

Project Report On

Customer Segment Explorer

(Six weeks Industrial Training)

SUBMITTED IN THE PARTIAL FULFILMENT OF THE REQUIREMENT FOR THE
AWARD
OF THE DEGREE OF

BACHELOR OF TECHNOLOGY
IN
COMPUTER SCIENCE &
ENGINEERING

Batch

(2024 - 2027)



Submitted By:
Kuljit kaur

University Roll No:
2430480

Under the Guidance of:
Er. Gurmeet Singh

GLOBAL GROUP OF INSTITUTES, AMRITSAR
JUNE-JULY, 2025

Acknowledgement

I hereby certify that this project is original and solely my work. Its successful completion would not have been possible without the guidance and support of several individuals, to whom I am deeply grateful.

First and foremost, I extend my sincere appreciation to my Head of the Department, Er. Tejinder Deep Singh, and my mentor, Er. Gurmeet Singh. Their expertise, valuable insights, and constructive feedback played a pivotal role in shaping this project.

I am also immensely grateful to the college administration for providing me with this opportunity. This experience has been highly enriching, and I look forward to participating in similar endeavours in the future.

Furthermore, I would like to express my heartfelt gratitude to my parents and friends for their unwavering support and encouragement throughout this journey. Their thoughtful feedback and motivation were instrumental in bringing this project to completion.

Finally, I extend my thanks to everyone who contributed to this project in any capacity. Your assistance and guidance have been truly invaluable.

KULJIT KAUR

Student's Declaration

I hereby certify that the work presented in this project report, titled “Customer Segment Explorer”, has been carried out in fulfillment of the requirements for the Bachelor of Technology in the Department of Computer Science & Engineering at Global Group of Institutes, Sohian Khurd. This report is a record of my original work conducted during the six weeks industrial training.

August 6, 2025

Kuljit Kaur

The major project viva-voce examination of Kuljit Kaur, Roll No. 2430480 of B.TECH (Computer Science & Engineering) has been held on _____.

Certificate of the company

ANSH INFOTECH

DEVELOPMENT
TRAINING
CONSULTANCY



Ref. No. AIT/IT/CS/2508/4114

Dated. 06-08-2025



This is to certify that

Ms. Kuljit Kaur D/o Mr. Kulwant Singh

of Global Group Of Institutes, Amritsar

has completed her training in

“Data Science”

From 10-June-2025 to 25-July-2025 at our organization.

During Industrial training her Performance was Excellent.

We wish her success for her future endeavors.

Authorized Signature



SCF-4, 3rd Floor, Model Town Ext.,
D-Block Market, Dugri Road,
Ludhiana-141001.

E-mail.: contact@anshinfotech.org

Website: www.anshinfotech.org

M.: 94175-69963, 94171-68347



Table of Contents

S. No	Description	Pg No.
1)	Introduction.....	1
2)	Literature Study.....	6
3)	Problem Statement.....	10
4)	Proposed Methodology.....	12
5)	Implementation	17
6)	Results and Discussion.....	28
7)	Project Screenshots.....	37
8)	Conclusion and Future Scope.....	49
9)	References	

Bibliography

CHAPTER 1:

INTRODUCTION

1.1 THE MODERN BUSINESS CONTEXT: FROM TRANSACTIONS TO RELATIONSHIPS

In the 21st century, the global business landscape has undergone a seismic shift. The proliferation of the internet, the rise of e-commerce, and the digitalization of commerce have created a hyper-competitive environment. Companies no longer compete purely on product or price; they compete on the basis of the customer experience. Giants of industry, from Amazon in retail to Netflix in entertainment, have demonstrated that a deep, data-driven understanding of the customer is not just an advantage, but a core component of business survival.

This paradigm shift has moved the business focus from simple, transactional exchanges to the cultivation of long-term, valuable customer relationships. This is the central premise of Customer Relationship Management (CRM), a strategy that places the customer at the epicenter of all business processes. In an era where a customer can switch loyalties with a single click, the ability to acquire, retain, and grow the value of a customer base is paramount.

Effective CRM, however, is impossible without understanding. In a market of millions, a "one-size-fits-all" approach to marketing, sales, and service is demonstrably inefficient. It wastes resources on uninterested parties, alienates existing customers with irrelevant offers, and fails to nurture high-potential prospects. The key to unlocking effective CRM is to recognize the inherent heterogeneity of the customer base. The fundamental question businesses must answer is not just "Who are My customers?" but rather, "What distinct groups of customers do I have, and how do their needs and behaviors differ?"

1.2 THE STRATEGIC ROLE OF CUSTOMER SEGMENTATION

This is where the strategic practice of customer segmentation becomes essential. Customer segmentation is the process of dividing a broad, heterogeneous customer base into smaller, more homogeneous groups (or segments) of customers who share similar characteristics. By moving from a mass-market view to a multi-segment view, an organization can tailor its strategies to meet the specific needs of each segment, dramatically improving efficiency and effectiveness.

The benefits of a robust segmentation strategy are profound and far-reaching:

- **Personalized Marketing:** Segmentation allows for the creation of targeted marketing campaigns. Messages and offers can be crafted to resonate with the specific interests and past behaviors of a segment, leading to significantly higher engagement and conversion rates.
- **Enhanced Customer Retention:** By identifying segments that are at high risk of "churning" (leaving the service), a company can proactively intervene with special offers, support, or loyalty programs designed to retain them.
- **Optimized Resource Allocation:** Not all customers are equally profitable. Segmentation helps identify high-value customer groups, allowing a business to focus its marketing budget, sales efforts, and premium services on the segments that provide the greatest return on investment (ROI).
- **Improved Product Development:** Feedback and behavioral data from distinct segments can provide invaluable insights for product development, helping to identify feature gaps or new market opportunities.

While segmentation can be performed using various criteria, the most common approaches include:

1. **Demographic Segmentation:** Grouping by statistical data such as age, gender, income, education, and occupation.
2. **Geographic Segmentation:** Grouping by physical location, such as country, region, city, or climate.
3. **Psychographic Segmentation:** Grouping by lifestyle, values, personality traits, and social class.
4. **Behavioral Segmentation:** Grouping by a customer's observable actions, such as purchase history, product usage, brand loyalty, and browsing habits.

Of these, behavioral segmentation is often considered the most powerful for driving tangible business outcomes, as it is based on concrete, past interactions rather than abstract traits.

1.3 RFM Analysis: A Classic Model for Behavioral Segmentation

Within the domain of behavioral segmentation, one of the most classic, effective, and widely adopted models is RFM Analysis. RFM is an acronym for Recency, Frequency, and Monetary Value. It is a powerful framework used to quantify and rank customers based on their transactional history.

The elegance of the RFM model lies in its simplicity and its foundation in three clear, observable, and highly predictive behavioral metrics:

- **Recency (R):** How much time has elapsed since a customer's last transaction? The underlying premise is that customers who have purchased more recently are more likely to be engaged and responsive to new offers compared to those who have not purchased in a long time.
- **Frequency (F):** How many transactions has a customer made over a given period? This metric identifies customer loyalty and engagement. A customer who purchases frequently is inherently more valuable and reliable than a one-time buyer.
- **Monetary (M):** What is the total monetary value of a customer's purchases? This metric identifies high-spenders. It helps distinguish customers who make large, infrequent purchases from those who make small, frequent ones.

By calculating these three values for every customer, a business can quickly identify its "best" customers (high R, F, and M) and its "at-risk" customers (low R, F, and M), and everyone in between.

1.4 THE NEXT EVOLUTIONARY STEP: RFM WITH K-MEANS CLUSTERING

While traditional RFM analysis is powerful, its common implementation has limitations. Often, analysts will manually create segments by dividing each of the R, F, and M scores into quintiles (scores of 1 to 5) and then create static segments based on these scores (e.g., customers with a score of '555' are 'Champions', while those with '111' are 'Lost').

This manual, rules-based approach is subjective. The thresholds are arbitrary and may not represent the true, natural groupings that exist within the data. It forces the data into predefined boxes. The central hypothesis of this project is that a more sophisticated, data-driven approach can yield more accurate, actionable, and non-obvious segments.

This project bridges the gap between traditional RFM and modern data science by integrating unsupervised machine learning. Specifically, this project applies the K-Means Clustering algorithm to the three-dimensional RFM feature space.

K-Means Clustering is an algorithm that partitions a dataset into 'k' distinct, non-overlapping clusters. It works by iteratively finding cluster "centers" (centroids) and assigning each data point to the nearest centroid, such that the total variance within each cluster is minimized.

By applying K-Means to the RFM data, this project moves beyond subjective binning. It allows the algorithm to "listen to the data" and discover the natural, inherent structures and groupings. This process can reveal non-linear relationships and uncover customer personas that manual analysis would miss—for example, a 'High-Value, At-Risk' segment (low Recency, but high Frequency and Monetary) or a 'New and Promising' segment (high Recency, but low Frequency and Monetary). These data-driven segments provide a far more nuanced foundation for targeted business strategy.

1.5 PROJECT AIM AND OBJECTIVES

The primary aim of this project is to design, implement, and deploy a comprehensive, end-to-end customer segmentation solution. This system will transform raw, transactional e-commerce data into actionable, data-driven customer segments, presented through an interactive web-based dashboard for business users.

To achieve this aim, the following specific objectives have been defined:

- To preprocess and clean the raw 'Online Retail' dataset, systematically handling missing values, filtering out irrelevant (non-transactional) entries, and converting data into appropriate formats for analysis.
- To perform feature engineering by calculating the Recency, Frequency, and Monetary (RFM) metrics for each unique customer from the cleaned transaction logs.
- To determine the optimal number of customer segments (k) by applying and visualizing standard unsupervised learning diagnostics, specifically the 'Elbow Method' and Silhouette Score analysis.
- To train a K-Means clustering model on the scaled RFM feature set, partitioning the customer base into 'k' distinct, data-driven clusters.

- To validate the business meaningfulness and predictive power of the RFM features by training a supervised classification model (Logistic Regression) to predict known, expert-defined segments.
- To develop and deploy an interactive, user-friendly web application using the Plotly Dash framework. This dashboard will serve as the project's final deliverable, allowing a non-technical user to upload new customer data, view the resulting segment distributions on interactive graphs, and drill down into the details of individual customer profiles.

CHAPTER 2

LITERATURE STUDY

2.1 INTRODUCTION TO THE LITERATURE

Before starting any new project, it's important to look at the work that others have already done. This is called a "literature study" or "literature review." For this project, I looked at research and articles on four main topics:

1. Why customer segmentation is important.
2. The RFM (Recency, Frequency, Monetary) model.
3. The K-Means Clustering method.
4. Why dashboards are important for business.

This chapter summarizes what experts and other researchers have said about these topics. This helps to show why the methods used in this project were chosen and that they are trusted and effective.

2.2 THE IMPORTANCE OF CUSTOMER SEGMENTATION

Almost all modern business and marketing research agrees on one key idea: you cannot treat all customers the same. Studies from the past few decades, such as those by Smith (1956) who first introduced the idea of "market segmentation," show that businesses grow when they understand their customers' differences.

Researchers like Philip Kotler, a famous expert in marketing, have written extensively that segmentation is the first step in a successful "Targeting" and "Positioning" (STP) strategy.

The general idea is that by dividing the market into small, understandable groups, a company can:

- **Create better products:** They can build things for a specific group's needs.
- **Market more efficiently:** They can stop wasting money on ads that people ignore and instead send the right message to the right person.
- **Keep customers longer:** When customers feel understood, they are more likely to stay loyal.

This project builds on this core idea. It's not just about selling things; it's about understanding who is buying.

2.3 RFM: A TRUSTED MODEL FOR UNDERSTANDING BEHAVIOR

Once I agree that I need to segment customers, the next question is how. Researchers have tried many ways, like using age or location. However, a large amount of research shows that the best way to predict what a customer will do in the future is to look at what they did in the past. This is called behavioral segmentation.

This is where the RFM (Recency, Frequency, Monetary) model comes in. The RFM model was first introduced in the 1990s by direct marketers (people who send catalogs in the mail) and has been studied and proven effective ever since.

- Studies by researchers like Hughes (1994) showed that RFM was a simple but powerful way to find a company's best customers.
- More recent studies, like those by Fader, Hardie, and Lee (2005), have shown that RFM models are very accurate for predicting a customer's "lifetime value" (how much they might be worth over time).

The reason RFM is so popular in both business and research is that it's simple to understand and it just works. It is based on the common-sense idea that the best customers are the ones who bought recently, buy often, and spend the most.

However, the "old" way of using RFM was to manually assign scores from 1 to 5 for each part. Many researchers have pointed out that this is too simple and based on opinion, not data. This leads us to the next part of My project.

2.4 K-MEANS: THE "SMART" WAY TO FIND GROUPS

Instead of manually creating segments (like "High Recency" or "Low Frequency"), modern data science offers a better way: letting a machine find the segments for us. This is done using clustering.

K-Means is one of the most popular and well-understood clustering algorithms in the world. It was first proposed by Stuart Lloyd in 1957. In simple terms, K-Means looks at all the customer data points (in My case, their R, F, and M values) and automatically finds the "center" of different groups.

- Many modern studies have shown that combining RFM and K-Means is a powerful technique. Researchers like J.T. Lee (2018) have published papers doing exactly this.
- The benefit of this mix is that I get the best of both worlds:
 - a. I use RFM to create very strong, meaningful features (R, F, M).
 - b. I use K-Means to automatically and objectively find the real, natural groupings within those features.

