

FOUR WEEK TRAINING REPORT

at

ThinkNEXT Technologies Private Limited

SUBMITTED IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE
AWARD OF THE COMPUTER SCIENCE AND ENGINEERING DEGREE

BACHELOR OF TECHNOLOGY
Computer Science and Engineering



JUNE-JULY

SUBMITTED BY:
Jashan
URN-2302560

Department of Computer Science and Engineering
Guru Nanak Dev Engineering College
Ludhiana, 141006

ThinkNEXT Technologies Private Limited



ThinkNEXT®
Innovation at every step...
ISO 9001:2015 Certified Company



Scan and verify your Certificate
Certificate ID:433375
Ref.No. TNT/C-25/17335

Certificate

This Certificate do hereby recognizes that

Jashan S/o Ramesh Kumar

has successfully completed Industrial Training Program

from 23rd June 2025 to 21st July 2025

in AI & ML using Python Grade A

For ThinkNEXT Technologies Pvt. Ltd.

Nagma Khan
Authorised Signatory

Member Training Division

ThinkNEXT Technologies Pvt. Ltd.

Munish Mittal
Director

Director

Corporate Office: S.C.F 113, Sector 65, Mohali

Outstanding Excellent Very Good Good Satisfactory
100-90% 89-80% 79-70% 69-60% 59-50%



ISO Certified
CII
Member of Confederation
of Indian Industry



Candidate's Declaration

GURU NANAK DEV ENGINEERING COLLEGE, LUDHIANA CANDIDATE'S DECLARATION

I, "**Harshpreet Singh**", hereby declare that I have undertaken one-month training at "**ThinkNEXT Technologies Private Limited**" during the period from 24 June 2025 to 21 July 2025 in partial fulfillment of the requirements for the award of the degree of **B.Tech (Computer science and Engineering)** at **GURU NANAK DEV ENGINEERING COLLEGE, LUDHIANA**. The work which is being presented in the training report, submitted to the Department of Computer science and Engineering at GURU NANAK DEV ENGINEERING COLLEGE, LUDHIANA, is an authentic record of the training work.

Jashan
URN: 2302560

The one-month industrial training Viva–Voce Examination of _____
has been held on _____ and accepted.

Signature of Internal Examiner

Signature of External Examiner

Abstract

This report summarizes the work undertaken during my four-week industrial training at ThinkNEXT Technologies Private Limited, focusing on **Artificial Intelligence and Machine Learning (AI/ML)** with practical exposure to real-world problem-solving using data-driven approaches. The training was designed to enhance my understanding of machine learning algorithms, data preprocessing, feature engineering, and model evaluation techniques.

During this period, I developed a project titled “**College Dropout Prediction System**”, which aims to predict the likelihood of students dropping out based on various academic, personal, and socio-economic factors. The system utilizes machine learning models such as **Logistic Regression, Decision Tree, and Random Forest** to analyze historical student data and classify students as either “likely to drop out” or “likely to continue.” The project also emphasizes data visualization, accuracy comparison among models, and interpretation of results to identify the most influential factors leading to dropout.

This application demonstrates the effective use of Python libraries such as **Pandas, NumPy, Scikit-learn, and Matplotlib**, providing valuable insights that educational institutions can use for early intervention and student retention strategies. Throughout the training, I gained hands-on experience in data preprocessing, model training, and evaluation, along with practical exposure to building predictive models from scratch. This experience strengthened my technical foundation in AI and ML, improved my analytical thinking, and enhanced my ability to apply machine learning concepts to real-world educational challenges.

Acknowledgment

I would like to express my sincere gratitude to **ThinkNEXT Technologies Private Limited** for providing me the opportunity to undergo a four-week industrial training in **Artificial Intelligence and Machine Learning (AI/ML)**. The experience significantly enhanced my understanding of machine learning applications and helped me gain practical knowledge of how predictive analytics systems are designed and implemented.

I am especially thankful to my mentor and the professionals at ThinkNEXT for their continuous guidance, valuable suggestions, and constructive feedback throughout the training period. Their support was instrumental in helping me complete my final project, **Student Dropout Prediction System**, which incorporated multiple features such as data visualization, machine learning models (Logistic Regression and Random Forest), real-time prediction capabilities, feature importance analysis, and an interactive web interface using Streamlit.

I also wish to acknowledge the support of the **Department of Computer Science and Engineering, Guru Nanak Dev Engineering College, Ludhiana**, and **I.K. Gujral Punjab Technical University** for incorporating this training as a part of the B.Tech curriculum, thereby providing students like me an opportunity to gain hands-on industry exposure in AI/ML technologies and strengthen our data science and analytics skills.

List of Tables

2.1	Python Data Types and Examples	14
2.2	Python & ML Tools/Functions and Their Usage	16
3.1	Project Features and Implementation Status	21

List of Figures

2.1	Screenshot of the main project dashboard, showing key visualizations.	18
2.2	The real-time prediction interface with user input fields.	19
2.3	The real-time prediction interface with user input fields report.	19
3.1	The main dashboard of the Student Dropout Prediction System.	21
3.2	Performance and accuracy score of model.	22
3.3	Feature on which dropout of student depend.	23

Contents

1	Introduction	8
1.1	Background	8
1.2	Objective of the Training	10
1.3	Overview of Artificial Intelligence & Machine Learning	11
1.3.1	Artificial Intelligence Overview	11
1.3.2	Machine Learning (ML) Overview	11
1.4	Importance of Artificial Intelligence and Machine Learning in Modern Technology	12
1.5	Scope of Training	12
2	Training Work Undertaken	14
2.1	Week 1 – Introduction to Python, Data Types, Functions & OOP Concepts	14
2.2	Week 2 – Collections, preprocessing, Manipulation & Mini Project	14
2.3	Week 3 – Machine Learning Algorithms, Model Training Evaluation	15
2.4	Week 4 – Model Integration, Dashboard Development & Deployment	16
2.5	Final Project Development – Student Dropout Prediction System	17
3	Results and Discussions	20
3.1	Overview of the Project	20
3.2	Implementation Results	20
3.2.1	Interactive Dashboard and Data Analysis	20
3.2.2	Machine Learning Models and Prediction System	22
3.2.3	Feature Analysis and Performance Metrics	23
3.3	Discussion and Observations	23
3.4	Summary	24
4	Conclusion and Future Scope	25
4.1	Conclusion	25
4.2	Future Scope	26
4.3	Reflection and Learning Outcome	26
References		27

