

Six weeks industrial training project report on

“Streamlytix”

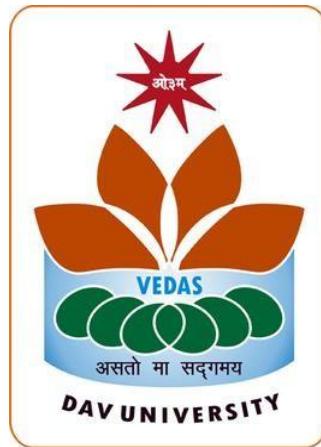
Submitted in the partial fulfilment of the requirement for the award of degree of

Bachelor of technology

In

Computer Science and Engineering

Batch (2022-2026)



Submitted to:

Mr. Om Dang

Assistant Professor (Department of CSE)

Submitted by:

Aniket Taneja

12200351

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
DAV UNIVERSITY JALANDHAR-PATHANKOT NATIONAL
HIGHWAY, NH 44, SARMASTPUR PUNJAB**

1440012

ACKNOWLEDGEMENT

I express my gratitude to all those who helped us in various stages of the development of this project. First, I would like to express my sincere gratitude to Dr. Manoj Kumar (Vice Chancellor), Dr. Rahul Hans (HOD), Mr. Om Dang (Assistant Professor) of DAV University for allowing me to undergo the summer training of 45 days at **O7 Services, Jalandhar**. I am also thankful to all faculty members of the Department of Computer Science and Engineering, for their true help, inspiration and for helping me for the preparation of the final report and presentation. Last but not least, I pay my sincere thanks and gratitude to all the Staff Members of the **O7 Services, Jalandhar** for their support and for making our training valuable and fruitful.

DECLARATION

I, **ANIKET TANEJA**, hereby declare that the work which is being presented in this project/training titled "**Streamlytix**" by me, in partial fulfilment of the requirements for the award of Bachelor of Technology (B.Tech) Degree in "Computer Science and Engineering" is an authentic record of my own work carried out under the guidance of Dr. Rahul Hans & Mr. Om Dang. To the best of my knowledge, the matter embodied in this report has not been submitted to any other University/ Institute for the award of any degree or diploma.

Aniket Taneja

12200351

Mr.Om Dang

Assistant Professor

(Department of CSE)

CERTIFICATE

Ref. No. OTS/61N24/336



Certificate

OF COMPLETION



THIS CERTIFICATE IS PROUDLY PRESENTED TO
ANIKET TANEJA

FOR SUCCESSFULLY COMPLETING TRAINING COURSE IN
SPECIALIZATION IN REACT JS & FIREBASE



FROM JUNE 2024 TO JULY 2024

ISO-9001:2015

Aniket
Instructor's Signature



Director's Signature

COMPANY PROFILE



O7 Services is an established company founded in 2015 and authorised by the government. They specialise in a variety of services including Web Development, Mobile Application Development, Custom Software Development, UI/UX Designing, Hosting services, Digital Marketing, Registration of Domain Names with various extensions, AMC & MMC Services, Bulk SMS and voice calls. O7 Services offers advanced IT solutions supporting the entire business cycle - from consulting to system development, deployment, quality assurance, and 24x7 support. With over 10 years of experience, the company aims to form long-lasting strategic partnerships with clients, offering affordable prices, timely delivery, and measurable business results. Their headquarters is located in Jalandhar with a branch office in Hoshiarpur. Some of the products developed by O7 Services include Vehicle Tracking System, Invoice Software, School Management System, Hospital Management System, Parents-Teacher Communication App, Fee Management system, Task Management System, Online Food Ordering App, Security App, Admission system, Inventory Software, and Car Servicing App. Additionally, O7 Services provides training programs including 6 Weeks/6 Months Industrial Training, Project-Based Training, Corporate Training, and Job-Oriented Courses Training, covering major IT trends such as Full Stack Development (MEAN/MERN), Flutter, Kotlin Android, Swift UI (iOS), Firebase, Python, Angular, React Js, Vue Js, Node Js, ASP.NET, .NET Core, PHP, Laravel, CodeIgniter, Software Testing, Cloud Computing, Blockchain, DevOps, Data Science, Artificial Intelligence, Machine Learning, UI/UX Designing, Digital Marketing, WordPress, Linux, CCNP, CCNA Security, Network Security, Cyber Security, MCSE, MCITP, Java, Spring, Hibernate, C/C++, Photoshop, Adobe Illustrator, Figma, CorelDraw and many more

+91- 8437365007, +91-181-5015007

E-Mail: enquiry@o7services.com , hr@o7services.com

www.o7services.com

DAV UNIVERSITY

TABLE OF CONTENT

S.NO.	CONTENT	PAGE NO.
1.	INTRODUCTION 1.1 Introduction 1.2 Project description 1.2.1 User Intercation Module 1.2.2 Sentiment Analysis Module 1.3 Problem definition 1.4 Existing system 1.5 Proposed system	1-4
2.	SYSTEM REQUIREMENTS 1.1 Hardware requirements 2.2 Software requirements	5
3.	FEASIBILITY STUDY 3.1 Introduction 3.2 Types of feasibility study	6-7
4.	SYSTEM ANALYSIS 4.1 Introduction 4.2 Key aspects	8
5.	TECHNOLOGIES USED 5.1 Python 5.1.1 Introduction 5.1.2 Features 5.2 Pandas 5.2.1 Introduction 5.2.2 Features 5.2.3 Advantages 5.3 Numpy 5.3.1 Introduction 5.3.2 Features 5.3.3 Advantages 5.4 Matplotlib 5.4.1 Introduction 5.4.2 Features 5.4.3 Advantages 5.5 Plotly Express 5.5.1 Introduction	9-33

	<p>5.5.2 Features 5.5.3 Advantages</p> <p>5.6 Seaborn 5.6.1 Introduction 5.6.2 Features 5.6.3 Advantages</p> <p>5.7 Dataframes 5.7.1 Introduction 5.7.2 Features 5.7.3 Advantages</p> <p>5.8 Datasets 5.8.1 Introduction 5.8.2 Features 5.8.3 Advantages</p> <p>5.9 Streamlit 5.9.1 Introduction 5.9.2 Features 5.9.3 Advantages</p> <p>5.10 Scikit Learn 5.10.1 Introduction 5.10.2 Features 5.10.3 Advantages</p>	
6.	<p>SOFTWARE PROCESS MODEL</p> <p>6.1 Waterfall model 6.1.1 Introduction 6.1.2 Why use SDLC waterfall model 6.1.3 Advantages of waterfall model 6.1.4 Disadvantages of waterfall model</p>	34-36
7.	<p>DESIGN</p> <p>7.1 Introduction 7.2 Important components of system design 7.3 Common system design principles</p>	37-38
8.	<p>DATA FLOW DIAGRAMS</p> <p>8.1 Introduction 8.2 Important rules 8.3 Some standard DFD symbols 8.4 Levels in DFDs 8.5 DFD level 0 - Context level diagram 8.6 DFD level 1 8.6.1 DFD for admin</p>	39-40

	8.6.2 DFD for user	
9.	SNAPSHOTS	41-45
10.	CONCLUSION AND FUTURE SCOPE 10.1 Conclusion 10.2 Future scope	46-47
11.	REFERENCES 11.1 References	48

Chapter 1

INTRODUCTION

1.1 Introduction

In the evolving landscape of digital content consumption, understanding user behavior and preferences through data-driven approaches has gained significant relevance. With the widespread availability of movies and web series on streaming platforms, users engage with large volumes of multimedia content on a regular basis. However, there is limited availability of recommendations and analytical feedback outside the streaming platforms which can only either be accessed by paying subscriptions or doesn't have all of the content at one place.

Streamlytix is a Python-based data science, machine learning application that aims to bridge this gap by providing personalized analysis and recommendations based on the streaming content marked as watched by individual users. Using an intuitive Streamlit-based interface, users can select and submit the movies and series they have watched. The application then performs analysis on the selected content, provide user with 5 similar movies/TV shows to watch next. These insights are presented through graphs, summaries, and visual reports etc.

By combining natural language processing, data visualization, machine learning algorithms and a user-centric design, Streamlytix introduces a novel approach to understanding streaming habits.

1.2 Project Description

The Streamlytix platform is developed using Python and includes several coordinated modules, each designed to handle a specific aspect of user interaction, data processing, and sentiment visualization.

1.2.1 User Interaction Module:

This module allows users to interact with the system via a web-based frontend developed using Streamlit. It includes the following key features:

- **Watched Content Input:** Users can browse or search from a curated list of movies and series (sourced from the TMDB dataset) and mark the titles they have watched.
- **Session Handling:** The application stores the user's selections temporarily during the session for personalized analysis.
- **Ease of Use:** The interface is minimal and user-friendly, with login and register functionality to save the user's data.

1.2.2 Analysis Module:

This module is responsible for analyzing the content that the user has marked as watched. It operates as follows:

- **Metadata Mapping:** Each title is matched with associated metadata.
- **Content Evaluation:** The content is evaluated based on the above meta data.
- **Mood Profiling:** The results of all watched items are aggregated to generate a user-specific preference profile, showing dominant trends in their viewing history.

1.2.3 Visualization And Reporting Module:

This module focuses on presenting the personalized insights in a clear and engaging manner through visual outputs:

- **Graph Generation:** The application uses data visualization libraries like Matplotlib, Plotly, and Seaborn to create informative bar charts, pie charts, and line graphs.
- **Behavioral Summary:** A final report is generated, summarizing the types of content watched, their sentiment distribution, and other behavioral observations based on the user's input.

- **Interpretation:** These visuals allow the user to reflect on their content preferences and emotional tendencies associated with their viewing habits.

1.2.4 Recommendations Module:

- **Metadata processing:** The metadata related to the title is analyzed, processed and relevant columns are combined to form a new column.
- **Vectorization :** The words are processed, vectorized and plotted in a graph.
- **Machine learning algorithms:** Machine learning algorithms are used to find similar movies/TV shows.

