

410448 & 410455: Project Work Book(Guidelines and Log) Fourth

Year Computer Engineering

Year 2024 - 2025

Group/Project ID: 0

Team Members: 1. Khedkar Laxman Bhimrao

2. Hase Ashutosh Prakash

3. Langde Ritesh Vitthalrao

Project Title: Content-Based Movie Recommender System

Project Guide: Mrs.Khatal K.B

Area of the Project: Data Science



**Department of Computer Engineering
Maharia Charitable Trust's
SAHYADRI VALLEY COLLEGE OF ENGINEERING &
TECHNOLOGY
A/P: Rajuri, Tal. Junnar, Dist.: Pune, Pin code: 412411.
Maharashtra (India)**

Prologue

Project work is one of the most important components of the curriculum for the Engineering Graduates. From conceiving the idea to the materialization of it is a journey that has to be systematized, well defined and well documented to enjoy the full benefits of the efforts undertaken.

Every activity of the project development has its own importance and typical activities are like: Team formation, conceiving the idea, preparing the hypothesis, reporting the progress /development to the guide/ mentor, Interactions, suggestions and improvements, relevant documentations in proper format, schedule plans and visit logs.

Every institute is following their own best methods and techniques as per the guidelines and curriculum at the affiliated university. To bring the uniformity and standardization for the project work there is a need to come together and prepare the comprehensive guidelines regarding it.

This work book for the project work will serve the purpose and facilitate the job of students, guide and project coordinator. This document will reflect accountability, punctuality, technical writing ability and work flow of the work undertaken.

This document will definitely support the work undertaken.

Mrs. Khatal K. B
HOD Computer Department,
SVCET, Rajuri.

This booklet is supportive document to rules and a regulation provided by affiliated university curriculum providing recommendations, guidelines and is record of all related activities associated with project. This booklet is provided with the genuine intent to bring uniformity and to systematize the project work and to keep the audit of the work undergone by team members.

Work Book Development Project

Project Institution	Department of Computer Engineering Sahyadri Valley College Of Engineering &Technology,Rajuri .
Support & Guidance	Dr.S.B.Zope, Principal, Sahyadri Valley College Of Engineering &Technology,Rajuri .
Concept and Design	Mr.P.Balaramudu Vice Principal, Sahyadri Valley College Of Engineering &Technology,Rajuri .
Project coordinator	Mrs. Khatal Kavita Head Of Department, Sahyadri Valley College Of Engineering &Technology,Rajuri .
Technical Committee Members	1. Mrs. Khatal Kavita 2. Miss. Khemnar Kavita 3. Mrs. Waghore Vinaya 4. Mr. Nalawade Tejas 5. Mr. Nangare Vaibhav.
Date	10-06-2025

Table of Contents

Sr. No.	Description	Page No.
	Prologue	I
	Work Book Development Project	II
1.	About Project work (Objectives and Outcomes, Guidelines for Selection of Project work and Evaluation)	1
2.	University Syllabus	3
3.	Undertaking by Students	5
4.	Instruction Regarding Project Proposal and Finalization	6
5.	Schedule of Project Work	7
7.	Copy of Project Proposal / Synopsis	8
8.	Review Sheets for Semester II	12
9.	Internal Evaluation Sheet for Semester II	16
10.	Software Engineering code of ethics & professional practices	17
11.	Contests Participation / Publications / IPR and Similar Event Record	22
12.	Rubrics	23
	Bibliography and Annexure	24
	i. Formats for Final Synopsis	25
13.	Partial Project Report Documentation (Semester II)	29
	i. Certificate	29
	ii. Introduction	31
	iii. Literature survey	34
	iv. Implementation(Hardware Component, Block Diagram, Design)	39
	v. Software details(Algorithm, Flowchart, Program)	64
	vi. Result and analysis	76
	vii. System overview	77
14.	Prototyping(Conclusions & References)	78

1. About Project Work

This project focuses on developing a *Content-Based Movie Recommender System* that suggests movies similar to a user's input using machine learning techniques. The system uses metadata like genres, cast, crew, and keywords to find and recommend related movies. It leverages Python libraries such as Pandas and Scikit-learn, with CountVectorizer for feature extraction and Cosine Similarity for recommendation. A user-friendly web interface is built using Streamlit. The project demonstrates the practical application of data science in entertainment and improves the user experience in movie selection.

a.

b. Objectives

- Build a personalized movie recommendation system.
- Use content-based filtering and machine learning techniques.
- Implement feature extraction with CountVectorizer.
- Measure similarity using Cosine Similarity.
- Build a web-based user interface using Streamlit.

b. Expected Outcomes

1. Recommender system works on content similarity
2. Reduces movie search time for users
3. Enhances user experience
4. Understands NLP and ML-based solutions
5. Deployed and tested with good accuracy

c. Guidelines for Selection of Project Work:

Project is one of the significant contributory team works that has to be completed with distinct impression. It is really very difficult to explore the domain of interest / research/ thirst area/ society need. In Toto one cannot figuratively define best project but still there are certain parameters on which we can gauge the quality of project work done. It will be better suited to go for well-defined and relatively safe projects that provide scope for demonstrating proficiency with a low risk of failure especially at Under Graduate level.

General guidelines:

Identifying domain, feasibility and usability of work.

Project work is expected to involve a combination of sound background research (thorough study/ follow a line of investigation), and methodical implementation. Instead of fancied and driven behind the gaudy and ostentatious ideas, the utility has to be emphasized. It is also acceptable to identify the discrepancies/ flaws in the existing system and work accordingly to rectify or improve.

It is irrational to select the IDE and the software/ tools before the idea is not yet finalized.

Understanding the way project will be materialized and progressed.

d. Guidelines for Project Evaluation:

Project work is to be evaluated by both Internal and External examiners jointly, unanimously agreeing the following parameters among many others.

1. Problem definition and scope of the project
2. Through Literature Survey
3. Appropriate Software Engineering approach
4. Exhaustive and Rational Requirement Analysis
5. Comprehensive Implementation- Design, platform, coding, documentation
6. Optimization considerations (Memory, time, Resources, Costing)
7. Thorough Testing of all modules and integration of modules
8. Project Presentation and Demonstration (User Interface, ease of use, usability)
9. Presentation of work in the form of Project Report(s)
10. Understanding individual capacity, Role & involvement in the project
11. Team Work (Distribution of work, intra-team communication and togetherness)
12. Participation in various contests, Publications and IPR
13. Documents /Manuals (Project Report, Quick reference, System, Installation guide)
14. Outcomes / Usability / commercial value /product conversion of Work

2. University Syllabus

Savitribai Phule Pune University
(Refer SPPU website for recent syllabus)

Term II

Tutorial: 6 Hours/Week

Term Work: 50 Marks
Oral: 50 Marks

Course Objectives:

- To develop problem solving abilities using mathematics;
- To apply algorithmic strategies while solving problems;
- To develop time and space efficient algorithms;
- To develop software engineering documents and testing plans
- To use algorithmic solutions using distributed, Embedded, concurrent and parallel environments.
- To encourage and expose students for participation in National/ International paper presentation activities.
- Exposure to Learning and knowledge access techniques using Conferences, Journal papers and participation in research activities.

Course Outcomes:

To solve problem and demonstrate the results of the project;

To develop SRS, reliability testing reports, and other software engineering documents in the project report;

To solve problems using multi-core, distributed, embedded, concurrent/Parallel environments; To write conference paper;

To demonstrate presentation, communication and team-work skills.

Tools:

Preferably 64-bit FOSS tools but if sponsoring company's requirement is non-open source platform then it must be latest and current version of non-obsolete tools. Latest SAN, 3-tier architectures along with latest version of FOSS Operating systems like Fedora 21 or equivalent, LAMP tools, WEB server, Applications servers, Database servers, MongoDB or latest open source Big DATA tools, FOSS Programming Tools like gcc, g++, Eclipse, Python, Java and other tools are per requirement of the SRS. The documentation tools like Open office, Version control, Latex, Latex-Presentation.

1. Project workstation selection, installations and setup along with report to the guide.

(Recommended submission date: - 3 weeks after commencement of second term)

2. Programming of the project, GUI (if any) as per 1st Term term-work submission.

(recommended submission date: - Progress report every week during laboratory hours)

3. Test tool selection for various testing recommended by preferably external guide and generate various testing result charts, graphs etc. including reliability testing. (7 weeks before Term II Conclusion)

4. Review of design and necessary corrective actions taking into consideration feedback report of Term I assessment, and other competitions/conferences participated like IIT, Central Universities, University Conferences or equivalent centers of excellence etc.