Chapter 1

Introduction

1.1 Background

In the rapidly evolving field of computer science, industrial training serves as a vital bridge between theoretical knowledge and practical implementation. While classroom learning primarily focuses on algorithms, statistics, and mathematical foundations, industry training introduces students to real-world machine learning applications, tools, and workflows. The true value of engineering education is realized when students apply their conceptual understanding to solve actual data-driven problems — and that is exactly what this industrial training aimed to achieve.

As a Computer Science and Engineering (CSE) student, my academic journey had already introduced me to the principles of data science, statistical analysis, and algorithmic thinking. However, the modern AI industry demands more than theoretical understanding. Companies seek data scientists and ML engineers who can work efficiently with frameworks, libraries, and tools that enhance model performance and deployment capabilities. Therefore, the 4-week industrial training at ThinkNEXT Technologies Pvt. Ltd., Mohali was an invaluable opportunity to gain industry-level exposure in machine learning applications using Python and its powerful ecosystem.

Python, along with libraries like scikit-learn, pandas, and Streamlit, has become the de-facto standard for developing machine learning applications. These tools allow data scientists to build sophisticated predictive models, perform complex data analysis, and create interactive web applications — all using a single technology stack. The combination of these tools has revolutionized how developers approach machine learning projects, making it easier to go from data analysis to production-ready applications.

During the training, I was introduced to the complete machine learning pipeline, starting from the basics of data preprocessing and gradually progressing to complex model development and web deployment techniques. I learned how feature engineering and selection serve as the core building blocks in ML, allowing developers to create robust and accurate prediction systems efficiently. The interactive nature of Streamlit made it easier to visualize real-time predictions, which significantly improved my understanding of model behavior and performance analysis.

Moreover, I explored various machine learning algorithms, particularly focusing on classification techniques that play a crucial role in predictive analytics. Understanding concepts like Logistic Regression, Random Forests, and evaluation metrics helped me appreciate how production-grade ML systems are developed in real-world projects. I also learned the importance of model interpretation and visualization, which enables stakeholders to understand and trust the predictions made by the system.

The training sessions at ThinkNEXT Technologies were structured in a way that balanced both theoretical understanding and practical implementation. Each module included demonstrations, hands-on practice, and guidance from experienced mentors. The instructors emphasized not just the technical aspects but also the business context and ethical considerations in AI/ML applications. Planning a machine learning application, choosing appropriate algorithms, and implementing efficient data processing pipelines.

One of the most valuable aspects of this industrial training was learning about real-world ML project workflows. From setting up development environments with Python and essential libraries to model optimization and performance tuning, I became familiar with tools and techniques used by professional data scientists. I also understood the importance of maintaining clean data pipelines, feature engineering practices, and reusable model components, which are critical in production ML systems.

In addition to learning core machine learning concepts, I also explored how different algorithms like Logistic Regression and Random Forest can be used to predict student outcomes. These techniques were particularly useful in developing my final project — "Student Dropout Prediction System", a web-based application with advanced features like real-time predictions, feature importance analysis, interactive visualizations, and comprehensive model performance metrics. The project served as a complete reflection of what I learned during the training and allowed me to integrate multiple AI/ML concepts into a single functional system.

The training experience also provided me with exposure to modern data visualization principles. I learned how effective visualization techniques — such as confusion matrices, feature importance plots, correlation analyses, and interactive charts — significantly enhance model interpretation. Streamlit's flexible components allowed me to create an intuitive and interactive user interface without complex web development.

Beyond technical learning, the training enhanced my analytical and debugging abilities. Each model development phase introduced challenges that encouraged me to think critically about data preprocessing, feature selection, and model evaluation. Handling imbalanced datasets, tuning hyperparameters, and optimizing model performance gave me the confidence to tackle real-world machine learning problems independently.

Furthermore, I learned the significance of version control and data version control systems like Git and DVC, which are essential for maintaining ML projects in production environments. Although the project was individual, I practiced maintaining versions using GitHub repositories to simulate real-world collaborative development and to track model iterations effectively.

Another key takeaway was understanding the machine learning project lifecycle — which includes data collection, preprocessing, feature engineering, model development, evaluation, and deployment. The mentors explained how each stage contributes to creating reliable, production-ready ML systems. For instance, in the evaluation phase, I learned to use various metrics and visualization techniques to ensure model robustness and reliability.

The environment at ThinkNEXT Technologies Pvt. Ltd. was highly conducive to learning. The mentors were supportive, encouraging experimentation with different algorithms, and motivating us to explore advanced concepts beyond basic machine learning. Their practical insights into industry trends, model deployment strategies, and ethical AI considerations helped me connect academic learning with actual career goals.

Overall, this 4-week industrial training transformed my understanding of AI/ML development. It taught me how to think like a data scientist — analyzing problems, design-

ing efficient models, and writing clean, reproducible code. The training also reinforced the importance of continuous learning, as AI technologies evolve rapidly in the industry.

In conclusion, this training emphasizes how essential it is for engineering students to experience hands-on industry exposure in AI/ML. Through this program, I gained not only technical expertise in Python, scikit-learn, and Streamlit but also a broader understanding of machine learning as a structured, iterative process. By the end of the training, I successfully transitioned from theoretical understanding to practical ML application development, marking a significant milestone in my journey as a future AI/ML engineer.

