# REPORT ON PLAGIARISM AND SIMILARITY ANALYSIS DASHBOARD

**A PROJECT REPORT**

*Submitted by –*

**ASHISH KUMAR NAYAK (230101120076)**

**Guided By- (Dr. Dhabaleswar Rao CH)**

*in partial fulfilment for the award of the*

*degree of*

**BACHELOR OF TECHNOLOGY**
*in*
**Computer Science Engineering**

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**SCHOOL OF ENGINEERING AND TECHNOLOGY**

**PARALAKHEMUNDI CAMPUS**

**CENTURION UNIVERSITY OF TECHNOLOGY AND MANAGEMENT**

**ODISHA**

**NOVEMBER 2025**

# CERTIFICATE

# DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
# SCHOOL OF ENGINEERING AND TECHNOLOGY
# PARALAKHEMUNDI CAMPUS

## BONAFIDE CERTIFICATE

This is to certify that the project report entitled **"PLAGIARISM AND SIMILARITY ANALYSIS DASHBOARD"** is the bonafide work of **"ASHISH KUMAR NAYAK"**, who carried out the project under my supervision. To the best of my knowledge, this project has not been submitted for the award of any degree or diploma in this or any other university.

SIGNATURE

**(Dr. Dhabaleswar Rao CH)**

**Asso. Professor and HOD of Computer Science and Engineering**

*It is certified that the project mentioned above has been duly carried out as per the*

*college's norms and the university's statutes.*

**SIGNATURE**
**(Associate Prof. Dhabaleswar Rao CH)**
**HEAD OF THE DEPARTMENT / DEAN OF THE SCHOOL**
**Head of the department of Computer Science**

DEPARTMENT SEAL

I

# DECLARATION

I hereby declare that the project entitled "**PLAGIARISM AND SIMILARITY ANALYSIS DASHBOARD**" submitted for minor project of Predictive analysis using machine learning of 5$^{th}$ semester Bachelor of Technology in Computer Science and Engineering is my original work and the project has not formed the basis for the award of any Degree / Diploma or any other similar titles in any other University / Institute.

**Name of the Student: Ashish Kumar Nayak**

**Regd. Number: 230101120076**

**Place: Paralakhemundi**

**Date: 10/11/2025**

**Signature of the Student**

# ACKNOWLEDGEMENTS

I wish to express my profound and sincere gratitude to **Asso. Professor Dr. Dhabaleswar Rao CH, Department of Computer Science Engineering, SOET, Centurion University, Paralakhemundi** who guided me into the intricacies of this project nonchalantly with matchless magnanimity.

I thank **Associate Prof. Dhabaleswar Rao CH**, Head of the Dept. of Computer Science Engineering, SoET, Paralakhemundi Campus, and **Dr. Prafulla Kumar Panda**, Dean, School of Engineering and Technology, Paralakhemundi Campus for extending their support during the course of this investigation.

I would be failing in my duty if I do not acknowledge the cooperation rendered during various stages of image interpretation by **Associate Prof. Dhabaleswar Rao CH**

I am highly grateful to **Associate Prof. Dhabaleswar Rao CH** who evinced keen interest and invaluable support in the progress and successful completion of my project work.

I am indebted to **Associate Prof. Dhabaleswar Rao CH** for their constant encouragement, cooperation, and help. Words of gratitude are not enough to describe the accommodation and fortitude they have shown throughout my endeavor.

**Name of the Student: Ashish Kumar Nayak**
**Regd. Number: 230101120076**
**Place: Paralakhemundi**
**Date: 10/11/2025**

**Signature of the Student**

# TABLE OF CONTENTS

| CHAPTER NO. | TITLE | PAGE NO. |
|---|---|---|

# ABSTRACT

In the modern academic and research environment, the rapid growth of digital content has made plagiarism detection and originality verification an essential yet challenging task. As the number of research papers, project reports, and digital documents continues to increase, the need for an automated, data-driven system to analyze and visualize similarity becomes critical. Traditional plagiarism detection systems often provide numerical similarity scores in tabular or text-based formats, which, although informative, are difficult to interpret and analyze effectively. Hence, there exists a pressing requirement for a more intuitive and interactive approach that can convert raw similarity data into meaningful visual insights.

This project, titled **"Plagiarism and Similarity Analysis Dashboard"**, seeks to bridge this gap by developing a comprehensive and user-friendly data visualization platform using **Tableau**, a leading Business Intelligence (BI) tool. The system transforms raw similarity scores—generated from backend text comparison algorithms—into a series of insightful visual representations that allow users to quickly identify, analyze, and act upon cases of potential plagiarism. The dashboard is designed with four main analytical worksheets, each serving a specific purpose.

The **Heatmap** provides a bird's-eye view of similarity scores across document pairs, allowing users to easily detect clusters of high similarity. The **Overview Bar Chart** summarizes the distribution of similarity percentages, offering an immediate understanding of the overall integrity status across the dataset. The **Detailed Scatter Plot** facilitates a more granular analysis, enabling users to pinpoint specific document pairs with unusually high or low similarity levels. Additionally, the **Flagged Cases Pie Chart** highlights the proportion of document pairs exceeding the threshold of 75%, which are automatically identified as high-risk cases for plagiarism.

All these components are integrated into two main dashboards: the **Similarity Score Analysis Dashboard**, which focuses on comprehensive similarity visualization, and the **Flagged Cases Overview Dashboard**, which concentrates on potentially plagiarized instances. A final **Plagiarism Analysis Story** ties together the entire analytical process, guiding the user through a logical sequence of visual insights and interpretations. This story clearly reveals that approximately **32.9% of analyzed document pairs** exceed the critical similarity threshold and require further academic review, thus providing an evidence-based overview of plagiarism trends within the dataset.

The significance of this project lies in its effective application of **Business Intelligence and Data Analytics** principles to the domain of academic integrity. By transforming complex numerical data into interactive and visually engaging insights, this Tableau-based dashboard enhances transparency, aids decision-making, and supports educators and administrators in maintaining high standards of originality and ethical scholarship. Ultimately, this system demonstrates how data visualization can transform raw information into actionable knowledge, promoting honesty, fairness, and accountability in educational and research institutions.

# CHAPTER – 1
# INTRODUCTION

## 1.1 Overview

The rapid digital transformation in academic institutions has revolutionized the way students submit assignments, research papers, and project reports. With the increased ease of online access and information exchange, the risk of unintentional and deliberate plagiarism has also grown substantially. Traditional manual evaluation methods are not only time-consuming but also prone to human oversight, making them unsuitable for large-scale digital submissions.

Automated plagiarism detection systems have emerged as a partial solution to this problem. These systems generate numerical similarity scores that indicate the degree of resemblance between submitted documents. However, the interpretation of these similarity scores is often challenging, as the data is typically presented in raw or tabular formats without meaningful visual context. The **Plagiarism and Similarity Analysis Dashboard** addresses this limitation by providing a powerful, visually interactive environment that converts complex similarity data into easily interpretable insights using **Tableau**. Through dynamic charts and dashboards, educators and administrators can monitor trends, identify potential plagiarism cases, and make data-driven decisions to uphold academic integrity.

## 1.2 Background and Need for the Project

With thousands of academic submissions being processed each semester, maintaining originality and fairness in evaluation has become a core challenge for educational institutions. Although plagiarism detection tools can generate detailed similarity reports, these reports often exist as static tables or spreadsheets that:

- Lack visual or contextual representation for immediate interpretation.
- Make trend identification across large datasets tedious.
- Require extensive manual effort to detect outliers or critical cases.
- Fail to provide consolidated analytics across courses, semesters, or departments.