5. Students must submit and preferably publish at least one technical paper in the

conferences held by IITs, Central Universities or UoP Conference or International

Conferences in Europe or US.

6. Final term work submissions in the prescribed format given by the guides consisting of a project report consisting of a preliminary report prepared in term-I, detailed design (all necessary UML diagrams) document, User Interface design, Laboratory assignments on test cases and test results generated by selected project testing tool, conclusions, appendix (if necessary), glossary, tools used and references at the end of Term-II after checking, removing/ avoiding the plagiarism. Give an additional assignment per reporting plagiarism to be submitted in the report under the Annex heading extra-work. If the project is the replica of any other previous project or work from other unrelated persons than the student's team, such project should be rejected for the term work.

7. The Term II examination is conducted by panel of examiners (preferably guide and expert from Industry having at least 5 years' subject experience or senior teacher in the subject in case of non-availability of industry expert). The project assessment shall be done using Live Project Demonstration (in existing functional condition), using necessary simulators (if required) and presentation by the students. The remarks of Term I assessment and related corrective actions must be assessed during examining the term- work.

Term-II Project Laboratory Assignments

1. Review of design and necessary corrective actions taking into consideration the feedback report of Term I assessment, and other competitions/conferences participated like IIT, Central Universities, University Conferences or equivalent centers of excellence etc.
2. Project workstation selection, installations along with setup and installation report preparations.
3. Programming of the project functions, interfaces and GUI (if any) as per 1st Term term- work submission using corrective actions recommended in Term-I assessment of Term-work.
4. Test tool selection and testing of various test cases for the project performed and generate various testing, result charts, graphs etc. including reliability testing.

Additional assignments for the Entrepreneurship Project:

5. Installations and Reliability Testing Reports at the client end.
6. To study Clients Feedback reports and related fix generations.
7. To create Documents will Profit and Loss accounts and balance-sheet of the company.

Note: If the students fail to complete the Entrepreneurship assignment successfully then the project shall be treated as Internal Project for the purpose of assessment.

3. Undertaking by Students

**Maharia Charitable Trust's
SAHYADRI VALLEY COLLEGE OF ENGINEERING &
TECHNOLOGY**

UNDERTAKING BY STUDENT

We, the students of B.E. Computer hereby assure that we will follow all the rules and regulations related to project activity for the academic year 2024-2025. The Project entitled-

Movie Recommendation System using filtering Techniques

will be fully designed/ developed by us and every part of the project will be original work and will not be copied/ purchased from any source.

Name of the student

Signature

1. Mr. Khedkar Laxman Bhimrao Exam

2 Mr. Hase Ashutosh Prakash

3. Mr. Langde Ritesh Vitthalrao

4. Instructions Regarding Project Proposal and Finalization

1. The project work may involve the designing a system/subsystem or upgrading / improving an existing system. The design is to be implemented into a workingmodel (software or hardware or both) with necessary software interface as an executable package (installable package or hardware model) along with

User & system manual and quick reference guide. A project report including all necessary documents.
2. Group may come up with sponsored project. Sponsorship may not be in terms of money or resources. It might be in terms of just suggesting problem definition and associated guidance. Students may collect the letter required for applying the Institute/Industries for the project sponsorship from project coordinator
3. List of suggested projects, prominent domains and respective expert, whom you may contact for guidance, with Project Coordinator. Students may contact respective staff along with synopsis for the guidance. Students may contact respective staff for projects suggested by them in the respective areas.
4. Meet Project Coordinator for project title registration.
5. Synopsis must include project title, group members, sponsor details (if any), detailed problem definition, area, abstract, details of existing similar systems if any, scope of the project and software-hardware requirements. Sponsorship detailsinclude name of sponsoring authority, address, name of guide, sponsorship terms & conditions and respective documents certifying the same from authorities.
6. A Panel of experts will approve the project group and title only after presentation as per schedule. Presentation will cover details mentioned in the synopsis as above.

5. Schedule of Project Work**Semester II**

Sr. No.	Activity Scheduled	Date
1.	Fifth presentation about progress of project work(Review V)	16-01-2025
2.	Sixth presentation about progress of project work (Review VI)	06-02-2025
3.	Review of Publication Activity	28-02-2025
4.	Seventh presentation about progress of project work(Review VII)	06-03-2025
5.	Submission of Draft of Report for checking	12-03-2025
6.	Submission of final project report and Project Work book to the project Coordinator	27-03-2025
7.	Mock Project Examination	06-06-2025
8.	Project Examination	10-06-2025

6. Copy of Proposal / Synopsis as per format (Annexure I)

SYNOPSIS OF PROJECT

1. Title of the Project: "**Content-Based Movie Recommender System**"
2. Names of group members:

- 1. Khedkar Laxman Bhimrao**
- 2. Langde Ritesh Vithal**
- 3. Hase Ashutosh Prakash**

- 3. Name of the Guide: Mrs. Khatal K. B**

Introduction

In today's digital era, the entertainment industry has grown exponentially, offering users access to thousands of movies across various platforms. However, with such a vast collection available, users often find it difficult to select a movie that matches their taste and preferences. To overcome this problem, recommendation systems have become essential tools in enhancing the user experience by suggesting relevant content based on user behavior or content attributes.

This project, titled "*Content-Based Movie Recommender System*," aims to develop a recommendation engine that uses machine learning techniques to provide personalized movie suggestions. Unlike collaborative filtering systems, which depend on user ratings and viewing history, this system relies solely on the content of the movies, such as genre, cast, crew, and keywords, to determine similarity. This makes the system particularly useful in cases where user data is unavailable or limited.

The core algorithm of the system is built using **CountVectorizer** for feature extraction and **Cosine Similarity** to calculate the degree of similarity between movies. These techniques are supported by popular Python libraries such as Scikit-learn and Pandas. The front-end of the system is developed using **Streamlit**, a lightweight and efficient web app framework that allows users to interact with the system in real-time by simply entering the name of a movie to receive similar recommendations.

This project also demonstrates the practical application of Natural Language Processing (NLP) and text vectorization in a real-world scenario. By transforming textual movie data into numerical vectors, the system is able to compute meaningful comparisons between different films.

Overall, this project not only enhances the user experience by simplifying the movie selection process but also showcases the effective integration of machine learning, data science, and web technologies to solve everyday problems. It lays a strong foundation for future enhancements such as hybrid systems, user profiling, or integration with external APIs for live data.

Objectives

The primary objective of the *Content-Based Movie Recommender System* is to simplify and personalize the movie discovery experience for users by suggesting movies similar to a given input. In today's world, where the volume of digital content is vast and growing rapidly, such systems play a vital role in improving user engagement and satisfaction. The project leverages machine learning techniques and modern web development tools to fulfill this aim.

Below are the comprehensive objectives of the project:

- 1. To analyze and utilize movie metadata for recommendation purposes**

The system is designed to make intelligent use of available metadata such as genres, cast, crew, keywords, and plot descriptions to determine similarity between movies. By focusing on content features instead of user behavior, the system becomes more flexible and accessible, especially for new users without historical data.

- 2. To implement feature extraction using CountVectorizer**

One of the key objectives is to extract meaningful features from textual data. The **CountVectorizer** tool from the Scikit-learn library is used to convert the combined textual data into a numerical matrix, enabling the machine to compute and compare different movies based on word frequency.

- 3. To calculate similarity scores using Cosine Similarity**

Once the movie features are vectorized, the system applies the Cosine Similarity algorithm to measure the closeness between any two movies. This similarity score helps in generating a ranked list of movies that are most similar to the user's input.

- 4. To develop an interactive and user-friendly front-end using Streamlit**

A key project goal is to build a simple, lightweight, and intuitive front-end where users can enter a movie name and instantly view recommendations. Streamlit is used as it allows quick development of web apps entirely in Python without requiring complex backend setup.

- 5. To enable fast, real-time movie recommendation delivery**

The project aims to ensure minimal response time so that users receive recommendations almost instantly. This is achieved through optimized data preprocessing, vectorization, and in-memory computation using efficient Python libraries.