1.2 Objective of the Training

The main objective of this 4-week industrial training was to gain practical knowledge of Artificial Intelligence (AI) and Machine Learning (ML) for developing real-world predictive models and to understand the modern data science workflow followed in the industry. While classroom education focuses on theoretical algorithms and formulas, industry training emphasizes data preprocessing, model building, and result interpretation, which are the real skills required in professional environments. At ThinkNEXT Technologies Pvt. Ltd., the training was systematically designed to achieve the following objectives:

- Understand the fundamentals of AI and ML, including supervised and unsupervised learning, classification, regression, clustering, and model evaluation metrics.
- Learn Python programming and essential data science libraries such as NumPy, Pandas, Matplotlib, and Scikit-learn for efficient data handling and visualization.
- Develop data preprocessing and feature engineering skills, including handling missing values, encoding categorical data, and scaling numerical features.
- Work on real-world datasets to analyze real use of different ML algorithms.
- Understand the workflow of model training, testing, and validation, and develop a complete project to apply all learned concepts — from planning to implementation.
- Enhance debugging, evaluation, and visualization skills to interpret results effectively and make data-driven decisions.
- Learn professional data science practices, including documentation, reproducibility, and version control using Git and GitHub.

In summary, the objective of this industrial training was not only to learn a new technology but also to develop the mindset of an industry-ready developer. It was about transforming theoretical learning into a structured, creative, and functional development process.

1.3 Overview of Artificial Intelligence & Machine Learning

1.3.1 Artificial Intelligence Overview

Artificial Intelligence (AI) is a branch of computer science focused on creating systems capable of performing tasks that normally require human intelligence. These include learning from data, reasoning, problem-solving, perception, and decision-making. AI combines principles from mathematics, statistics, and computer science to enable machines to mimic human-like thinking and behavior.

Key features of AI include:

- **Learning Ability:** Systems can learn and improve automatically through experience and data analysis.
- **Automation:** Enables computers to perform repetitive and complex tasks without manual intervention.
- **Reasoning and Decision-Making:** Helps systems make predictions and logical decisions based on patterns and data.
- **Natural Language Processing (NLP):** Allows machines to understand and respond to human language.
- **Computer Vision:** Enables machines to interpret and process visual information like images and videos.

1.3.2 Machine Learning (ML) Overview

Machine Learning (ML) is a subset of AI that focuses on building algorithms that allow computers to learn from data and make predictions or decisions without being explicitly programmed. It uses statistical techniques to recognize patterns and improve model performance over time.

Key advantages of ML include:

- **Data-Driven Decision Making:** Learns patterns from past data to predict future outcomes.
- **Automation of Analytical Tasks:** Reduces human effort by automating processes like classification and prediction.
- **Wide Range of Algorithms:** Supports supervised, unsupervised, and reinforcement learning techniques.
- **Scalability:** Can handle large datasets efficiently for real-world applications.
- **Continuous Improvement:** Model accuracy enhances with more data and re-training.

During the training, I learned how AI and ML complement each other — AI provides the conceptual foundation for building intelligent systems, while ML supplies the algorithms and techniques to make those systems learn and improve autonomously.

1.4 Importance of Artificial Intelligence and Machine Learning in Modern Technology

In today's data-driven world, organizations rely heavily on Artificial Intelligence (AI) and Machine Learning (ML) to make faster, smarter, and more accurate decisions. With the growing volume of data and the need for automation, AI and ML have become essential technologies across industries such as healthcare, education, finance, and transportation. They empower systems to analyze patterns, predict outcomes, and adapt intelligently — reducing human effort and increasing efficiency.

AI and ML enable companies to derive insights from data and automate complex tasks that were once considered manual and time-consuming. From recommendation engines to fraud detection and predictive analytics, these technologies are transforming how businesses operate and innovate.

Some of the major reasons for AI and ML growing importance include:

- **Data-Driven Insights:** Help organizations make evidence-based decisions and identify hidden trends.
- **Automation:** Reduces repetitive tasks and increases overall productivity.
- **High Accuracy:** Improves prediction and classification results through continuous learning.
- **Wide Applicability:** Used in various domains such as education, healthcare, retail, and cybersecurity.
- **Community & Open Source Tools:** Libraries like TensorFlow, Scikit-learn, and PyTorch are continuously enhanced by global developers and researchers.
- **Future-Ready Technology:** Acts as the foundation for emerging fields like deep learning, robotics, and intelligent systems.

In my training project, College Dropout Prediction System, I experienced how machine learning algorithms can effectively analyze student data and predict dropout risks with high accuracy. This demonstrated the real-world potential of AI and ML in solving educational challenges and supporting better decision-making through predictive analytics.

1.5 Scope of Training

The scope of this 4-week industrial training program was extensive, covering both theoretical understanding and practical implementation of Artificial Intelligence (AI) and Machine Learning (ML) concepts. It was designed to transform a beginner into a capable data science enthusiast, equipped to build intelligent predictive systems using Python and modern ML tools.

The training began with the fundamentals of machine learning, gradually moving toward data preprocessing, feature engineering, model building, and performance evaluation. By the end of the course, I had gained the confidence to independently design, train, and test predictive models on real-world datasets.

The scope included:

- Understanding the fundamentals of **AI and ML algorithms** such as classification, regression, and clustering.
- Learning Python programming and essential libraries like **pandas, NumPy, and scikit-learn**.
- Performing **data cleaning, encoding, and scaling** to prepare datasets for model training.
- Implementing and comparing multiple ML models such as **Logistic Regression and Random Forest Classifier**.
- Visualizing data and model performance using Matplotlib, Seaborn, and Plotly.
- Working on mini-projects and a final project — **Student Dropout Prediction System** — to apply learned concepts.
- Gaining exposure to **industry workflows, debugging techniques, version control, and documentation practices**.

This training not only enhanced my technical knowledge but also strengthened my analytical thinking, problem-solving ability, and project management skills. It laid a strong foundation for my future career in Artificial Intelligence and Data Science, and gave me a clear vision of how AI/ML technologies are transforming modern decision-making systems.