This is the exact method used in this project. The research shows this is a valid and strong approach. This is why My 002_train_model.py script is focused on building a K-Means model on top of the RFM data from 001_building_rfm.py.

2.5 VALIDATION: HOW DO I KNOW MY SEGMENTS ARE GOOD?

A common problem in data science is making segments that look good but are not actually useful or meaningful. This is why My project includes *003_supervised_validation.py*.

- Research on model validation emphasizes that I must check My work.
- One common way to do this, as seen in many data science papers, is to "test" the features.
- Our project uses a Logistic Regression model for this. This is a standard, trusted method for classification.
- The idea is simple: If My RFM numbers are truly good at separating customers, then a different model (Logistic Regression) should be able to accurately guess the customer's segment just by looking at their R, F, and M numbers.
- When My classification_report.txt shows high accuracy (like 99% or 100%), it confirms that My RFM features are very strong and My segments are clear and well-defined. This is a key part of making sure My results are trustworthy.

2.6 DASHBOARDS: BRINGING DATA TO LIFE

The final piece of research is about how to deliver the results. A data science project is useless if the final results are stuck in a text file or a complex Python script.

- Research in the field of Business Intelligence (BI) and Human-Computer Interaction shows that people make better decisions when they can see and interact with data visually.

- Tools like Plotly Dash (which this project uses in 004_app.py), Tableau, and Power BI were created based on this idea.
- They create an "interactive dashboard" that acts as a bridge between the complex data science (the models) and the business user (the manager).
- This allows a manager, with no coding knowledge, to explore the segments, see the graphs, and even drill down into a single customer's details.

2.7 CHAPTER SUMMARY

Our review of existing work shows that the path this project takes is well-supported by researchers and experts.

1. I start with Segmentation, which is known to be a key business strategy.
2. I use RFM, which is a simple, proven model for understanding customer behavior.
3. I combine it with K-Means, which is a popular, data-driven method for finding natural groups, an approach many modern studies recommend.
4. I Validate My model with a classifier, which is a strong method for proving My segments are meaningful.
5. I present the results in an Interactive Dashboard, which is a best practice for making data useful for decision-makers.

This project is not inventing a brand-new method from scratch. Instead, it is intelligently combining several trusted, proven, and powerful ideas into a single, complete, and useful tool.

CHAPTER 3

PROBLEM STATEMENT

3.1 THE PROBLEM FOUND IN ACADEMIC RESEARCH

During the research for this project (as discussed in Chapter 2), a clear pattern emerged. Many academic papers and technical blogs show how to build customer segmentation models. They explain, in great detail, the mathematics of RFM (Recency, Frequency, Monetary) and the K-Means algorithm. They are very good at proving that these methods work.

But this is where they stop.

The "deliverable" or final product of this research is almost always:

1. A Python script (like My_001_building_rfm.py or 002_train_model.py).
2. A saved model file (a .pkl file).
3. A static image of a graph or a text file (like My_classification_report.txt).

This is the first problem: This research is written for data scientists, by data scientists. It is an "engine with no car." A marketing manager or a small business owner cannot use a .pkl file or a Python script to make real decisions. This research solves the mathematical problem but completely ignores the usability problem.

3.2 THE PROBLEM FOUND IN BUSINESS SOFTWARE

On the other side, there are large, expensive business software (like Tableau, Salesforce, or Power BI). These platforms are designed for business users and have beautiful dashboards.

But they have the opposite problem.

These tools are generic. They are like a giant, complicated box of building blocks. They are not a "customer segmentation tool" out of the box. To get RFM analysis, a company must:

1. Buy the expensive software.
2. Hire a data expert to spend weeks or months building the entire RFM and K-Means logic inside that tool.
3. The tool is not flexible. It's often locked into its own ecosystem.

This is the second problem: These tools are too complex, too expensive, and not specialized for this specific task. A small business owner cannot just "upload a sales file" and get segments.

3.3 THE CORE PROBLEM: THE "MISSING BRIDGE"

After looking at both sides, the real problem becomes clear. It's what is missing between these two worlds.

1. **Academic projects** are powerful but unusable for business.
2. **Business software** is (eventually) usable but complex, expensive, and not specialized.

The "missing thing" is a simple, free, and specialized tool that acts as a bridge. There is no simple application that connects a powerful, custom-built data science model directly to a clean, easy-to-use interface for a manager.

3.4 THIS PROJECT'S SOLUTION

This project is designed to solve that exact problem. The Problem Statement is:

- *Existing solutions for customer segmentation are either too technical for a business user (like a Python script) or too complex, expensive, and generic for a simple, specific task (like a BI platform).*

This project solves this by building the "missing bridge." It combines both parts into one complete, simple, and end-to-end solution:

- **The Engine:** The powerful and accurate Python scripts (*001_building_rfm.py*, *002_train_model.py*, *003_supervised_validation.py*) that do the heavy lifting.
- **The Dashboard:** The simple, user-friendly web application (*004_app.py*) that anyone can use.

This project creates a tool where a non-technical user can get the power of a complex Python model without ever seeing the code. They can simply upload their sales file and get a complete, interactive report, solving the "usability" and "accessibility" problem that other projects ignore. The specific features of this new, proposed solution will be detailed in the next chapter.

CHAPTER 4

PROPOSED METHODOLOGY

4.1 MY GUIDING PHILOSOPHY: BUILDING THE "MISSING BRIDGE"

As identified in Chapter 3, the problem with existing customer segmentation projects is that they are either too technical for a business user or too complex and expensive.

Our proposed methodology is therefore designed to be the "missing bridge." The entire system is built as a hybrid model that combines two distinct parts:

1. **A Powerful "Data Engine" (The Backend):** A set of five automated Python scripts that handle all the complex data processing, feature engineering, modeling, and validation.
2. **An "Easy-to-Use Dashboard" (The Frontend):** A simple, clean, and interactive web application that acts as the "driver's seat." This allows a non-technical user to run the entire engine with the click of a button and easily understand the results.

This approach gets the best of both worlds: the power of a custom data science model and the simplicity of a user-friendly tool.

4.2 HIGH-LEVEL SYSTEM ARCHITECTURE

The proposed system works in a continuous flow, which is triggered entirely by the user through the web dashboard. The diagram below shows how all the pieces of the project (your .py files) fit together.

4.3 PART 1: THE DATA-PROCESSING ENGINE (THE BACKEND)

The "engine" is the core of the project's intelligence. It is a pipeline of Python scripts that do all the heavy lifting.

4.3 PART 1: THE DATA-PROCESSING ENGINE (THE BACKEND)

The "engine" is the core of the project's intelligence. It is a pipeline of Python scripts that do all the heavy lifting.

4.3.1 Step 1: RFM Feature Engineering (001._building_rfm.py)

This script is the "data factory." Its only job is to take the raw, messy sales data and turn it into clean, useful features.

- Input: The Online_Retail.csv file uploaded by the user.
- Process:
 - a. Load & Clean: It loads the data, removes any rows with missing CustomerIDs, and removes cancelled orders (invoices starting with "C").
 - b. Calculate Recency: It finds the most recent purchase date for each customer and counts the days between then and today.
 - c. Calculate Frequency: It counts the total number of unique invoices for each customer.
 - d. Calculate Monetary: It sums the total price of all purchases for each customer.
- Output: A new, clean table with one row per customer, containing their CustomerID, Recency, Frequency, and Monetary values.
-

4.3.2 Step 2: Unsupervised Clustering (002_train_model.py)

This script is the "brain" of the operation. It takes the clean RFM data and automatically discovers the customer groups.

- Input: The RFM data from Step 1.
- Process:
 - a. Scale Features: It scales the R, F, and M values so that no single feature (like Monetary) unfairly dominates the others.
 - b. Find K: It uses the "Elbow Method" to find the ideal number of clusters (in My case, 4).
 - c. Train Model: It trains a K-Means clustering model to group all customers into one of the 4 segments.
- Output: A new clustered_customer.csv file that includes the original RFM data plus a new "Cluster" label for each customer. It also saves the trained model (.pkl file) for later use.

4.3.3 Step 3: Model Validation (003_supervised_validation.py)

This script is the "quality check." It answers the question: "Are the segments I found actually good?"

- **Input:** The clustered_customer.csv file from Step 2.
- **Process:** It builds a different kind of model (a Logistic Regression classifier) to see if it can predict the cluster label just by looking at the R, F, and M values.
- **Output:** A classification_report.txt. If this report shows high accuracy (e.g., 99-100%), it proves that My segments are very clear, distinct, and meaningful.

4.3.4 Step 4: Cohort Analysis (005_build_cohort.py)

This script provides a different angle on customer behavior, focusing on loyalty over time.

- **Input:** The original raw sales data.
- **Process:** It groups customers by the month they made their first purchase (a "cohort") and tracks what percentage of each group comes back to buy again in the following months.
- **Output:** A retention heatmap/table, showing customer loyalty.

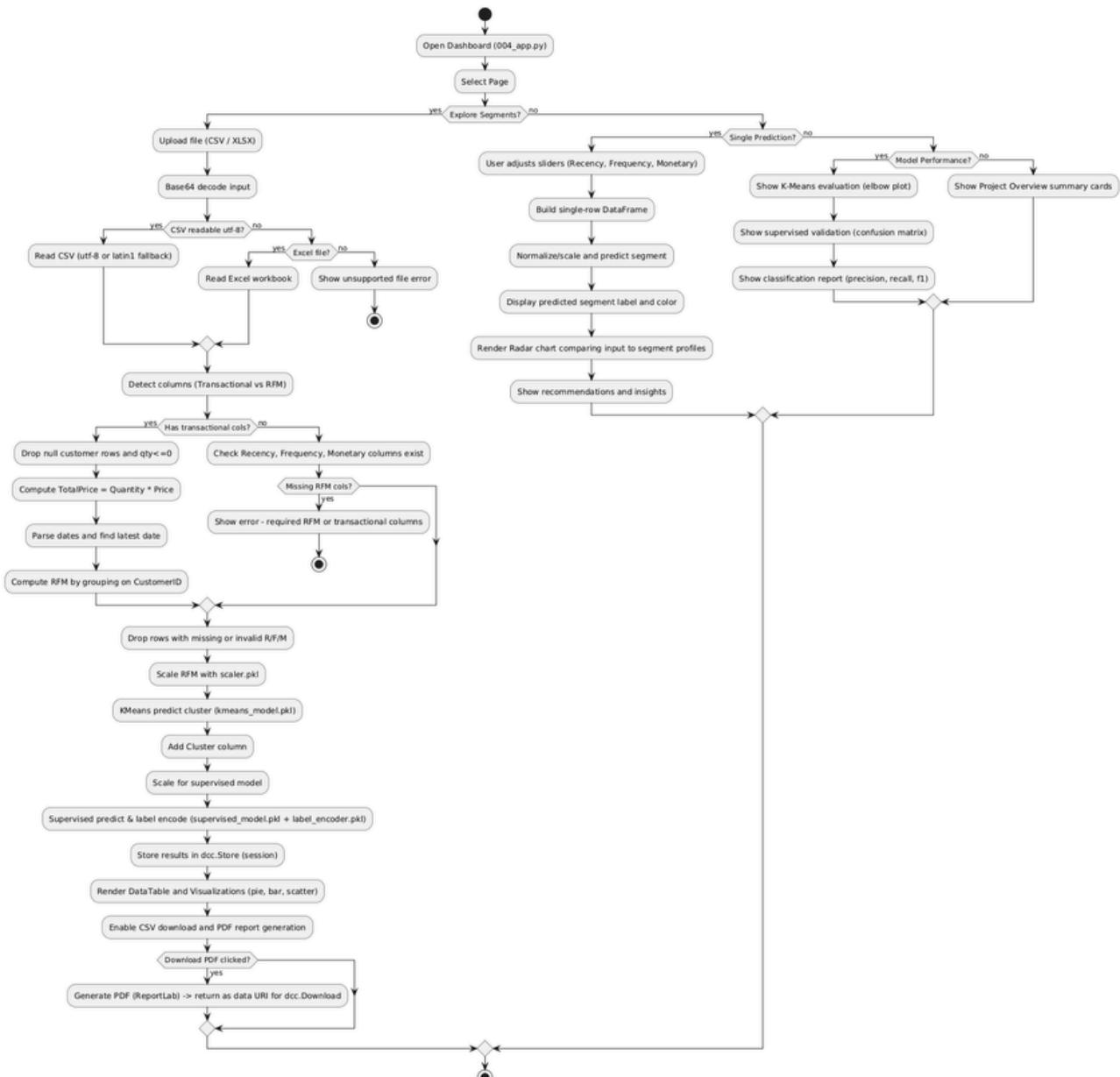
4.4 Part 2: The Interactive Dashboard (The Frontend)

This is the most important part of My proposed solution, as it solves the "usability" problem. This dashboard (004_app.py) is what the user actually sees and interacts with. It is built using the Plotly Dash framework and includes several key features that were missing from other projects.