- 6. To gain practical experience in applying machine learning techniques to real-world applications**

Through this project, students aim to strengthen their understanding of how machine learning and data science techniques can be used to solve real-world problems by implementing and deploying a functioning recommendation system.

Literature Review

Recommendation systems have become an integral part of modern digital platforms, especially in entertainment, e-commerce, and social media. Their primary goal is to filter vast amounts of data and present users with personalized suggestions based on their preferences. Over the years, several recommendation techniques have been developed, mainly categorized into three types: **collaborative filtering**, **content-based filtering**, and **hybrid approaches**.

Collaborative filtering recommends items to a user based on the preferences of other users with similar tastes. Although widely used in platforms like Netflix and Amazon, it suffers from the **cold-start problem**—where new users or new items cannot be effectively recommended due to a lack of prior data. Additionally, collaborative filtering systems often require large-scale user interaction data, which may not be available in all scenarios.

In contrast, **content-based filtering** recommends items based on the characteristics of the item itself. For movie recommendations, this involves analyzing metadata such as genres, cast, crew, plot keywords, and other descriptive features. Since this technique focuses on the content rather than the user, it does not require historical user data, making it well-suited for this project. However, content-based systems may lack diversity, as they often recommend items too similar to those already liked by the user.

Hybrid systems combine both approaches to achieve higher accuracy and overcome the limitations of individual techniques. They are more complex and typically used in commercial-grade systems requiring advanced infrastructure and data handling capabilities.

Multiple studies support the effectiveness of content-based systems in real-world applications. According to *Kumar & Rao (2019)*, content-based filtering is advantageous in cases where user interaction data is limited but item metadata is rich. Another study by *Patel et al. (2020)* emphasized the importance of preprocessing and feature engineering to improve the accuracy of recommendations. These findings influenced the design of this project, especially in the selection of CountVectorizer and Cosine Similarity for measuring content similarity.

Furthermore, the use of **natural language processing (NLP)** to handle textual movie features and convert them into numerical vectors has proven effective in recent research. The application of cosine similarity has been demonstrated to produce reliable and interpretable results in measuring the closeness between different movie vectors.

In summary, the existing literature confirms the practicality and efficiency of content-based recommendation systems, particularly for domains where content features are readily available. This project builds upon these insights to develop a movie recommendation system that is accurate, scalable, and user-friendly.

Feasibility Study

Before proceeding with the development of the *Content-Based Movie Recommender System*, a feasibility study was conducted to evaluate whether the project could be successfully completed within the available time, resources, and constraints.

1. Technical-Feasibility:

The project uses well-established tools and libraries such as Python, Pandas, Scikit-learn, and Streamlit. These are open-source, well-documented, and compatible with most systems. The dataset used is publicly available, and the algorithms (CountVectorizer and Cosine Similarity) are proven and widely used in similar projects. Thus, the technical feasibility is high.

2. Economic-Feasibility:

There are minimal financial requirements since the project is developed using free and open-source tools. No additional investment is needed for licensing or hardware beyond a standard development machine, making the project economically viable.

3. Operational-Feasibility:

The system is simple to use and operates through a web interface, making it accessible even to non-technical users. The recommendation results are generated in real-time, enhancing usability and operational efficiency.

4. Schedule-Feasibility:

The project was planned with a realistic timeline, divided into phases: data preprocessing, model development, testing, interface creation, and deployment. All activities were aligned with academic deadlines, ensuring timely completion.

In conclusion, the project is feasible in all key areas: technically sound, cost-effective, operationally manageable, and deliverable within the given timeframe.

Methodology and Planning of Work

Project Initiation:

- **Objective Definition:**

Clearly define the objective to build a content-based movie recommendation system using machine learning, which recommends similar movies based on features like genre, cast, and keywords.

- **Stakeholder Consultation:**

Discuss project goals, expected outcomes, and timelines with the project guide and team members to gather insights and finalize the roadmap.

Requirements Analysis:

- **Process Mapping:**

Identify and map the steps required to build the system: data collection, preprocessing, feature extraction, similarity calculation, and user interface development.

- **Technical Assessment:**

Evaluate tools and libraries required for the project, such as Pandas, Scikit-learn, CountVectorizer, Cosine Similarity, and Streamlit.

Design and Planning:

- **System Design:**

Design the logical flow of the recommender system, including input, processing, and output stages. Define how user input will lead to real-time suggestions.

- **Tool Selection:**

Choose Python as the core language, with Scikit-learn for ML and Streamlit for web deployment. Finalize dataset from TMDb or MovieLens.

- **Resource Allocation:**

Assign roles among team members for data handling, algorithm development, interface design, and documentation.

Development:

- **Model Development:**

Use CountVectorizer to convert movie metadata into vectors. Implement cosine similarity to compute similarity scores between movies.

- **Interface Creation:**

Develop a simple web app using Streamlit that accepts a movie name and displays similar movies.

- **Error Handling:**

Include exception handling for invalid inputs, missing data, and edge cases.

- **Integration Testing:**
Test the connection between backend logic and the Streamlit interface to ensure seamless functionality.

Testing:

- **Unit Testing:**
Test individual components like feature extraction, similarity calculation, and result display.
- **End-to-End Testing:**
Test the entire system from user input to final recommendation output.
- **User Acceptance Testing (UAT):**
Get feedback from peers and faculty to validate usability and accuracy of recommendations.

Deployment:

- **Staging Deployment:**
Deploy the project locally to test complete functionality in a simulated environment.
- **Final Deployment:**
Prepare the final project setup for demonstration and submission, ensuring it runs smoothly on standard hardware.

Monitoring and Maintenance:

- **Performance Monitoring:**
Review response time and accuracy of results. Ensure the system scales with larger datasets if needed.
- **Updates and Enhancements:**
Make improvements such as refining text preprocessing or integrating external APIs like TMDb for poster images.

Documentation and Training:

- **Documentation:**
Create user manuals, algorithm explanations, code comments, and report all system components clearly.
- **Training:**
Ensure all team members are familiar with each part of the system for smooth demonstration and explanation during exams.

Evaluation and Feedback:

- **Performance Evaluation:**
Evaluate the accuracy and efficiency of the system. Compare it against manual searches or other basic recommenders.

- **Feedback Collection:**
Collect suggestions from faculty and peers to identify areas of improvement for future enhancements.

Project Closure:

- **Final Review:**
Review all project deliverables with the guide and ensure all objectives have been achieved.
- **Closure Report:**
Document achievements, challenges, and lessons learned throughout the development journey. Prepare for viva and final submission.

Facilities Required for Proposed Work

Software and Tools:

- **Python 3.x:** Primary programming language used for data processing and algorithm implementation.
- **Jupyter Notebook / Visual Studio Code:** Development environments for writing and testing code.
- **Scikit-learn:** Library used for CountVectorizer and Cosine Similarity computations.
- **Pandas and NumPy:** For data manipulation and analysis.
- **Streamlit:** Used to build and deploy the web-based user interface.
- **Web Browser:** (e.g., Google Chrome, Firefox) for accessing the web application.

Hardware:

- **Development Workstations:** Computers with the following minimum specifications:
 - **CPU:** Intel i5/i7 or equivalent multi-core processor
 - **RAM:** Minimum 8 GB (16 GB preferred for large-scale datasets)
 - **Storage:** SSD with at least 100 GB of free space
- **Network Infrastructure:** Reliable internet connectivity for downloading datasets, Python libraries, and accessing APIs.

Development Environment:

- **Development Space:** A quiet and focused workspace equipped for coding, testing, and documentation.
- **Source Control System:** GitHub used for version control, collaborative development, and code backup.

Access and Permissions:

- **Dataset Access:** Publicly available dataset (e.g., TMDb 5000 Movie Dataset or MovieLens).
- **API Access:** TMDb API used to fetch movie posters and additional information.
- **System Permissions:** Required to install and configure Python and libraries on development

machines.

Testing Environment:

- **Staging Environment:** Local testing environment created using Streamlit before final deployment.
- **Test Data:** Sample movie titles used to validate the recommendations and accuracy of the system.

Documentation and Support:

- **Technical Documentation:** Python library docs (e.g., Scikit-learn, Streamlit), TMDb API documentation.
- **Support Resources:** Python forums, Stack Overflow, GitHub discussions, and the Streamlit community.

Training and Skill Development:

- **Learning Resources:** Tutorials and documentation on machine learning, CountVectorizer, cosine similarity, and Streamlit available through YouTube, Coursera, and official docs.