Chapter 2

Training Work Undertaken

2.1 Week 1 – Introduction to Python, Data Types, Functions & OOP Concepts

The first week of training served as the foundation for the entire **AI and Machine Learning** journey. We began by exploring **Python**, one of the most popular programming languages for data science and ML development. I learned about variables, data types, control structures, and functions, which helped in writing clean and reusable code for analytical tasks. The trainer emphasized the importance of writing efficient and readable code using Python's built-in libraries and modular programming practices. The week also introduced **Object-Oriented Programming (OOP)** concepts such as

classes, objects, constructors, inheritance, and polymorphism. I implemented small programs that simulated real-world data-driven problems like student data management and basic statistical calculations. In addition, we covered **exception handling** using try-except-finally blocks, which helped manage runtime errors effectively during data operations. By the end of the week, I could write well-structured Python programs

and gained confidence in logic building, debugging, and file handling — laying a strong foundation for the upcoming machine learning modules.

Table 2.1: Python Data Types and Examples

Data Type	Description	Example
int	Whole numbers (integer values)	<code>age = 21</code>
float	Decimal numbers (floating-point)	<code>price = 99.99</code>
str	Textual data	<code>name = "Harshpreet"</code>
bool	Boolean values: True or False	<code>isActive = True</code>
list	Ordered collection of elements	<code>marks = [85, 90, 88]</code>

2.2 Week 2 – Collections, preprocessing, Manipulation & Mini Project

The second week of training focused on more advanced aspects of data preprocessing and manipulation, which are essential for preparing datasets for machine learning. I learned

to use Pandas and NumPy, two core Python libraries for efficient data handling and numerical computation. Through hands-on practice, I explored data structures such as Series and DataFrames, performing operations like data filtering, grouping, merging, and aggregation.

We focused on handling missing values, outliers, and categorical variables, and applied transformations such as normalization and encoding to make the data model-ready. These skills are crucial for maintaining data consistency and accuracy in AI/ML projects.

The week also covered file handling in Python — reading and writing CSV files, managing directories, and automating data import/export processes. I practiced working with multiple datasets and learned how to handle exceptions while loading or saving data files.

To consolidate all the concepts learned, we were assigned a Mini Project, where I performed a small-scale Student Data Analysis task. The project included:

- Cleaning and organizing raw student data.
- Handling missing or inconsistent entries.
- Performing descriptive statistics and visualizing key features using Matplotlib.
- Writing results to new output files using Pandas.

This mini project significantly improved my confidence in working with real datasets and prepared me for the next phase — applying Machine Learning algorithms for predictive modeling in Week 3.

2.3 Week 3 – Machine Learning Algorithms, Model Training Evaluation

The third week marked the transition from data preprocessing to machine learning model development. I was introduced to supervised learning algorithms and the process of building predictive models for real-world applications.

We started with an overview of the machine learning workflow, which includes:

- Data preprocessing and feature engineering
- Model selection and training
- Performance evaluation and tuning

I learned about classification algorithms, which are essential for predicting categorical outcomes such as student dropout status. The trainer explained Logistic Regression and Random Forest Classifier in detail, including their mathematical foundation, parameters, and practical use cases.

Next, I explored model training and evaluation techniques:

- Splitting data into training and testing sets (80-20 split)
- Training models using cross-validation to avoid overfitting
- Calculating performance metrics like accuracy, precision, recall, and F1-score

- Hyperparameter tuning to optimize model performance

Additionally, I learned how to implement real-time predictions and generate probability scores to classify students into Low, Medium, and High risk categories.

By the end of Week 3, I was able to build and evaluate machine learning models, interpret results using visualization tools, and apply these models to make data-driven predictions for student dropout risk. This week laid the groundwork for integrating the models into a fully interactive application in Week 4.

Table 2.2: Python & ML Tools/Functions and Their Usage

Tool / Function	Purpose	Example Usage
pandas.DataFrame	Store and manipulate tabular data	df = pd.DataFrame(data)
pandas.read_csv()	Load CSV datasets	df = pd.read_csv("students.csv")
numpy.array	Create numerical arrays for computation	arr = np.array([1, 2, 3])
train_test_split	Split data into training and testing sets	X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
LogisticRegression()	Build classification model	model = LogisticRegression()
RandomForestClassifier()	Build ensemble classification model	rf = RandomForestClassifier()
matplotlib.pyplot.plot()	Visualize data trends	plt.plot(x, y)
seaborn.heatmap()	Visualize correlations	sns.heatmap(df.corr(), annot=True)

2.4 Week 4 – Model Integration, Dashboard Development & Deployment

The fourth week focused on making the machine learning project interactive, dynamic, and deployable. Key learnings and tasks included:

- **User Input Handling:** Created input forms in Streamlit to accept student parameters (age, grades, study time, family support, etc.). Implemented real-time validation to ensure accurate and clean data collection.
- **Dashboard and Visualization:** Built an interactive dashboard to display data insights using Matplotlib, Seaborn, and Plotly. Visualized feature correlations, class distribution, and risk levels with dynamic charts and graphs.
- **Model Integration:** Integrated trained Logistic Regression and Random Forest models into the web application. Enabled real-time predictions of student dropout risk with probability scores and classification into Low, Medium, or High risk categories.

- **Notifications and Feedback:** Added interactive elements to provide instant feedback to users about predicted risks. Displayed actionable recommendations for each student based on model output.
- **API Handling and File Management:** Loaded datasets dynamically from CSV files and implemented automated preprocessing pipelines. Learned how backend computations interact with the front-end dashboard for smooth data flow.
- **Deployment:** Prepared the application for Streamlit Cloud deployment with minimal dependencies. Ensured the app was scalable, responsive, and user-friendly for educational administrators.

By the end of Week 4, I successfully deployed a fully functional, interactive ML-based application. This final project, Student Dropout Prediction System, integrated all prior learning — from data handling and model training to real-time predictions and dashboard visualization — demonstrating practical application of AI/ML in educational analytics.