1.3 Problem Definition

Existing streaming and movie platforms focus primarily on content discovery, watchlist management. They lack the capability to provide user-specific feedback about the content that has already been consumed and provide recommendations based on the content. This creates a gap for viewers who wish to understand about their streaming habits and want to find something similar to watch.

Some key limitations of current systems include:

- Absence of any analysis for watched content.
- No support for generating personalized summaries.
- Lack of meaningful visual reports showing content-based breakdown.
- Limited user interaction beyond marking content as watched or giving ratings
- Limited recommendations for the content watched.
- Needing subscription for using the platforms.

These limitations prevent users from gaining deeper insights into their viewing behavior.

1.4 Existing System

Content platforms like TMDB, IMDb, and watchlist trackers allow users to rate or mark media as watched. However, they do not provide personalized analysis or personalized recommendation. Limitations in the existing systems include:

- No categorization of consumed content.
- No visualization of user trends.
- Minimal analytical reporting capabilities.
- Lack of a dedicated interface for personalized recommendations or feedback.

These systems serve functional purposes but do not offer the added value of data-driven insights based on content consumption.

1.5 Proposed System

Streamlytix is designed to overcome the limitations of current content tracking systems by offering a platform that analyzes and interprets watched content and provide recommendations. The focus is on providing each user with a unique, personalized analysis report based on their own streaming data and providing user with more watch options similar what they have already watched. The key objectives of the proposed system are as follows:

- Enable users to mark watched content from a curated TMDB dataset using a user-friendly interface.
- Perform analysis on the selected content.
- Generate visual reports summarizing the user's content consumption patterns.
- Provide personalized recommendations based on users watch history.
- 100% free to use and no-ads on the platform with interactive UI.

The system offers an innovative way for users to reflect on their multimedia engagement and gain a better understanding of their content-based emotional footprint.

Chapter 2

SYSTEM REQUIREMENTS

2.1 Hardware Requirements

- **Processor:** Intel Core i3 or higher (e.g., Intel Core i5, Intel Core i7)
- **Ram:** 4 GB
- **SSD:** 256GB

2.2 Software Requirements

- **Front End:** Streamlit
- **DB Tool:** Cloud based PostgreSQL DB (neonDB)
- **Browser:** Mozilla Firefox/Chrome/Edge or any other relevant browser
- **OS:** Windows operating system/Linux
- **Text Editor:** Visual Studio, Jupyter Notebook, Google COLAB

Chapter 3

FEASIBILITY STUDY

3.1 Introduction

A feasibility study is an analysis and evaluation of the practicality, viability, and potential success of a proposed project or business practice. It involves assessing various technical, economic, legal, operational, and scheduling considerations, to determine whether the project is feasible and worth pursuing. The goal of the feasibility study is to provide stakeholders with insights into potential risks, costs, and benefits, helping them to decide whether they should proceed with the project or not. This study is the foundation of any project. It helps answer the question "Is the project feasible or not?"

3.2 Types Of Feasibility Study

There are five types of feasibility studies and those types are as follows:

- **Technical Feasibility:** Technical feasibility assesses the technical complexity of an expert system, including its potential implementation using different techniques and tools. An important aspect of technical feasibility involves selecting the appropriate development shell because the choice of the shell impacts the system's quality and is necessary for its success. While the ideal characteristics of an expert system shell may vary based on the task and domain requirements, it should possess the flexibility to incorporate expert reasoning effectively. Moreover, a user-friendly interface in the shell promotes increased usage by end users, contributing to the system's overall success.
- **Economic Feasibility:** The process of collecting information about a proposed project and then analysing it is necessary to determine the feasibility of implementation, mitigate risks, and check the project's profitability. Understanding the project's potential success or loss in comparison to the local market and its demands is important. Moreover, evaluating the company's capacity to sustain profitability over a specified period is necessary. This detailed examination is the economic feasibility study.

- **Legal Feasibility:** A legal feasibility assessment involves examining a proposed restructuring plan for potential legal issues and formulating an implementation strategy with a focus on tax considerations. Tax implications are often a critical factor in an International Business Restructuring (IBR) project, and legal teams play a key role in assessing and addressing these concerns. Once the restructuring plan is well-defined, legal teams conduct a thorough review for any potential red flags. The findings from this analysis are incorporated back into the plan, enhancing its effectiveness. In certain cases, the legal team may also conduct traditional legal due diligence as part of the overall feasibility evaluation. An example of a feasibility assessment framework in project management is TELOS, which encompasses five areas of consideration.
- **Operational Feasibility:** Operational feasibility involves assessing the extent to which a service meets requirements and the ease of operating and maintaining the product post-deployment. This analysis includes evaluating the usability of the product and determining whether the suggested solution proposed by the software development team is acceptable. Other operational aspects, such as the practicality of implementing the suggested solution, are also considered in this evaluation.
- **Scheduling Feasibility:** This evaluation is important for the success of a project since timely completion of the projects is necessary. Scheduling feasibility involves estimating the time required to complete the project. Once these aspects are thoroughly examined, the feasibility analysis helps pinpoint potential constraints that the proposed project may encounter, encompassing internal project constraints such as technical, technology, budget, resource, etc., internal constraints of an organisation, i.e., financial, marketing, export, etc., and external constraints including logistics, environment, laws, and regulations.
- **Environmental Feasibility:** An environmental feasibility study tests the practicality of a project. It checks the project's environmental and social impact and potential challenges and risks associated with it.

Chapter 4

SYSTEM ANALYSIS

4.1 Introduction

In the areas of science, information technology, and knowledge, the difficulty of systems is of much importance. As systems became more complicated, the traditional method of problem-solving became inefficient. System analysis is to examine a business problem, identify its objectives and requirements, and then design the most optimal solution to fulfil those needs.

4.2 Key Aspects

It is the very first step in any system development and the critical phase where developers come together to understand the problem, needs, and objectives of the project.

Some of the key aspects of system analysis are:

1. **Problem Identification:** It involves identifying the issues that the system is aiming to address. Whether it is automating a business process, improving data management, or improving the user experience, understanding the problem is the first and most important step.
2. **Requirements Gathering:** Once the problem is identified, the next step is to gather and write down the requirements. This involves communicating with the customer and developer to gather information about how the system is to be designed.
3. **Feasibility study:** Before going into development, it is important to check the feasibility of the project. This includes the evaluation of technical, operational, and financial aspects to determine the feasibility of the proposed solution.
4. **Analysis and modelling:** To get a deep insight into the system, analysts develop various models, such as Data Flow Diagrams(DFD), Use Cases, and Entity-Relationship(ER) diagrams. These models help the customer to visualise the system and its interactions.

Chapter 5

TECHNOLOGIES USED

5.1 Python

5.1.1 Introduction

Python is a powerful, high-level programming language known for its readability and simplicity. It follows the object-oriented programming paradigm, which means it's organized around objects rather than actions, making it intuitive and efficient for developers.

Python's design philosophy emphasizes code readability and simplicity, allowing developers to write clear, logical code for small—and large-scale projects. As a high-level language, Python abstracts away much of the complexity involved in programming, enabling developers to focus on solving problems rather than worrying about underlying technical details.

Python is at the core of many technologies and applications we use daily. For instance, YouTube uses it for video processing and search engines to handle vast amounts of data.

5.1.2 Features

1. Python is versatile and flexible

Python is a general-purpose language, which means it can be used to create a wide variety of applications. From web development to data analysis, from artificial intelligence to scientific computing, Python's versatility is unmatched.

2. Python is simple and easy to learn

Python's simple and clean syntax makes it an ideal language for beginners. Its commands are English-based, and its straightforward layout helps new programmers understand code easily. This simplicity also makes Python suitable for rapid development and prototyping, reducing the time it takes from concept to implementation.

3. Python is open-source

Python's open-source nature has led to the development of a vast ecosystem of libraries and frameworks. Whether you need tools for web development (Django, Flask), data analysis (pandas, NumPy), machine learning (TensorFlow, scikit-learn), or any other task, Python has a library for it.

4. Python has strong community support

Python boasts a large and active community of developers contributing to its continuous improvement. This community support means that there are countless tutorials, forums, and documentation available to help newcomers and experienced developers alike.

5. Python is used everywhere

Python's widespread use in various industries makes it a valuable skill for developers. Companies across the globe, from tech giants like Google and Facebook to financial institutions like JP Morgan Chase, rely on Python for their technological solutions.

6. Python is in continuous evolution

Python is continually evolving to meet the needs of modern developers. Recent versions, such as Python 3.10 and 3.11, have introduced significant performance improvements and new features, keeping the language relevant and efficient.



Figure 5.1

5.2 Pandas

5.2.1 Introduction

pandas is a data manipulation package in Python for tabular data. That is, data in the form of rows and columns, also known as DataFrames. Intuitively, you can think of a DataFrame as an Excel sheet.

pandas' functionality includes data transformations, like sorting rows and taking subsets, to calculating summary statistics such as the mean, reshaping DataFrames, and joining DataFrames together. pandas works well with other popular Python data science packages, often called the PyData ecosystem, including

5.2.2 Features

Pandas is used throughout the data analysis workflow. With pandas, you can:

- Import datasets from databases, spreadsheets, comma-separated values (CSV) files, and more.
- Clean datasets, for example, by dealing with missing values.
- Tidy datasets by reshaping their structure into a suitable format for analysis.
- Aggregate data by calculating summary statistics such as the mean of columns, correlation between them, and more.
- Visualize datasets and uncover insights.

5.2.3 Advantages

- **Made for Python:** Python is the world's most popular language for machine learning and data science.
- **Less verbose per unit operations:** Code written in pandas is less verbose, requiring fewer lines of code to get the desired output.

- **Intuitive view of data:** pandas offers exceptionally intuitive data representation that facilitates easier data understanding and analysis.
- **Extensive feature set:** It supports an extensive set of operations from exploratory data analysis, dealing with missing values, calculating statistics, visualizing univariate and bivariate data, and much more.
- **Works with large data:** pandas handles large data sets with ease. It offers speed and efficiency while working with datasets of the order of millions of records and hundreds of columns, depending on the machine.