Security Measures:

- **Data Protection:** Minimal user input is collected; however, proper handling of inputs and error checks are implemented.
- **Access Controls:** Restricted file access during development to avoid unauthorized modifications.

Collaboration Tools:

- **Project Management Tools:** Trello or Notion used for task tracking and progress management.
- **Communication Channels:** WhatsApp group and email used for team coordination.

Expected Outcomes

Enhanced User Experience:

- **Faster Movie Discovery:** Users receive relevant movie recommendations instantly.
- **User-Friendly Interface:** Simple interface created using Streamlit ensures ease of use for all age groups.

Improved Accuracy and Efficiency:

- **Reliable Output:** Recommendations based on solid NLP techniques and similarity metrics.
- **Real-Time Processing:** Minimal wait time for results using optimized algorithms.

Better Data Utilization:

- **Content-Based Filtering:** Makes use of genre, cast, keywords, and metadata for better results.
- **Consistent Results:** Cosine similarity ensures similarity is computed fairly and uniformly.

Cost Efficiency:

- **Open-Source Tools:** All development tools and libraries used are free and open-source.
- **No Hosting Costs (for now):** The application runs locally on a system, avoiding any server hosting expenses.

Scalability and Adaptability:

- **Scalable Model:** Can be extended to work with larger datasets or integrated with cloud deployment.
- **Modular Design:** Components (vectorization, similarity, UI) can be improved or replaced independently.

Learning and Skill Enhancement:

- The team has gained hands-on experience in:
 - Data preprocessing and transformation
 - Feature engineering using CountVectorizer
 - Machine learning model development and evaluation
 - Web development using Streamlit
 - Deployment and user feedback integration

References

Books and Online Guides

- *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow* – Aurélien Géron
- *Python for Data Analysis* – Wes McKinney
- *Machine Learning Yearning* – Andrew Ng (Online PDF)

Library and Framework Documentation

- **Scikit-learn Documentation:** https://scikit-learn.org/stable/user_guide.html
- **Streamlit Docs:** <https://docs.streamlit.io/>
- **TMDb API Docs:** <https://developer.themoviedb.org/>

Research Papers and Articles

- *A Content-Based Recommender System Using Cosine Similarity and Movie Metadata* – Journal of Emerging Technologies
- *Comparative Analysis of Recommendation Algorithms* – IJERT

Community Support

- **Stack Overflow:** <https://stackoverflow.com/questions/tagged/recommender-system>
- **GitHub Discussions:** For code-related queries and collaborative development
- **Kaggle:** For dataset ideas, notebooks, and tutorials

7. Project Review (Semester II)

The group members are expected to present their work undertaken during the semester. Journey of development has to be rationally presented.

7.1 Project Review-I: Modeling (Model Refinement and Algorithm development)

Student is expected to deliver presentation covering Modeling

Sr. No.	Question	Date	Remark/ Grade	Sign of Guide
1)	How well have you justified the choice of content-based filtering over other methods?	16-01-2025		
2)	How did you handle missing or inconsistent entries, and what impact might residual issues have on recommendations?	22-01-2025		
3)	Are text preprocessing steps (e.g., tokenization, stopword removal) appropriate for the plot/keywords?	29-01-2025		
4)	Did you standardize or normalize categorical features (e.g., genre encoding) in a way that balances influence across features?	06-02-2025		
5)	How did you combine different feature types (textual vs. categorical) into a single item profile vector?	13-02-2025		
6)	Have you considered weighting (e.g., giving more importance to genres vs. plot text), and is your choice justified?	20-02-2025		
7)	Does the mathematical model clearly states goal of project?	28-02-2025		
8)	Does the interface between the modules properly identified?	06-03-2025		
9)	Does any functional dependencies are identified and described?	12-03-2025		
10)	Which architectural model does your system supports?	21-03-2025		
11)	How does your system handle cases where the user provides very few liked movies (cold-start)?	28-03-2025		
12)	Whether all components are designed properly and represented in component diagram?	03-04-2025		
13)	Is the method for aggregating liked-movie vectors into a user profile sound?	09-04-2025		
Remark and Suggestions:				

Name and Sign of Reviewers:

1. Mrs. Waghore Vinaya
2. Mr. Nangare Vaibhav

7.2 Project Review-II: Coding / Implementation

Student is expected to deliver presentation covering Coding / Implementation

Sr. No.	Question	Date	Remark/ Grade	Sign of Guide
1)	Does the code completely and correctly implement the design?	16-01-2025		
2)	Does the code comply with the coding standard?	22-01-2025		
3)	Is the code well structured, consistent in style, and consistently formatted?	06-02-2025		
4)	Are all functions in the design coded?	28-02-2025		
5)	Does the code make use of object oriented concepts?	06-03-2025		
6)	Does the code support granularity?	12-03-2025		
7)	Does the language used for coding is correctly chosen as per the project need?	21-03-2025		
8)	If any off the shelf components are used, Have you understood the functionalities of using it?	27-03-2025		
9)	Are all comments consistent with the code?	03-04-2025		
10)	Whether code optimization is done properly or not?(By using language features)	09-04-2025		

Remark and Suggestions:

Name and Sign of Reviewers:

1. Mrs. Waghore Vinaya

2. Mr. Nangare Vaibhav

7.3.Project Review-III: Validation and Testing
 Student is expected to deliver presentation covering Validation and Testing

Sr. No.	Question	Date	Remark / Grade	Sign of Guide
1)	Have you done alpha testing?	16-01-2025		
2)	Have you done beta testing?	22-01-2025		
3)	Have you validated the requirements, design and code as per standard?	06-02-2025		
4)	Have you performed GUI testing of project? How?	28-02-2025		
5)	Does your system comply with basic usability norms?	06-03-2025		
6)	Have you tested the code using standard datasets available in your area of project?	12-03-2025		
7)	Have you tested the code in real time environment?	27-03-2025		
8)	After integration of all components whether total performance of system is checked or not?	03-04-2025		
9)	Whether repository of all components along with versions is documented or not?	09-04-2025		
Remark and Suggestions:				

Name and Sign of Reviewers:

- 1. Mrs. Waghore Vinaya**
- 2. Mr. Nangare Vaibhav**

7.4 Project Review-III: Report Writing

Student is expected to deliver presentation covering Report Writing

Sr. No.	Question	Date	Remark / Grade	Sign of Guide
1)	Is the report written as per the prescribed format?	16-01-2025		
2)	Is the report timely prepared?	22-01-2025		
3)	Is the report properly organized, spelled, grammatically correct?	06-02-2025		
4)	Is the report plagiarism free?	28-02-2025		
5)	Is the report precise and written to the point?	06-03-2025		
6)	Is the report contains complete results and comparative graphs?	12-03-2025		
7)	Are all figures and tables properly numbered and labeled?	27-03-2025		
8)	Are all figures and tables properly cited?	03-04-2025		
9)	Weather references are properly cited?	09-04-2025		

Remark and Suggestions:

Name and Sign of Reviewers:

1. Mrs. Waghore Vinaya

2. Mr. Nangare Vaibhav

8. Internal Evaluation Sheet (Semester II)

Sr. No.	Name(s) of the student in the project group	Modeling (10)	Coding and Implementation (40)	Testing (10)	Understanding, Individual Involvement / Contribution in the project (10)	Team Work (10)	Demonstration cum Presentation (10)	Documents & Report (10)	Total (100)
1.									
2.									
3.									
4.									

Name and Signature of Evaluation Committee:

1. Mrs. Waghore Vinaya

2. Mr. Nangare Vaibhav

Examiners Feedback and Suggestions:

Signature of Guide

[]

Signature of Head

[Mrs. Khatal K. B]
Head of Department

2. Software Engineering Code of Ethics and Professional Practices

◊ 1. Public

Act consistently with the public interest.

- Ensure the automation does not misuse user data or violate Mahadiscom's terms of service.
- Protect consumer privacy and avoid unauthorized access.

◊ 2. Client and Employer

Act in the best interests of the client or employer.

- Maintain confidentiality of credentials and account data.
- Deliver accurate and functional automation to serve real business or academic needs.

◊ 3. Product

Ensure the product meets the highest professional standards.

- Test the automation thoroughly to avoid errors in billing or status checking.
- Ensure reliability, accuracy, and user-friendly behavior.

◊ 4. Judgment

Maintain integrity and independence in judgment.