2.5 Final Project Development – Student Dropout Prediction System

The final week of the training was dedicated to the culmination of all learned concepts: the development and deployment of the **Student Dropout Prediction System**. This project served as a practical application of the entire machine learning workflow, from data handling to creating an interactive, user-facing web application.

Project Overview

The Student Dropout Prediction System is a powerful web application designed to identify students at risk of dropping out. By leveraging machine learning models, the system provides educational institutions with real-time predictions and actionable insights, enabling timely interventions. The entire application was built using Python, with Streamlit for the frontend, and scikit-learn for the backend modeling.

Key Features Developed

During this week, the following key features were implemented to create a comprehensive and user-friendly tool:

- **Interactive Dashboard:** A multi-page dashboard was built using Streamlit to allow users to navigate between different sections like Data Overview, Model Performance, Prediction, and Feature Analysis.
- **Dual Model Integration:** Both the trained **Logistic Regression** and **Random Forest Classifier** models were integrated into the backend, allowing users to see predictions from both for comparison.
- **Real-Time Prediction System:** An input form was created with over 30 student parameters. The system processes this input in real-time to generate a dropout risk assessment, complete with probability scores and a clear classification (Low, Medium, or High risk).

- **Dynamic Visualizations:** Using Plotly and Matplotlib, interactive charts were developed to display feature importance, correlation analysis, and model performance metrics, helping users understand the key factors influencing predictions.

Development and Deployment Process

The development followed a structured, module-by-module approach:

1. **UI Design:** The user interface was crafted with Streamlit's widgets to ensure an intuitive and responsive experience for administrators.
2. **Backend Integration:** The pre-trained models and the data processing pipeline (including one-hot encoding and scaling) were connected to the frontend input forms.
3. **Visualization:** Performance metrics (Accuracy, Precision, Recall) and feature analysis charts were integrated into dedicated pages on the dashboard.
4. **Deployment:** Finally, the application was prepared for deployment on Streamlit Cloud, with all dependencies managed to ensure smooth and scalable operation.

This final project successfully combined data science principles with web development, resulting in a practical tool that demonstrates the business impact of machine learning in the education sector.

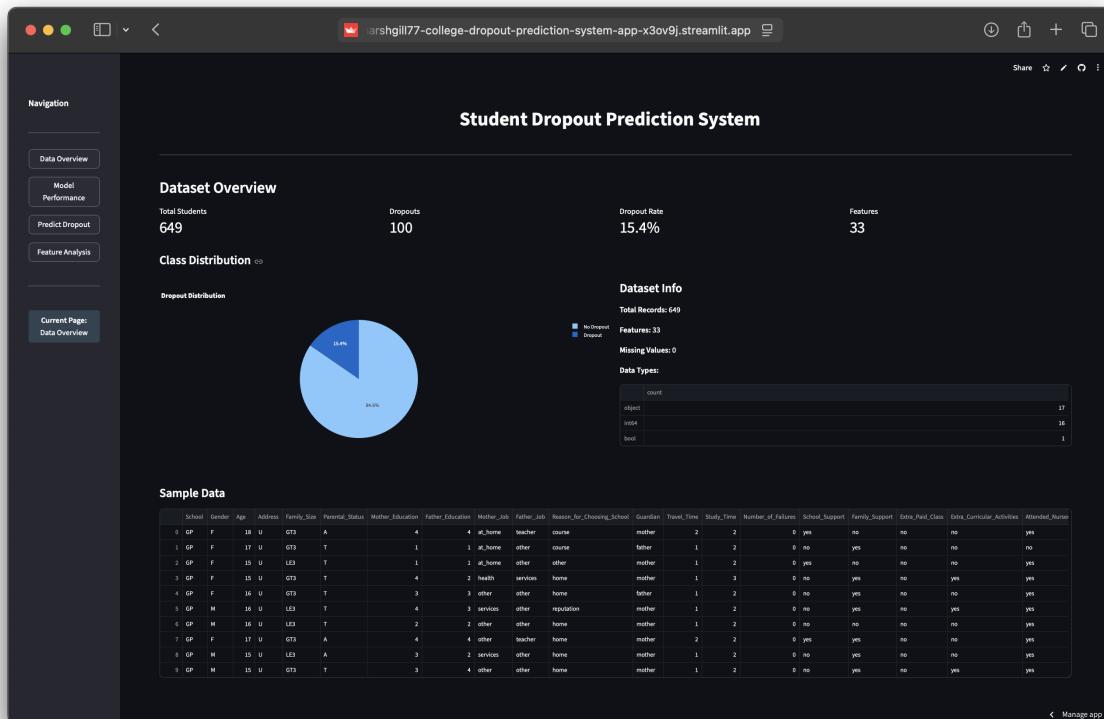


Figure 2.1: Screenshot of the main project dashboard, showing key visualizations.

Predict Student Dropout

Enter the student information below to predict their dropout risk. The system will use both Logistic Regression and Random Forest models for prediction.

Basic Information

School: GP

Gender: F

Age: 18

Address: U

Family Size: LE3

Parental Status: T

Academic Performance

Travel Time (1-4): 2

Study Time (1-4): 2

Number of Failures: 0

Parental Education

Education levels: 0=none, 1=primary, 2=5th-9th grade, 3=secondary, 4=higher

Mother Education: 0

Father Education: 0

Support Systems

School Support: yes

Family Support: yes

Extra Paid Classes: yes

Extra Curricular Activities: yes

Attended Nursery: yes

Wants Higher Education: yes

[Manage app](#)

Figure 2.2: The real-time prediction interface with user input fields.

Predict Dropout Risk

Prediction Results

Logistic Regression Model

LOW RISK - Student will likely STAY

Probability of dropping out: 0.6%

Probability of staying: 99.4%

Random Forest Model

LOW RISK - Student will likely STAY

Probability of dropping out: 4.0%

Probability of staying: 96.0%

Combined Prediction

Overall Risk Level: LOW

Average dropout probability: 2.3%

Recommendations

Student appears to be on track for success. Continue current support.

