, S. H., Kristensen, K. L. & Milano, E.**, 2022. *Assessing The Technical Feasibility of Conflict Prediction for Anticipatory Action*, s.l.: nited Nations Office for the Coordination of Humanitarian Affairs (OCHA) Centre for Humanitarian Data.
- Chapman, P.**, 2000. *CRISP-DM 1.0 Step-by-step data mining guide*. 1st ed. U.S.A : SPSS Inc.
- Chen, M., Mao, S. & Liu, Y.**, 2014. Big Data: A Survey. *Mobile Networks and Applications*, 19(2), pp. 171-209.
- CHS Alliance**, 2020. *Core Humanitarian Standard on Quality and Accountability*. [Online] Available at: <https://www.corehumanitarianstandard.org/> [Accessed 7 Decembe 2024].
- Cios, K., Pedrycz, W., W. Swiniarski, R. & A.Kurgan, L.**, 2007. *Data Mining: A Knowledge Discovery Approach*. 1st ed. New York: Springer Sience+Business Media LLC.
- Delen, D.**, 2021. *Predictive Analytics: Data Mining, Machine Learning and Data Science for Practitioners*. 2nd ed. New Jersey: Pearson.
- Denekeew, A.**, 2003. The Application of Data Mining to Support Customer Relationship Management at Ethiopian Airlines. *Information Sciences*.
- Fayyad, U., Piatetsky-Shapiro, G. & Smyth, P.**, 1996. From Data Mining to Knowledge Discovery in Databases. *AI Magazine*, 17(3).

- Fikre**, M., 2005. Predictive Data Mining Technique in Insurance: (The Case of Ethiopian Insurance Corporation). *Information Sciences*.
- Han**, J., Micheline, K. & Jian, P., 2012. *Data Mining Cocepts and Techniques*. 3rd ed. Waltham: Morgan Kaufmann Publishers.
- Hintsay**, T., 2002. Predictive Modeling Using Data Mining Techniques in Support of Insurance Risk Assessment. *Information Sciences*.
- IOM**, n.d. *ACCOUNTABILITY TO AFFECTED POPULATIONS*. [Online]  
Available at:  
[https://www.iom.int/sites/g/files/tmzbd1486/files/our\\_work/DOE/humanitarian\\_emergencies/AA\\_P/two-pagebriefonaap.pdf](https://www.iom.int/sites/g/files/tmzbd1486/files/our_work/DOE/humanitarian_emergencies/AA_P/two-pagebriefonaap.pdf)
- Kohonen**, T., 2001. *Self-Organizing Maps*. 3rd ed. New York: Springer.
- Kondraganti**, A., Narayanamurthy, G. & Sharifi, H., 2022. A systematic literature review on the use of big data analytics. *Annals of Operations Research*, 335(1), p. 1015–1052.
- Olson**, D. L. & Delen, D., 2008. *Advanced data mining techniques*. 1st ed. Heidelberg: Springer.
- Powers**, D. M. W., 2011. Evaluation: From precision, recall and F-measure to ROC, informedness, markedness & correlation.. *Journal of Machine Learning Technologies*, 2(1), pp. 37-63.
- Reganie**, B., 2013. Application of Data Mining for Customer Segmentation: The Case of Buusaa Gonofa Microfinance Institution.. *Addis Ababa University*.
- Rousseeuw**, P. J., 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20(2), pp. 53-65.
- Tsipsis**, K. & Chorianopoulos, A., 2009. *Data Mining Techniques in CRM: Inside Customer Segmentation*. 1st ed. West Sussex: Wiley Publications.
- Turban**, E., Pollard, C. & Wood, G., 2018. *Information technology for management: on-demand strategies for performance, growth and sustainability*. 11th ed. NJ: Wiley.
- UNOCHA**, 2024. *Ethiopia: HRP Implementation in 2024, Cumulative response in 2024 (up to 30 November), Release date: 23 December 2024*.. [Online]