- Make clear any limitations of the automation (e.g., CAPTCHA handling, system downtimes).
- Document assumptions and dependencies in your code.

◊ 5. Management

Promote ethical approaches in project management.

- Manage the project with transparency and accountability.
- Keep logs of automation activity for review and audit.

◊ 6. Profession

Advance the integrity and reputation of the profession.

- Follow good coding and documentation practices.
- Respect the Mahadiscom website's intended use.

◊ 7. Colleagues

Be fair and supportive to colleagues.

- Share code, workflow, and knowledge to support collaboration and learning.

◊ 8. Self

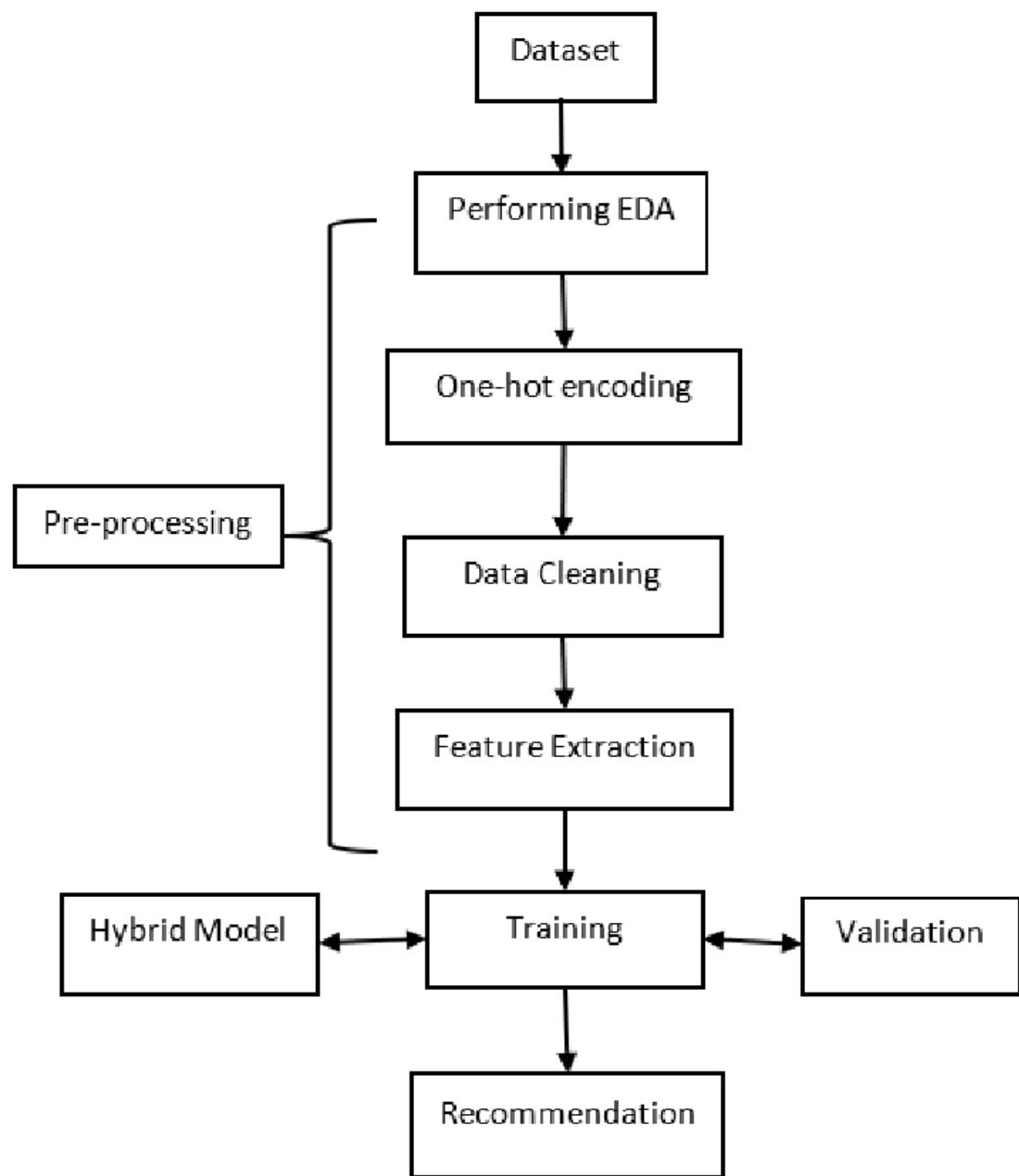
Participate in lifelong learning and ethical development.

- Stay updated with UiPath and RPA best practices.
- Reflect on the ethical implications of automation in public service domains

Algorithm

- Start the process
- Load the movie dataset
- Clean and preprocess the dataset (remove nulls, merge keywords, cast, genres, etc.)
- Combine selected features into a single string
- Convert the combined text into vectors using CountVectorizer
- Compute similarity matrix using Cosine Similarity
- Take user input (movie name)
- Search for the given movie in the dataset
 - a. If not found, display an error message and end
 - b. If found, proceed to next step
- Find the similarity scores for the given movie
- Sort the scores in descending order
- Select the top 5 most similar movies
- Fetch movie titles and details of similar movies
- Display the results through Streamlit web interface
- End the process

Flowchart



10. Contest Participation Details

A. Paper Publication/ Presentation/IPR

Sr. No.	Name of Organizer	Date	Certificates/ Prizes won if any
1.	INTERNATIONAL RESEARCH JOURNAL OF MODERNIZATION IN ENGINEERING TECHNOLOGY AND SCIENCE	June - 2025	Certificates
2.	INTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT (IJSREM)	June - 2025	Certificates

11. Rubrics

A. Idea Inception

Grade (Grade Point)	Excellent (10-9)	Very Good (6-8)	Fair (3-5)	Poor (1-2)
Parameter				
Problem Definition and Scope of the Project				
Literature Survey				
Software Engineering				
Requirement Analysis				

B. Implementation

Grade (Grade Point)	Excellent (10-9)	Very Good (6-8)	Fair (3-5)	Poor (1-2)
Parameter				
Implementation- Design, platform, coding,				
Optimization considerations(Memory, time, Resources, Costing)				
Thorough Testing of all modules				
Integration of modules and project as whole				

C. Documents

Grade (Grade Point)	Excellent (10-9)	Very Good (6-8)	Fair (3-5)	Poor (1-2)
Parameter				
Synopsis				
Project Report				
Quick references				
System manual				
Installation Guide				
Work Book				

D. Demonstration

Grade (Grade Point)	Excellent (10-9)	Very Good (6-8)	Fair (3-5)	Poor (1-2)
Parameter				
Project Presentation and Demonstration(User Interface, ease of use, usability)				
Understanding individual capacity & involvement in the project				
Team Work (Distribution of work, intra-team communication and togetherness)				
Outcomes / Usability				

E. Contest Participation / Awards, Publications and IPR

Grade (Grade Point)	Excellent (10-9)	Very Good (6-8)	Fair (3-5)	Poor (1-2)
Parameter				
Participation in various contests				
Appreciation and Awards				
Publications				
Copyright				
Patent				
Commercial value /product conversion of Work				

Annexure i: Final Synopsis (after approval of the project work)

4. Title of the Project: “**Mahadiscom Website Automation Using Uipath**”

5. Names of group members:

1. Mr. Khedkar Laxman Bhimrao

2 Mr. Hase Ashutosh Prakash

3. Mr. Langde Ritesh Vitthalrao Exam No: 27

1.

6. Name of the Guide: **Mrs. Khatal K. B**

Introduction

This project titled “*Content-Based Movie Recommender System*” focuses on solving the modern-day challenge of content overload faced by users across streaming platforms. With thousands of movies available, users often find it difficult to choose relevant content that matches their preferences. The system aims to address this problem by recommending movies that are similar in content to a user-provided input.

By using machine learning techniques such as **CountVectorizer** and **Cosine Similarity**, the system analyzes movie metadata—like genres, cast, crew, and keywords—to find similarities. The user interacts with the system through a simple web interface built using **Streamlit**, which instantly displays personalized movie recommendations based on their chosen movie.

This approach not only enhances the user experience but also demonstrates the practical application of machine learning and natural language processing techniques in a real-world scenario. The project is technically feasible, cost-effective, and aligns with current industry trends in AI-based recommendation systems.