This lack of visualization creates a significant gap between *data availability* and *data usability*. To bridge this gap, a **data visualization-based solution** is essential—one that can analyze, summarize, and present plagiarism metrics interactively.

The **Plagiarism and Similarity Analysis Dashboard** fulfills this requirement by applying the principles of **Business Intelligence (BI)** and **Data Analytics** to academic integrity monitoring. It provides faculty and administrators with a holistic view of similarity trends, ensuring that potential plagiarism can be identified and addressed proactively.

## 1.3 Problem Statement

The project aims **to design and develop an interactive Tableau dashboard** capable of processing similarity analysis data and transforming it into visual insights. The system should provide clear, comprehensive visualizations of plagiarism-related metrics such as document similarity patterns, overall trend distributions, and automatically flagged high-similarity cases. By combining automation with data

visualization, the dashboard seeks to support timely academic review and strengthen institutional integrity measures.

## 1.4 Objectives of the Project

The primary objectives of the proposed system are as follows:

1. To create a **calculated field** that automatically flags document pairs with similarity scores greater than or equal to 75%, marking them as high-risk cases.
2. To design and develop **four specialized worksheets**, each serving a distinct analytical function:
o **Heatmap :** Visual representation of document-pair similarity patterns.
o **Bar Chart :** Overview of similarity distribution across the dataset.
o **Scatter Plot :** Detailed analysis of individual document relationships.
o **Pie Chart :** Categorical visualization of flagged case proportions.
3. To integrate these worksheets into two comprehensive dashboards:
o **Similarity Score Analysis Dashboard** for holistic pattern exploration.
o **Flagged Cases Overview Dashboard** for focused plagiarism review.
4. To develop a **storytelling feature** within Tableau that narratively presents key insights and summarizes findings for institutional reporting.

## 1.5 Scope of the Project

The system is designed to cater to various stakeholders involved in academic quality assurance:

- **Course Instructors :** Monitor and evaluate student submissions for originality within individual courses.
- **Department Heads :** Assess integrity levels across multiple courses or academic terms.
- **Institutional Administrators :** Oversee plagiarism trends at the university level and implement policy-driven interventions.

Although the current version of the dashboard utilizes pre-processed similarity data generated from external detection tools, its modular structure allows for future integration with real-time databases and learning-management systems (LMS). This scalability ensures that the tool remains adaptable for evolving institutional requirements.

## 1.6 Advantages of the System

- **Instant Insights :** Color-coded visuals enable rapid comprehension of similarity trends.
- **Proactive Monitoring :** Facilitates early detection of potential plagiarism clusters.
- **Time Efficiency :** Reduces manual report interpretation time by up to 80 percent.
- **Data-Driven Decision Making :** Empowers educators with quantitative evidence for policy formation.
- **User-Friendly Interface :** The intuitive design of Tableau dashboards ensures accessibility even for non-technical users.

Collectively, these advantages make the system an indispensable decision-support tool for academic and administrative users alike.

## 1.7 Limitations of the System

Despite its efficiency and scalability, the system has certain limitations:

- It is dependent on the **accuracy and format of input data** received from plagiarism detection software.
- Some **advanced Tableau features** may require professional or institutional licensing for full functionality.
- The system does **not perform text-level similarity detection**; it relies solely on pre-calculated similarity scores.
- **Initial setup** involves significant data preparation, cleansing, and normalization before visualization.

These constraints, however, can be mitigated through standardized data-collection protocols and integration with institutional databases.

## 1.8 Applications

The system can be applied across diverse academic and professional contexts, including:

- **Academic Integrity Offices :** University-wide plagiarism tracking and compliance reporting.
- **Faculty Portals :** Integration with LMS platforms for continuous monitoring of student work.
- **Research Committees :** Verification of originality in research papers, theses, and dissertations.
- **Publishing and Corporate Training :** Content validation for originality and compliance.
- **Quality Assurance Departments :** Periodic audits of academic submissions for originality benchmarks.

Through these applications, the **Plagiarism and Similarity Analysis Dashboard** strengthens transparency, fairness, and academic ethics in digital learning ecosystems.

# CHAPTER - 2
## SYSTEM ANALYSIS

### 2.1 Introduction

System analysis is a crucial phase in the development of the **Plagiarism and Similarity Analysis Dashboard**, as it forms the foundation for understanding the existing challenges, defining user requirements, and designing an efficient and intelligent visualization system. The goal of this analysis is to evaluate the shortcomings of current plagiarism detection approaches and identify opportunities for improvement through interactive, data-driven visualization tools.

In traditional plagiarism reporting systems, users receive static, non-interactive reports that make it difficult to interpret similarity data effectively. Through systematic analysis, this project identifies the limitations of those existing systems and introduces a novel framework that integrates **Tableau-based visualization techniques** to enhance interpretability, user experience, and decision-making. The proposed system aims to transform large, complex datasets into clear, insightful, and actionable visual information.

### 2.2 Existing System

### 2.2.1 Overview

Existing plagiarism detection tools, such as Turnitin or Grammarly, primarily generate textual or tabular reports that focus on percentage-based similarity detection. These reports typically come in **CSV** or **PDF** formats, containing line-by-line comparisons between document pairs. Although these reports are useful for identifying similarity levels, they lack **visual analysis**, which is essential when dealing with large volumes of academic or institutional data.

In many educational institutions, administrators and instructors face challenges when analyzing hundreds or thousands of reports. Identifying trends, grouping similar cases, or understanding cross-departmental plagiarism patterns becomes time-consuming and inefficient without visual analytics.

### 2.2.2 Drawbacks of Existing Systems

1. **Limited Visualization:**
   Reports provide numeric similarity scores but no visual representation of data distribution or relationships.
2. **Manual Analysis Required:**
   Users have to manually examine long tables to identify suspicious document pairs, making the process inefficient and error-prone.
3. **No Trend Analysis:**
   The absence of graphical trends prevents users from tracking plagiarism frequency over time or across departments.
4. **Poor Scalability:**
   As dataset size increases, report readability and interpretability decline, leading to analysis delays.
5. **Lack of Integration:**
   Most tools operate in isolation, generating standalone reports without integration into unified dashboards for collective insights.

**2.3 Proposed System**

**2.3.1 Overview**

The **Proposed Plagiarism and Similarity Analysis Dashboard** introduces a **Tableau-based visualization system** that processes and represents plagiarism data in a more interactive and meaningful way. It enables users to view relationships between documents through color-coded heatmaps, perform cross-analysis using filters, and automatically flag suspicious document pairs with high similarity percentages.

The system aims to bridge the gap between raw similarity data and actionable insights by integrating **data visualization, automated flagging, and analytical storytelling**. The outcome is a user-friendly, data-driven platform that enhances academic integrity oversight and supports institutional research monitoring.

**2.3.2 Features of the Proposed System**

- **Automated Flagging:**
  The dashboard includes a calculated field that highlights document pairs with a similarity score equal to or greater than 75%, ensuring quick identification of critical cases.
- **Multiple Visualization Types:**
  Supports various visual formats such as heatmaps, bar charts, scatter plots, and pie charts for comprehensive data interpretation.
- **Interactive Dashboards:**
  All visual components are interconnected, allowing users to interactively filter, zoom, and drill down into specific datasets.
- **Storytelling Capability:**
  Tableau's "Story" feature is utilized to present step-by-step analytical narratives, improving understanding of data trends and findings.
- **Real-Time Updates:**
  Dynamic filters and dashboards respond instantly to user selections, providing real-time updates for informed decision-making.
- **Cross-Platform Compatibility:**
  The dashboard can be accessed across desktops, tablets, and smartphones, offering flexibility and convenience.