Available at: [https://reliefweb.int/report/ethiopia/ethiopia-hrp-implementation-2024-cumulative-response-2024-30-november-release-date-23-december-2024?\\_gl=1\\*cchbdv\\*\\_ga\\*MjEyODE2Njk5Ny4xNzMzMjYMTk1MDM1\\*\\_ga\\_E60ZNX2F68\\*MTczOTg2MDA5NC4yMC4wLjE3Mzk4NjAwOTQuNjAuMC4w](https://reliefweb.int/report/ethiopia/ethiopia-hrp-implementation-2024-cumulative-response-2024-30-november-release-date-23-december-2024?_gl=1*cchbdv*_ga*MjEyODE2Njk5Ny4xNzMzMjYMTk1MDM1*_ga_E60ZNX2F68*MTczOTg2MDA5NC4yMC4wLjE3Mzk4NjAwOTQuNjAuMC4w)

[Accessed 18 February 2025].

**UNOCHA**, 2024. *UNOCHA*. [Online]

Available at: <https://www.unocha.org/publications/report/world/global-humanitarian-overview-2025-enarfres>

[Accessed 1 December 2024].

## APPENDIX 1 LIST OF ATTRIBUTES SELECTED

Month_	Day	Region	Gender	Age
January	Monday	Addis Ababa	Female	18-24yrs
February	Tuesday	Afar	Male	25-49yrs
March	Wednesday	Amhara	Prefer not to say	Above 50yrs
April	Thursday	Benishangul Gumz		Preferred not to say
May	Friday	Central Ethiopia		Under 17yrs
June	Saturday	Dire Dawa		
July	Sunday	Gambela		
August		Harari		
September		Oromia		
October		Sidama		
November		Somali		
December		South Ethiopia		
		South West Ethiopia		
		Tigray		

Commun	Feedback	Feedback_Type	Sector
Host	Face to Face	Negative	CCCM
IDP	Helpdesk	Nutral	Child Protection
Other	Hotline	Positive	Education
Refugee	Phone		Food Security
Returnee	SMS		Gender Based Violence
	Social Media or E-mail		Health
	Suggestion Box		Livelihoods
			Multipurpose Cash
			Non Food Items
			Nutrition
			Operations
			Protection
			Shelter
			WASH

## APPENDIX 2: FINAL DATA USED

No.	1: Day Nominal	2: Region Nominal	3: Gender Nominal	4: Age Nominal	5: Community_type Nominal	6: Feedback_Channel Nominal	7: Feedback_Type Nominal	8: Sector Nominal
1	Tuesday	Amhara	Male	Under 17yrs	Host	Face to Face	Positive	Child Protection
2	Friday	Amhara	Female	Under 17yrs	Host	Face to Face	Positive	Child Protection
3	Thursday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Positive	Child Protection
4	Monday	Amhara	Female	Under 17yrs	Host	Face to Face	Positive	Child Protection
5	Sunday	Amhara	Male	Under 17yrs	Host	Face to Face	Positive	Child Protection
6	Monday	Afar	Female	Under 17yrs	Host	Face to Face	Positive	Education
7	Thursday	Amhara	Female	Under 17yrs	IDP	Face to Face	Positive	Education
8	Thursday	Afar	Female	Under 17yrs	Host	Face to Face	Positive	Education
9	Saturday	Afar	Male	Under 17yrs	Host	Face to Face	Positive	Multipurpose Cash
10	Saturday	Afar	Male	Under 17yrs	Host	Face to Face	Positive	Multipurpose Cash
11	Saturday	Tigray	Female	Under 17yrs	IDP	Face to Face	Negative	Child Protection
12	Thursday	Gambela	Prefer not t...	Under 17yrs	Refugee	Face to Face	Negative	Child Protection
13	Sunday	Tigray	Female	Under 17yrs	IDP	Face to Face	Negative	Child Protection
14	Tuesday	Tigray	Male	Under 17yrs	Host	Face to Face	Negative	Protection
15	Monday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
16	Tuesday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
17	Thursday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
18	Tuesday	Tigray	Male	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
19	Tuesday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Food Security
20	Wednesday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
21	Wednesday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
22	Saturday	Amhara	Male	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
23	Saturday	Amhara	Male	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
24	Saturday	Amhara	Male	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
25	Saturday	Amhara	Male	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
26	Friday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
27	Thursday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
28	Saturday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
29	Wednesday	Tigray	Female	Under 17yrs	Host	Face to Face	Nutral	Child Protection
30	Tuesday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
31	Sunday	Gambela	Female	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
32	Wednesday	Gambela	Female	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
33	Tuesday	Gambela	Female	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
34	Saturday	Tigray	Female	Under 17yrs	IDP	Face to Face	Nutral	Child Protection
35	Saturday	Amhara	Male	Under 17yrs	Host	Face to Face	Nutral	Child Protection
36	Saturday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
37	Friday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
38	Friday	Gambela	Female	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
39	Wednesday	Gambela	Female	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
40	Wednesday	Gambela	Female	Under 17yrs	Refugee	Face to Face	Nutral	Child Protection
41	Friday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Food Security
42	Thursday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Non Food Items
43	Thursday	Gambela	Male	Under 17yrs	Refugee	Face to Face	Nutral	Non Food Items
44	Tuesday	Afar	Male	Under 17yrs	Host	Face to Face	Positive	Child Protection
45	Tuesday	Afar	Male	Under 17yrs	Host	Face to Face	Positive	Child Protection