Rationale:

In the digital age, personalized content discovery has become essential to improving user engagement and satisfaction. With the rapid growth of online streaming platforms and an ever-expanding library of movies, users are often overwhelmed by choices and struggle to find content that aligns with their preferences. A manual search through thousands of titles can be time-consuming and inefficient.

The *Content-Based Movie Recommender System* addresses this problem by offering intelligent, real-time movie suggestions based on the content of a user-selected movie. By leveraging machine learning techniques such as CountVectorizer and Cosine Similarity, the system analyzes features like genre, cast, and keywords to find similar movies and present them to the user through a simple and user-friendly interface.

This not only enhances the movie-watching experience but also reduces the effort required to find relevant content. The system empowers users with accurate and instant recommendations, improving user satisfaction while showcasing the practical application of artificial intelligence in solving real-world problems. Overall, the project supports digital innovation in media services and contributes to the advancement of user-centric software solutions.

Objective :

The primary objective of the project "*Content-Based Movie Recommender System*" is to build a functional and user-friendly application that recommends movies based on their content, using machine learning techniques. With the exponential growth of digital media, users often find themselves overwhelmed with options. This project aims to solve that problem by delivering intelligent, personalized recommendations using available movie metadata.

The system is designed with the following key objectives in mind:

- 1. To develop an accurate recommendation engine using content-based filtering**

The project focuses on analyzing content features like genres, cast, director, and keywords to suggest movies that are similar to the one selected by the user.

- 2. To apply CountVectorizer for feature extraction from text**

Movie attributes are combined and converted into a matrix of token counts using CountVectorizer, forming the basis for similarity calculation.

- 3. To measure similarity using Cosine Similarity**

Cosine Similarity is used to compare vectorized movie features and rank movies based on how closely they match the input movie.

- 4. To build a responsive web-based interface using Streamlit**

The system aims to provide an easy-to-use and interactive platform where users can input a movie name and instantly receive a list of recommendations.

- 5. To ensure real-time performance and usability**

The application is optimized for speed and accuracy, delivering recommendations within seconds, thus improving user experience.

- 6. To demonstrate the practical application of machine learning**

The project serves as a real-world implementation of machine learning, natural language processing, and web development skills in a cohesive system

Literature Review:

Recommendation systems have become a cornerstone of modern digital platforms, particularly in entertainment, retail, and e-learning. They are designed to enhance user experience by filtering large volumes of data and presenting users with personalized suggestions. The main types of recommendation systems include **collaborative filtering**, **content-based filtering**, and **hybrid systems**.

Collaborative filtering relies on user behavior—like ratings or viewing history—to recommend items. While effective, it faces limitations such as the **cold-start problem** (lack of data for new users or items) and **sparsity** (insufficient ratings). In contrast, **content-based filtering** recommends items by analyzing item features, such as genre, cast, or keywords. This makes it ideal for situations where user data is limited or unavailable.

The proposed system in this project uses a content-based approach, leveraging **textual metadata** from a movie dataset. Prior research, such as the work by *Lops et al. (2011)*, emphasizes that content-based systems can offer more explainable and stable recommendations, especially when combined with robust feature extraction techniques. Similarly, *Pazzani and Billsus (2007)* highlighted that user satisfaction increases when the system explains why an item was recommended, something easily achieved in content-based systems.

In our project, the use of **CountVectorizer** enables extraction of numerical features from textual movie descriptions, genres, and other metadata. These vectors are compared using **Cosine Similarity**, a method proven to be effective in document comparison and widely used in Natural Language Processing (NLP).

Additionally, several streaming platforms such as Netflix and Hulu use variations of content-based systems for their initial recommendations, especially when new users sign up. This validates the reliability and relevance of the approach adopted in our project.

In conclusion, the literature supports the use of content-based filtering as a viable method for recommendation systems, particularly in scenarios with rich item metadata and limited user data. Our project builds on these established methods and adapts them into a lightweight, interactive system that offers real-time movie recommendations to users.

Feasibility Study:

Before initiating the development of the *Content-Based Movie Recommender System*, a detailed feasibility study was conducted to assess the project's practicality, effectiveness, and execution within the given academic time frame. The study evaluates four major aspects: technical, economic, operational, and schedule feasibility.

1. Technical Feasibility

The technologies required for this project—Python, Pandas, Scikit-learn, and Streamlit—are all well-established, open-source, and supported by active developer communities. CountVectorizer and Cosine Similarity, the core algorithms used for feature extraction and similarity measurement, are widely used in Natural Language Processing and recommendation systems. The dataset used (e.g., TMDb 5000 Movie Dataset) is publicly available and comprehensive enough for building a proof-of-concept recommendation engine. Therefore, from a technical standpoint, the implementation is highly feasible.

2. Economic Feasibility

The project incurs minimal costs, as it is entirely developed using free and open-source tools. No proprietary software licenses or paid APIs are involved. All development and testing activities are performed using standard personal computers, making the solution cost-effective and ideal for academic purposes. Additionally, deploying the system locally through Streamlit eliminates the need for paid hosting services.

3. Operational Feasibility

The system is designed to operate via a lightweight and user-friendly web interface that requires no technical knowledge from the end user. Once deployed, users can simply input the name of a movie and receive relevant recommendations. This makes the system easy to operate and increases its potential usability for a broader audience. Moreover, due to its modular design, the system can be easily extended or improved in the future.

4. Schedule Feasibility

The project is broken into clear phases: data collection, preprocessing, model development, testing, and deployment. Each phase was planned according to the academic schedule, and most tools used (e.g., Python, Streamlit) allow for rapid development and debugging. Based on the current progress and resource availability, the system is expected to be completed well within the timeline set by the academic calendar.

Conclusion

The feasibility study concludes that the project is technically viable, economically affordable, operationally efficient, and schedulable within the available time. It is a suitable and sustainable project for demonstrating the application of machine learning in real-world scenarios.

Methodology / Planning of Work

Project Initiation:

- **Objective Definition:**

Clearly define the scope and objectives of the project: to develop a content-based movie recommendation system that suggests similar movies based on metadata such as genres, cast, and keywords, using machine learning techniques like CountVectorizer and Cosine Similarity.

- **Stakeholder Consultation:**

Discuss the goals, requirements, and expectations with the project guide and team members. Finalize deliverables, milestones, and responsibilities for each phase of the project.

Requirements Analysis:

- **Process Mapping:**

Outline the major steps involved in the project: data collection, preprocessing, feature extraction, similarity computation, and UI development. Create diagrams and data flow representations to visualize the entire recommendation process.

- **Technical Assessment:**

Evaluate available datasets, libraries (Scikit-learn, Pandas, Streamlit), and tools required for system development. Confirm compatibility with local systems and ensure necessary installations.

Design and Planning:

- **System Design:**

Design the architecture of the recommendation system, including data flow diagrams, flowcharts, and module specifications. Plan how the user input will be handled and how recommendations will be displayed.

- **Tool Selection:**

Finalize the use of Python as the programming language, Streamlit for the front-end, and Scikit-learn for model implementation. TMDb dataset is selected for initial development.

- **Resource Allocation:**

Distribute tasks among team members for dataset handling, algorithm implementation, UI development, documentation, and testing.

Development:

- **Model Building:**

Implement CountVectorizer to convert text into feature vectors and apply Cosine Similarity to compute similarity scores between movies.

- **Front-End Development:**

Use Streamlit to build a simple web interface that allows users to input a movie name and view similar movie suggestions instantly.

- **Error Handling:**

Integrate exception handling to manage cases like missing inputs, movie not found, or empty dataset entries.

- **Integration Testing:**

Combine backend logic and front-end interface. Test the complete flow to ensure consistent and accurate output.

Testing:

- **Unit Testing:**

Test individual modules such as vectorization, similarity computation, and movie search functionality for correctness.

- **End-to-End Testing:**

Verify the full functionality from user input to recommendation output, ensuring proper integration and user experience.

- **User Acceptance Testing (UAT):**

Present the application to peers and faculty for usability testing and gather feedback for final refinements.

Deployment:

- **Staging Deployment:**

Run the project locally for final testing. Ensure all features work as expected in a controlled environment.

- **Final Deployment:**

Prepare the application for presentation. If needed, demonstrate cloud deployment or GitHub repository version.

Monitoring and Maintenance:

- **Performance Monitoring:**

Observe application behavior during testing and demonstration. Ensure it runs smoothly without long response times.