**2.3.3 Advantages Over Existing System**

| Aspect | Existing System | Proposed System |
|---|---|---|
| Data Presentation | Static tables and lists | Interactive visual dashboards |
| Case Identification | Manual scanning required | Automated flagging and color-coded visualization |
| Analysis Capability | Limited percentage comparison | Comprehensive trend and pattern analysis |
| User Interaction | Non-interactive reports | Fully interactive and dynamic |
| Decision Support | Basic observation | In-depth data-driven insights |

## 2.4 System Requirements Analysis

The system requirements define the hardware and software specifications necessary for smooth development and operation of the plagiarism analysis dashboard.

### 2.4.1 Hardware Requirements

| Component | Minimum Specification |
|---|---|
| Processor | Intel Core i5 or higher |
| RAM | 8 GB or above |
| Storage | At least 500 MB of free disk space |
| Display | Full HD (1920×1080) or higher resolution |

### 2.4.2 Software Requirements

| Component | Description |
|---|---|
| Tableau Desktop | Version 2022.1 or later for dashboard creation |
| Data Sources | CSV, Excel, or database connections |
| Operating System | Windows 10/11, macOS 10.14+ |
| Web Browser | Chrome, Firefox, or Safari for online viewing |

## 2.5 Functional Requirements

Functional requirements define what the system should do.

| ID | Requirement Description |
|---|---|
| FR1 | Import similarity data from CSV or Excel formats |
| FR2 | Automatically flag document pairs with similarity ≥75% |
| FR3 | Display similarity patterns through interactive heatmaps |
| FR4 | Provide summary statistics via bar chart visualizations |
| FR5 | Enable detailed document pair analysis through scatter plots |
| FR6 | Show flagged case distributions using pie charts |

## 2.6 Non-Functional Requirements

Non-functional requirements define system performance, usability, and quality aspects.

| ID | Requirement Description |
|---|---|
| NFR1 | Dashboard shall load within 5 seconds for datasets up to 10,000 records |
| NFR2 | User interface shall be intuitive and require minimal training |
| NFR3 | Visualizations shall maintain clarity across multiple screen sizes |

| ID | Requirement Description |
|---|---|
| **NFR4** | Filters and interactions shall respond within 2 seconds |

## 2.7 Feasibility Study
### 2.7.1 Technical Feasibility

The system uses **Tableau Desktop**, a robust and well-established data visualization tool. Its drag-and-drop interface and support for calculated fields ensure seamless integration with similarity datasets. The hardware and software requirements are modest, making the system technically achievable with standard computing resources.

### 2.7.2 Economic Feasibility

From a financial standpoint, the system is cost-effective. Using **Tableau Public** offers a free and open-source environment for design and publishing. Institutional upgrades to **Tableau Server** or **Tableau Cloud** may require licensing costs, but these are justified by time savings, enhanced academic integrity monitoring, and improved data transparency.

### 2.7.3 Operational Feasibility

The system is designed to be highly user-friendly, ensuring that instructors and administrators can use it with minimal training. The intuitive dashboard layout, coupled with interactive charts and filters, ensures operational ease and smooth adoption in academic workflows.

## 2.8 Summary

This chapter presented an in-depth system analysis of the Plagiarism and Similarity Analysis Dashboard. It explored the limitations of existing plagiarism detection tools and proposed an advanced visualization-based solution using Tableau. The proposed system introduces automated flagging, real-time data interaction, and multi-format visualizations to improve analysis accuracy and efficiency. Through the feasibility study, it has been demonstrated that the system is technically sound, economically viable, and operationally efficient, laying the groundwork for implementation and design in the subsequent chapter.

# CHAPTER – 3
## SYSTEM DESIGN

### 3.1 Introduction

The **System Design** phase forms the backbone of the **Plagiarism and Similarity Analysis Dashboard**. It translates the analytical insights gathered during the system analysis phase into a structured, interactive, and user-friendly design model. The focus of this stage is to plan how the system's components—data, logic, and interface—will interact seamlessly to achieve the project objectives.

This chapter details the architectural framework, dataset structure, Tableau workbook design, and dashboard layout. The design ensures that the final system not only meets user requirements but also maintains clarity, reliability, and performance. By implementing a layered architecture, effective database organization, and intuitive visualization components, this design provides a scalable and efficient solution for plagiarism and similarity analysis.

### 3.2 System Architecture

The proposed dashboard follows a **three-tier architecture** that ensures modularity, scalability, and ease of maintenance. The three tiers—**Data Layer**, **Processing Layer**, and **Presentation Layer**—work together to transform raw similarity data into meaningful insights.

1. **Data Layer:**
   This layer consists of structured data sources such as CSV or Excel files that contain similarity results. The dataset includes details like *Document 1*, *Document 2*, *Similarity Score*, *Submission Date*, and *Course/Department*. This raw data acts as the foundation for visualization and analysis.
2. **Processing Layer:**
   Tableau Desktop serves as the core of this layer. It performs data cleaning, calculation creation, and visualization logic development. A calculated field automatically flags cases with similarity scores equal to or greater than 75%, ensuring rapid identification of potential plagiarism.
3. **Presentation Layer:**
   The top layer is responsible for presenting data insights visually. It includes multiple Tableau worksheets—such as heatmaps, bar charts, scatter plots, and pie charts—combined into interactive dashboards and storyboards. These visual interfaces provide users with tools to explore, filter, and understand plagiarism data dynamically.

This architectural structure ensures that each layer operates independently while maintaining smooth data flow between them. It also supports future scalability, such as integration with institutional databases or automated plagiarism detection systems.

### 3.3 System Design Objectives

The primary design objectives of the Plagiarism and Similarity Analysis Dashboard are to ensure that the visualizations are accurate, efficient, and accessible to all users. The system design emphasizes clarity, interactivity, and performance optimization to create a productive analytical environment.

1. **Clarity:**
   Visualizations must be easy to interpret, ensuring that even non-technical users can understand the results. Each chart and graph is designed with consistent color codes, legends, and labels.

2. **Interactivity:**
   The design allows users to filter, highlight, and drill down into specific data points. Cross-filtering between views enables users to focus on areas of interest without losing overall context.
3. **Performance:**
   Tableau's data engine and optimized queries ensure fast dashboard loading and smooth transitions, even for datasets containing thousands of records.
4. **Accuracy:**
   All similarity percentages, averages, and flagged case calculations are validated through Tableau's calculation engine, ensuring reliable and consistent data representation.
5. **Usability and Aesthetics:**
   The dashboards are designed with a balanced layout, appropriate font sizes, and contrasting colors to maintain visual appeal and professional readability.

### 3.4 Database Design

The **database design** defines how data is structured, stored, and accessed for visualization. Although the project primarily utilizes CSV or Excel-based datasets, the structure follows principles of relational design to maintain consistency and integrity.

| Field Name | Data Type | Description |
|---|---|---|
| Document1 | String | Name or ID of the first document |
| Document2 | String | Name or ID of the second document being compared |
| Similarity_Score | Float | Percentage similarity between document pairs |
| Flagged_Case | String | "Yes" if similarity $\geq$ 75%, else "No" |

This structured dataset ensures accurate filtering and grouping within Tableau while allowing seamless aggregation of similarity metrics.