Figure 5.2

5.3 NumPy

5.3.1 Introduction

NumPy(Numerical Python) is a fundamental library for Python numerical computing. It provides efficient multi-dimensional array objects and various mathematical functions for handling large datasets making it a critical tool for professionals in fields that require heavy computation.

5.3.2 Features

NumPy has various features that make it popular over lists.

- **N-Dimensional Arrays:** NumPy's core feature is ndarray, a powerful N-dimensional array object that supports homogeneous data types.
- **Arrays with High Performance:** Arrays are stored in contiguous memory locations, enabling faster computations than Python lists (Please see Numpy Array vs Python List for details).
- **Broadcasting:** This allows element-wise computations between arrays of different shapes. It simplifies operations on arrays of various shapes by automatically aligning their dimensions without creating new data.
- **Vectorization:** Eliminates the need for explicit Python loops by applying operations directly on entire arrays.
- **Linear algebra:** NumPy contains routines for linear algebra operations, such as matrix multiplication, decompositions, and determinants.

5.3.3 Advantages

Following are the main advantage of NumPy:

- **Efficient array operations:** NumPy arrays are highly optimized for numerical computations. They allow for element-wise operations such as addition,

subtraction, multiplication, and division, as well as more complex operations like dot products and matrix inversions. The underlying C code of NumPy makes these operations much faster than equivalent Python code.

- **Memory efficiency:** NumPy arrays are stored in contiguous blocks of memory, which makes them more efficient in terms of memory usage. They also provide a way to view and manipulate memory blocks as arrays of different shapes and sizes, without the need for copying data. This makes NumPy ideal for working with large datasets.
- **Broadcasting:** NumPy allows for broadcasting, which is a powerful feature that enables the use of arrays with different shapes and sizes in arithmetic operations. Broadcasting avoids the need for explicitly resizing or reshaping arrays, which can be a time-consuming process.
- **Universal functions:** NumPy provides a set of universal functions (ufuncs) that operate element-wise on arrays, such as sin, cos, and exp. These functions are highly optimized and can operate on entire arrays without the need for a loop. They are also compatible with other NumPy functions, making it easy to chain operations together.
- **Interoperability with other libraries:** NumPy is designed to work seamlessly with other scientific computing libraries in Python, such as SciPy, Matplotlib, and Pandas. This allows for a wide range of advanced computations and data visualization tasks.



Figure 5.3

5.4 Matplotlib

5.4.1 Introduction

Matplotlib is a popular 2D plotting library in Python. It was created by John D. Hunter in 2003 with the specific purpose of generating publication-quality plots in Python. The motivation behind its development was to create a plotting tool that could match the plotting capabilities offered by MATLAB, a widely-used software for numerical computing and visualization. By aiming to replicate MATLAB's functionalities, Matplotlib aimed to provide a familiar and powerful plotting solution for Python users, enabling them to create professional-looking plots for various applications, including scientific research, data analysis, and data visualization.

This library allows you to create high-quality, customizable visualizations. Matplotlib provides a wide range of plotting functions and customization options, making it suitable for creating static, animated, and interactive visualizations in Python. These visualizations make data easier to understand through visual representations.

5.4.2 Features

- **Matplotlib's plotting possibilities:** Matplotlib provides users with a wide variety of plot types, including line plots, scatter plots, bar charts, histograms, pie charts, contour plots, 3D plots, and more. This versatility enables users to cater to a diverse range of data visualization needs.
- **Customization Options:** With Matplotlib, users have the flexibility to customize every aspect of their plots, including colors, line styles, markers, labels, titles, axes, legends, gridlines, and more. This level of customization empowers users to create visually stunning and highly informative plots that align with their specific requirements and visual aesthetics.
- **Publication-Quality Plots:** Matplotlib excels in producing publication-quality plots with high-resolution output. Users have control over the resolution, dimensions, and formats of their plots, enabling the creation of professional-looking visualizations that are suitable for academic papers, reports, presentations, and publications.

- **Integration with Python ecosystem:** Matplotlib's strong integration with the Python ecosystem is a notable advantage, as it seamlessly integrates with popular libraries such as NumPy and Pandas, facilitating efficient data manipulation and analysis prior to visualization. Moreover, its compatibility with Jupyter notebooks enhances the interactive exploration and presentation of data, providing a cohesive environment for data analysis and visualization workflows.
- **Interactive Exploration:** Matplotlib can be used in interactive environments like Jupyter notebooks, enabling users to interactively explore and manipulate plots, zoom in/out, and dynamically visualize data.
- **Extensive Documentation and Resources:** Matplotlib stands out with its comprehensive gallery of examples, extensive documentation, tutorials, and strong community support, providing users with valuable resources to learn, troubleshoot, and find inspiration for their visualizations.

5.4.3 Advantages

- **Wide range of plots:** We can create a line plot, a scatter plot, a bar plot, a histogram, a 3D plot, and a polar plot through the matplotlib library. In the code below, we use plt.tight_layout() to adjust the layout of subplots to avoid overlapping. These plots help us visualize data in two-dimension and three-dimension as well.
- **Publication-quality plot:** Using Matplotlib, we can create a publication-quality plot. The code creates sample data for the x-axis and calculates corresponding y-axis values using the sine function. We can customize a plot with a blue line and solid line style.
- **Integration with NumPy and Pandas:** We can integrate NumPy and Pandas, allowing users to create plots directly from NumPy arrays or Pandas DataFrames. If users are working with different libraries, they will prefer using Matplotlib for feasibility. In the code below, we create the data through NumPy and introduce Pandas DataFrame. From this, we plotted the graph.
- **Support for LaTeX and math text:** Matplotlib supports LaTeX and math text, allowing users to display mathematical equations and symbols in various aspects

of their plots, such as axis labels, titles, and annotations. This feature is particularly helpful when dealing with scientific or mathematical data visualizations.

- **Customization options:** Matplotlib offers extensive customization options, empowering users to tailor their visualizations to specific requirements. In the code below, we are customizing the plot with linestyle= '-'.
- **Interactive plots:** Matplotlib's animation module allows users to create interactive and dynamic visualizations, perfect for showcasing time-based data or data with changing characteristics. The resulting interactive plot showcases the sine wave shifting horizontally, and the title dynamically changes to indicate the current iteration.
- **Rich ecosystem:** Matplotlib's plt.bar() function is utilized to create a bar plot, showcasing a salary comparison for each individual in the DataFrame.



Figure 5.4

5.5 Plotly Express

5.5.1 Introduction

The plotly.express module (usually imported as px) contains functions that can create entire figures at once, and is referred to as Plotly Express or PX. Plotly Express is a built-in part of the plotly library, and is the recommended starting point for creating most

common figures. Every Plotly Express function uses graph objects internally and returns a `plotly.graph_objects.Figure` instance. Throughout the plotly documentation, you will find the Plotly Express way of building figures at the top of any applicable page, followed by a section on how to use graph objects to build similar figures. Any figure created in a single function call with Plotly Express could be created using graph objects alone, but with between 5 and 100 times more code.

Plotly Express provides more than 30 functions for creating different types of figures. The API for these functions was carefully designed to be as consistent and easy to learn as possible, making it easy to switch from a scatter plot to a bar chart to a histogram to a sunburst chart throughout a data exploration session. Scroll down for a gallery of Plotly Express plots, each made in a single function call.

5.5.2 Features

- **A single entry point into plotly:** just import `plotly.express` as `px` and get access to all the plotting functions, plus built-in demo datasets under `px.data` and built-in color scales and sequences under `px.color`. Every `PX` function returns a `plotly.graph_objects.Figure` object, so you can edit it using all the same methods like `update_layout` and `add_trace`.
- **Sensible, Overridable Defaults:** `PX` functions will infer sensible defaults wherever possible, and will always let you override them.
- **Flexible Input Formats:** `PX` functions accept input in a variety of formats, from lists and dicts to long-form or wide-form `DataFrames` to `numpy` arrays and `xarrays` to `GeoPandas GeoDataFrames`.
- **Automatic Trace and Layout configuration:** `PX` functions will create one trace per animation frame for each unique combination of data values mapped to discrete color, symbol, line-dash, facet-row and/or facet-column. Traces' `legendgroup` and `showlegend` attributes are set such that only one legend item appears per unique combination of discrete color, symbol and/or line-dash. Traces are automatically linked to a correctly-configured subplot of the appropriate type.

- **Automatic Figure Labelling:** PX functions label axes, legends and colorbars based in the input DataFrame or xarray, and provide extra control with the labels argument.
- **Automatic Hover Labels:** PX functions populate the hover-label using the labels mentioned above, and provide extra control with the hover_name and hover_data arguments.
- **Styling Control:** PX functions read styling information from the default figure template, and support commonly-needed cosmetic controls like category_orders and color_discrete_map to precisely control categorical variables.
- **Uniform Color Handling:** PX functions automatically switch between continuous and categorical color based on the input type.
- **Faceting:** the 2D-cartesian plotting functions support row, column and wrapped facetting with facet_row, facet_col and facet_col_wrap arguments.
- **Marginal Plots:** the 2D-cartesian plotting functions support marginal distribution plots with the marginal, marginal_x and marginal_y arguments.
- **A Pandas backend:** the 2D-cartesian plotting functions are available as a Pandas plotting backend so you can call them via df.plot().
- **Trendlines:** px.scatter supports built-in trendlines with accessible model output.
- **Animations:** many PX functions support simple animation support via the animation_frame and animation_group arguments.
- **Automatic WebGL switching:** for sufficiently large scatter plots, PX will automatically use WebGL for hardware-accelerated rendering.

5.5.3 Advantages

Plotly Express offers several advantages for data visualization in Python:

- **Simplicity and Ease of Use:** Plotly Express provides a high-level, concise syntax that simplifies the creation of a wide range of common plot types. It abstracts away

much of the complexity found in the core Plotly Graph Objects library, making it accessible for users with varying levels of experience.