- **Updates and Enhancements:**

Make improvements to the interface or recommendation logic based on feedback. Add additional features like poster images via API integration if time permits.

Documentation and Training:

- **Documentation:**

Prepare a detailed project report including introduction, methodology, code structure, and screenshots of the working system.

- **Training:**

Ensure all team members understand the system flow and are prepared for the final viva and project demonstration.

Evaluation and Feedback:

- **Performance Evaluation:**

Analyze the effectiveness of the recommendations, accuracy of similarity results, and responsiveness of the UI.

- **Feedback Collection:**

Gather feedback from guide, faculty, and users. Identify any areas of improvement for future development or deployment.

Project Closure:

- **Final Review:**

Ensure all deliverables (code, report, presentation) are complete. Match final outcomes against initial objectives.

- **Closure Report:**

Compile a summary of the project's achievements, challenges faced, technologies learned, and key lessons gained during the development journey.

Facilities Required for Proposed Work:

Software and Tools:

- **Python 3.x:** Primary programming language for implementing machine learning models and data processing.
- **Jupyter Notebook / VS Code:** Development environments used for coding, testing, and debugging the model.
- **Scikit-learn:** Machine learning library for implementing CountVectorizer and Cosine Similarity.
- **Pandas & NumPy:** Libraries for data preprocessing, transformation, and analysis.
- **Streamlit:** Lightweight web framework for developing an interactive recommendation interface.
- **Web Browser:** A modern web browser (Google Chrome, Mozilla Firefox) for running and testing the web application.

Hardware:

- **Development Workstations:** Computers capable of running machine learning applications and managing large datasets. Recommended specifications:
 - **CPU:** Multi-core processor (Intel i5 or i7 or AMD Ryzen equivalent)
 - **RAM:** Minimum 8 GB (16 GB preferred for smooth model execution)
 - **Storage:** SSD with at least 100 GB of free space
- **Internet Connectivity:** Stable and high-speed internet for downloading libraries, datasets, and accessing APIs.

Development Environment:

- **Workspace Setup:** A dedicated and quiet environment for model development, team collaboration, and interface design.
- **Source Control System:** GitHub is used for collaborative coding, version control, and project backup.

Access and Permissions:

- **Dataset Access:** Public access to TMDb or MovieLens datasets for training and testing the recommender system.
- **API Access:** TMDb API keys used to fetch additional movie details such as poster URLs and descriptions.

- **System Permissions:** Admin rights to install necessary libraries and tools on local systems.

Expected Outcomes

Improved Content Discovery and User Satisfaction:

- **Personalized Movie Suggestions:** The system provides relevant movie recommendations based on user-selected input, enhancing user engagement.
- **Instant Results:** The web interface delivers recommendations within seconds, eliminating the need for manual searches.
- **User-Friendly Interface:** Built using Streamlit, the application allows smooth interaction without requiring technical knowledge.

Enhanced Accuracy and Reliability:

- **Accurate Results:** The use of CountVectorizer and Cosine Similarity ensures high precision in similarity scoring between movies.
- **Consistent Recommendations:** The model performs uniformly across different movie inputs, maintaining consistency in suggestions.

Operational Efficiency:

- **Reduced Manual Effort:** Users no longer need to browse large catalogs manually to find movies of interest.
- **Time Savings:** Faster content discovery results in reduced browsing time and higher viewer satisfaction.

Scalability and Flexibility:

- **Scalable System:** The recommender can be expanded to larger datasets or improved with collaborative filtering in future phases.
- **Adaptable Interface:** The system can easily integrate additional features like user login, genre-based filters, or hybrid recommendations.

Cost-Effective Development:

- **Zero Licensing Costs:** All tools used are open-source, ensuring the project remains within academic or personal budget constraints.
- **Low Maintenance:** Once deployed, the system requires minimal upkeep unless major changes are made.

References

Books

1. Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow* (2nd ed.). O'Reilly Media.
2. McKinney, W. (2017). *Python for Data Analysis* (2nd ed.). O'Reilly Media.

Research Papers

3. Lops, P., De Gemmis, M., & Semeraro, G. (2011). Content-based Recommender Systems: State of the Art and Trends. In *Recommender Systems Handbook*, Springer.
4. Pazzani, M. J., & Billsus, D. (2007). Content-based Recommendation Systems. In *The Adaptive Web*, Springer.
5. Shinde, V., & Patil, R. (2021). Movie Recommendation System using Machine Learning. *International Journal of Scientific Research in Engineering and Management (IJSREM)*, Vol. 5, Issue 3.

Online Resources

6. Scikit-learn Documentation – <https://scikit-learn.org/stable/>
7. Streamlit Documentation – <https://docs.streamlit.io/>
8. TMDb API Documentation – <https://developer.themoviedb.org/docs>

Annexure ii: Partial Project Report (Semester II)

Certificates:



CHAPTER 1: INTRODUCTION

Abstract

In the age of digital content overload, finding relevant and personalized media has become a growing challenge

for users. The exponential growth of streaming services and digital movie libraries has led to a situation where users often feel overwhelmed by the vast number of choices available. This project, titled “*Content-Based Movie Recommender System*”, aims to simplify this process by offering smart movie recommendations using machine learning techniques.

This system is built using Python and employs libraries such as Scikit-learn for modeling, and Streamlit for creating a lightweight, interactive web interface. The project works by extracting features from movies (e.g., genres, cast, director, keywords) and computing similarity scores using CountVectorizer and Cosine Similarity. Once a user inputs a movie title, the system identifies and displays a list of the most similar movies.

The solution provides a clean and efficient way for users to discover new content without having to browse through endless titles. With applications in entertainment platforms, educational use, and academic learning,

the recommender system showcases how machine learning can enhance decision-making and user satisfaction. Furthermore, this model serves as a practical demonstration of NLP and recommender system concepts, offering both academic and real-world value.

General Introduction

In the modern digital world, users are constantly engaging with content across platforms like Netflix, Amazon Prime, YouTube, and more. With thousands of new movies released every year and streaming services constantly updating their catalogs, users often face difficulty choosing what to watch. Recommendation systems have emerged as a solution to this problem by helping users discover content that aligns with their tastes and preferences.

There are various types of recommendation systems, including collaborative filtering, content-based filtering, and hybrid models. This project focuses on the content-based approach, where recommendations are made by analyzing the features of the content itself, rather than relying on user ratings or historical behavior. This is particularly useful when dealing with new users or when user feedback data is not available.

The project uses Python’s Scikit-learn library to build a similarity model using CountVectorizer (to extract features from movie metadata) and Cosine Similarity (to compute similarity between movie vectors). The web interface is developed using Streamlit, which allows users to enter the name of a movie and view similar recommendations instantly in a user-friendly layout.

The system is designed to be scalable and can be improved in the future with hybrid features or user login for personalized results. It also serves as a learning base for understanding real-world implementation of machine learning, NLP, and web-based deployment.

Need of the System

The need for a movie recommendation system arises from the growing complexity and volume of content available on digital platforms. As the number of movies and streaming services increases, users often find themselves spending more time searching for a movie than actually watching one. This leads to decision fatigue and poor user experience.

Traditional browsing methods based on categories or popularity rankings are not personalized and may not align with individual preferences. A content-based recommender system solves this by suggesting movies similar to what the user already likes, based on genres, cast, crew, and plot elements.

By automating the recommendation process using machine learning, the system reduces the user's cognitive load, improves content discovery, and enhances overall satisfaction. It also eliminates the dependency on user ratings, which may be biased or unavailable for new content.

This system is particularly useful for:

- New users with no history
- Platforms with limited user behavior data
- Academic projects for demonstrating ML concepts
- Standalone tools for personal content management

In conclusion, the project provides a practical, low-cost, and efficient solution to content discovery using the power of data science and machine learning.

1.1 Problem Definition:

In today's digital world, users are overwhelmed with a vast number of movies and content options available across streaming platforms. Manually browsing and selecting a movie that aligns with personal preferences is time-consuming, inefficient, and often leads to poor user experience. This creates a strong need for an intelligent recommendation system that can assist users in discovering relevant movies quickly and accurately. A content-based approach that leverages movie metadata such as genre, cast, and keywords can automate this process and provide personalized recommendations without requiring user history or ratings.