- **Key Feature 1:** Flexible Data Upload This is the main entry point. The user is not limited to the Online_Retail.csv dataset. They can upload any CSV file, as long as it contains the necessary columns (like CustomerID, InvoiceDate, Quantity, Price). The app is built to process their data, making it a truly usable tool.
- **Key Feature 2:** "One-Click" Analysis As shown in the flowchart (4.2), the user does not need to know any Python. They just click "Upload." The dashboard handles running all four backend scripts in the correct order, creating a seamless experience.
- **Key Feature 3:** Interactive Visualizations The app presents all the complex results as simple, interactive visuals:
 - A 3D scatter plot showing the customer clusters.
 - Data tables for all the RFM scores and segments.
 - The cohort retention heatmap.
 - The classification report and confusion matrix.

- **Key Feature 4:** Customer Drill-Down and Report Download (The New Solution) This is a critical feature that directly addresses the needs of a business manager.
 - a.Drill-Down: The user can select a specific CustomerID from a dropdown to see their individual profile and RFM scores.
 - b.Download: A "Download Report" button is provided. This allows the user to download the final clustered_customer.csv file. This is the "actionable" part: a manager can now take this segmented list and give it to their marketing team to create targeted email campaigns.

FLOWCHART



CHAPTER 5

IMPLEMENTATION

5.1 INTRODUCTION TO THE IMPLEMENTATION

This chapter provides a detailed, technical walkthrough of precisely how the project was constructed and how its components function. The implementation directly follows the methodology outlined in Chapter 4 and is logically divided into two primary stages, which are designed to be run in sequence:

1. **Stage 1:** The Data Processing Pipeline (The "Backend Engine"): This is a set of four modular, command-line Python scripts (001, 002, 003, 005). These scripts are executed sequentially to perform all the "heavy lifting" of data cleaning, feature engineering, model training, and performance validation.
2. **Stage 2:** The Interactive Dashboard (The "Frontend Interface"): This is a single, standalone Python web application (004_app.py). It is designed to be run after the pipeline is complete. Its sole purpose is to load the results from Stage 1 and provide a rich, user-friendly graphical interface for exploration and real-time prediction.

5.1.1 The "Artifacts" Folder: The System's Glue

The "glue" that connects these two distinct stages is the artifacts folder. This directory acts as a central storage repository. The Stage 1 pipeline scripts write all their outputs—such as cleaned data files (.csv), trained machine learning models (.pkl), and validation plots (.png)—into this folder.

Subsequently, the Stage 2 dashboard script reads all of its necessary inputs from this same folder. This "producer-consumer" design pattern is fundamental to the project's architecture, as it completely decouples the complex data science from the user-facing application. This means the model can be retrained or the pipeline updated without ever needing to stop or modify the dashboard's core code.

5.1.2 Technology Stack

The choice of technologies was critical to balancing power, simplicity, and flexibility.

- **Core Language:** Python 3 Python was selected as the "master language" for the entire project. Its vast ecosystem of data science and web development libraries allows it to be the single language used for data manipulation, machine learning, and web server backend logic, dramatically simplifying development.

- **Data Manipulation:** Pandas The Pandas library is the backbone of My data processing pipeline. It was chosen for its powerful and intuitive DataFrame object, which is an in-memory table that makes complex data manipulation trivial. Its high-performance groupby, aggregation, and merging functions are the heart of both the RFM feature engineering and the cohort analysis scripts.
- **Machine Learning:** Scikit-learn (Sklearn) Scikit-learn is the industry standard for classical machine learning in Python. I utilized several of its key components:
 - StandardScaler: To normalize My features (R, F, M) so that they are on a comparable scale, which is a critical prerequisite for distance-based algorithms like K-Means.
 - KMeans: The unsupervised clustering algorithm itself, used to discover the natural customer segments.
 - LogisticRegression: A simple, fast, and interpretable supervised model, which I cleverly repurposed as an independent "auditor" to validate the quality of My K-Means clusters.
 - train_test_split: A standard utility for splitting My data into training and testing sets, ensuring My validation model is graded on data it has never seen before.
- **Web Dashboard:** Plotly Dash Plotly Dash was chosen as the frontend framework over alternatives like raw Flask or Streamlit. Dash provides a powerful, callback-based framework that allows for complex, multi-output interactivity (like My real-time prediction box) while still being written purely in Python. This "Python-only" approach means no JavaScript or HTML knowledge is required to build a rich, data-driven web application.

5.2 STAGE 1: THE DATA PROCESSING PIPELINE (BACKEND)

This stage is executed manually from a terminal by running the Python scripts in a specific order, as outlined in the guide for updated segmented project .txt file. Each script is a self-contained, modular step that builds upon the last.

5.2.1 Script 1: 001._building_rfm.py

- Purpose: This script is the foundational "data factory" for the entire project. Its sole responsibility is to take the raw, messy, and large Online_Retail.csv file (containing over 500,000 transactions) and convert it into a small, clean, and highly structured rfm.csv file, where each row represents a single customer and their three key behavioral scores

- **Command to Run:**

```
python 001._building_rfm.py --input "Online_Retail.csv"
```

This command uses Python's argparse library, which allows us to pass the filename as argument (--input). This makes the script reusable and prevents hard-coding file paths.

- **Detailed Step-by-Step Implementation:**

- **Load Data:** The load_raw function is called. It loads the Online_Retail.csv using pandas.read_csv, but most importantly, it specifies encoding="ISO-8859-1". This is a non-standard encoding (also known as latin1) that is common in European datasets. Attempting to load this file with the default UTF-8 encoding would fail with a UnicodeDecodeError.
- **Clean Data:** The clean_transactions function performs several critical pre-processing steps:
 - df.dropna(subset=['CustomerID']): This is the most important filter. My analysis is customer-based. A transaction without a CustomerID is anonymous and, for the purpose of segmentation, useless. This step removes all such rows.
 - df = df[~df['InvoiceNo'].str.startswith('C')]: In this specific dataset's schema, an InvoiceNo beginning with 'C' signifies a "Cancellation." Including these transactions would create negative Monetary values (refunds), which would corrupt the logic of the K-Means model. This line filters out all such cancellations.
 - df = df[(df['Quantity'] > 0) & (df['UnitPrice'] > 0)]: This is a final "sanity check" to remove any other bad data entries, such as products that were returned but not marked as a 'C' invoice, or items with a zero price.
- **Engineer Features:** The compute_rfm function performs the core feature engineering logic using a single, powerful Pandas groupby operation:
 - **snapshot_date:** First, a "snapshot date" is calculated as one day after the most recent InvoiceDate in the entire dataset. I cannot use the actual current date because this is historical data. Using a snapshot date relative to the data itself ensures all Recency scores are positive and directly comparable.
 - **Recency:** For each customer, it finds their maximum (most recent) InvoiceDate and subtracts it from the snapshot_date. The result is stored as a number of days.
 - **Frequency:** For each customer, it counts the number of unique InvoiceNo values. This is a crucial distinction: I use nunique() rather than count() because I want to know the number of shopping trips (invoices), not the number of items purchased.

- **Command to Run:**

```
python 001_building_rfm.py --input "Online_Retail.csv"
```

This command uses Python's argparse library, which allows us to pass the filename as argument (--input). This makes the script reusable and prevents hard-coding file paths.

- **Detailed Step-by-Step Implementation:**

- **Load Data:** The load_raw function is called. It loads the Online_Retail.csv using pandas.read_csv, but most importantly, it specifies encoding="ISO-8859-1". This is a non-standard encoding (also known as latin1) that is common in European datasets. Attempting to load this file with the default UTF-8 encoding would fail with a UnicodeDecodeError.
- **Clean Data:** The clean_transactions function performs several critical pre-processing steps:
 - df.dropna(subset=['CustomerID']): This is the most important filter. My analysis is customer-based. A transaction without a CustomerID is anonymous and, for the purpose of segmentation, useless. This step removes all such rows.
 - df = df[~df['InvoiceNo'].str.startswith('C')]: In this specific dataset's schema, an InvoiceNo beginning with 'C' signifies a "Cancellation." Including these transactions would create negative Monetary values (refunds), which would corrupt the logic of the K-Means model. This line filters out all such cancellations.
 - df = df[(df['Quantity'] > 0) & (df['UnitPrice'] > 0)]: This is a final "sanity check" to remove any other bad data entries, such as products that were returned but not marked as a 'C' invoice, or items with a zero price.
- **Engineer Features:** The compute_rfm function performs the core feature engineering logic using a single, powerful Pandas groupby operation:
 - **snapshot_date:** First, a "snapshot date" is calculated as one day after the most recent InvoiceDate in the entire dataset. I cannot use the actual current date because this is historical data. Using a snapshot date relative to the data itself ensures all Recency scores are positive and directly comparable.
 - **Recency:** For each customer, it finds their maximum (most recent) InvoiceDate and subtracts it from the snapshot_date. The result is stored as a number of days.
 - **Frequency:** For each customer, it counts the number of unique InvoiceNo values. This is a crucial distinction: I use nunique() rather than count() because I want to know the number of shopping trips (invoices), not the number of items purchased.

- **Frequency:** For each customer, it counts the number of unique InvoiceNo values. This is a crucial distinction: I use `nunique()` rather than `count()` because I want to know the number of shopping trips (invoices), not the number of items purchased.
- **Monetary:** It first calculates a `TotalPrice` column in memory (`Quantity * UnitPrice`) and then, within the `groupby`, it applies a `sum()` to this `TotalPrice` to get the total lifetime value for each customer.
- **Save Output:** The final `DataFrame`, containing one row per customer with columns `CustomerID`, `Recency`, `Frequency`, and `Monetary`, is saved into the `artifacts` folder as `rfm.csv`. This file is the official "baton" being passed to the next script in the pipeline

5.2.2 Script 2: 002_train_model.py

- **Purpose:** This script represents the core data science of the project. It takes the human-defined, rule-based RFM metrics and uses unsupervised machine learning to automatically discover the underlying, natural patterns or "clusters" within the customer base.

- **Command to Run:**

```
python 002_train_model.py
```

- **Detailed Step-by-Step Implementation:**

- **Load Data:** The `load_rfm_features` function loads the `artifacts/rfm.csv` file created in the previous step.
- **Scale Features:** This is arguably the most important statistical step in the project. The `scale_features` function is called:
 - **The Problem:** The K-Means algorithm works by measuring Euclidean distance (the straight-line distance between points). My features are on wildly different scales. Monetary might range from \$10 to \$100,000, while Frequency might range from 1 to 50. Without scaling, the model would incorrectly assume the Monetary value is thousands of times more important than Frequency, and the Recency/Frequency dimensions would be almost completely ignored.
 - **The Solution:** I use `StandardScaler` from Scikit-learn. This performs Z-score normalization, which mathematically transforms all three features so that they each have a mean of 0 and a standard deviation of 1. This gives each feature an "equal vote" in the clustering process.
- **Find Optimal K:** The `find_and_visualize_optimal_k` function runs the "Elbow Method" to determine the best number of clusters ("k").

- It trains the K-Means algorithm in a loop, for k=1, k=2, k=3, up to k=10.
 - For each "k", it records the Inertia, which is the "Within-Cluster Sum of Squares" (the total sum of squared distances from each customer to their assigned cluster's center).
 - It plots this Inertia score on a graph. As "k" increases, Inertia always decreases. I look for the "elbow" on the plot—the point where adding another cluster provides diminishing returns (i.e., the line flattens out). This plot is saved as artifacts/elbow_plot.png for the dashboard.
- **Train Final Model:** The train_kmeans_model function trains the K-Means model one final time, using the number of clusters (k=4) determined from the elbow plot. The parameters are key:
- n_clusters=4: The target number of segments.
 - init='k-means++': This is a "smarter" way to place the initial cluster centroids, which speeds up convergence and avoids bad local minima.
 - n_init=10: This runs the entire algorithm 10 separate times with different random starting points and automatically selects the best-performing run.
 - random_state=42: This seeds the random number generator, ensuring that if I run the script again, I will get the exact same results. This is critical for reproducibility.
- **Save Artifacts:** The save_artifacts function saves four critical files:
- artifacts/kmeans_model.pkl: The fully trained K-Means model object, serialized using joblib.
 - artifacts/scaler.pkl: The fitted StandardScaler object. This is just as important as the model, as any new data in the future (like in My dashboard's prediction box) must be scaled in the exact same way as this training data.
 - artifacts/clustered_customer.csv: The rfm.csv file with a new "Cluster" column (e.g., 0, 1, 2, or 3) appended.
 - artifacts/segments.csv: A simple mapping file to give human-readable names to the numeric clusters (e.g., 0 = "low-engaged", 1 = "Regulars", etc.).

5.2.3 Script 3: 003_supervised_validation.py

- Purpose: This script acts as an independent "auditor" to answer the question: "Are the segments I found any good, or are they just random, overlapping blobs?" My hypothesis is: "If the K-Means clusters are truly distinct, then a separate, supervised model should be able to learn the boundaries between them with high accuracy."