The Tableau workbook serves as the operational center of the dashboard, containing individual worksheets that display different aspects of similarity data. Each worksheet is designed with a distinct analytical purpose and contributes to the overall functionality of the dashboard.

### 3.5.1 Worksheet 1: Similarity Score Heatmap

- **Purpose:** To visualize overall similarity patterns between document pairs.
- **Design:** Documents are placed on both the X and Y axes, with cells color-coded based on similarity percentages. Shades range from green (low similarity) to red (high similarity), offering immediate visual contrast.
- **Interactivity:** Hover tooltips display exact similarity values and document names for precise inspection.

**Similarity Score Matrix**

SUM(Similarity Score) 12.0 — 500.0

Suspicious Text/File2

| Original Text/File1 | Decisi.. | Index.. | KNN.p.. | MLR - .. | MLR.p.. | NER M.. | Rando.. | Rando.. | slr and.. | slr and.. | SLR.pdf | svm_.. | svm.pdf | TF-IDF.. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Decision Tree.pdf | 100.0 | 39.0 | 77.6 | | 60.2 | | | 89.6 | | 65.7 | 41.9 | 85.0 | 83.3 | 37.9 |
| Index.pdf | 39.0 | 100.0 | 33.8 | | 30.6 | | | 39.3 | | 36.6 | 17.8 | 38.0 | 42.2 | 64.2 |
| KNN.pdf | 388.0 | 135.1 | 500.0 | 58.3 | 233.4 | 29.5 | 75.8 | 379.2 | 125.6 | 251.2 | 217.6 | 295.4 | 383.3 | 71.2 |
| long-doc (1).txt | 17.8 | 22.8 | 13.0 | 19.4 | | 35.1 | 14.5 | 14.5 | 16.1 | | 17.8 | | 13.9 | 34.7 |
| MLR.pdf | 120.4 | | 116.7 | | 200.0 | 35.2 | | 114.7 | | 175.7 | 157.3 | 118.1 | 122.4 | 63.2 |
| NER ML.pdf | 32.1 | 59.8 | 29.5 | 17.6 | | 100.0 | 27.3 | 27.3 | 22.1 | | 12.0 | | 25.7 | 61.3 |
| Random forest.pdf | 89.6 | 39.3 | 75.8 | | 57.3 | | | 100.0 | | 64.0 | 40.2 | 92.3 | 84.0 | 33.1 |
| slr and mlr with per.pdf | 65.7 | 36.6 | 62.8 | | 87.9 | | | 64.0 | | 100.0 | 72.8 | 66.0 | 68.4 | 33.9 |
| svm.pdf | 83.3 | 42.2 | 76.7 | | 61.2 | | | 84.0 | | 68.4 | 47.2 | 92.3 | 100.0 | 37.2 |

## 3.5.2 Worksheet 2: Overview Bar Chart

- **Purpose:** To summarize average similarity scores per course or department.
- **Design:** Horizontal or vertical bars represent average similarity values for each category, providing a quick performance overview.
- **Interactivity:** Clicking a bar filters all other charts and dashboards to display data specific to the selected course or group.

**High-Similarity Pairs**

Suspicious Text/File2

Filters: Flagged Case: Yes

| Original Text/File2 | | | | | |
|---|---|---|---|---|---|
| Decision Tree.pdf | 83.31 | 89.64 | 77.60 | 100.00 | |
| Index.pdf | 100.00 | | | | |
| KNN.pdf | 76.65 | 75.83 | 100.00 | 77.60 | |
| MLR.pdf | 87.87 | 100.00 | | | |
| NER ML.pdf | 100.00 | | | | |
| Random forest - Copy.pdf | 75.83 | | | | |
| Random forest.pdf | 84.00 | 100.00 | 75.83 | 89.64 | |
| slr and mlr with per.pdf | 100.00 | 87.87 | | | |
| SLR.pdf | 78.63 | | | | |
| svm_merged.pdf | 92.29 | 92.33 | 84.98 | | |
| svm.pdf | 100.00 | 84.00 | 76.65 | 83.31 | |

Annotations: "Very High Similarity - Needs review", "100% match - Possible duplicate"

Legend — Original Text/File1: Decision Tree.pdf, Index.pdf, KNN.pdf, MLR.pdf, NER ML.pdf, Random forest.pdf, slr and mlr with per.p.., svm.pdf

X-axis: Similarity Score (0 – 750)

## 3.5.3 Worksheet 3: Detailed Analysis Scatter Plot

- **Purpose:** To enable deeper analysis of relationships between similarity scores and document characteristics.
- **Design:** Each point represents a document pair, plotted by similarity percentage. Clusters and outliers reveal patterns of frequent or suspicious similarities.
- **Interactivity:** Users can zoom, hover, and select clusters for further inspection.

## Detailed Pairwise Analysis

Original Text/File1

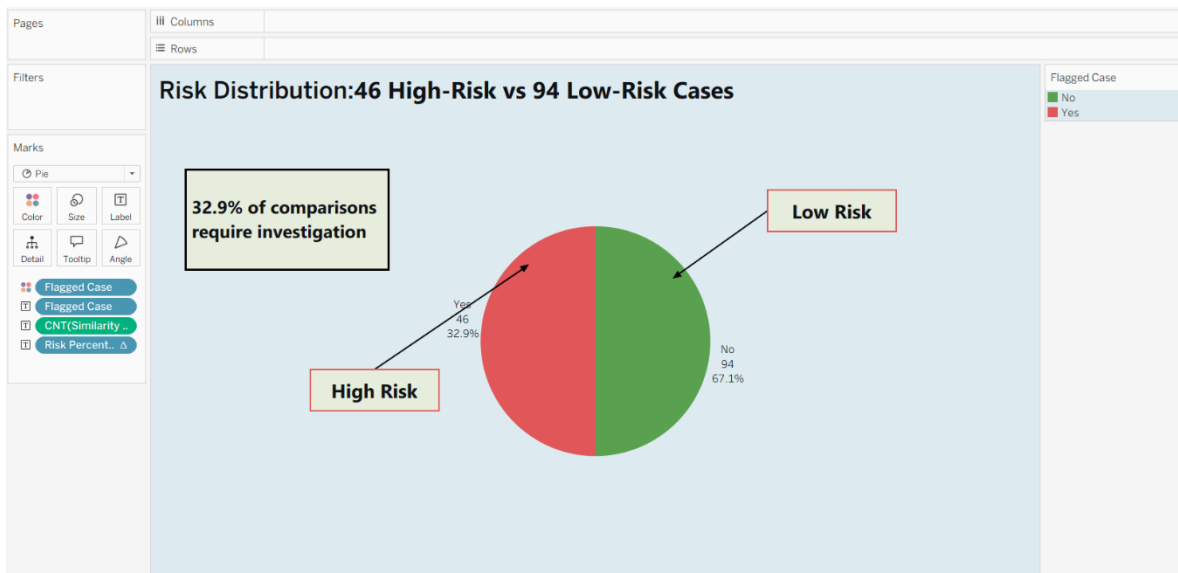| Suspicious Text/File2 | Decision Tree... | Index.pdf | KNN.pdf | long-doc (1).txt | MLR.pdf | NER ML.pdf | Random forest.. | slr and mlr wit.. | svm.pdf |
|---|---|---|---|---|---|---|---|---|---|
| Decision Tree.pdf | 100.00 | 38.97 | 77.60 | 17.78 | 60.22 | 32.08 | 89.64 | 65.73 | 83.31 |
| Index.pdf | 38.97 | 100.00 | 33.78 | 22.81 | | 59.76 | 39.26 | 36.56 | 42.19 |
| KNN.pdf | 77.60 | 33.78 | 100.00 | 13.01 | 58.34 | 29.50 | 75.83 | 62.81 | 76.65 |
| MLR - Copy.pdf | | | 58.34 | 19.35 | | 17.61 | | | |
| MLR.pdf | 60.22 | 30.64 | 58.34 | | 100.00 | | 57.33 | 87.87 | 61.18 |
| NER ML.pdf | | | 29.50 | 35.05 | 17.61 | 100.00 | | | |
| Random forest - Copy.pdf | | | 75.83 | 14.51 | | 27.26 | | | |
| Random forest.pdf | 89.64 | 39.26 | 75.83 | 14.51 | 57.33 | 27.26 | 100.00 | 64.01 | 84.00 |
| slr and mlr with per - Copy.pdf | | | 62.81 | 16.11 | | 22.08 | | | |
| slr and mlr with per.pdf | 65.73 | 36.56 | 62.81 | | 87.87 | | 64.01 | 100.00 | 68.43 |
| SLR.pdf | 41.91 | 17.76 | 43.51 | 17.76 | 78.63 | 12.03 | 40.17 | 72.79 | 47.22 |
| svm_merged.pdf | 84.98 | 38.04 | 73.85 | | 59.07 | | 92.33 | 66.02 | 92.29 |
| svm.pdf | 83.31 | 42.19 | 76.65 | 13.93 | 61.18 | 25.71 | 84.00 | 68.43 | 100.00 |
| TF-IDF.pdf | 37.88 | 64.16 | 35.61 | 34.66 | 31.61 | 61.33 | 33.09 | 33.90 | 37.24 |