[Manage app](#)

Figure 2.3: The real-time prediction interface with user input fields report.

Chapter 3

Results and Discussions

3.1 Overview of the Project

The Student Dropout Prediction System is a comprehensive machine learning web application developed during my 4-week industrial training. The application provides a versatile environment for predicting student dropout risk, allowing educational institutions to intervene proactively. Built with Python, Streamlit, and scikit-learn, the system operates as a powerful analytical tool without requiring complex installations.

This project is designed to address the critical challenge of student retention. Unlike static reports, this system offers an interactive dashboard for real-time predictions. It not only predicts dropout likelihood but also provides insights into the key factors influencing these predictions, such as academic performance, social factors, and personal characteristics. The application features a multi-page interface for data overview, model performance analysis, real-time prediction, and feature importance visualization.

The app's modular structure allows smooth navigation between its different analytical functions. Its data-driven approach ensures that predictions are based on historical patterns, providing accuracy and reliability. The integration of two different ML models (Logistic Regression and Random Forest) allows for comparative analysis, making the system a unique and efficient tool for educational administrators.

3.2 Implementation Results

The Student Dropout Prediction System has been implemented as a multi-functional analytical tool, with several integrated modules designed for seamless usability. Each module demonstrates the practical application of Python libraries, machine learning concepts, and web deployment with Streamlit. Below is a detailed description of the app's features and their implementation status.

3.2.1 Interactive Dashboard and Data Analysis

A key component of the system is its interactive and user-friendly dashboard built with Streamlit. It allows administrators to explore the dataset and understand its underlying patterns without writing any code. Key features include:

- A comprehensive data overview page showing the raw dataset.
- Visualizations for class distribution (Dropout vs. Enrolled) to check for data imbalance.

Table 3.1: Project Features and Implementation Status

Feature	Description	Status
Interactive Dashboard	Real-time data visualization and responsive UI with multi-page navigation.	Completed
Data Analysis	Comprehensive dataset overview, class distribution, and feature correlation studies.	Completed
Machine Learning Models	Integration of Logistic Regression and Random Forest classifiers for prediction.	Completed
Real-Time Prediction	A form to input student parameters and get instant risk assessment with probability scores.	Completed
Feature Importance	Visualization of the most influential factors affecting dropout predictions.	Completed
Deployment	Application prepared for cloud deployment using Streamlit Cloud.	Completed

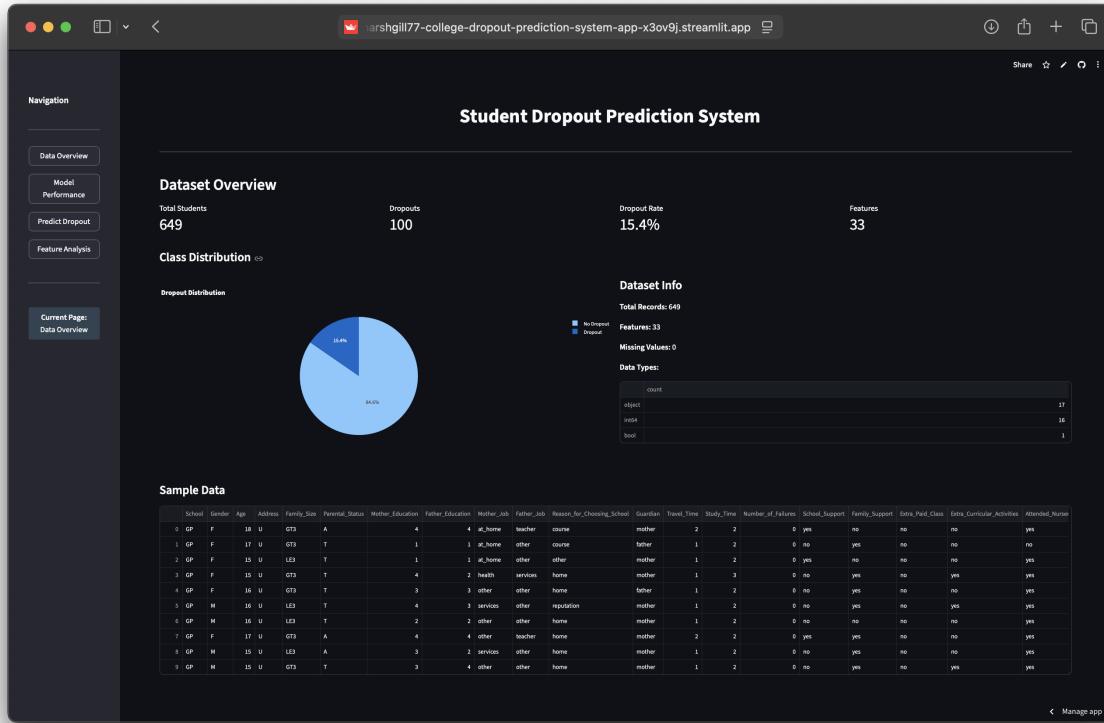


Figure 3.1: The main dashboard of the Student Dropout Prediction System.

- Interactive correlation heatmaps to study the relationships between different student attributes.

The dashboard interface is designed for simplicity and efficiency, utilizing Streamlit's widgets for dynamic rendering of charts and data tables. Users can quickly navigate between different analytical pages using a sidebar menu, providing instant access to all functions.

3.2.2 Machine Learning Models and Prediction System

The core of the project relies on two robust classification models from scikit-learn to predict dropout risk. Key features include:

- **Logistic Regression:** A reliable linear model providing baseline performance with an **accuracy of approximately 93%**.
- **Random Forest Classifier:** A more complex ensemble model that captures non-linear relationships, achieving a higher accuracy of **around 99%**.
- **Real-Time Prediction:** A dedicated page with an input form for over 30 student parameters (e.g., previous grades, study time, family support).
- **Risk Classification:** The system processes the inputs and provides a prediction, classifying the student's risk as Low, Medium, or High, along with a probability score.