- **Rapid Visualization and Exploratory Data Analysis (EDA):** Its straightforward syntax allows for quick generation of plots, making it ideal for exploratory data analysis where rapid iteration and visualization are crucial for understanding data patterns.
- **Built-in Interactivity:** Plotly Express automatically creates interactive plots, enabling users to zoom, pan, hover over data points for detailed information, and interact with legends to toggle data visibility. This enhances data exploration and understanding.
- **Automatic Handling of Plot Components:** It automates many aspects of plot configuration, such as axis labels, titles, and legends, based on the input data, reducing the need for manual customization in many cases.
- **Integration with Pandas DataFrames:** Plotly Express works seamlessly with Pandas DataFrames, directly accepting them as input for plotting functions, which simplifies the visualization workflow for data analysts and scientists.
- **Support for Various Plot Types:** It supports a wide array of visualizations, including scatter plots, line charts, bar charts, histograms, box plots, and more, offering versatility for different data visualization needs.
- **Seamless Integration with Dash:** Plotly Express integrates well with Dash, Plotly's framework for building analytical web applications.



Figure 5.5

5.6 Seaborn

5.6.1 Introduction

Seaborn is a library for creating statistical graphics in Python. It is based on Matplotlib and integrates with Pandas data structures.

This library is as powerful as Matplotlib but brings simplicity and unique features. It allows for quick data exploration and understanding.

Complete data frames can be captured, and internal functions for semantic mapping and statistical aggregation allow you to convert data into graphical visualizations.

Seaborn abstracts away all the complexity of Matplotlib. However, it is still possible to create graphics that meet all your needs and requirements.

5.6.2 Features

- **High-Level API for Statistical Graphics:** Seaborn provides a simplified, high-level interface for creating complex statistical plots, requiring less code compared to Matplotlib for similar visualizations.
- **Integration with Pandas Data Structures:** It seamlessly integrates with Pandas DataFrames, allowing direct use of DataFrame labels and structures for plotting, simplifying data preparation.
- **Built-in Themes and Color Palettes:** Seaborn offers a variety of aesthetically pleasing default themes and color palettes, enhancing the visual appeal of plots with minimal effort and providing easy customization options.
- **Specialized Statistical Plots:** It includes functions for a wide range of statistical plots, such as:
 - **Distribution Plots:** Histograms, KDE plots, and rug plots for visualizing data distributions.

- **Relational Plots:** Scatter plots and line plots for showing relationships between variables.
- **Categorical Plots:** Bar plots, count plots, box plots, violin plots, and strip plots for visualizing categorical data.
- **Regression Plots:** Functions for visualizing linear regression models and their fits.
- **Faceting for Complex Data:** Seaborn supports creating grids of plots (faceting) based on categorical variables, enabling easy comparison across different subsets of data.
- **Pair Plots:** The pairplot function generates a matrix of relationships between all numerical variables in a dataset, including histograms or KDE plots on the diagonal and scatter plots for off-diagonal relationships.
- **Customization Options:** While offering opinionated defaults, Seaborn also provides extensive customization options for fine-tuning plot aesthetics, including labels, titles, and figure sizes.
- **Time Series Plotting:** It includes functionalities well-suited for visualizing time-based data and patterns over time.

5.6.3 Advantages

- **High-Level Interface and Simplified Syntax:** Seaborn provides a high-level API built on Matplotlib, abstracting away much of the low-level detail and boilerplate code required for complex plots. This results in a more concise and intuitive syntax, allowing users to create elaborate visualizations with fewer lines of code.
- **Enhanced Statistical Visualization:** It specializes in statistical plotting, offering specialized functions for visualizing relationships, distributions, and categorical data. This includes built-in capabilities for statistical estimation and aggregation within plot functions, providing immediate insights into data characteristics.

- **Attractive Aesthetics and Built-in Themes:** Seaborn comes with visually appealing default styles and color palettes, enhancing the aesthetic quality of plots without requiring extensive manual customization. It also offers built-in themes for consistent and professional-looking visualizations.
- **Seamless Integration with Pandas:** Seaborn is designed to work seamlessly with Pandas DataFrames, a widely used library for data manipulation in Python. This integration facilitates direct visualization of data structures without the need for extensive data preparation.
- **Support for Complex Plot Types:** While excelling at simpler plots, Seaborn also supports the creation of more complex visualizations like multi-panel figures, violin plots, pair plots, and heatmaps, enabling comprehensive data exploration.
- **Customization Options:** Despite its high-level interface, Seaborn provides flexibility for customizing various aspects of plots, allowing users to tailor visualizations to specific needs and preferences.



Figure 5.6

5.7 Dataframes

5.7.1 Introduction

A DataFrame is a data structure that organizes data into a 2-dimensional table of rows and columns, much like a spreadsheet. DataFrames are one of the most common data structures used in modern data analytics because they are a flexible and intuitive way of storing and working with data.

Every DataFrame contains a blueprint, known as a schema, that defines the name and data type of each column. Spark DataFrames can contain universal data types like StringType and IntegerType, as well as data types that are specific to Spark, such as StructType. Missing or incomplete values are stored as null values in the DataFrame.

A simple analogy is that a DataFrame is like a spreadsheet with named columns. However, the difference between them is that while a spreadsheet sits on one computer in one specific location, a DataFrame can span thousands of computers. In this way, DataFrames make it possible to do analytics on big data, using distributed computing clusters.

The reason for putting the data on more than one computer should be intuitive: either the data is too large to fit on one machine or it would simply take too long to perform that computation on one machine.

5.7.2 Features

- **Two-dimensional structure:** DataFrames are organized as a table with rows and columns, similar to a spreadsheet or SQL table.
- **Labeled axes:** Each row and column has a unique label, allowing for easy access and manipulation of data using these labels.
- **Heterogeneous data:** Columns in a DataFrame can hold different data types (e.g., integers, floats, strings, booleans).
- **Size mutability:** The size of a DataFrame can be changed (rows or columns can be added or removed).

- **Efficient data manipulation:** Pandas provides a wide range of functions for data selection, filtering, cleaning, aggregation, and transformation.
- **Integration with other tools:** DataFrames work well with other Python libraries for data visualization, machine learning, and more.
- **Handles missing data:** Pandas DataFrames can handle missing or incomplete data (represented as NaN or Null).
- **Arithmetic operations:** Allows for arithmetic operations (addition, subtraction, multiplication, etc.) to be performed on rows and columns.
- **Index and column labels:** Rows are indexed (labeled) by default with integers, but can be customized with other labels. Columns also have labels.
- **Input flexibility:** Accepts various input formats like dictionaries, lists, NumPy arrays, and other DataFrames to create a DataFrame.

5.7.3 Advantages

1. High-Level Abstraction and Ease of Use:

- DataFrames provide a user-friendly, high-level API, simplifying complex data operations and making them accessible to users with varying levels of programming experience.
- They offer a SQL-like interface, allowing users to query and manipulate data using familiar syntax, making it easier to work with large datasets.
- Intuitive methods for filtering, sorting, aggregating, merging, reshaping, and transforming data make data manipulation tasks straightforward.

2. Optimized Query Execution:

- Catalyst Optimizer, a core component of DataFrames, optimizes query execution by analyzing and transforming logical plans into efficient physical plans.
- This optimization leads to faster query processing, especially for complex data manipulations and large datasets.

3. Schema Enforcement and Data Integrity:

- DataFrames have a built-in schema, which enforces data types and structure, ensuring data consistency and integrity.
- This feature is crucial for working with structured and semi-structured data, preventing errors caused by inconsistent data types or formats.

4. Compatibility and Integration:

- DataFrames are compatible with various data sources and file formats, including CSV, JSON, SQL databases, and more.
- They integrate seamlessly with business intelligence (BI) tools like Tableau and Power BI, facilitating data visualization and reporting.
- DataFrames also integrate well with other programming languages like Java, Scala, R, and Python, making them versatile for different project needs.

5. Scalability and Performance:

- DataFrames can handle large datasets that might be difficult to load into spreadsheets or other data structures.
- They can be used with distributed programming tools like Apache Spark, allowing them to scale to handle petabytes of data on large clusters.
- In-memory columnar data format, like Apache Arrow, allows for faster data transfer and processing, further enhancing performance.

6. Flexibility and Data Manipulation:

- DataFrames allow for efficient data alignment based on labels, simplifying operations on data with different index/column labels or missing values.
- They offer a wide range of functions and methods for manipulating data, including filtering, sorting, aggregation, and joining.
- DataFrames also support operations on entire columns, making it easier to perform calculations and transformations on specific data subsets.

7. Comparison with other data structures:

- DataFrames offer advantages over RDDs in Apache Spark, particularly when working with structured data, due to their optimized query execution and schema enforcement.
- While Spark RDDs are beneficial for unstructured data, DataFrames are generally preferred for structured data analysis due to their performance and ease of use.

5.8 Datasets

5.8.1 Introduction

In Python, "datasets" refer to structured collections of data, typically organized in a way that facilitates analysis, processing, and use in applications like machine learning or data visualization. While the term "dataset" itself is a general concept in data science, in Python, it often relates to specific implementations and libraries designed for handling and manipulating these data collections.

5.8.2 Features

- **Numerical Features:** These may include numerical values such as height, weight, and so on. These may be continuous over an interval, or discrete variables.
- **Categorical Features:** These include multiple classes/ categories, such as gender, colour, and so on.
- **Metadata:** Includes a general description of a dataset. Generally in very large datasets, having an idea/ description of the dataset when it's transferred to a new developer will save a lot of time and improve efficiency.
- **Size of the Data:** It refers to the number of entries and features it contains in the file containing the Dataset.
- **Formatting of Data:** The datasets available online are available in several formats. Some of them are JSON (JavaScript Object Notation), CSV (Comma Separated Value), XML (eXtensible Markup Language), DataFrame, and Excel Files (xlsx or xlsm). For particularly large datasets, especially involving images for disease detection, while downloading the files from the internet, it comes in zip files which will be needed to extract in the system to individual components.
- **Target Variable:** It is the feature whose values/attributes are referred to to get outputs from the other features with machine learning techniques.

- **Data Entries:** These refer to the individual values of data present in the Dataset. They play a huge role in data analysis.