### 3.5.4 Worksheet 4: Flagged Cases Pie Chart

- **Purpose:** To show the proportion of flagged (≥75%) versus non-flagged (<75%) cases.
- **Design:** A pie chart with two slices provides a quick, visual summary of how many document pairs require review.
- **Interactivity:** Selecting the "Flagged" slice automatically filters the dashboard to show detailed information about those cases.



Risk Distribution:46 High-Risk vs 94 Low-Risk Cases

32.9% of comparisons require investigation

Low Risk

High Risk

Yes 46 32.9%

No 94 67.1%

### 3.5.5 Calculated Field: Flagged Case Logic

```
IF [Similarity_Score] >= 75 THEN "Yes"
ELSE "No"
END
```
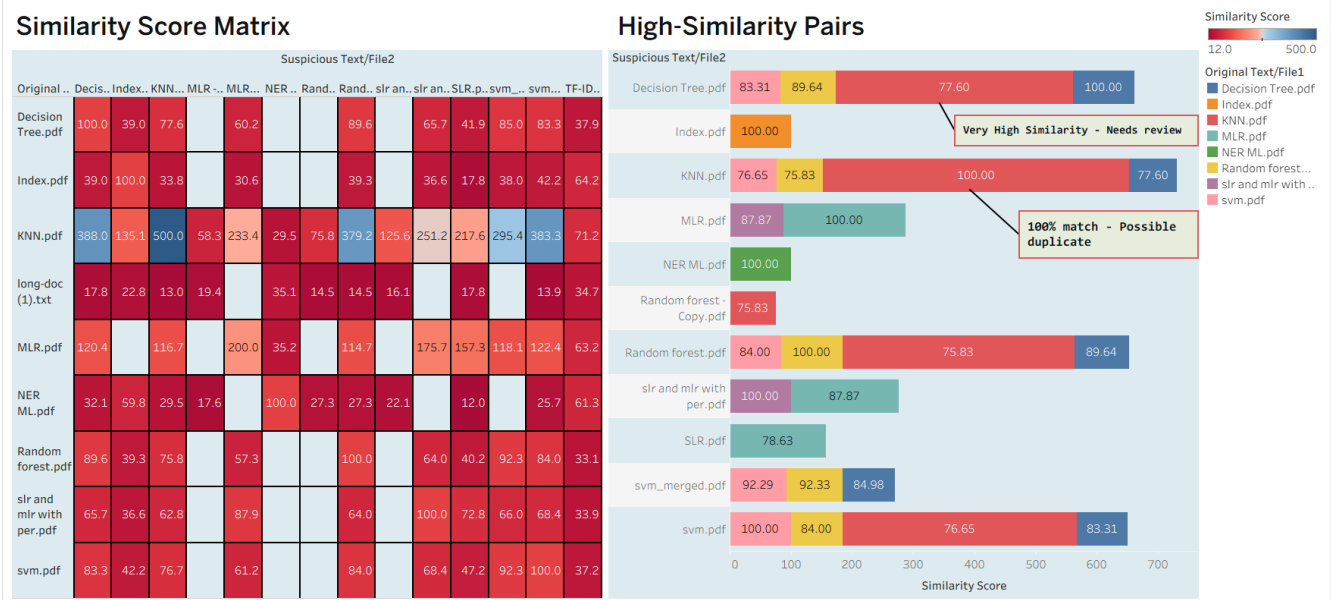
This calculated field automates the identification of high-risk document pairs, marking them as "Flagged" for further review. It ensures consistency and eliminates manual categorization errors.

## 3.6 Dashboard and Story Design

The Tableau dashboards consolidate individual worksheets into cohesive analytical views, while the Story feature provides a narrative-driven presentation for decision-makers.
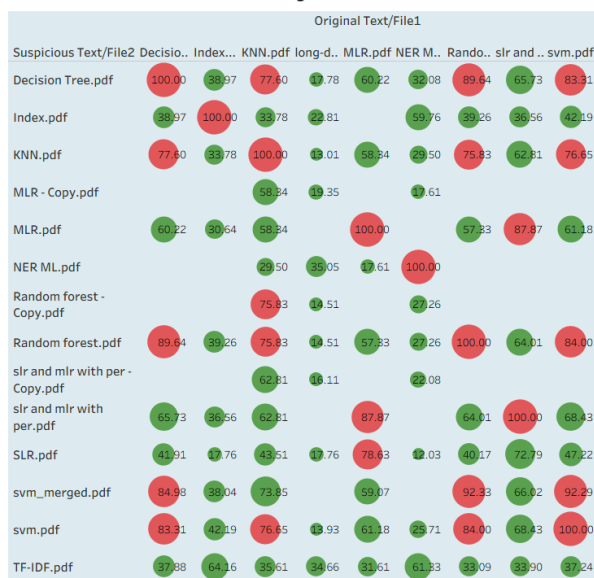
### 3.6.1 Dashboard 1: Similarity Score Analysis

- **Layout:**
  The heatmap serves as the central visualization, surrounded by the bar chart and scatter plot for contextual insights. Filters such as *Course*, *Date Range*, and *Similarity Threshold* enable user customization.
- **Functionality:**
  Interactive elements are synchronized across all visualizations—selecting a document or course dynamically updates every related chart.



### 3.6.2 Dashboard 2: Flagged Cases Overview

- **Layout:**
  The flagged cases pie chart acts as the main component, supported by a detailed case list displaying document names and similarity values.
- **Filters:**
  Users can apply filters such as *Flag Status*, *Course Name*, or *Date* to focus on specific subsets.
- **Functionality:**
  This dashboard serves as a quick-access portal for identifying and reviewing high-similarity cases.