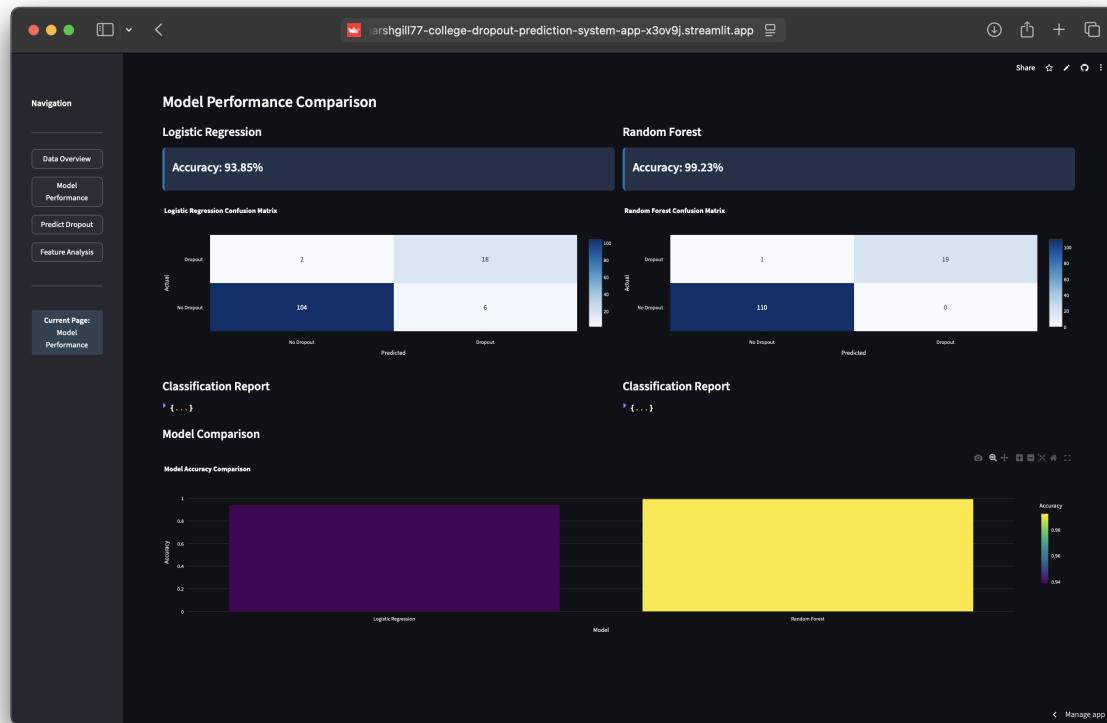


Figure 3.2: Performance and accuracy score of model.

This dual-model approach allows administrators to compare outputs and gain more confidence in the final assessment, making the tool both powerful and transparent.

3.2.3 Feature Analysis and Performance Metrics

To enhance usability and trust in the model's predictions, the system includes features for model interpretability:

- **Feature Importance:** A dedicated page displays a chart of the most important factors that the Random Forest model uses to make predictions. This helps identify key drivers of student dropout.
- **Key Predictive Factors:** The analysis consistently highlights academic performance (grades, failures), social factors (family support), and personal characteristics (age, health) as top predictors.
- **Performance Metrics:** The dashboard includes a section detailing the performance of each model, with metrics such as Accuracy, Precision, Recall, and F1-Score, ensuring transparency about the system's reliability.

These features provide a deeper understanding of why a prediction is made, moving the system from a "black box" to an insightful decision-support tool.

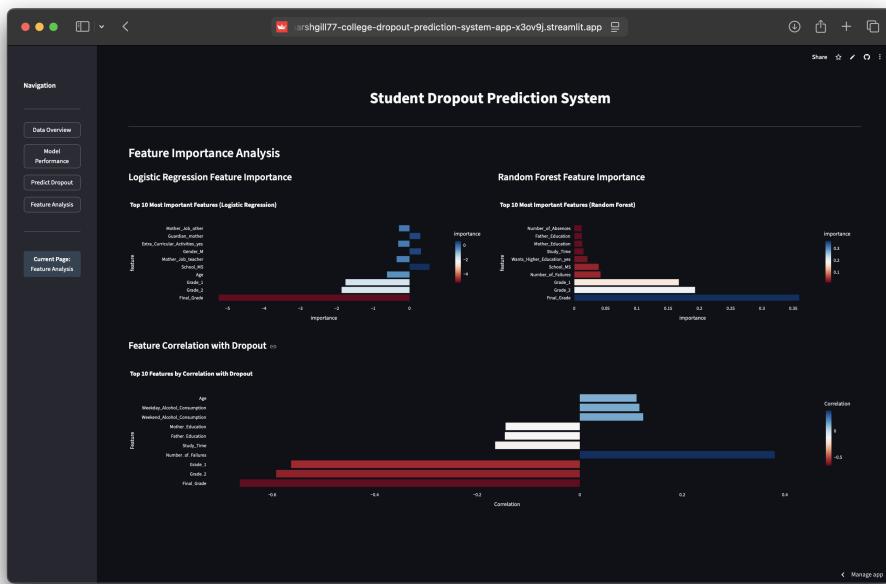


Figure 3.3: Feature on which dropout of student depend.

3.3 Discussion and Observations

The development of the Student Dropout Prediction System provided practical insights into creating end-to-end machine learning applications. Key observations include:

- Streamlit's framework allowed for rapid UI development and iteration, making it ideal for creating data-centric web applications with minimal front-end code.
- The scikit-learn pipeline was critical for creating a reproducible and clean workflow for data preprocessing, training, and prediction.

- Model interpretation features, like feature importance plots, were crucial for building trust and making the application's outputs actionable for end-users.
- The performance difference between Logistic Regression and Random Forest highlighted the trade-offs between model simplicity and predictive power.