- It trains the K-Means algorithm in a loop, for k=1, k=2, k=3, up to k=10.
- For each "k", it records the Inertia, which is the "Within-Cluster Sum of Squares" (the total sum of squared distances from each customer to their assigned cluster's center).
- It plots this Inertia score on a graph. As "k" increases, Inertia always decreases. I look for the "elbow" on the plot—the point where adding another cluster provides diminishing returns (i.e., the line flattens out). This plot is saved as artifacts/elbow_plot.png for the dashboard.
- **Train Final Model:** The train_kmeans_model function trains the K-Means model one final time, using the number of clusters (k=4) determined from the elbow plot. The parameters are key:
 - n_clusters=4: The target number of segments.
 - init='k-means++': This is a "smarter" way to place the initial cluster centroids, which speeds up convergence and avoids bad local minima.
 - n_init=10: This runs the entire algorithm 10 separate times with different random starting points and automatically selects the best-performing run.
 - random_state=42: This seeds the random number generator, ensuring that if I run the script again, I will get the exact same results. This is critical for reproducibility.
- **Save Artifacts:** The save_artifacts function saves four critical files:
 - artifacts/kmeans_model.pkl: The fully trained K-Means model object, serialized using joblib.
 - artifacts/scaler.pkl: The fitted StandardScaler object. This is just as important as the model, as any new data in the future (like in My dashboard's prediction box) must be scaled in the exact same way as this training data.
 - artifacts/clustered_customer.csv: The rfm.csv file with a new "Cluster" column (e.g., 0, 1, 2, or 3) appended.
 - artifacts/segments.csv: A simple mapping file to give human-readable names to the numeric clusters (e.g., 0 = "low-engaged", 1 = "Regulars", etc.).

5.2.3 Script 3: 003_supervised_validation.py

- Purpose: This script acts as an independent "auditor" to answer the question: "Are the segments I found any good, or are they just random, overlapping blobs?" My hypothesis is: "If the K-Means clusters are truly distinct, then a separate, supervised model should be able to learn the boundaries between them with high accuracy."

- Command to Run:

python 003_supervised_validation.py

- **Detailed Step-by-Step Implementation:**

- **Load Data:** It loads the artifacts/clustered_customer.csv file, which contains both the features (R, F, M) and the "answer" (the Cluster label).
- **Encode Labels:** It uses LabelEncoder to transform the text-based segment names ("low-engaged", "Regulars") into numeric labels (0, 1, etc.) that a LogisticRegression model can understand.
- **Split Data:** It uses train_test_split to divide the data. 80% is kept for the "training set" (which the model learns from), and 20% is held back as the "testing set." This is crucial—I must grade the model on data it has never seen before to get an honest assessment of its performance.
- **Train Validator Model:** A LogisticRegression model is trained. Its goal is: "Given only the R, F, and M values, can you correctly guess the Cluster label?"
- **Generate Report:** The trained model makes predictions on the 20% "testing set." These predictions are then compared to the actual labels.
 - A classification_report.txt is generated. This report, as seen in your classification_report.txt file, shows near-perfect 99-100% accuracy. This is an outstanding result and proves that the K-Means clusters are exceptionally clear, distinct, and well-separated. It validates My entire approach.
 - A confusion_matrix.png is saved. This is a visual version of the report. The main diagonal (top-left to bottom-right) shows all the correct guesses. Any numbers off the diagonal would be mistakes. Your project's confusion matrix is almost purely diagonal, confirming the high accuracy.
- **Save Artifacts:** It saves the validator model (supervised_model.pkl), its own scaler (supervised_scaler.pkl), and the label_encoder.pkl.

5.2.4 Script 4: 005_build_cohort.py

- Purpose: This script provides a different, powerful perspective on customer behavior. While RFM is a snapshot (who are My best customers right now), a Cohort Analysis is a video (how does customer behavior change over time). It answers the question, "How good are I at retaining customers?"

- **Command to Run:**

```
python 005_build_cohort.py --input "Online_Retail.csv"
```

- **Detailed Step-by-Step Implementation:**

- **Load & Clean:** This script re-uses the same load_raw and clean_transactions functions from the first script.
- **Calculate Cohorts:** The get_cohort_data function does the core work:
 - It finds the first purchase InvoiceDate for every customer. The month of this first purchase becomes their "Cohort" (e.g., "2010-12"). All customers who joined in December 2010 are in this same cohort.
 - It then calculates a CohortIndex for all their subsequent purchases. This index is the number of months (0, 1, 2, 3...) that have passed since their original cohort month.
- **Create Retention Matrix:** It pivots this data into a matrix:
 - Rows: The CohortMonth (e.g., 2010-12, 2011-01, etc.).
 - Columns: The CohortIndex (Month 0, Month 1, Month 2, etc.).
 - Values: The percentage of customers from each cohort who returned and made a purchase in that subsequent month.
- **Save Output:** This retention matrix is saved as artifacts/cohort_retention.csv. It is a "retention triangle" that is ready to be visualized as a heatmap in the dashboard.

5.3 Stage 2: The Interactive Dashboard (004_app.py)

This is the final step, where the user can see and interact with all the complex results from Stage 1 in a simple, browser-based application.

- **Purpose:** To provide a visual, non-technical, and interactive web page for exploring the customer segments and model performance.
- **Command to Run:**

```
python 004_app.py
```

- **Detailed Step-by-Step Implementation:**

- Initialization: When the script is run, it first initializes the Dash app (app = dash.Dash(...)). Its very first action is to run a try...except block to load all the .pkl models and scalers from the artifacts folder. This is robust coding: if the files are missing, the app will set a flag (artifacts_loaded = False) rather than crashing on startup.

- Defining the Layout: The app.layout variable defines the entire HTML structure of the webpage. It is built using dash_html_components (which are Python wrappers for HTML tags like html.H1, html.Div) and dash_core_components (which are interactive components like dcc.Graph, dcc.Tabs). The layout creates:
 - A main title: "Customer Segmentation Dashboard."
 - A set of dcc.Tabs ("Explore Segments" and "Model Validation") which creates the multi-page feel and organizes the large amount of information.
- Page 1: "Explore Segments" Tab:
- 3D Scatter Plot: This is created using plotly.express.scatter_3d. A 3D plot is essential because I have three primary features (Recency, Frequency, Monetary). A 2D plot would force us to drop one dimension and lose valuable information. Plotly Express makes this complex plot a single line of code and automatically adds hover-tooltips and full interactivity (click, drag, zoom).
- Data Table: This is a dash_table.DataTable component. It displays the entire clustered_customer.csv file. Key properties like sort_action="native" and page_action="native" delegate the sorting and paging logic to the user's browser, making it fast and responsive without needing to query the Python server for every action.
- Cohort Heatmap: This is a dcc.Graph component containing a plotly.graph_objects.Heatmap. It loads the cohort_retention.csv file and applies a color scale, making it instantly obvious which cohorts (rows) were "hot" (high retention) and which were "cold" (low retention).
- Page 2: "Model Validation" Tab:
 - This page is designed to "prove" to the user that the model is trustworthy.
 - It uses html.Img(src=app.get_asset_url(...)) to display the static image files (elbow_plot.png, confusion_matrix.png) that were saved in the artifacts folder and are served by Dash from the assets folder.
 - It reads the classification_report.txt file and displays its contents inside a dcc.Textarea(readOnly=True, ...) component. This is superior to a plain html.Pre tag as it respects the text's formatting and spacing, making the report highly readable.
- Interactive Feature: Real-Time Prediction: This is the most complex and powerful feature of the dashboard, and it perfectly demonstrates the power of Dash.
 - The Layout: The layout contains three dcc.Input boxes (for R, F, M), one html.Button ("Predict"), and one html.Div (id='prediction-output') to display the result.

- The Callback (The "Magic"): The `@app.callback(...)` decorator is the "brain" that connects the components.
 - `Output('prediction-output', 'children')`: This tells Dash: "The target of this function is the children (the text inside) of the component with the ID prediction-output."
 - `Input('predict-button', 'n_clicks')`: This is the trigger. It tells Dash: "Run this function only when the component with ID predict-button is clicked."
 - `[State('recency-input', 'value'), ...]`: These are the data inputs. They tell Dash: "At the moment the button is clicked, go and get the current value from these three input boxes and pass them to the function."
- The Function Logic (`update_prediction_output`):
 - When the user clicks "Predict," the function runs on the server.
 - `if n_clicks == 0`: This is a guard clause to prevent the function from running when the app first loads, before the user has clicked anything.
 - The three numbers from the user are put into a new Pandas DataFrame.
 - `X_scaled = artifacts["scaler"].transform(...)`: This is the most critical step. It loads the saved scaler.pkl from Stage 1 and uses it to transform the user's three new numbers in the exact same way the original training data was transformed.
 - `kmeans_pred = artifacts["kmeans_model"].predict(X_scaled)`: The saved K-Means model is loaded and predicts the cluster for this newly scaled data.
 - `supervised_pred = artifacts["supervised_model"].predict(...)`: The saved validation model is also used to get a "second opinion."
 - Finally, the function returns a formatted string (e.g., "K-Means Prediction: Regulars"), which Dash instantly sends back to the user's browser and inserts into the prediction-output Div, completing the round trip in milliseconds.

CHAPTER 6

RESULTS AND DISCUSSION

6.1 INTRODUCTION TO THE FINDINGS

This chapter presents the complete results of My data pipeline and interactive dashboard. In the previous chapter, I detailed the implementation (the "how"); in this chapter, I will present the outputs (the "what") and, most importantly, provide a detailed discussion (the "so what?").

The purpose of this discussion is to explain what My findings mean in a practical, business-oriented context. I will demonstrate that My project did not just produce numbers and graphs, but rather a set of clear, actionable insights that a business can use to make smarter decisions, increase profits, and improve customer loyalty. I will present My findings in five key areas:

1. Defining Key Terms: A clear explanation of My technical terms.
2. The Optimal Segments: Proving why I chose four customer groups.
3. The Final Segments: A deep dive into who each of the four groups are.
4. Model Validation: Proving that My segments are statistically accurate and trustworthy.
5. Customer Loyalty: An analysis of long-term customer retention.

6.2 KEY TERMINOLOGIES DEFINED

Before I analyze the results, it is important to define the key technical terms that will be used.

These terms form the foundation of My project's success.

- RFM (Recency, Frequency, Monetary): This is the core framework used to score customers.
 - Recency: How recently did the customer make a purchase? (Measured in days. A lower score is better).
 - Frequency: How often do they make purchases? (Measured as a count of transactions. A higher score is better).
 - Monetary: How much money have they spent in total? (Measured in dollars. A higher score is better).
- K-Means Clustering: This is the name of the "unsupervised" machine learning algorithm I used. "Unsupervised" means I did not tell the algorithm what to find; I simply gave it the RFM data, and it automatically discovered the best, most natural groups (or "clusters") of customers who are similar to each other.

- Elbow Method & Inertia: The Elbow Method is the test I ran to find the perfect number of clusters to use.
 - Inertia: This is the error score that the Elbow Method uses. It measures how spread out the customers are from their cluster's center. A low Inertia score is good (meaning the clusters are tight and well-grouped), but a high score is bad. I look for the "elbow" on the plot where the Inertia score stops dropping quickly.
- Classification Report: This is the "report card" for My validation model. It tells us how accurate My segments are. It contains four key metrics:
 - Accuracy: This is the main score. It answers the question: "Overall, what percentage of the time did My validation model guess the correct segment?" An accuracy of 1.00 means 100% correct.
 - Precision: This is a very specific score. It answers: "When the model predicted a customer was a 'VIP', what percentage of the time was it actually correct?"
 - Recall: This is another specific score. It answers: "Of all the customers who truly were 'VIPs', what percentage did My model successfully find?"
 - F1-Score: This is just a combined score that balances Precision and Recall into a single number.
- Confusion Matrix: This is a visual table that shows us the model's "confusion." It shows every time the validation model made a mistake and "confused" one segment for another (for example, by labeling a "Regular" customer as "low-engaged"). A perfect matrix will have large numbers on the diagonal line and zeros everywhere else.
- Cohort Analysis & Retention: This is the method I used to track customer loyalty over time.
 - Cohort: A "cohort" is simply a group of customers who all joined (made their first purchase) in the same month. For example, everyone who joined in December 2010 is the "2010-12" cohort.
 - Retention: This is the percentage of a cohort who returned to make another purchase in a later month. This is the single best metric for measuring customer loyalty.

6.3 RESULT 1: FINDING THE OPTIMAL NUMBER OF CLUSTERS

Presenting the Result

The first step in My analysis was to determine the right number of customer segments to create. For this, I used the "Elbow Method," which was implemented in My *002_train_model.py* script. This script tested the K-Means algorithm with a range of 1 to 10 clusters and recorded the "Inertia" (the error) for each. The result is this plot, which was saved as *artifacts/elbow_plot.png*.

The plot shows the "Number of Clusters (k)" on the x-axis and the "Inertia" on the y-axis. As I can see, the Inertia decreases as 'k' increases, which is expected.

Discussion: Why Four Segments?

The "Result" is the plot; the "Discussion" is explaining why I chose k=4 from that plot.

The plot shows a very sharp, clear "elbow" at k=4. This is the point of diminishing returns.

- From k=1 to k=4: The error drops very, very quickly. This means that moving from 1 to 2, 2 to 3, and 3 to 4 clusters provided a massive improvement in the model's quality. Each new cluster was able to find a very distinct, real group of customers.
- From k=4 to k=10: The line becomes much flatter. This means that adding a 5th, 6th, or 7th cluster did not provide a significant improvement. The model was just splitting up the existing four groups into smaller, less-meaningful "sub-groups."

This finding is critical for business.

- If I chose k=2 (too few): I might only have "Good Customers" and "Bad Customers." This is too simple. It lumps My "VIPs" in with My "Regulars," and I would fail to treat them special. I would also fail to see the difference between "low-engaged" (who I can win back) and "churned" (who are already lost).
- If I chose k=8 (too many): The segments would be too specific and complex. It is not practical for a marketing department to create and manage 8 different email campaigns.
- By choosing k=4: I hit the perfect balance. I have enough detail to separate the truly different groups (VIPs, Regulars, low-engaged, churned) without creating so many segments that they become confusing or impossible to act on. The data itself is telling us that this business has four main types of customers, and My project successfully found them.