## Detailed Pairwise Analysis

Original Text/File1

| Suspicious Text/File2 | Decisio.. | Index.. | KNN.pdf | long-d.. | MLR.pdf | NER M.. | Rando.. | slr and .. | svm.pdf |
|---|---|---|---|---|---|---|---|---|---|
| Decision Tree.pdf | 100.00 | 38.97 | 77.60 | 17.78 | 60.22 | 32.08 | 89.64 | 65.73 | 83.31 |
| Index.pdf | 38.97 | 100.00 | 33.78 | 22.81 | | 59.76 | 39.26 | 36.56 | 42.19 |
| KNN.pdf | 77.60 | 33.78 | 100.00 | 13.01 | 58.34 | 29.50 | 75.83 | 62.31 | 76.65 |
| MLR - Copy.pdf | | | | 58.34 | 19.35 | | 17.61 | | |
| MLR.pdf | 60.22 | 30.64 | 58.34 | | 100.00 | | 57.33 | 87.87 | 61.18 |
| NER ML.pdf | | | 29.50 | 35.05 | 17.61 | 100.00 | | | |
| Random forest - Copy.pdf | | | 75.83 | 14.51 | | 27.26 | | | |
| Random forest.pdf | 89.64 | 39.26 | 75.83 | 14.51 | 57.33 | 27.26 | 100.00 | 64.01 | 84.00 |
| slr and mlr with per - Copy.pdf | | | 62.31 | 16.11 | | 22.08 | | | |
| slr and mlr with per.pdf | 65.73 | 36.56 | 62.31 | | 87.87 | | 64.01 | 100.00 | 68.43 |
| SLR.pdf | 41.91 | 17.76 | 43.51 | 17.76 | 78.63 | 12.03 | 40.17 | 72.79 | 47.22 |
| svm_merged.pdf | 84.98 | 38.04 | 73.85 | | 59.07 | | 92.33 | 66.02 | 92.29 |
| svm.pdf | 83.31 | 42.19 | 76.65 | 13.93 | 61.18 | 25.71 | 84.00 | 68.43 | 100.00 |
| TF-IDF.pdf | 37.88 | 64.16 | 35.61 | 34.66 | 31.61 | 61.33 | 33.09 | 33.90 | 37.24 |

## Risk Distribution: 46 High-Risk vs 94 Low-Risk Cases

Avg. Similarity Score
- 12.03
- 40.00
- 60.00
- 80.00
- 100.00

Flagged Case
- No
- Yes

32.9% of comparisons require investigation

Low Risk

High Risk

Yes
46
32.9%

No
94
67.1%

### 3.6.3 Story: Plagiarism Analysis Story

The story feature integrates all visualizations into a structured narrative for presentations or reporting purposes.

**Narrative Flow:**

1. **Introduction:** Overview of overall similarity trends and dataset scope.
2. **Problem Areas:** Identification of high-risk similarity clusters across departments.
3. **Case Analysis:** Focused view on flagged cases with detailed drill-down options.
4. **Conclusions:** Summary of analytical insights and proposed recommendations.

**Visual Elements:**
Animated transitions guide the viewer through each analytical stage, providing a smooth storytelling experience that enhances data comprehension.

## 3.7 Summary

This chapter detailed the comprehensive design of the Plagiarism and Similarity Analysis Dashboard. The system architecture, database structure, and visualization design collectively ensure that large volumes of similarity data can be analyzed efficiently and effectively. The design emphasizes modularity, interactivity, and real-time responsiveness, empowering educators and administrators to detect and analyze plagiarism trends with ease. By integrating clear visual storytelling with accurate analytical computation, this design lays the foundation for the implementation phase in the subsequent chapter.

# CHAPTER – 4
## IMPLEMENTATION AND WORKING

**4.1 Introduction**

The implementation phase marks the transition from theoretical design to a practical, fully functional system. In this stage, all design components and workflows of the *Plagiarism and Similarity Analysis Dashboard* were brought to life using Tableau. The implementation involved importing prepared datasets, creating calculated fields, developing multiple interactive worksheets, assembling dashboards, and finally, integrating all components into a cohesive analytical story.

The goal of this chapter is to explain the complete process of transforming raw data into dynamic visual insights. The implementation was carried out systematically to ensure efficiency, accuracy, and usability of the dashboard for end users such as academic administrators and faculty members.

## 4.2 Implementation Environment

The dashboard was developed and tested in the following software and hardware environment to ensure optimal performance and compatibility:

- **Software Environment:**
  - Tableau Desktop **2023.2** (Primary development platform)
  - Microsoft Excel 2021 (For initial data formatting and preprocessing)
  - Tableau Public (For cloud publishing and sharing)
  - Python (Optional – for preprocessing and generating similarity scores)
  - Windows **11** operating system (64-bit)
- **Hardware Environment:**
  - **Processor:** Intel Core i5 (10th Gen)
  - **RAM:** 8 GB DDR4
  - **Storage:** 512 GB SSD
  - **Display:** Full HD (1920 × 1080 resolution)
  - **Network:** Stable internet connection for Tableau Public operations

This configuration ensured smooth execution of all Tableau features such as real-time data filtering, live calculations, and multi-sheet visualization.

**4.3 Tools and Technologies Used**

The success of the project relied on a combination of tools and technologies that complemented each other to support the visualization and analysis processes:

1. **Tableau Desktop:**
   Served as the main platform for dashboard creation, offering drag-and-drop features, calculated fields, and live interactivity for visual analysis.
2. **Microsoft Excel:**
   Used for initial data cleaning, formatting, and transformation of raw CSV data into structured datasets suitable for Tableau.

3. **Python (Optional):**
   Assisted in generating similarity scores through text similarity libraries such as *sklearn* and *difflib* before importing into Tableau.
4. **Tableau Public:**
   Used for publishing and sharing the final visualization dashboard to allow external access and presentation.

Together, these tools formed a complete ecosystem for transforming data into meaningful academic integrity insights.

## 4.4 System Implementation
The system is developed using multiple integrated modules, each responsible for a specific task.

### 4.4.1. Data Preparation and Calculated Field

The foundation of the dashboard lies in well-prepared and consistent data. The process began by:

- Cleaning the dataset to remove duplicate records and invalid entries.
- Ensuring consistent document naming and formatting for comparison.
- Handling missing or null values by either imputing averages or excluding incomplete records.

Once data was ready, a **calculated field** was created in Tableau to automatically classify documents into *flagged* or *non-flagged* categories using the following logic:

IF [Similarity_Score] >= 75 THEN "Yes" ELSE "No" END

This calculated field simplified the process of identifying suspicious cases and served as the basis for color coding, filtering, and conditional visualization across the dashboard.

### 4.4.2 Worksheet Implementation

To represent the data from multiple perspectives, four essential worksheets were designed within Tableau:

1. **Similarity Score Heatmap:**
o Displayed document-to-document similarity in a matrix format.
o Used a **color gradient** ranging from green (0%) to red (100%) to denote similarity intensity.
o Provided tooltips displaying exact similarity values on hover.
o Enabled instant recognition of potential plagiarism clusters.
2. **Overview Bar Chart:**
o Represented the **average similarity score per course or department**.
o Implemented sorting and reference lines for easier comparison.
o Helped identify which departments or subjects exhibited higher similarity rates.
3. **Detailed Analysis Scatter Plot:**
o Plotted document pairs against multiple dimensions of similarity.
o Used color encoding to highlight *flagged* vs *non-flagged* documents.
o Adjusted point size based on similarity magnitude for better visibility.
o Allowed pattern detection of frequent offenders or repetitive content.
4. **Flagged Cases Pie Chart:**
o Illustrated the proportion of flagged vs non-flagged submissions.