5.8.3 Properties

- **Center of data:** This refers to the "middle" value of the data, often measured by mean, median, or mode. It helps understand where most of the data points are concentrated.
- **Skewness of data:** This indicates how symmetrical the data distribution is. A perfectly symmetrical distribution (like a normal distribution) has a skewness of 0. Positive skewness means the data is clustered towards the left, while negative skewness means it's clustered towards the right.
- **Spread among data members:** This describes how much the data points vary from the center. Common measures include standard deviation or variance, which quantify how far individual points deviate from the average.
- **Presence of outliers:** These are data points that fall significantly outside the overall pattern. Identifying outliers can be important as they might influence analysis results and require further investigation.
- **Correlation among the data:** This refers to the strength and direction of relationships between different variables in the dataset. A positive correlation indicates values in one variable tend to increase as the other does, while a negative correlation suggests they move in opposite directions. No correlation means there's no linear relationship between the variables.
- **Type of probability distribution that the data follows:** Understanding the distribution (e.g., normal, uniform, binomial) helps us predict how likely it is to find certain values within the data and choose appropriate statistical methods for analysis.

5.9 Streamlit

5.9.1 Introduction

Streamlit is an open-source Python library used to create custom web applications for machine learning and data science projects. It enables developers and analysts to convert data scripts into interactive web interfaces with minimal effort and without requiring advanced knowledge of frontend development. Streamlit is particularly popular for building data-driven dashboards, visual analytics tools, and quick prototypes for model deployment.

In the context of Streamlytix, Streamlit plays a central role as the user interface framework, allowing users to select movies and series, submit their watched list, and view sentiment analysis results in an interactive and intuitive environment.

5.9.2 Features

- **Rapid Development:** Streamlit allows developers to build functional web applications using simple Python scripts, significantly reducing development time.
- **Pythonic Syntax:** The entire application layout and logic are written in Python, eliminating the need for separate HTML, CSS, or JavaScript code.
- **Built-in Widgets:** Offers a wide range of interactive components such as buttons, sliders, checkboxes, radio buttons, and file uploaders to capture user input.
- **Auto-refresh Mechanism:** Automatically updates the UI when the underlying Python script is modified or re-run, ensuring a dynamic and responsive user experience.
- **Seamless Integration with Data Libraries:** Works well with popular Python libraries such as Pandas, NumPy, Matplotlib, Plotly, and Seaborn for data manipulation and visualization.

- **Deployment Ready:** Streamlit applications can be deployed easily using cloud platforms like Streamlit Community Cloud, Heroku, or AWS with minimal configuration.

5.9.3 Advantages

- **Ease of Use:** Enables non-web developers, such as data scientists and analysts, to create interactive applications using familiar Python code.
- **Real-time Interaction:** Allows users to interact with data in real-time through widgets, making data exploration more engaging.
- **Minimal Overhead:** Does not require complex web frameworks or boilerplate code; the entire app can be built within a single Python file.
- **Open Source and Actively Maintained:** Backed by an active community and frequently updated with new features and improvements.
- **Effective Visualization:** Supports seamless embedding of charts, tables, and images, making it ideal for presenting analytical results.

In Streamlytix, Streamlit is used to render the user interface, accept user input (watched list), and dynamically display sentiment analysis results, including mood graphs and behavioral summaries. Its simplicity and power make it the ideal choice for building an interactive, data-driven web application without compromising on performance or user experience.



Figure 5.7

5.10 Scikit-Learn

5.10.1 Introduction

Scikit-learn is an open-source Python library widely used for machine learning and data analysis. It provides simple and efficient tools for data preprocessing, model training, evaluation, and deployment. Built on top of libraries like NumPy, SciPy, and Matplotlib, Scikit-learn offers a consistent interface for implementing a wide range of supervised and unsupervised learning algorithms.

In the context of Streamlytix, Scikit-learn serves as the core library for applying machine learning to the data. It is used to train and evaluate machine learning models that analyze user preferences and generate recommendations from the watched content.

5.10.2 Features

- **Comprehensive Algorithm Library:** Includes a wide range of machine learning algorithms such as regression, classification, clustering, and dimensionality reduction.
- **Preprocessing Utilities:** Provides tools for data cleaning, normalization, feature extraction, and transformation to prepare datasets for modeling.
- **Model Evaluation and Validation:** Offers various metrics and validation techniques, including cross-validation and confusion matrices, to assess model performance.
- **Pipeline Support:** Allows creation of end-to-end machine learning workflows that combine preprocessing steps and estimators for reproducible results.

5.10.3 Advantages

- **Ease of Implementation:** Simplifies machine learning tasks through high-level APIs and minimal code requirements.
- **Consistency:** Provides a unified interface for all algorithms, ensuring smooth transitions between models during experimentation

- **Performance:** Optimized for efficiency and scalability using underlying NumPy and Cython operations.
- **Community Support:** Actively maintained by a large open-source community with regular updates and new features.
- **Interoperability:** Easily integrates with other Python libraries, enabling the development of complex machine learning pipelines.



Figure 5.8

Chapter 6

SOFTWARE PROCESS MODEL

6.1 Waterfall model

6.1.1 Introduction

Winston Royce introduced the Waterfall Model in 1970. This model has five phases: Requirements analysis and specification, design, implementation, and unit testing, integration and system testing, and operation and maintenance. The steps always follow in this order and do not overlap. The developer must complete every phase before the next phase begins. This model is named "Waterfall Model", because its diagrammatic representation resembles a cascade of waterfalls.

1. **Requirements analysis and specification phase:** The aim of this phase is to understand the exact requirements of the customer and to document them properly. Both the customer and the software developer work together so as to document all the functions, performance, and interfacing requirements of the software. It describes the "what" of the system to be produced and not "how." In this phase, a large document called Software Requirement Specification (SRS) document is created which contains a detailed description of what the system will do in the common language.
2. **Design Phase:** This phase aims to transform the requirements gathered in the SRS into a suitable form which permits further coding in a programming language. It defines the overall software architecture together with high level and detailed design. All this work is documented as a Software Design Document (SDD).
3. **Implementation and unit testing:** During this phase, design is implemented. If the SDD is complete, the implementation or coding phase proceeds smoothly, because all the information needed by software developers is contained in the SDD.

During testing, the code is thoroughly examined and modified. Small modules are tested in isolation initially. After that these modules are tested by writing some

overhead code to check the interaction between these modules and the flow of intermediate output.

4. **Integration and System Testing:** This phase is highly crucial as the quality of the end product is determined by the effectiveness of the testing carried out. The better output will lead to satisfied customers, lower maintenance costs, and accurate results. Unit testing determines the efficiency of individual modules. However, in this phase, the modules are tested for their interactions with each other and with the system.
5. **Operation and maintenance phase:** Maintenance is the task performed by every user once the software has been delivered to the customer, installed, and operational.

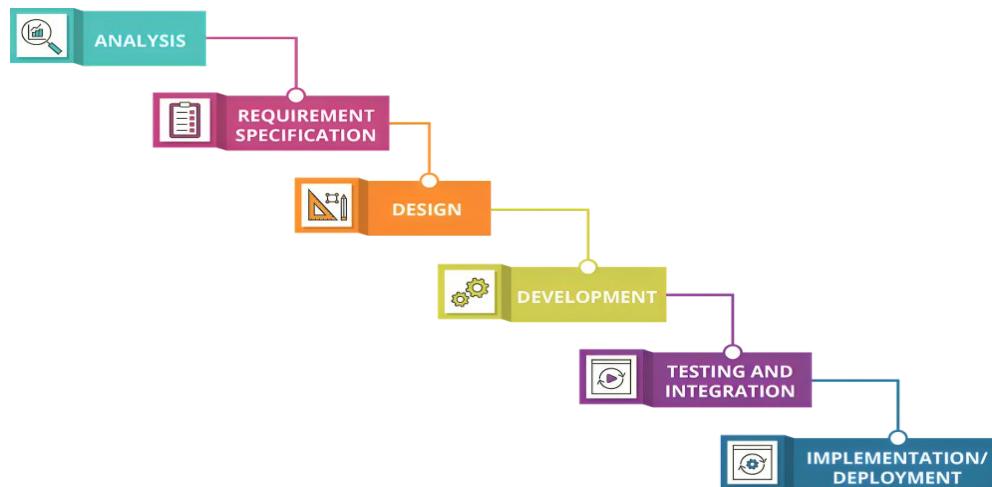


Figure 6.1

6.1.2 Why use SDLC Waterfall model

Some Circumstances where the use of the Waterfall model is most suited are:

- When the requirements are constant and not changed regularly.
- A project is short

- The situation is calm
- Where the tools and technology used is consistent and is not changing
- When resources are well prepared and are available to use.

6.1.3 Advantages of Waterfall model

1. This model is simple to implement and the number of resources that are required for it is minimal.
2. The requirements are simple and explicitly declared; they remain unchanged during the entire project development.
3. The start and end points for each phase are fixed, which makes it easy to cover progress.
4. The release date for the complete product, as well as its final cost, can be determined before development.
5. It gives easy control and clarity for the customer due to a strict reporting system.

6.1.4 Disadvantages of Waterfall model

1. In this model, the risk factor is higher, so this model is not suitable for more significant and complex projects.
2. This model cannot accept the changes in requirements during development.
3. It becomes tough to go back to the phase. For example, if the application has now shifted to the coding phase, and there is a change in requirement, It becomes tough to go back and change it.
4. Since the testing is done at a later stage, it does not allow identifying the challenges and risks in the earlier phase, so the risk reduction strategy is difficult to prepare.

Chapter 7

DESIGN

7.1 Introduction

The full process of thinking, developing, and creating the architecture, structure, and components of a system to meet specified criteria, functions, or purposes is known as system design. It specifies the structure, behaviour, interfaces, and interactions between the system's different pieces.

The translation of user needs and limitations into a blueprint or Model that specifies how different components or modules of the system will work together cohesively to accomplish the intended results is included in this design phase. While developing the system's development and implementation framework, system designers must consider elements such as performance, scalability, dependability, security, and usability.