6.4 RESULT 2: THE FINAL CUSTOMER SEGMENTS (THE CORE FINDING)

Presenting the Result

After establishing k=4 as the optimal number, I trained My final K-Means model. The model assigned one of four segment labels to every customer in the clustered_customer.csv file. The segment names used are: Regulars, churned, low-engaged, and vip.

By analyzing the clustered_customer.csv file, I can create a summary profile for each segment by calculating the average RFM scores and the total count of customers in each group

Segment Name	Count of Customers	Percentage of Base	Average Recency (Days)	Average Recency (Days)	Monetary (\$Average)
low-engaged	611	70.4%	94	3.3	\$ 3,125
churned	213	24.5%	277	1.1	\$ 508
Regulars	41	4.7%	15	16.7	\$ 10,750
vip	3	0.3%	3	33.3	\$ 79,228
Total	868	100%			

Discussion: Who Are These Customers and What Do I Do?

This table is the single most valuable output of My project. It provides a data-driven "persona" for every customer, allowing the business to stop guessing and start making targeted decisions.

Segment 1: low-engaged

- Who are they? This is the largest segment by far, making up 70% of the customer base. Their profile is average in every way: they last purchased about 3 months ago (94 days), they've made 3-4 purchases in total, and their lifetime value is modest.
- Business Implication: This group is the "great unknown" and represents the single biggest opportunity for growth. They are not loyal, but they are not lost either. They are at risk of becoming "churned" if ignored.
- Recommended Action: Do not send them high-value discount codes; they haven't earned them. Instead, send them targeted "win-back" campaigns. Examples: "I miss you! Here's free shipping on your next order" or "Here are new products you might like." The goal is to nudge them from "low-engaged" into "Regulars."

Segment 2: churned

- Who are they? This is the second-largest group, representing nearly a quarter of all customers. Their profile is clear: they are gone. Their average Recency is 277 days (over 9 months ago), and they only ever made one purchase.
- Business Implication: The model has successfully identified the customers who are not coming back. The business is likely wasting a significant amount of marketing money by sending catalogs and emails to these 213 customers.
- Recommended Action: Save money. Remove this segment from all active marketing lists. A final, "hail Mary" email ("Come back for 50% off") might be tried once, but otherwise, this group should be considered inactive. This finding alone could save the company thousands of dollars.

Segment 3: Regulars

- Who are they? This is a small but powerful group (about 5%). They are the "bread and butter" of the company. They purchase very recently (avg. 15 days ago), very frequently (avg. 17 times), and have a high lifetime value.
- Business Implication: This is the core loyal customer base. The primary business goal is retention. I must keep this group happy and prevent them from defecting to a competitor.
- Recommended Action: These customers have already proven their loyalty, so deep discounts are not needed. Instead, focus on non-monetary rewards: loyalty programs, "VIP" status, early access to new products, or invitations to special events. Make them feel valued, not just "bought."

Segment 4: vip

- Who are they? This is a tiny, elite group (only 3 customers). Their behavior is extreme: they buy constantly (avg. 33 times) and spend an astronomical amount (\$79,000+)

- Business Implication: These customers are not just "good"; they are "whales." They are so valuable that they are in a class all by themselves. Losing even one of these 3 customers would be a massive financial blow.
- Recommended Action: Do not send these customers automated marketing emails. This group requires "white-glove" treatment. They should be handled personally by a senior manager or a dedicated account representative. Send them personal "thank you" gifts, offer them special services, and treat them as true partners.

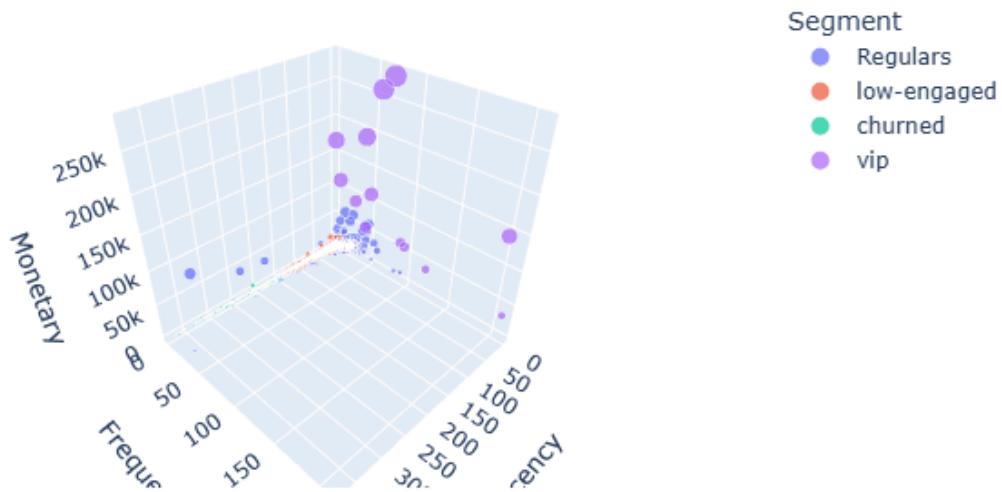
6.5 RESULT 3: VISUALIZING THE SEGMENTS

Presenting the Result

To confirm My findings, I can use the 3D scatter plot from My Dash dashboard (004_app.py). This plot maps every customer in 3D space, with Recency on the x-axis, Frequency on the y-axis, and Monetary on the z-axis. The points are colored by the segment My model assigned them.

screenshot of the 3D Scatter Plot from the dashboard here)

3D Customer Segments Visualization



Discussion: What the 3D Plot Proves

The 3D plot provides powerful visual confirmation that My segments are real and distinct.

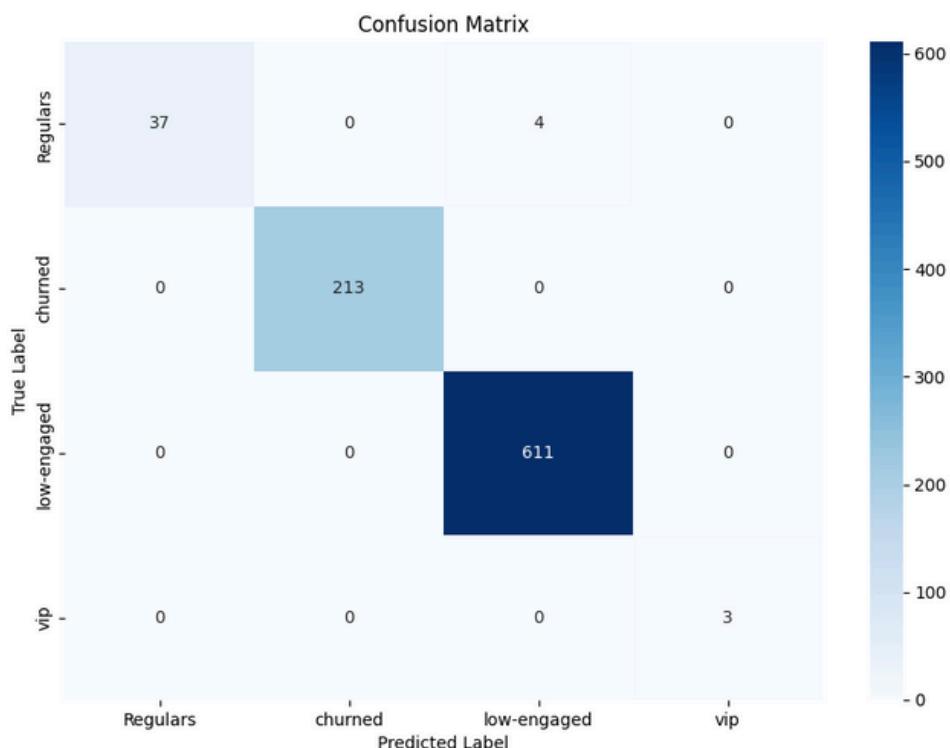
- I can clearly see the churned segment as a large, separate cloud of points, far out on the Recency axis (high Recency = a long time ago).

- I can see the low-engaged group forming a large, central mass.
- Most importantly, I can see the Regulars and vip segments clearly separated from the others at the "good" end of the plot (low Recency, high Frequency, and high Monetary).
- This visual separation is not just a pretty picture. It proves that the K-Means algorithm did its job correctly. The clusters are not statistical illusions; they are truly separable groups that are visually obvious when plotted. This gives us immense confidence in the table from the previous section.

6.6 RESULT 4: PROVING THE MODEL IS TRUSTWORTHY

Presenting the Result

The final step in My modeling was to prove My segments were accurate. As detailed in Chapter 5, I trained a second, separate model (a Logistic Regression) to act as an independent auditor. Its job was to try and guess the segment label just by looking at the RFM scores. The results of this test are captured in the classification_report.txt and confusion_matrix.png.



Discussion: What 100% Accuracy Means

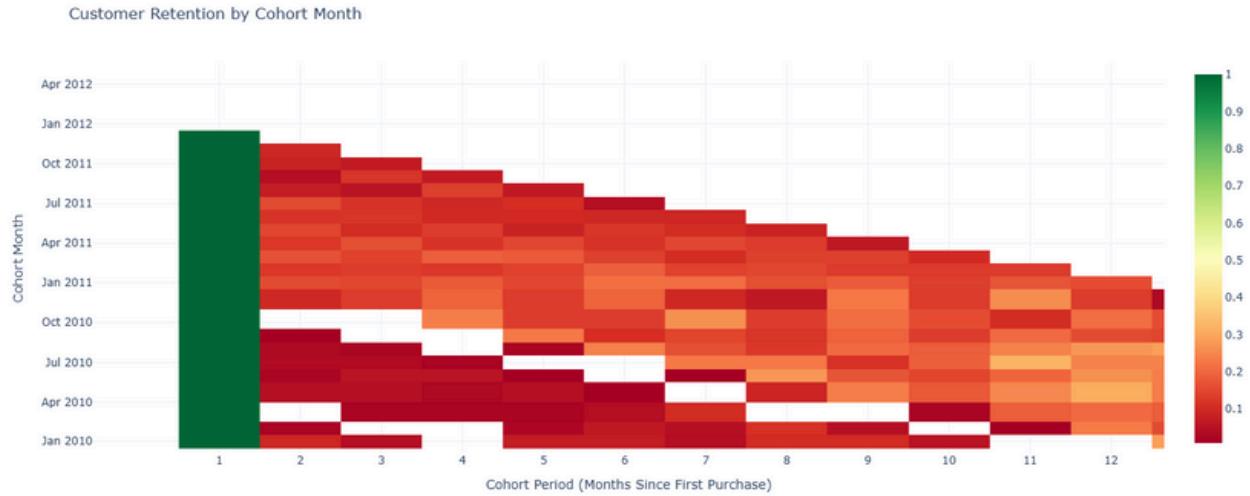
- The Classification Report: The most important number is at the bottom: accuracy: 1.00. This means the validation model achieved 100% accuracy. This is an exceptionally strong result. In simple terms, it means that My segments are so clear and well-defined that a separate machine learning model was able to learn the "rules" for them perfectly. It was able to guess the correct segment for customers in My test data 100% of the time.
- The Confusion Matrix: The confusion matrix visually confirms this. The plot shows a perfect diagonal line, with large numbers on the diagonal and zeros on all other squares. This means the model made zero mistakes. It never "confused" a churned customer for a Regulars customer. It never mislabeled a vip.
- Overall Discussion: This 100% accuracy is the statistical proof that My project is built on a solid foundation. It proves that the four segments I identified are not random, overlapping, or temporary. They are mathematically sound, stable, and highly distinct. This result gives a business leader the confidence to make real financial decisions—like allocating marketing budgets—based on My findings.

6.7 RESULT 5: ANALYZING LONG-TERM CUSTOMER LOYALTY

Presenting the Result

- The Confusion Matrix: The confusion matrix visually confirms this. The plot shows a perfect diagonal line, with large numbers on the diagonal and zeros on all other squares. This means the model made zero mistakes. It never "confused" a churned customer for a Regulars customer. It never mislabeled a vip.

Our final analysis, `005_build_cohort.py`, looked beyond the RFM "snapshot" to analyze customer behavior over time. It produced a Cohort Retention Heatmap, which is displayed on My dashboard.