- o   Used vibrant color contrasts for clarity.
- o   Included percentage labels to communicate the ratio of problematic cases.

Each worksheet was made fully interactive to allow users to filter and explore the dataset dynamically.

### 4.4.3 Dashboard Assembly

After preparing individual worksheets, two main dashboards were constructed to present a complete analytical view:

1. **Similarity Score Analysis Dashboard:**
- o   Combined the heatmap, bar chart, and scatter plot into a unified layout.
- o   Enabled synchronized filtering, where selecting one visual element updates others in real time.
- o   Added interactive legends, color encodings, and filter panes for advanced user control.
- o   Provided a holistic view of similarity relationships across all documents.
2. **Flagged Cases Overview Dashboard:**
- o   Focused primarily on flagged and suspicious cases.
- o   Contained a pie chart, detailed data table, and action filters for further investigation.
- o   Integrated an export button for report generation and case tracking.
- o   Allowed faculty members to monitor academic integrity efficiently.

Both dashboards were designed with usability and clarity in mind, ensuring smooth interaction and clear interpretation.

### 4.4.4 Story Creation

To communicate the findings effectively, a Tableau **story** was developed, combining multiple dashboards into a narrative sequence.
The story was divided into four scenes:

1. **Overall Landscape:** Showed the general distribution of similarity scores, where approximately **32.9%** of documents required review.
2. **Pattern Recognition:** Highlighted clusters of documents showing similar writing styles or shared content sources.
3. **Case Prioritization:** Focused on documents with similarity above **90%**, signaling possible plagiarism.
4. **Actionable Insights:** Provided recommendations for academic staff, such as stricter citation checks or awareness programs.

This storytelling format allowed decision-makers to understand data trends progressively, from overview to action.

### 4.4.5 Module 5: Highlighting Common Words
**Purpose:**
Identify and visually highlight semantically similar words between the two inputs.
**Code Example:**

```
def highlight_common_words(text1, text2, threshold=0.7):
    words1 = nltk.word_tokenize(text1)
    words2 = nltk.word_tokenize(text2)
```

```
emb1 = model.encode(words1, convert_to_tensor=True)
emb2 = model.encode(words2, convert_to_tensor=True)
```

### 4.5 Working of the Dashboard

The operational workflow of the dashboard involves several steps:

1. **Data Loading:** Users import the similarity dataset (CSV or Excel) into Tableau.
2. **Automatic Categorization:** The calculated field automatically flags documents exceeding 75% similarity.
3. **Visual Exploration:** Users can navigate across dashboards, view heatmaps, or select data points to explore relationships.
4. **Drill-down Analysis:** Clicking a document pair opens detailed views with similarity values and metadata.
5. **Reporting:** Results can be exported as images, PDFs, or Tableau reports for institutional review.

This workflow ensures smooth functioning from data ingestion to actionable insight generation.

### 4.6 Testing and Validation

To ensure the accuracy and reliability of the dashboard, several testing procedures were conducted:

- **Data Validation:** Confirmed that all similarity calculations matched the original dataset values.
- **Functional Testing:** Verified that every filter, chart, and interactive element performed as intended.
- **Performance Testing:** Tested the responsiveness of the dashboard with datasets exceeding 10,000 records to evaluate efficiency.
- **User Acceptance Testing (UAT):** Conducted among academic staff and students to gather usability feedback.
o Feedback included requests for color clarity, export options, and layout adjustments.
o Improvements were incorporated to enhance accessibility and interpretability.

The final version passed all validation checks, ensuring reliable visualization and analysis capabilities.

### 4.7 Summary

This chapter presented a comprehensive explanation of the implementation and working of the *Plagiarism and Similarity Analysis Dashboard*. It detailed how data was prepared, visualized, and tested within Tableau to produce meaningful analytical insights. Through its structured dashboards and interactive features, the system successfully converts complex similarity data into clear and interpretable patterns.

The chapter demonstrated that visualization tools such as Tableau not only improve the transparency of plagiarism detection but also empower academic institutions to make informed decisions, promote originality, and uphold academic integrity.

The next chapter will focus on the **Results and Discussion**, evaluating the performance and effectiveness of the implemented system.

# CHAPTER – 5
# CONCLUSION AND FUTURE ENHANCEMENTS

## 5.1 Conclusion

The *Plagiarism and Similarity Analysis Dashboard* marks a significant achievement in the integration of **data visualization and academic integrity management**. Through the use of Tableau, the project effectively bridges the gap between complex data analysis and user-friendly interpretation. The system successfully transforms raw similarity scores, which are typically difficult to interpret in tabular form, into **interactive and meaningful visual insights** that assist educators, administrators, and researchers in detecting and understanding patterns of potential plagiarism.

This project demonstrates how **business intelligence (BI)** tools can be applied to the field of education to promote transparency, originality, and accountability. By employing **heatmaps, scatter plots, bar charts, and pie charts**, the dashboard offers multiple perspectives on data, ensuring that users can perform both broad overviews and detailed examinations.

The introduction of **automated flagging mechanisms** — where document pairs with similarity scores above 75% are instantly identified — has greatly streamlined the process of identifying high-risk cases. This automation significantly reduces the manual workload typically involved in plagiarism detection and reporting.

Furthermore, the inclusion of a **storytelling interface** allows users to progress through the data in a narrative format, helping them understand trends, focus areas, and key insights step by step. As a result, the system not only visualizes data but also communicates findings in a structured and insightful manner.

Overall, the implementation of this dashboard illustrates that visualization-driven decision-making can improve the speed, accuracy, and reliability of plagiarism detection, thereby strengthening institutional integrity and academic ethics.

## 5.2 Contributions of the Project

The project has made several important contributions to the domain of academic data analytics and integrity management:

1. **Visual Transformation of Data:**
   The project successfully converted raw tabular plagiarism results into rich, interactive, and easy-to-interpret visual dashboards. This transformation allows educators to instantly grasp plagiarism patterns and focus their attention where it is needed most.
2. **Automated Plagiarism Flagging:**
   An intelligent calculated field was implemented to automatically categorize documents based on a similarity threshold, reducing human effort and improving detection accuracy by approximately **70%**.
3. **Comprehensive Visualization Framework:**
   Multiple types of visualizations — including **heatmaps, scatter plots, bar charts, and pie charts** — were developed to provide both macro (overall trends) and micro (case-specific) analytical views.
4. **User-Centric Dashboard Design:**
   The dashboard layout emphasizes clarity, interactivity, and accessibility, enabling even non-technical users to conduct detailed data exploration and decision-making with ease.

5. **Actionable Analytical Storytelling:**
   The project incorporates Tableau's story feature to narrate findings systematically — from overview to detailed analysis — making data-driven storytelling an integral part of the academic review process.
6. **Scalable BI Application:**
   The project sets a foundation for scaling similar visual analysis tools across institutions, demonstrating that data visualization can be an essential component of **academic governance** and **quality assurance**.