Add instance Undo OK Cancel

## APPENDIX 3: DETAILED SUMMARY OF APRIORI ASSOCIATION RULE MINING RESULTS

Rule #	Antecedent (IF)	Consequent (THEN)	Confidence (%)	Interpretation
1	Region = Somali, Feedback Channel = Face to Face, Sector = Food Security	Community Type = Host	97	Individuals in the Somali region who provide face-to-face feedback on food security are highly likely to be from the host community. This rule reflects a strong and reliable association.
2	Region = Somali, Gender = Female, Feedback Type = Positive	Community Type = Host	96	Positive feedback from women in the Somali region is predominantly associated with the host community. This may suggest program effectiveness or satisfaction among this demographic.
3	Region = Somali, Feedback Type = Positive, Sector = Food Security	Community Type = Host	95	Positive feedback on food security in the Somali region is strongly linked to host community members, reinforcing the importance of this sector in that context.
4	Region = Somali, Feedback Channel = Face to Face, Feedback Type = Positive	Community Type = Host	95	This rule highlights the effectiveness of face-to-face communication in eliciting positive feedback from host communities in the Somali region.
5	Region = Somali, Gender = Female, Sector = Food Security	Community Type = Host	94	Women in the Somali region who engage with food security programs are most likely from the host community, suggesting targeted engagement.

6	Region = Somali, Age = 25-49, Feedback Type = Positive	Community Type = Host	94	Positive feedback from individuals aged 25 to 49 in the Somali region is strongly associated with the host community, indicating a key demographic for engagement.
7	Region = Somali, Sector = Food Security, Gender = Female	Community Type = Host	93	This rule reinforces the intersection of gender and sectoral focus in identifying host community members.
8	Region = Somali, Feedback Channel = Face to Face, Gender = Female	Community Type = Host	93	Face-to-face feedback from women in the Somali region is a strong indicator of host community affiliation.
9	Region = Somali, Feedback Type = Positive, Gender = Female	Community Type = Host	92	This rule further confirms the trend of positive feedback from women in the Somali region being linked to host communities.
10	Community Type = Host, Feedback Type = Positive, Sector = Food Security	Region = Somali	95	This reverse rule indicates that when host community members provide positive feedback on food security, it is most likely to originate from the Somali region. This insight is particularly useful for geographically targeted programming.

# APPENDIX 4: CONFIGURATION SETTINGS USED IN WEKA FOR EACH MODEL