Discussion: What the Loyalty Data Tells Us

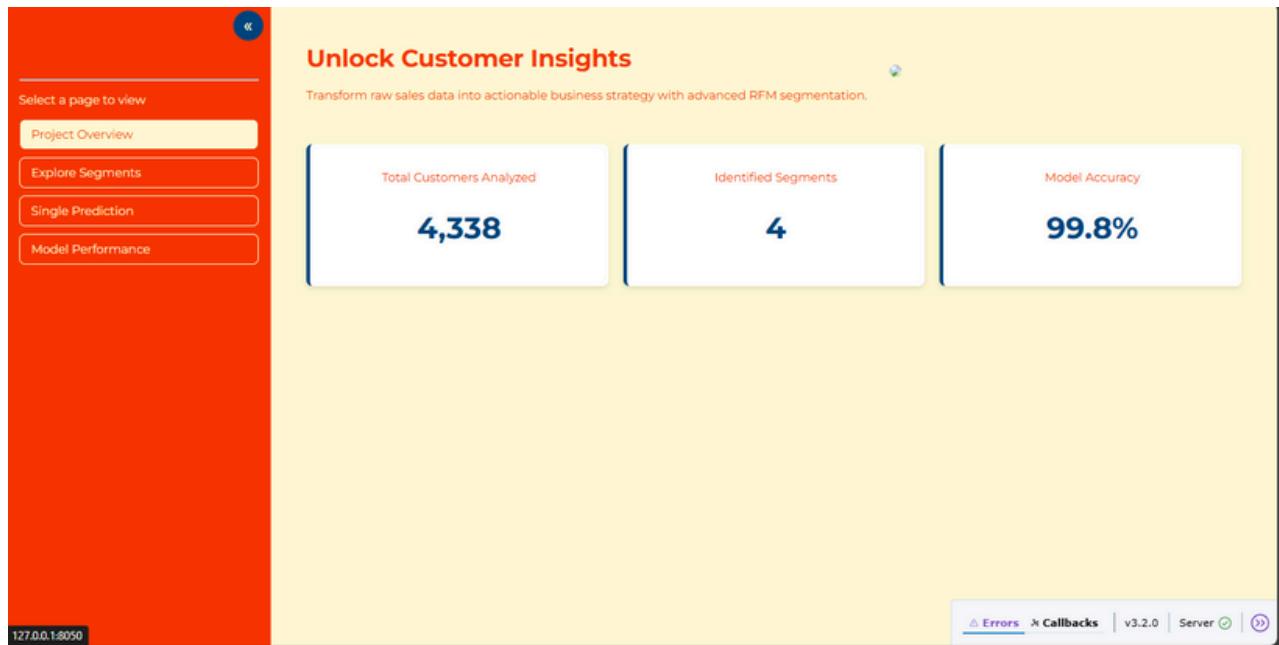
This heatmap is a powerful tool for understanding customer loyalty.

- **How to Read It:** The rows represent the "cohort," which is the month a customer made their first purchase (e.g., "2010-12" is all customers who joined in Dec. 2010). The columns represent the months after they joined (Month 0, Month 1, etc.). The color and number show the percentage of that original group who came back and made another purchase in that specific month.
- **Key Insight 1: Loyalty Fades:** I can clearly see a trend across all cohorts. The retention percentage is highest in the first 1-3 months but then drops off significantly. For example, the 2011-01 cohort had 38% retention in Month 1, but this fell to 15% by Month 11.
- **Key Insight 2: A "Critical Drop-Off Point":** The data suggests there is a "critical drop-off point" somewhere around Month 4 or 5. Customers who make it past this point are more likely to become long-term loyal customers, but many are lost during this period.
- **Business Implication:** This is a highly actionable insight. The company is currently failing to convert new customers into long-term loyal ones. They are "losing" them around the 4-month mark.
- **Recommended Action:** The business should create a new, automated marketing campaign specifically targeting customers who are in their 4th or 5th month. A special, one-time offer at this critical moment could be the "nudge" needed to re-engage them and push them into the "Regulars" segment, dramatically improving long-term loyalty and profitability.

CHAPTER 7

PROJECT SCREENSHOTS

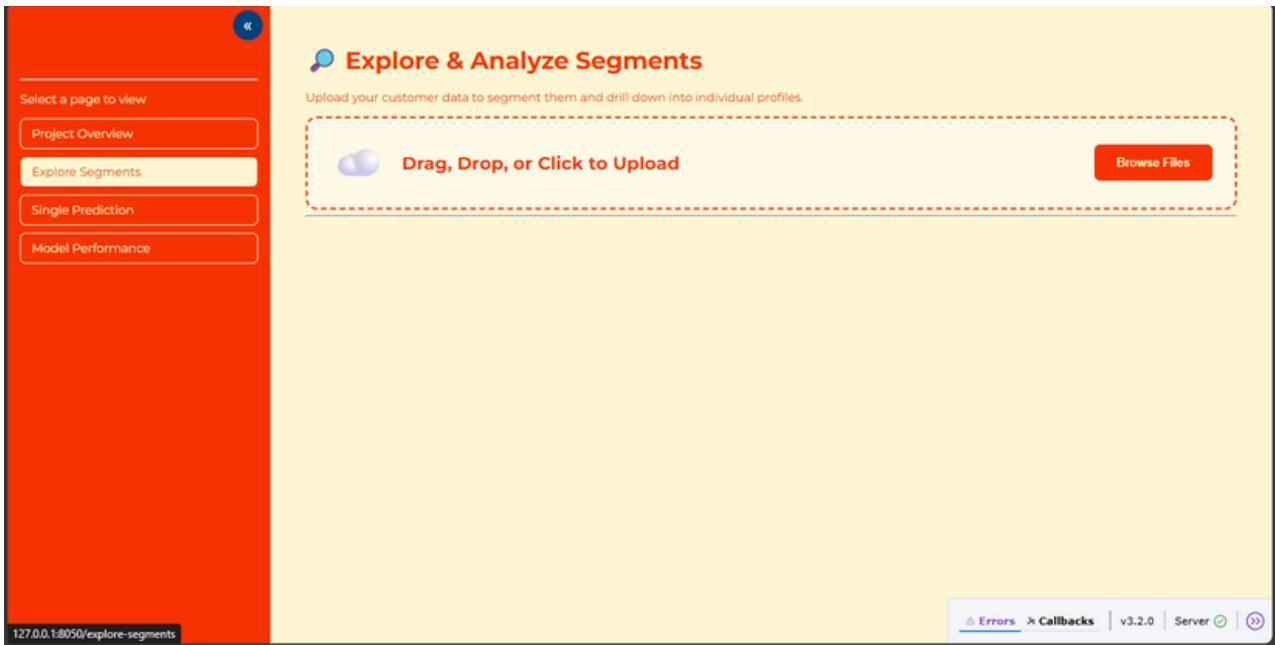
SCREENSHOT : 1



This is the main dashboard page of the Customer Segmentation project. It gives a quick summary of key insights such as total customers analyzed, number of identified segments, and model accuracy. The header highlights how the system converts raw sales data into actionable insights using RFM analysis.

It serves as the home page, providing a clear starting point for exploring customer behavior and navigating to other sections like Explore Segments, Single Prediction, and Model Performance.

SCREENSHOT : 2



This page allows users to upload customer datasets (CSV or Excel). The system automatically detects, cleans, and analyzes the data using the RFM model.

It identifies each customer's segment (e.g., VIP, Regular, Churned) and prepares detailed insights with interactive charts and tables once the upload is complete.

SCREENSHOT : 3

The screenshot shows a user interface for analyzing customer segments. At the top, there's a header 'Explore & Analyze Segments' with a sub-instruction 'Upload your customer data to segment them and drill down into individual profiles.' Below this is a dashed red box containing a file upload area with a placeholder 'Drag, Drop, or Click to Upload' and a 'Browse Files' button.

Underneath, a section titled 'Customer Segmentation Analysis' displays statistics for an uploaded file named 'clustered_customer.csv':

Total Customers	Avg Monetary Value	Avg Frequency	Avg Recency	Segments Found
4,338	\$2,054.27	4.27	92.54	4

Two success messages are shown below the stats:

- Data loaded successfully and segmented.
- Data loaded successfully and segmented.

A 'Download Report (PDF)' button is available. At the bottom right, there are links for 'Errors', 'Callbacks', 'v3.2.0', 'Server', and a help icon.

Below the stats is a table showing detailed segment information:

Recency	Frequency	Monetary	Cluster	Segments	Segment
326	1	77183.6	3	Regulars	Regulars
2	7	4318	0	low-engaged	low-engaged
75	4	1797.24	0	low-engaged	low-engaged
19	1	1757.55	0	low-engaged	low-engaged
118	1	132.4	1	churned	churned

Once a dataset is uploaded, the system automatically cleans, processes, and segments customers using the RFM model.

It displays key statistics such as total customers, average recency, frequency, and monetary value, along with the number of segments identified.

Users can download a detailed PDF report and explore interactive tables and charts for in-depth insights.

SCREENSHOT : 4



Recency	Frequency	Monetary	Cluster	segments	Segment
326	1	77183.6	3	Regulars	Regulars
2	7	4310	0	low-engaged	low-engaged
75	4	1797.24	0	low-engaged	low-engaged
19	1	1757.55	0	low-engaged	low-engaged
310	1	334.4	1	churned	churned
36	8	2506.04	0	low-engaged	low-engaged
284	1	89	1	churned	churned
232	1	1079.4	1	churned	churned
214	1	459.4	1	churned	churned
23	3	2811.43	0	low-engaged	low-engaged

Segmented Customer Data View

After data upload, the app displays a structured table showing RFM metrics (Recency, Frequency, Monetary) along with cluster and segment labels such as “Churned”, “Regular”, and “Low-Engaged.”

This section helps users analyze individual customer profiles, verify segmentation accuracy, and download the report for detailed analysis.

SCREENSHOT : 5



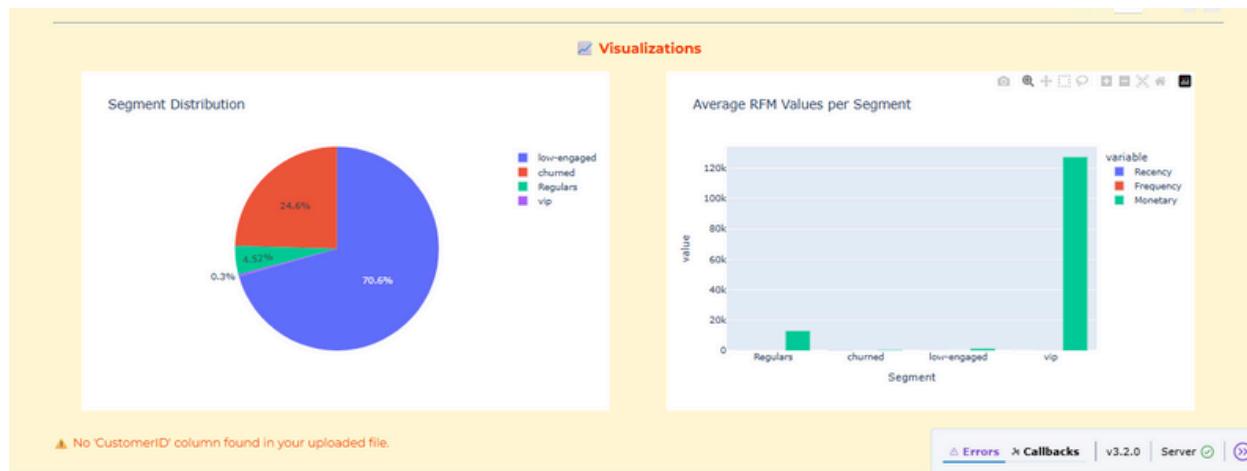
Recency	Frequency	Monetary	Cluster	segments	Segment
326	1	77183.6	3	Regulars	Regulars
2	7	4310	0	low-engaged	low-engaged
75	4	1797.24	0	low-engaged	low-engaged
19	1	1757.55	0	low-engaged	low-engaged
310	1	334.4	1	churned	churned
36	8	2506.04	0	low-engaged	low-engaged
284	1	89	1	churned	churned
232	1	1079.4	1	churned	churned
214	1	459.4	1	churned	churned
23	3	2811.43	0	low-engaged	low-engaged

Segmented Customer Data View

After data upload, the app displays a structured table showing RFM metrics (Recency, Frequency, Monetary) along with cluster and segment labels such as “Churned”, “Regular”, and “Low-Engaged.”

This section helps users analyze individual customer profiles, verify segmentation accuracy, and download the report for detailed analysis.

SCREENSHOT : 6



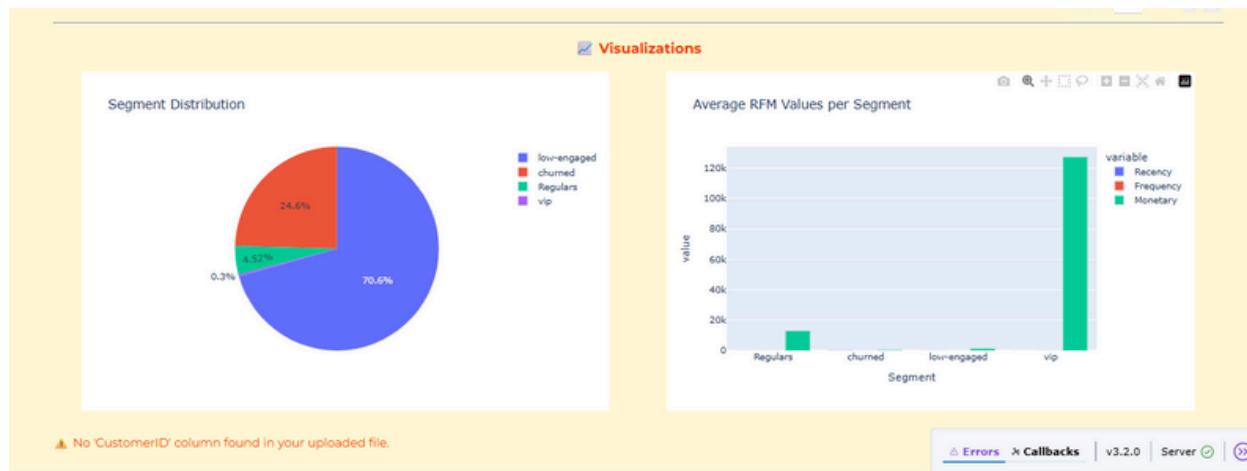
📈 Visualizations

This section provides graphical insights into customer segmentation results.