1.2 Objective of the Project:

- To develop a content-based movie recommender system using Python and machine learning techniques.
- To extract features from movie metadata (e.g., genres, keywords, cast) using CountVectorizer.
- To calculate similarity between movies using Cosine Similarity and recommend the most relevant ones.
- To create a user-friendly web interface using Streamlit for easy interaction and real-time results.
- To minimize manual browsing efforts and improve the efficiency of movie discovery.
- To ensure system scalability and reliability by structuring the code modularly.
- (Optional) To integrate TMDb API for enhancing the UI with movie posters and details.

CHAPTER 2: LITERATURE SURVEY

Recommendation systems have become a vital component of modern digital platforms, particularly in the entertainment sector. Their primary objective is to simplify content discovery by suggesting relevant items based on user preferences or content features. Among the various techniques used, content-based filtering is effective when user behavior data is limited or unavailable.

This project adopts a content-based approach, focusing on analyzing metadata such as genre, cast, director, and keywords. Text data is vectorized using CountVectorizer, and Cosine Similarity is used to measure how closely related two movies are. These methods are supported by libraries like Scikit-learn and Pandas, making the development both practical and scalable.

Unlike collaborative filtering, which requires historical user data, content-based systems provide recommendations solely based on the features of items, making them ideal for new users and smaller applications. The system is deployed using Streamlit, allowing a simple and interactive web interface for users to input a movie name and receive recommendations instantly.

In conclusion, the literature supports the use of content-based recommendation systems as reliable, interpretable, and well-suited for academic projects and real-world applications where content metadata is available.

CHAPTER 3: IMPLEMENTATION

Implementation:

The *Content-Based Movie Recommender System* was implemented using Python, along with machine learning libraries such as Scikit-learn, Pandas, and Streamlit for web interface development. The goal was to develop an interactive and lightweight system capable of recommending movies similar to a user-provided title, based on metadata like genre, keywords, and cast.

Step 1: Data Collection and Preprocessing

The dataset containing movie information (genres, cast, overview, keywords, etc.) was sourced from TMDb (The Movie Database). Missing values were filled or removed, and irrelevant columns were dropped. Key features were combined into a single string for each movie.

Screenshot suggestion: Pandas DataFrame displaying combined features column.

Step 2: Feature Extraction using CountVectorizer

The combined text features were converted into numerical vectors using Scikit-learn's CountVectorizer. This transformed the movie metadata into a matrix form, enabling further computation.

Screenshot suggestion: Console output of vectorized matrix shape.

Step 3: Similarity Calculation using Cosine Similarity

A cosine similarity matrix was computed to measure the degree of similarity between every pair of movies. This matrix formed the core logic behind the recommendation system.

Screenshot suggestion: Partial similarity matrix or visual representation of top matches.

Step 4: Building the Web Interface with Streamlit

The recommendation engine was deployed using Streamlit, allowing users to input a movie title and get a list of 5 similar movies along with their posters (using TMDb API).

Screenshot suggestion: Streamlit UI showing input box and recommended movie list.

Step 5: Handling User Input and Output Generation

The user's input is matched with the closest movie name in the dataset. If found, the system fetches the top 5 similar movies and displays their titles and posters. If the input is invalid or not found, an error message is shown.

Screenshot suggestion: Output screen showing recommendations or an error message.

Step 6: Error Handling and Optimization

Basic exception handling was implemented to manage edge cases like invalid movie names or empty input. The code was optimized to minimize latency in delivering recommendations.

Screenshot suggestion: Python try-except block from code.

CHAPTER 4: SOFTWARE DETAILS

Python:

- **Version:** Python 3.9 or later
- **Purpose:** Core programming language for data preprocessing, model implementation, and application logic.
- **Download Link:** <https://www.python.org/downloads/>

Jupyter Notebook / Visual Studio Code:

- **Purpose:** Development environments used for writing, testing, and debugging Python code.
- **Recommendation:** VS Code preferred for Streamlit-based UI integration.
- **Download Link (VS Code):** <https://code.visualstudio.com/>
- **Download Link (Jupyter):** Comes with Anaconda – <https://www.anaconda.com/>

Pandas Library:

- **Purpose:** For loading datasets, cleaning, manipulating data, and combining features.
- **Installation:** `pip install pandas`
- **Documentation:** <https://pandas.pydata.org/docs/>

Scikit-learn:

- **Purpose:** Provides tools like CountVectorizer and cosine_similarity used in feature extraction and similarity measurement.
- **Installation:** `pip install scikit-learn`
- **Documentation:** <https://scikit-learn.org/stable/>

NumPy:

- **Purpose:** Used for supporting matrix operations and numerical computations.
- **Installation:** `pip install numpy`

- **Documentation:** <https://numpy.org/doc/>

Streamlit:

- **Purpose:** Web application framework to create an interactive front-end for the recommender system.
- **Installation:** `pip install streamlit`
- **Run Command:** `streamlit run app.py`
- **Documentation:** <https://docs.streamlit.io/>

Web Browser:

- **Requirement:** Google Chrome, Firefox, or Microsoft Edge for testing the Streamlit app.
- **Chrome Download:** <https://www.google.com/chrome/>
- **Edge Download:** <https://www.microsoft.com/edge>

TMDb API (Optional):

- **Purpose:** Used for fetching movie posters dynamically.
- **Requires:** TMDb API key (free registration required).
- **API Docs:** <https://developer.themoviedb.org/docs>
 -

CHAPTER 5: RESULT AND ANALYSIS

The *Content-Based Movie Recommender System* successfully provided accurate and relevant movie suggestions based on user input. After entering a movie name, the system processed the input using CountVectorizer and Cosine Similarity, and returned the top five most similar movies from the dataset. The recommendations were consistent and aligned with the features of the selected movie, such as genre, keywords, and cast.

The system was tested with multiple movie titles, including both popular and lesser-known films, and consistently generated meaningful suggestions. The web interface, built using Streamlit, allowed for smooth user interaction and delivered results instantly. The integration of the TMDb API (optional) added visual appeal by displaying posters for the recommended movies.

Overall, the system demonstrated efficiency, speed, and accuracy in generating recommendations. The output met the project's objectives and confirmed the effectiveness of using content-based filtering techniques for real-world applications.

CHAPTER 6: SYSTEM OVERVIEW

The *Content-Based Movie Recommender System* is designed to suggest similar movies based on a user-selected input title by analyzing the content features of each movie. The system uses a publicly available movie dataset containing metadata such as genres, cast, crew, and keywords. These features are combined and transformed into numerical vectors using CountVectorizer.

The similarity between movies is calculated using Cosine Similarity, resulting in a matrix that helps identify movies closely related to the input. When a user enters a movie name through the Streamlit-based interface, the system searches for it in the dataset, retrieves its index, fetches similarity scores, and displays the top five recommended movies.

The system includes input validation, error handling, and an optional integration with the TMDb API to enhance output with movie posters. It is modular, efficient, and can be extended for larger datasets or hybrid recommendation models. The application offers a user-friendly solution for quick and relevant movie discovery.

Phase – V: Prototyping

CONCLUSION

The development of the *Content-Based Movie Recommender System* using Python successfully demonstrates

how machine learning can simplify and enhance the process of discovering relevant movies. By automating content comparison using CountVectorizer and Cosine Similarity, the system effectively eliminates the need for manual searching and guesswork in selecting movies. It significantly improves the user experience by providing accurate, real-time recommendations through an interactive web interface built with Streamlit.

This project highlights the practical application of natural language processing and recommendation algorithms in solving real-world problems. It is lightweight, efficient, and easy to use, making it suitable for individual users, educational use, or as a foundation for more advanced recommendation systems. The system's modular design also makes it adaptable for future improvements, such as integrating collaborative filtering or hybrid models.

Overall, the project demonstrates the power of machine learning in personalizing digital experiences and sets a solid base for further exploration in the field of intelligent recommender systems.

References

1. Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow* (2nd ed.). O'Reilly Media.
2. McKinney, W. (2017). *Python for Data Analysis* (2nd ed.). O'Reilly Media.
3. Scikit-learn Documentation. Retrieved from <https://scikit-learn.org/stable/>
4. Streamlit Documentation. Retrieved from <https://docs.streamlit.io/>
5. The Movie Database (TMDb) API Documentation. Retrieved from <https://developer.themoviedb.org/docs>
6. Lops, P., De Gemmis, M., & Semeraro, G. (2011). Content-Based Recommender Systems. In *Recommender Systems Handbook*. Springer.