Overall, the combination of data analysis, predictive modeling, and an interactive user interface makes this project a comprehensive tool for educational institutions.

3.4 Summary

In conclusion, the Student Dropout Prediction System successfully demonstrates the practical implementation of all concepts learned during the 4-week training:

- Python programming fundamentals and data manipulation with Pandas and NumPy.
- Machine learning model implementation with scikit-learn.
- Data visualization with Matplotlib, Seaborn, and Plotly.
- Web application development and deployment with Streamlit.

The project is a fully functional, data-driven web application that provides real-time predictions, interactive data analysis, and model performance evaluation. This chapter highlighted the implementation, results, and observations of the project, demonstrating how theoretical training translates into a functional and practical machine learning solution.

Chapter 4

Conclusion and Future Scope

4.1 Conclusion

The 4-week industrial training at ThinkNEXT Technologies Pvt. Ltd., Mohali provided me with a thorough understanding of the end-to-end machine learning workflow, combining both theoretical learning and practical implementation. The highlight of the training was the design and development of the Student Dropout Prediction System, a fully functional, data-driven web application.

Through this training, I was able to apply several core concepts of data science:

- **Python for Data Science:** I strengthened my understanding of Python and its core data science libraries, including Pandas for data manipulation, NumPy for numerical operations, and Matplotlib/Seaborn for visualization. Writing clean, efficient code was essential for building a reliable data pipeline.
- **Machine Learning Development:** The training enhanced my ability to build, train, and evaluate predictive models using scikit-learn. I learned about data pre-processing, feature engineering, model selection, and hyperparameter tuning. Implementing both Logistic Regression and Random Forest models helped me understand the practical trade-offs between different algorithms.
- **Interactive Web Application:** The project relies on Streamlit to provide an interactive and user-friendly interface. I learned to build multi-page dashboards, create data visualizations, and deploy a standalone web application, making the machine learning model accessible to non-technical users.
- **Model Interpretation:** Incorporating features like feature importance plots was crucial for making the model's predictions understandable and trustworthy. This introduced me to the important field of explainable AI (XAI).

Overall, this training has bridged the gap between academic theory and practical data science. I gained hands-on experience in creating a complex, interactive machine learning application from scratch, managing both the backend logic (data pipeline, model training) and the frontend design (UI layout, interactive charts). The development of this system also helped me improve my problem-solving skills, learn to debug ML models efficiently, and manage a project from conception to deployment. The project simulates a real-world data science workflow, making me confident in my ability to undertake ML projects independently in the future.

4.2 Future Scope

Although the Student Dropout Prediction System is fully functional, there is significant scope for future enhancements to make it even more versatile and robust. Some potential improvements include:

- **Cloud Deployment and Integration:** Integrating the system with cloud platforms like AWS or Azure would allow for greater scalability. Creating a REST API would enable the prediction service to be integrated directly into existing university Student Information Systems (SIS).
- **Automated Model Retraining:** Implementing an automated pipeline (e.g., using Kubeflow or Airflow) to periodically retrain the model with new student data would ensure the system's accuracy is maintained over time and prevents model drift.
- **Advanced ML and Deep Learning Models:** The current system can be enhanced by incorporating more advanced algorithms like Gradient Boosting (XGBoost, LightGBM) or even deep learning models (e.g., TabNet, or LSTMs if sequential data becomes available) to potentially improve predictive accuracy.
- **Enhanced Explainability (XAI):** Future versions could include more advanced model explanation techniques, such as SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations), to provide detailed reasons for individual student predictions.
- **Batch Prediction Functionality:** Adding a feature for administrators to upload a CSV file of multiple students and receive batch predictions would significantly improve the system's efficiency for institutional use.
- **Prescriptive Analytics:** Expanding the system from predictive to prescriptive analytics. Instead of just flagging at-risk students, the system could suggest specific, personalized intervention strategies based on the student's key risk factors.
- **User Authentication and Roles:** Adding user login/signup functionality would allow for personalized dashboards and role-based access for different stakeholders (e.g., administrators, counselors, faculty).

4.3 Reflection and Learning Outcome

Developing the Student Dropout Prediction System allowed me to experience the complete machine learning project lifecycle:

- Conceptualization and problem definition under mentor guidance.
- Data exploration, cleaning, and preprocessing using Pandas.
- Feature engineering and model selection with scikit-learn.
- UI/UX design and interactive dashboard development using Streamlit.

- Testing, debugging, and model evaluation for performance refinement.

This experience has strengthened my technical, analytical, and problem-solving skills, preparing me for real-world data science projects. I also learned to manage time effectively, plan development steps, and prioritize features, which are essential skills in the technology industry. Furthermore, I now have practical exposure to the Python data science ecosystem, interactive web development, and the end-to-end process of building and deploying a machine learning application, which forms a strong foundation for a career in AI and ML.

Bibliography

- [1] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Pearson, 4th edition, 2020.
- [2] Leo Breiman. "Random Forests". *Machine Learning*, 45(1):5–32, 2001.
- [3] F. Pedregosa et al. "Scikit-learn: Machine Learning in Python". *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [4] Wes McKinney. "Data Structures for Statistical Computing in Python". In *Proceedings of the 9th Python in Science Conference*, pages 51–56, 2010.
- [5] Charles R. Harris et al. "Array programming with NumPy". *Nature*, 585(7825):357–362, 2020.
- [6] Streamlit Inc. *Streamlit Documentation*. Available at: <https://docs.streamlit.io/>, 2023.
- [7] David W. Hosmer Jr., Stanley Lemeshow, and Rodney X. Sturdivant. *Applied Logistic Regression*. John Wiley & Sons, 3rd edition, 2013.
- [8] Ryan S.J.d. Baker and Kalina Yacef. "The State of Educational Data Mining in 2009: A Review and Future Visions". *Journal of Educational Data Mining*, 1(1):3-17, 2009.
- [9] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [10] Plotly Technologies Inc. *Plotly Python Graphing Library*. Available at: <https://plotly.com/python/>, 2023.