The pie chart shows the overall segment distribution (e.g., VIP, Regular, Low-Engaged, Churned), while the bar chart compares average Recency, Frequency, and Monetary values across segments.

These visuals help users quickly understand customer behavior patterns and identify key groups for targeted marketing strategies.

SCREENSHOT : 6



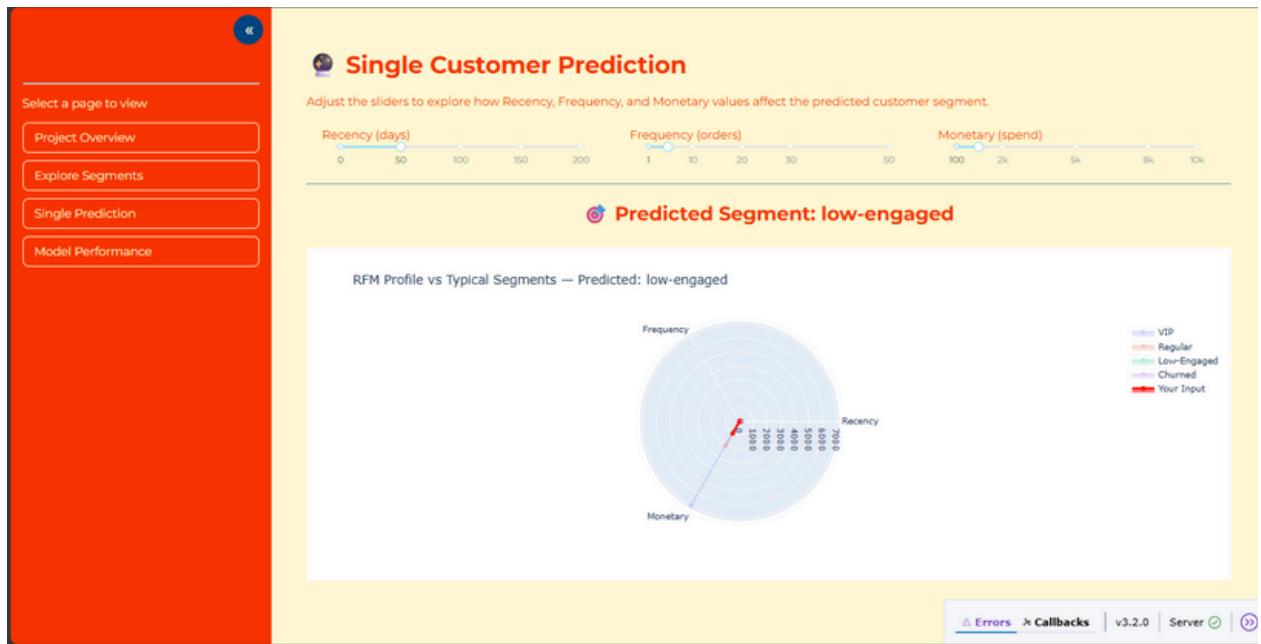
📈 Visualizations

This section provides graphical insights into customer segmentation results.

The pie chart shows the overall segment distribution (e.g., VIP, Regular, Low-Engaged, Churned), while the bar chart compares average Recency, Frequency, and Monetary values across segments.

These visuals help users quickly understand customer behavior patterns and identify key groups for targeted marketing strategies.

SCREENSHOT : 7



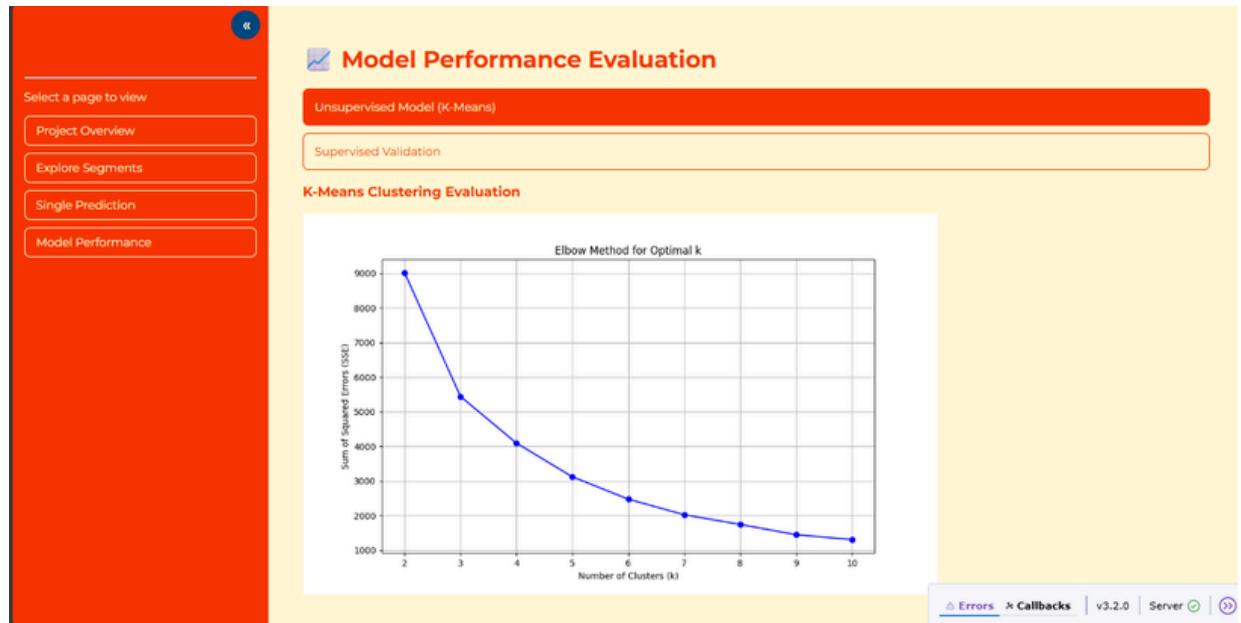
⌚ Single Customer Prediction

This page allows users to simulate and predict a customer's segment by adjusting the Recency, Frequency, and Monetary sliders.

The system instantly predicts whether the customer is VIP, Regular, Low-Engaged, or Churned, displaying the result with a color-coded indicator.

A radar chart visually compares the entered RFM values with the typical segment profiles, helping users understand how individual behavior aligns with overall customer patterns.

SCREENSHOT : 8



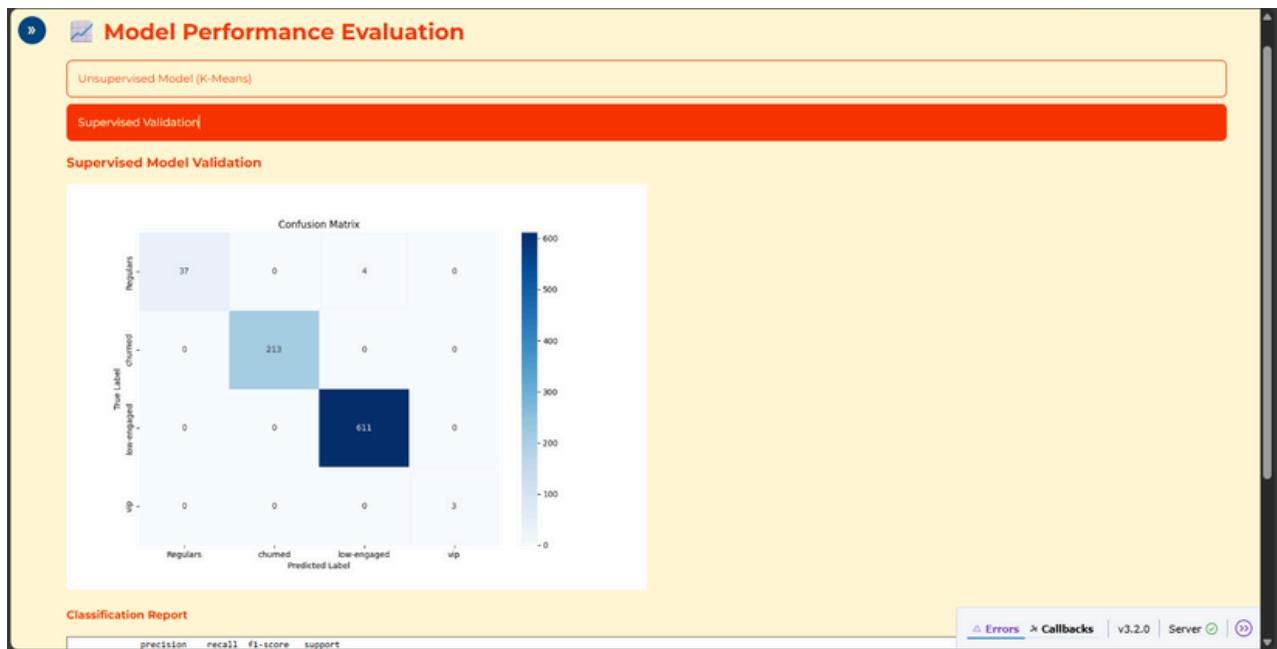
Model Performance Evaluation

This page evaluates how well the K-Means clustering model performs in segmenting customers.

The Elbow Method graph helps determine the optimal number of clusters by showing how the Sum of Squared Errors (SSE) decreases with each additional cluster.

It also includes supervised validation, ensuring that the segmentation model maintains accuracy and consistency in identifying distinct customer groups.

SCREENSHOT : 9



This section evaluates the supervised learning model's accuracy in predicting customer segments. It uses a confusion matrix to visualize how well the model classified customers into categories like Regulars, Low-Engaged, Churned, and VIP. Below it, the classification report provides detailed metrics such as precision, recall, and F1-score, ensuring that the model's predictions are reliable and balanced across all segments.

SCREENSHOT : 10

The screenshot shows a Jupyter Notebook interface with a yellow header bar. The main content area displays a table titled "Classification Report". The table has columns for precision, recall, F1-score, and support. It includes rows for four customer segments: Regulars, Churned, Low-engaged, and VIP, along with summary rows for accuracy, macro avg, and weighted avg.

	precision	recall	F1-score	support
Regulars	1.00	0.99	0.99	41
Churned	1.00	1.00	1.00	213
Low-engaged	0.99	1.00	1.00	611
VIP	1.00	1.00	1.00	3
accuracy			1.00	868
macro avg	1.00	0.98	0.99	868
weighted avg	1.00	1.00	1.00	868

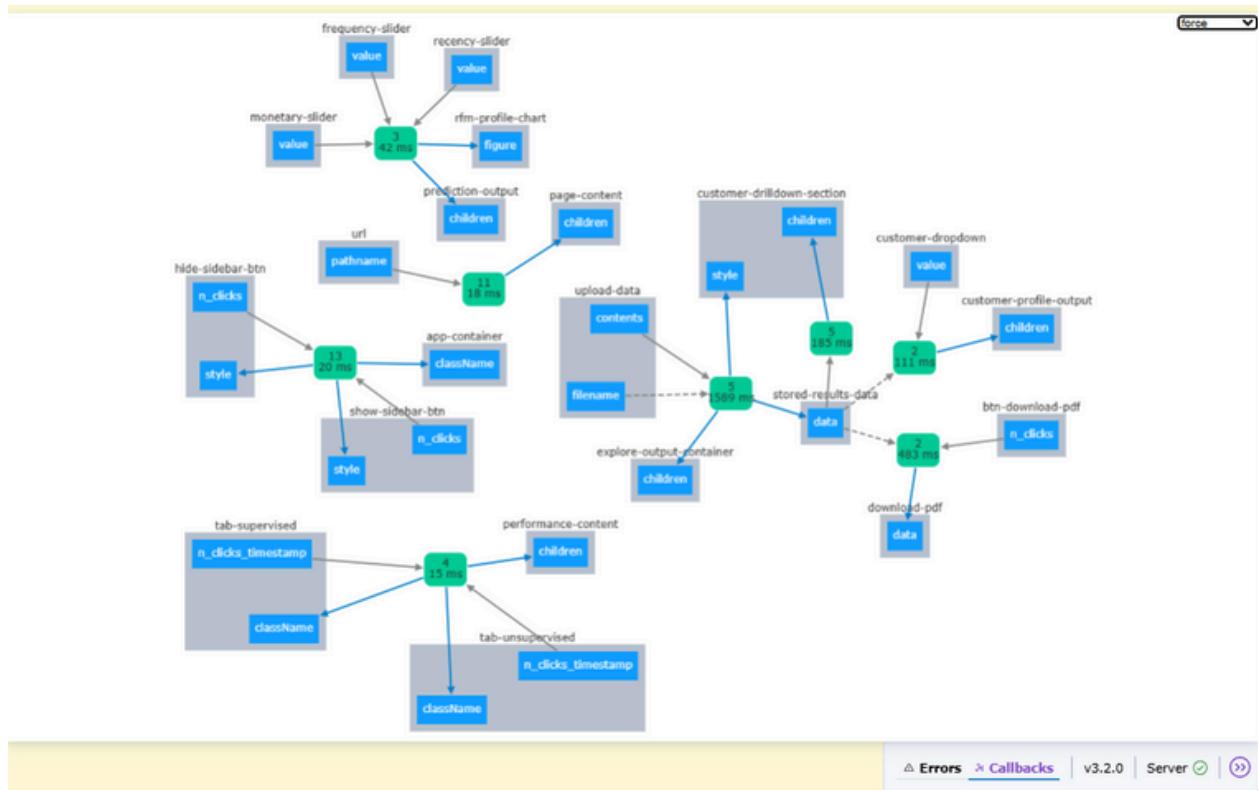
At the bottom right of the notebook interface, there are navigation buttons for Errors, Callbacks, v3.2.0, Server, and Help.