7.2 Important Components Of System Design

1. Understanding and documenting the system's needs, goals, and restrictions through stakeholder interactions.
2. Architectural Design: Creating the system's high-level structure and components and describing the links and interactions between these aspects.
3. Elaborating on the architectural design by describing precise requirements for individual components, algorithms, databases, interfaces, etc.
4. Building prototypes or models to validate the design, uncover any flaws, and ensure the system's functioning matches the requirements is known as prototyping and testing.
5. Creating the system based on the completed design, including coding, integration, and testing, to guarantee it functions properly.

7.3 Common System Design Principles

- Meeting needs: Understanding and fulfilling the needs of stakeholders is the fundamental purpose of system design. This entails collecting and converting user demands, functions, restrictions, and objectives into a system design that completely meets these factors.
- Scalability refers to the capacity of a system to grow and meet new workloads or user demands without major change or loss of performance. The system's scalability guarantees it can adapt to changing demands and surroundings.
- Building trustworthy, robust, resilient systems that work consistently and accurately under various scenarios. Systems should be designed to reduce the likelihood of failure and elegantly recover from faults or mistakes.
- Performance optimization is the process of ensuring that the system operates efficiently and satisfies performance standards. To accomplish ideal execution levels, speed, response times, asset use, and complete framework throughput should be streamlined.
- Client Experience and Convenience: developing user-friendly, intuitive, and simple technologies. While planning, it should be important to give clients connection points and communications that address their issues and upgrade their involvement in the framework.
- Safety efforts are set up to protect the framework against unapproved access, information breaks, and other potential risks. Data integrity, confidentiality, and system availability are all security goals.

Chapter 8

DATA FLOW DIAGRAMS

8.1 Introduction

A Data Flow Diagram (DFD) is a traditional visual representation of the information flows within a system. A neat and clear DFD can depict the right amount of the system requirement graphically. It can be manual, automated, or a combination of both.

It shows how data enters and leaves the system, what changes the information, and where data is stored.

The objective of a DFD is to show the scope and boundaries of a system as a whole. It may be used as a communication tool between a system analyst and any person who plays a part in the order that acts as a starting point for redesigning a system. The DFD is also called a data flow graph or bubble chart.

8.2 Important Rules

1. All names should be unique. This makes it easier to refer to elements in the DFD.
2. Remember that DFD is not a flow chart. Arrows is a flowchart that represents the order of events; arrows in DFD represent flowing data. A DFD does not involve any order of events.
3. Suppress logical decisions. If we ever have the urge to draw a diamond-shaped box in a DFD, suppress that urge! A diamond-shaped box is used in flow charts to represent decision points with multiple exists paths of which the only one is taken. This implies an ordering of events, which makes no sense in a DFD.
4. Do not become bogged down with details. Defer error conditions and error handling until the end of the analysis.

8.3 Some Standard DFD Symbols

- **Entities** - Entities are source and destination of information data. Entities are represented by rectangles with their respective names.



Figure 8.1

- **Process** - Activities and action taken on the data are represented by Circle or Round-edged rectangles.



Figure 8.2

- **Data Storage** - There are two variants of data storage - it can either be represented as a rectangle with absence of both smaller sides or as an open-sided rectangle with only one side missing.

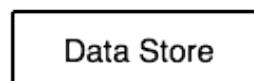


Figure 8.3

- **Data Flow** - Movement of data is shown by pointed arrows.