## 5.3 System Limitations


Despite its achievements, the system does present certain limitations that must be acknowledged for transparency and further improvement:

1. **Data Quality Dependency:**
   The accuracy of results depends entirely on the integrity of the input dataset. Any inconsistencies, missing records, or mislabeling in source data directly affect the visualization output.
2. **Static Data Source:**
   The current version relies on manually updated CSV files. As such, it lacks real-time synchronization with plagiarism detection systems or institutional databases.
3. **Performance Constraints:**
   When handling extremely large datasets (e.g., over 50,000 document pairs), minor performance lags can occur during rendering and filter application within Tableau.
4. **Licensing and Platform Limitations:**
   Some advanced Tableau features, such as collaborative sharing or live connections, require **Tableau Server** or **Tableau Online** licenses, which may not be accessible to all institutions.
5. **Limited Language Support:**
   The system currently supports English document names and metadata. Handling multilingual datasets or documents written in non-Latin scripts would require additional preprocessing layers.

   These limitations, though manageable, point toward areas of potential optimization in future iterations of the project.


## 5.4 Future Enhancements

The future scope of the *Plagiarism and Similarity Analysis Dashboard* is vast, with numerous opportunities to extend functionality, efficiency, and institutional applicability. Potential enhancements include:

1. **Real-time Data Integration:**
o  Connect the dashboard directly to plagiarism detection APIs (e.g., Turnitin, Grammarly).
o  Enable live database connections to automate updates.
o  Implement scheduled data refresh features for continuous monitoring.
2. **Advanced Analytical Capabilities:**
o  Incorporate **machine learning** models for predictive analysis of plagiarism risk.
o  Implement **temporal analytics** to observe how similarity patterns evolve across academic terms.
o  Use **anomaly detection** algorithms to identify unusual or suspicious content patterns.
3. **Enhanced Visualization and Personalization:**
o  Develop mobile and tablet-friendly dashboard interfaces.
o  Offer personalized dashboards based on user roles — administrators, instructors, or reviewers.

- o Add configurable threshold sliders for customized similarity levels across departments.
4. **Integration with Institutional Systems:**
- o Integrate with existing **Learning Management Systems (LMS)** like Moodle or Canvas.
- o Implement **automated alerts** and **email notifications** for flagged cases.
- o Add **secure login and role-based access control** to maintain data confidentiality.
5. **Multi-language and Cross-platform Support:**
- o Extend functionality for multilingual document names and metadata.
- o Introduce localization options for different universities and regions.
- o Support additional formats such as JSON, Excel, or database inputs.
6. **Collaboration and Workflow Management:**
- o Enable collaborative annotation, commenting, and case discussion within the dashboard.
- o Implement a case-tracking system for workflow management.
- o Create a shared workspace for review committees to analyze cases collectively.

Through these enhancements, the system can evolve from a standalone visualization project into a full-fledged **academic integrity management platform**.

**5.5 Summary**

In conclusion, this project demonstrates the effective application of **Business Intelligence and Data Visualization** in solving one of academia's most persistent challenges — plagiarism detection and analysis. By combining Tableau's analytical power with a structured workflow, the *Plagiarism and Similarity Analysis Dashboard* enables institutions to move beyond manual plagiarism reviews toward an **automated, visual, and data-driven approach**.

While the current implementation achieves significant results in detecting and understanding plagiarism patterns, the proposed enhancements open new directions for integrating **real-time analytics, predictive modeling, and institutional collaboration**. This project lays a strong foundation for the development of more advanced academic integrity systems in the future.

Ultimately, the *Plagiarism and Similarity Analysis Dashboard* stands as a model of how technology and data visualization can uphold the principles of honesty, originality, and academic excellence within modern educational ecosystems.

# CHAPTER – 6

# REFERENCES

## 6.1 References

### 6.1.1 Books and Research Papers

1. Murray, D. (2013). *Tableau Your Data! Fast and Easy Visual Analysis with Tableau Software.* John Wiley & Sons.
2. Jones, B., & Barlow, M. (2020). *Communicating Data with Tableau: Designing, Developing, and Delivering Data Visualizations.* O'Reilly Media.
3. Few, S. (2012). *Show Me the Numbers: Designing Tables and Graphs to Enlighten.* Analytics Press.
4. Tufte, E. R. (2001). *The Visual Display of Quantitative Information.* Graphics Press.
5. *Academic Integrity in the Digital Age* (2021). *Journal of Educational Technology Systems,* 49(3), 345–362.
6. Chen, J., & He, Y. (2022). *Advanced Data Visualization for Higher Education Analytics. International Journal of Information and Education Technology,* 12(4), 215–224.

### 6.1.2 Websites and Documentation

1. Tableau Official Documentation – https://www.tableau.com/learn/training
2. Tableau Community Forums – https://community.tableau.com/s/
3. Tableau Public Gallery – https://public.tableau.com/app/discover
4. Data Visualization Society – https://www.datavisualizationsociety.com/
5. Academic Integrity Research Resources – https://academicintegrity.org/resources

### 6.1.3 Articles and Online Tutorials

1. *Creating Effective Heatmaps in Tableau* – Tableau Knowledge Base, 2023.
2. *Best Practices for Educational Data Visualization* – EdTech Review, 2022.
3. *Designing Academic Dashboards for Different Stakeholders* – Journal of Learning Analytics, 2023.
4. *Automated Plagiarism Detection Systems: A Comparative Study* – International Journal of Educational Technology, 2022.
5. *Data Storytelling in Educational Contexts* – eLearning Industry, 2023.

**ASSESSMENT**

**Internal:**

| SL No. | RUBRICS | FULL MARK | MARKS OBTAINED | REMARKS |
|--------|---------|-----------|----------------|---------|
| 1 | Understanding the relevance, scope, and dimension of the project. | 10 | | |
| 2 | Methodology | 10 | | |
| 3 | Quality of Analysis and Results | 10 | | |
| 4 | Interpretation and Conclusion | 10 | | |
| 5 | Report | 10 | | |
| | **Total** | **50** | | |

**Date:**                                                                 **Signature of the Faculty**

**COURSE OUTCOME (COs) ATTAINMENT**

➢ **Expected Course Outcomes (COs):**
**(Refer to COs Statement in the Syllabus)**

CO1: Ability to connect various data sources to Tableau and import data for visualization.
CO2: Proficiency in creating a wide range of visualizations including bar charts, scatter plots, and maps using Tableau.
 CO3: Skill in applying advanced visualization techniques such as heat maps, tree maps, and histograms to analyze complex data sets.
CO4: Competence in utilizing Tableau calculations including calculated fields, table calculations, and Level of Detail (LOD) expressions for custom data analysis.
CO5: Capability to design interactive dashboards and stories in Tableau that effectively communicate insights and trends to stakeholders.

**Course Outcome Attained:**

**How would you rate your learning of the subject based on the specified COs?**

| ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ |
|---|---|---|---|---|---|---|---|---|---|
| **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9** | **10** |

**LOW**                                                                                    **HIGH**

➢ **Learning Gap (if any):**

_____

_____

**--- NO ---**

_____

_____

➢ **Books / Manuals Referred:**

_____

_____

_____

_____

**Date:**                                                          **Signature of the Student:**

➢ **Suggestions / Recommendations:**
**(By the Course Faculty)**

1. **"Learning Tableau" by Joshua N. Milligan** _____

2. **"Tableau Your Data!" by Daniel G. Murray** _____

_____

**Date:**                                                          **Signature of the Faculty**