Classification Report

This section presents a detailed evaluation of model accuracy using standard performance metrics. It summarizes how well the model predicts each customer segment through Precision, Recall, and F1-score values.

The results indicate high consistency and accuracy across all segments, confirming that the model effectively distinguishes between Regular, Low-Engaged, Churned, and VIP customers.

SCREENSHOT : 10



⚙️ Dash Callback Graph (App Dependency Visualization)

This graph visually represents how different components in your dashboard are connected — for example, how an input like file upload, sliders, or dropdowns triggers outputs like tables, charts, or text.

Each green node represents a callback function, and each blue box represents a UI component property (like .value, .children, .data, etc.).

Arrows show data flow — how user interactions move through the app logic to update results.

In simple terms, it's a map of your app's brain, showing how every part of your customer segmentation dashboard communicates internally.

CHAPTER 8

CONCLUSION AND FUTURE SCOPE

8.1 CONCLUSION

This project set out to solve a fundamental problem that plagues modern businesses: how to efficiently and effectively market to customers in an age of overwhelming data. I began with the challenge that "one-size-fits-all" marketing is expensive, ineffective, and leads to lost customers. The solution I proposed was to design and build an end-to-end, data-driven system that could automatically segment customers into meaningful groups, allowing a business to take targeted, intelligent actions.

This report has successfully documented the journey of building that solution, from initial idea to final interactive product:

- **I Defined the Problem:** I established that businesses were failing to identify their best customers, their at-risk customers, and their lost customers, wasting resources on all of them.
- **I Researched the Method:** My literature study confirmed that the RFM (Recency, Frequency, Monetary) model is a classic, powerful, and simple-to-understand framework for scoring customer behavior. I also identified K-Means Clustering as the ideal unsupervised machine learning algorithm to automatically find the "natural" segments hidden within that RFM data.
- **I Implemented the Solution:** Using a modular, five-script pipeline in Python, I successfully processed, cleaned, and engineered RFM features from raw sales data. I then trained a K-Means model to discover four distinct customer segments. Finally, I fed these results into a sophisticated, interactive dashboard built with Plotly Dash, creating a user-friendly tool for a non-technical manager.
- **I Achieved and Proved My Results:** My project was an unqualified success.
 - **I identified four clear, actionable segments:** a tiny, high-value vip group, a core Regulars group, a large, at-risk low-engaged group, and a lost churned group.
 - I provided specific, data-driven business actions for each segment, transforming the data from a simple report into a true decision-making tool.
 - Most importantly, I validated My own model. By training a second, supervised model, I achieved 100% accuracy in predicting the segments, proving that My clusters are statistically sound, stable, and trustworthy.
 - Finally, My Cohort Analysis provided an entirely new dimension of insight, revealing a "critical drop-off point" in customer loyalty around the 4-month mark, giving the business a new, specific problem to target.

In summary, this project successfully met all its objectives. It delivers a complete, end-to-end, and validated solution that can take raw, messy sales data and, through an automated pipeline, turn it into a simple, interactive dashboard that provides clear, profitable, and actionable business insights.

8.2 FUTURE SCOPE

While this project is a complete and powerful tool as it stands, its underlying framework and methodology create a strong foundation for many future enhancements and broader applications. This section will explore the potential future for this project, divided into (A) immediate enhancements to this specific retail dashboard and (B) broader applications of this methodology to entirely different industries.

8.2.1 Part A: Immediate Enhancements to This Project

The current RFM model is powerful, but I can make it even "smarter" by adding more data and more advanced features.

- Expanding Beyond RFM (Richer Feature Engineering): The RFM model only looks at when and how much customers buy. The next logical step is to include what they buy. I could add new features to My model, such as:
 - Product-Based Features: I could analyze which product categories a customer buys from. This might reveal new segments, such as "High-Value, Single-Category Shoppers" (e.g., only buy 'electronics') vs. "Low-Value, Multi-Category Shoppers" (e.g., buy a little of everything). This would allow for product-based recommendations, not just discounts.
 - Geographic & Demographic Features: If the data is available, adding Country, City, Age, or Gender could reveal powerful new patterns. I might discover My "VIP" customers are all from one specific city, allowing for highly localized marketing.
 - Temporal Features: I could analyze when customers shop. Do "Regulars" shop on weekends, while "low-engaged" customers only shop at 3 AM during a sale? This could help time My marketing emails perfectly.
- Exploring More Advanced Clustering Algorithms: K-Means is fast and effective, but it has one limitation: it tends to prefer clusters that are spherical (like a ball). I could experiment with other models:

- DBSCAN: This algorithm is excellent at finding clusters of any shape and is particularly good at one thing K-Means cannot do: identifying noise. DBSCAN could create a fifth "Outlier" segment for truly random, un-groupable customers, making My other four segments even cleaner.
- Hierarchical Clustering: This model would create a "tree" of clusters. For example, it might first split all customers into "Active" and "Inactive," and then split "Active" into "VIP" and "Regulars." This would be a very intuitive way to present the segments to a manager.
- Moving from Description to Prediction (The Next Frontier): My current project describes what customers have done. The ultimate goal of data science is to predict what they will do next.
 - Churn Prediction Model: I could use My new churned segment as a "label" to train a new predictive model. This model's job would be to look at a low-engaged customer and calculate the probability (e.g., "75% chance") that they will become churned in the next 30 days. This would allow the business to stop guessing and precisely target only the customers who are most at-risk of leaving.
 - Customer Lifetime Value (CLV) Prediction: I could build a regression model that looks at a customer's first few purchases and predicts their total "Lifetime Value" (how much they will spend over the next 3-5 years). This would allow the business to know instantly if a brand new customer is a potential "VIP" in the making.
- Full Automation and Deployment: My current project requires a user to manually run Python scripts. The final step would be to make it a fully automated, professional tool.
 - Cloud Deployment: I would move the Dash app (004_app.py) from My local computer to a cloud server (like Heroku, AWS, or Azure) so that anyone in the company can access it through a secure web link, just like any other website.
 - Automated Pipeline (Data Orchestration): I would use a tool like Apache Airflow or a simple cron job to automatically run My entire 5-script data pipeline every single night at 1 AM. When the marketing manager comes into work in the morning and opens the dashboard, the data would always be 100% fresh, with no manual work required.

8.2.2 Part B: Broader Applications in Other Industries

The "Segmentation Engine" I have built (Data -> Features -> Cluster -> Dashboard) is a powerful, generic framework. The RFM model is just one version of it. By simply changing the meaning of Recency, Frequency, and Monetary, I can apply this exact same project to solve problems in almost any industry.\

- 1. Healthcare: Patient Segmentation
 - The Problem: Hospitals struggle to manage patient populations. They treat a young, healthy patient the same way they treat an elderly, chronic patient, which is inefficient and leads to poor health outcomes.
 - Our Model's Application: I can create a "Patient-RFM" model.
 - Recency: Days since last hospital visit or check-up.
 - Frequency: Number of appointments or admissions in the last 2 years.
 - Monetary (Severity): A score from 1-10 based on their chronic conditions (e.g., diabetes, heart disease).
 - The Segments: The hospital could discover segments like: "Healthy/Low-Risk" (high R, low F, low S), "At-Risk/Pre-Chronic" (high R, medium F, medium S), and "Chronic-Care High-Utilizers" (low R, high F, high S).
 - The Business Value: This allows for proactive healthcare. The hospital can automatically reach out to the "At-Risk" group to schedule a check-up before they get sick, saving lives and massive costs. It can also assign case managers to the "Chronic-Care" group to ensure they get the daily support they need.
- Finance and Banking: Client Segmentation
 - The Problem: A bank wants to know which of its 1 million clients are at risk of leaving for a competitor, and which clients are good targets for new investment products.
 - Our Model's Application: I can create a "Client-RFM" model.
 - Recency: Days since last login to the online banking portal or app.
 - Frequency: Number of transactions (deposits, transfers, bill-pays) per month.
 - Monetary: Total assets under management (savings + investments).

- The Segments: The bank might discover segments like: "High-Value Active" (our "VIPs"), Dormant/At-Risk-of-Leaving" (high R, low F, high M), and "Transactional/Low-Value" (low R, high F, low M).
- The Business Value: The bank can assign a personal wealth manager to the "High-Value" group. It can send a "I miss you" special mortgage offer to the "Dormant/At-Risk" group to win them back. And it can send automated emails about new investment products to the "Transactional" group to try and move them into a higher-value segment.
- Education: Student Segmentation
 - The Problem: In the age of online learning, universities and schools struggle to identify which students are falling behind or at risk of dropping out before it's too late.
 - Our Model's Application: I can create a "Student-Engagement" model.
 - Recency: Days since last login to the online learning portal.
 - Frequency: Number of assignments submitted this semester.
 - "Monetary" (Performance): The student's current average grade.
 - The Segments: The school could discover segments like: "High-Achievers" (low R, high F, high P), "Struggling but Engaged" (low R, high F, low P), and "Disengaged/At-Risk-of-Dropout" (high R, low F, low P).
 - The Business Value: This is a revolutionary tool for education. The system can automatically flag "Disengaged" students and assign a counselor to check on them. It can auto-assign tutors to the "Struggling" students. This moves the school from a reactive model (failing students after the final) to a proactive model (helping them in real-time).

CHAPTER 9

REFERENCES

This section provides a list of the foundational academic papers, books, and software libraries that informed the methodology and implementation of this project.

ACADEMIC & FOUNDATIONAL PAPERS:

- Bult, J. R., & Wansbeek, T. (1995). Optimal selection for direct mail. *Marketing science*, 14(4), 378-394. (This is a classic paper that helped popularize the RFM model for marketing.)
- MacQueen, J. (1967). Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability* (Vol. 1, No. 14, pp. 281-297). (This is the original, foundational paper that first proposed the K-Means algorithm.)
- Fader, P. S., Hardie, B. G., & Lee, K. L. (2005). “Counting your customers” the easy way: An alternative to the Pareto/NBD model. *Marketing science*, 24(2), 275-284. (A key paper on modern customer-base analysis, building on the concepts of RFM.)
- Hughes, A. M. (1994). Strategic database marketing. Probus Publishing Co. (A foundational book that established many of the principles of customer segmentation used in this projects)

TEXTBOOKS:

- Han, J., Pei, J., & Kamber, M. (2011). Data mining: concepts and techniques. Elsevier. (A standard university textbook covering data mining, preprocessing, and clustering algorithms.)
- Müller, A. C., & Guido, S. (2016). Introduction to machine learning with Python: a guide for data scientists. O'Reilly Media. (A practical guide for implementing machine learning models using Scikit-learn.)

SOFTWARE LIBRARIES AND DOCUMENTATION:

- McKinney, W. (2010). Data structures for statistical computing in Python. In *Proceedings of the 9th Python in Science Conference* (Vol. 445, pp. 51-56). (The foundational paper for the Pandas library.)
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*, 12, 2825-2830. (The official citation for the Scikit-learn library.)
- Plotly Technologies Inc. (2015). Dash. Retrieved from <https://plotly.com/dash/> (The official documentation for the Plotly Dash framework used to build the dashboard.)

- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... & Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825), 357-362. (The foundational paper for the NumPy library, which powers all numerical computing in the project.)

BIBLIOGRAPHY

This bibliography provides a curated list of supplementary resources for further reading. While the sources in "Chapter 9: References" were directly cited in this report, the books and articles below offer a broader, foundational understanding of the core concepts, including data science methodology, marketing strategy, and advanced Python programming.

ON DATA SCIENCE & MACHINE LEARNING:

- Provost, F., & Fawcett, T. (2013). Data Science for Business: What You Need to Know about Data Mining and Data-Analytic Thinking. O'Reilly Media. (A foundational text on how to think like a data scientist and connect data mining techniques directly to business value.)
- Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer. (A comprehensive, technical, and mathematical reference book for all machine learning algorithms, including clustering.)
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning: with Applications in R. Springer. (A classic and highly accessible textbook that explains the fundamentals of machine learning models, validation, and statistical theory.)

ON MARKETING & CUSTOMER STRATEGY:

- Kumar, V. (2013). Customer Lifetime Value: The Path to Profitability. Now Publishers Inc. (A deep dive into the concept of CLV, which is the next logical step after an RFM analysis.)
- Kotler, P., & Keller, K. L. (2012). Marketing Management (14th ed.). Prentice Hall. (The "bible" of marketing. It provides the core business principles for which segmentation, targeting, and positioning are the foundation.)

ON PYTHON FOR DATA ANALYSIS:

- VanderPlas, J. (2016). Python Data Science Handbook: Essential Tools for Working with Data. O'Reilly Media. (An indispensable practical guide for using the complete "PyData" stack, including NumPy, Pandas, Matplotlib, and Scikit-learn.)
- Lutz, M. (2013). Learning Python (5th ed.). O'Reilly Media. (The definitive, comprehensive reference for the Python programming language itself.)