Figure 8.4

Chapter 9

SNAPSHOTS

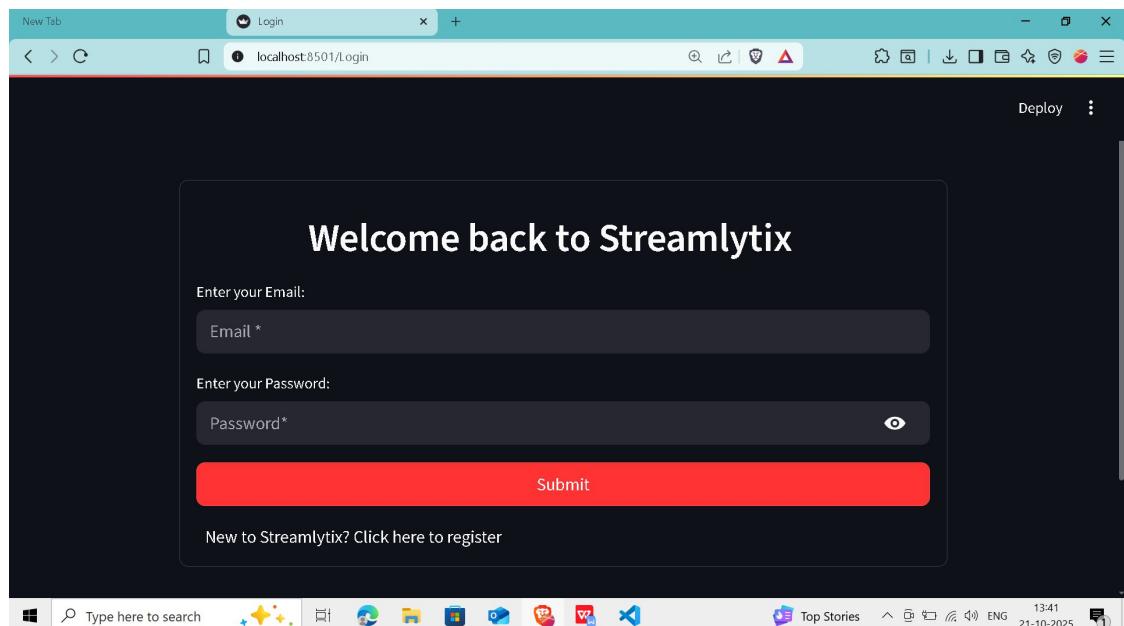


Figure 9.1 Introduction page

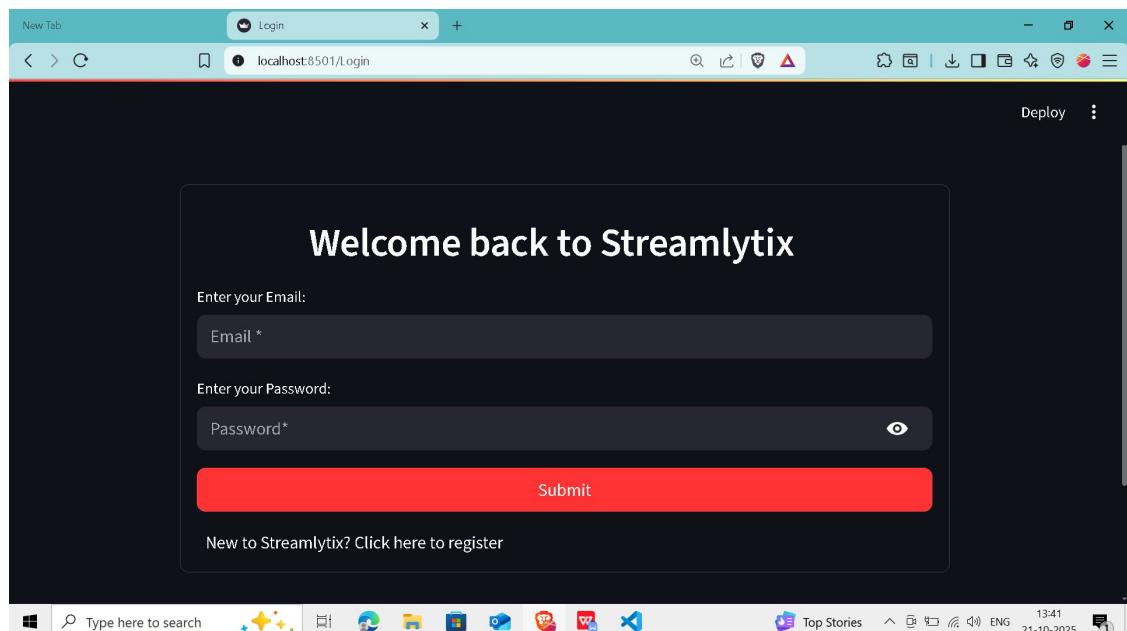


Figure 9.2 Login page

Streamlytix

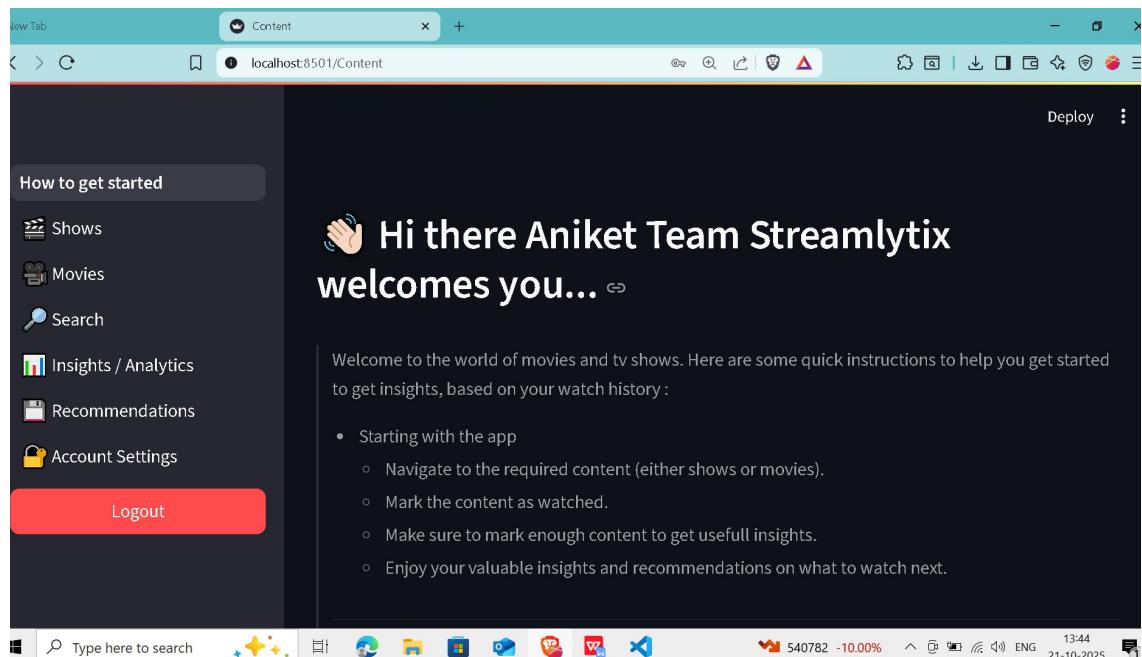


Figure 9.3 Home page

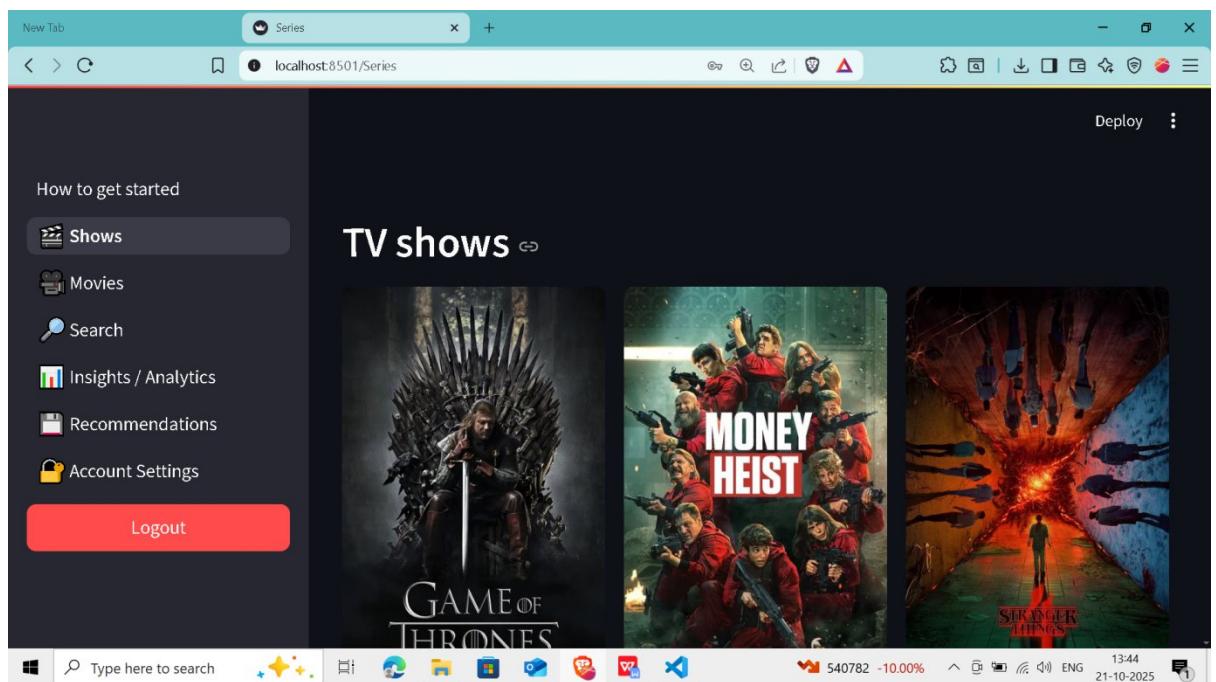


Figure 9.4 Shows page

Streamlytix

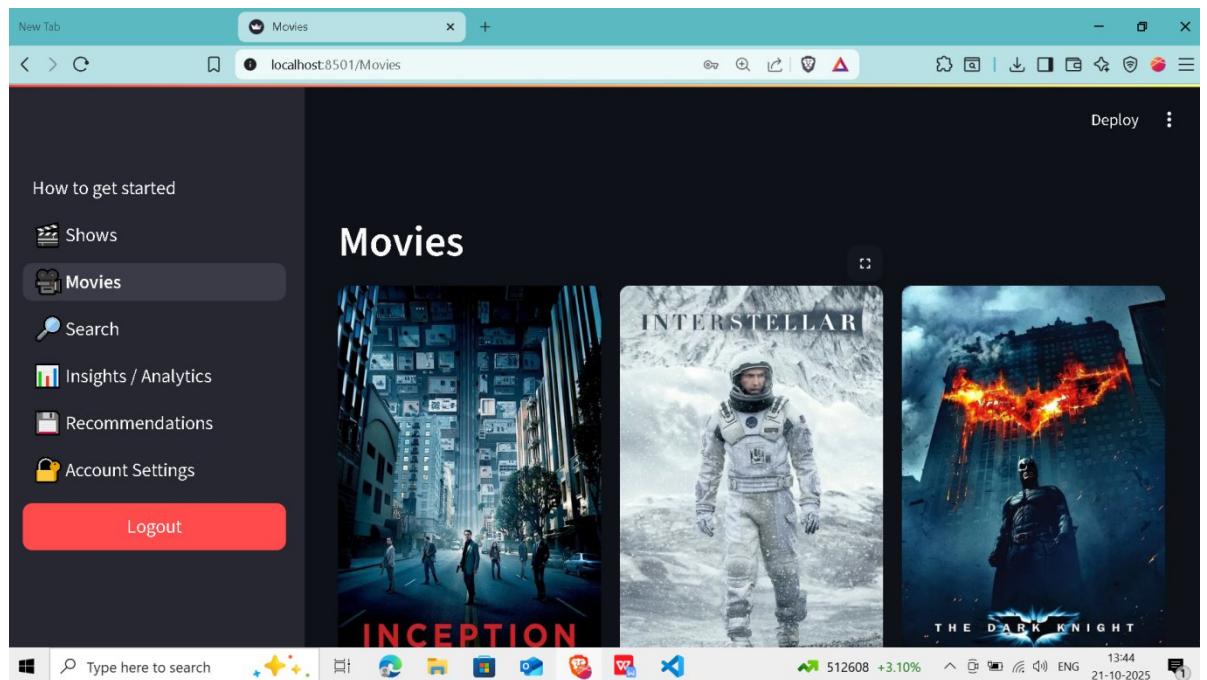


Figure 9.5 Movies page

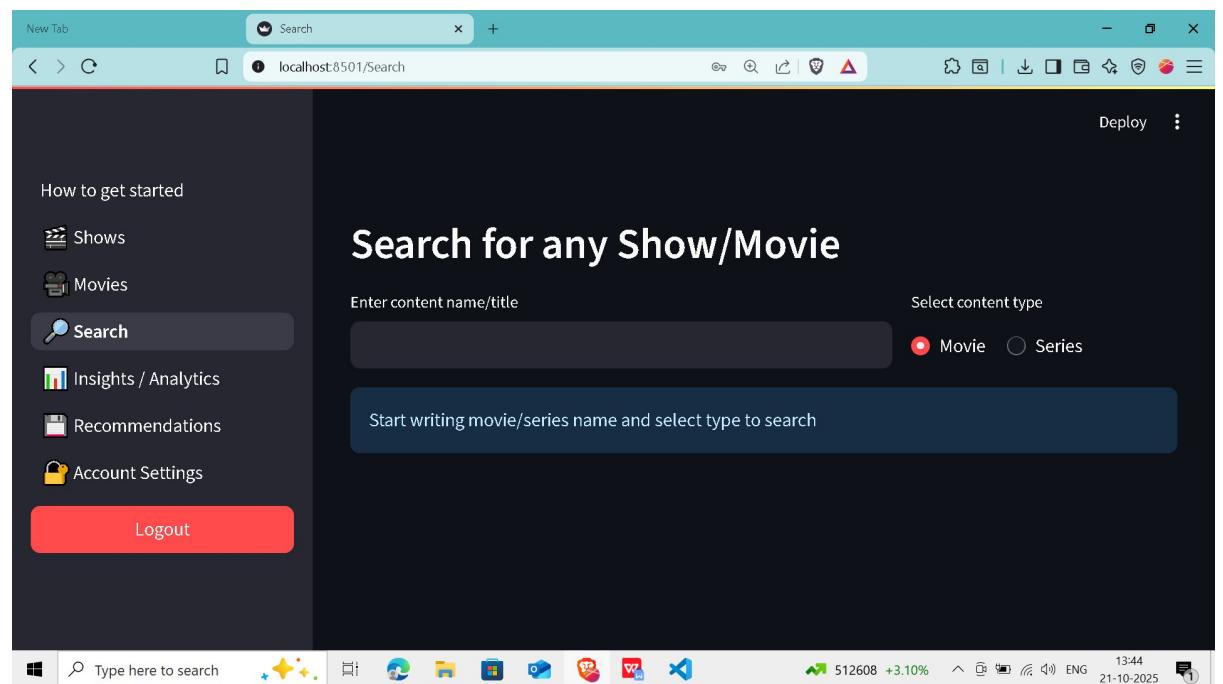


Figure 9.6 Search page

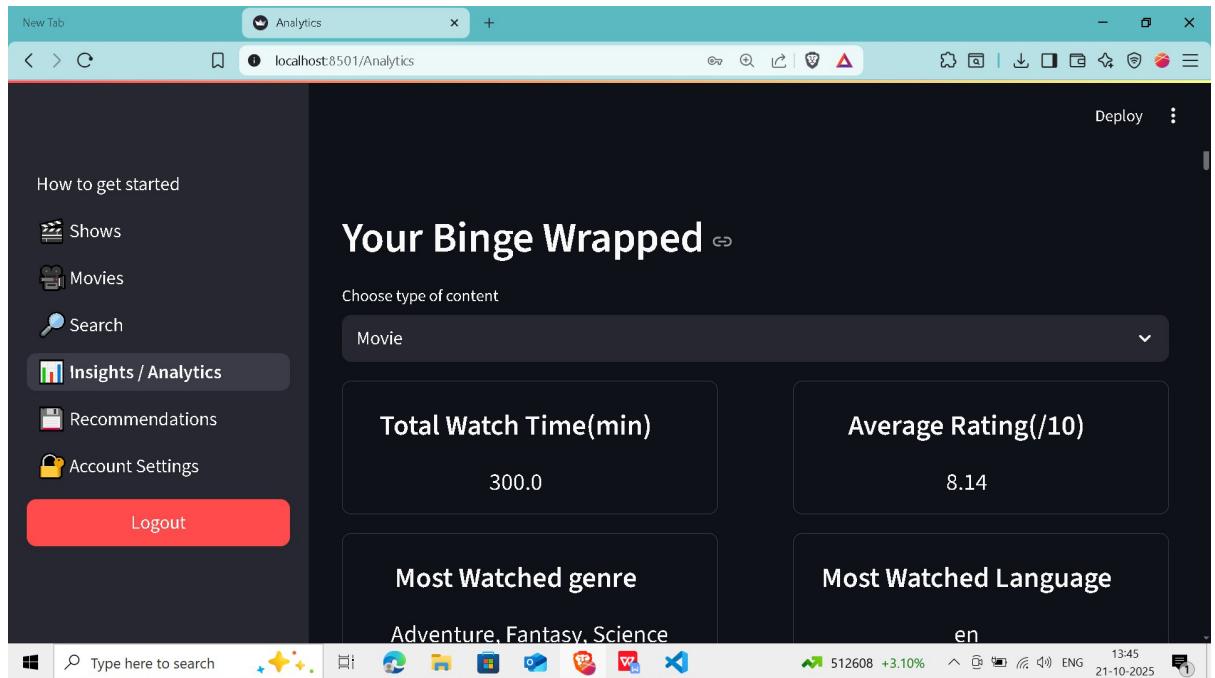


Figure 9.7 Analysis page

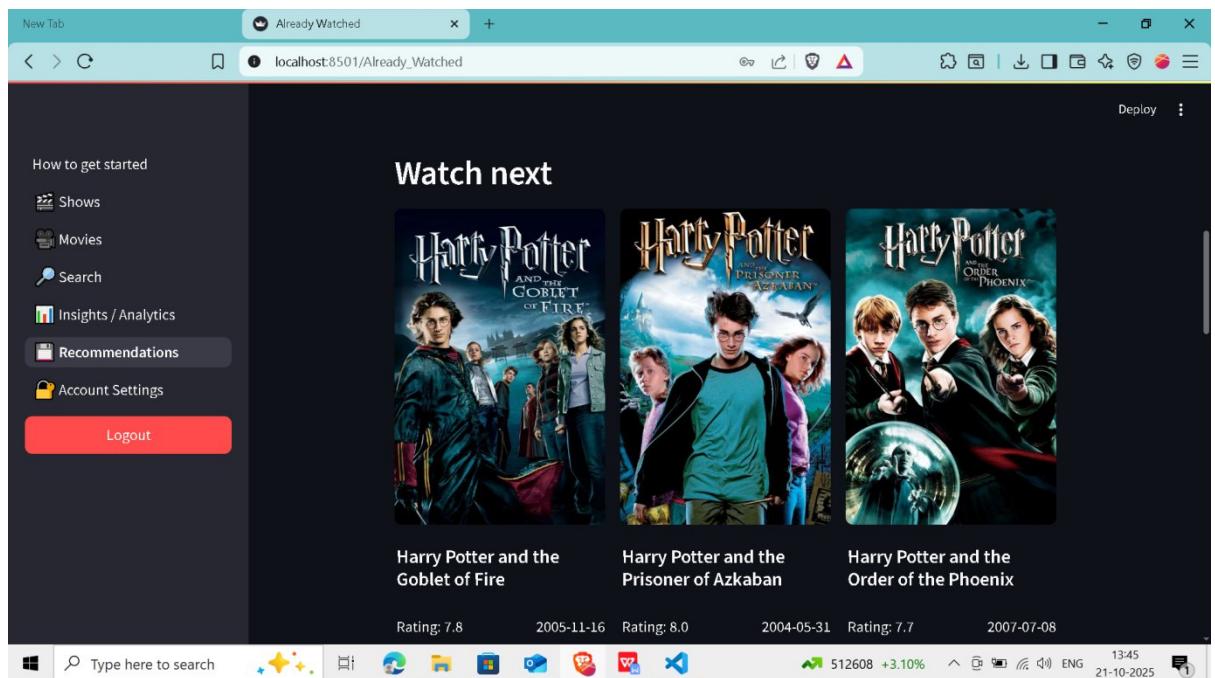


Figure 9.8 Recommendation page

Streamlytix

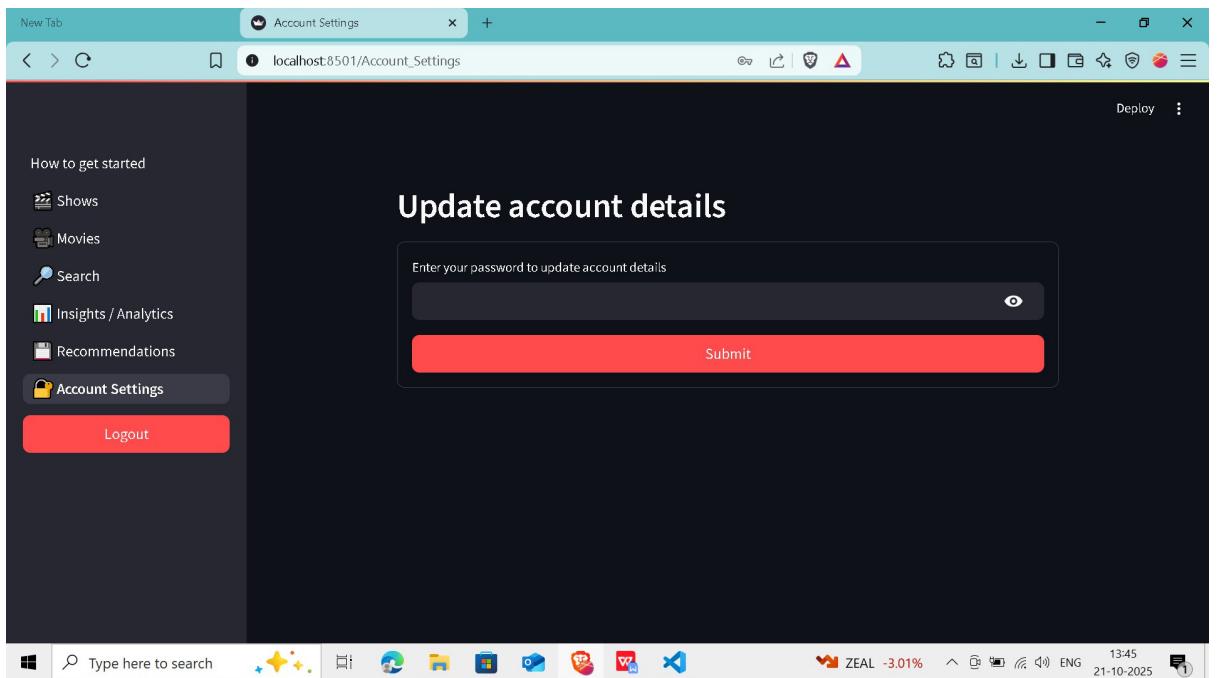


Figure 9.9 Account settings page

Chapter-10

CONCLUSION AND FUTURE SCOPE

10.1 Conclusion

The development of Streamlytix represents a significant step toward integrating data science, machine learning into everyday user experiences. By leveraging analysis and data visualization techniques, the application provides users with personalized insights into their content consumption . Unlike traditional media tracking systems that focus only on watchlists or content ratings, Streamlytix enables users to explore the recommendations for the movies and series they consume.

The system is built using Python and incorporates a rich dataset from TMDB, offering users an intuitive Streamlit-based interface to mark their watched content. It then applies analysis to this content and presents the results in the form of comprehensive graphs and summaries. These insights help users identify their viewing history, reflect on their preferences, and gain a clearer understanding of their media consumption patterns.

Throughout the development of Streamlytix, efforts were made to ensure modularity, scalability, and ease of use. Each module—from user interaction to evaluation to recommendation—was designed with a clear focus on delivering accurate, user-specific analysis in a simple yet effective format. The project successfully demonstrates the practical application of natural language processing, sentiment analysis, and data visualization in building a system that is both informative and engaging.

In conclusion, Streamlytix serves not only as a media tracking platform but also as a personal analytics tool that brings a new dimension to the way users understand their interaction with digital content.

10.2 Future scope

While the current version of Streamlytix provides essential features for content analysis, there are several areas where the system can be enhanced to increase its functionality and user engagement in future iterations:

1. **Similar user based recommendations:** Implementing machine learning algorithms to find similar users and recommend the content based on their watch patterns to provide diversified recommendations.
2. **Advanced Sentiment Analysis using Reviews or Synopses:** Currently, sentiment mapping is based on predefined labels or metadata. In future versions, integrating natural language processing techniques to analyze movie reviews, plot summaries, or user-generated content can provide more accurate and dynamic sentiment evaluations.
3. **Genre-based Emotional Profiling:** Enhancing the system to provide genre-wise sentiment breakdowns could help users understand how their mood varies with genres such as comedy, drama, thriller, or horror.
4. **Recommendation Engine Based on Emotional Patterns:** A future addition could include suggesting content that aligns with or contrasts the user's emotional consumption history, based on their past sentiment profiles.
5. **Integration with APIs for Real-time Content Updates:** Connecting the system directly with TMDB or other streaming platforms via APIs would allow for real-time data fetching, keeping the content library updated and relevant.
6. **Mobile Application Version:** A mobile-friendly version or standalone app can make the system more accessible and enhance usability, especially for users who prefer interacting with such tools on smartphones.
7. **Comparison Features:** Adding functionality to compare emotional viewing patterns between different time periods (e.g., monthly or annually) or across different users (in anonymized form) can make the tool more interactive and insightful.

Chapter 11

REFERENCES

11.1 References

- <https://www.python.org/doc/essays/blurb/>
- https://www.w3schools.com/python/python_intro.asp
- <https://www.geeksforgeeks.org/python/python-introduction-matplotlib/>
- <https://www.geeksforgeeks.org/python/introduction-to-seaborn-python/>
- https://www.tutorialspoint.com/scikit_learn/scikit_learn_introduction.htm