

**NANNAPANENI VENKATA RAO COLLEGE
OF ENGINEERING & TECHNOLOGY**
(Approved by AICTE and Affiliated to JNTUK- Kakinada)
**DEPARTMENT OF
COMPUTER SCIENCE & ENGINEERING**



Long Term Internship

(A.Y: 2024-2025)

Design and Developed by



**ANDHRA PRADESH STATE
COUNCIL OF HIGHER
EDUCATION**

(A STATUTORY BODY OF GOVERNMENT OF ANDHRA PRADESH)

NANNAPANENI VENKAT RAO
COLLEGE OF ENGINEERING & TECHNOLOGY

(Approved by AICTE, Affiliated to JNTUK- Kakinada)

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



CERTIFICATE

This is to Certify that the project Titled **“Student Performance Prediction”** is a bonafide work carried out by **MADDU. RAJESH (217T1A0542)**, in partial fulfilment of the Requirements for the award of Bachelor Of Technology degree in the Department of **Computer Science & Engineering** at **Nannapaneni Venkat Rao College of Engineering & Technology** during the Academic year 2024-2025. The result embedded in this report have not been submitted to any other University for the award of any Degree.

Signature of the Co-Ordinator

Signature of the HOD

External Examiner

Ints/2025/337400



CERTIFICATE OF INTERNSHIP

This Internship program certificate is proudly Awarded to

MADDU RAJESH (217T1A0542)

*For his/her outstanding completion of the internship program on **Python with Machine Learning** at corporate innovative technologies private limited from
1st, DECEMBER 2024 to 31st, MARCH 2025*

S.K. Khairing
CEO



S. Chaitanya
Program Manager

CONTENTS

TITLES	PAGE.NO
1.Abstract	1
2.Company Profile - INTS	2
3. Internship Objectives	3
4. Internship Roles and Responsibilities	5
5. Technologies and Tools Used	7
6. Machine Learning Overview	9
7. Algorithms Implemented	11
8. Project: Student Performance Prediction	14
9. System	17
10. Internship Experience & Learnings	21
11. Future Enhancements	22
12. Daily Work Log (Selected Days)	26
13. Conclusion	29
14. References	30
15. Annexure	31

1. Abstract

This report presents the internship experience on "Python with Machine Learning" at INTS. The focus was on understanding core ML algorithms and implementing a project titled Student Performance Prediction. The report includes theoretical background, tools used, methodologies, results, and future scope. The internship provided an excellent platform to bridge academic knowledge with industrial applications, reinforcing concepts through hands-on project work and collaboration with professionals. It also emphasized developing a solution-oriented mindset.

2. Company Profile - INTS

INTS is a technology firm specializing in Artificial Intelligence, Data Science, and Software Engineering solutions. It provides end-to-end solutions for clients using state-of-the-art ML frameworks and real-time systems. INTS focuses on developing customized technology platforms to optimize business operations through AI integration. The company values innovation and offers diverse learning opportunities for interns. The internship program at INTS is structured to enhance practical knowledge and professional development, encouraging students to solve real-world problems under expert guidance.

3. Internship Objectives and Learning Outcomes

1. Gaining Hands-on Experience with Python and ML Libraries

One of the primary objectives of this internship is to develop a strong foundational understanding of Python programming in the context of machine learning. Python has emerged as the dominant language in data science due to its simplicity, flexibility, and robust ecosystem of libraries. Through practical assignments and guided projects, I aim to deepen my familiarity with Python libraries such as:

- a) **NumPy** and **Pandas** for data manipulation.
- b) **Matplotlib** and **Seaborn** for data visualization.
- c) **Scikit-learn** for implementing core machine learning models.
- d) **TensorFlow** and **Keras** for deep learning techniques (if applicable).

Hands-on experience with these tools will provide me with the necessary skills to build and deploy machine learning models efficiently. By engaging in real-world coding tasks, I expect to enhance my programming logic, debugging skills, and ability to write clean, modular code.

2. Learning to Clean, Analyze, and Visualize Data

Data preprocessing is a critical step in the machine learning pipeline. In this internship, I aim to acquire in-depth knowledge and experience in handling raw datasets, which often include missing values, outliers, and inconsistent formats. By working with diverse datasets, I will develop expertise in:

- a) Identifying and treating missing or erroneous data.
- b) Feature engineering and transformation.
- c) Data normalization and standardization techniques.
- d) Exploring data through statistical summaries and visualization.

Visualization, in particular, plays a crucial role in understanding data distribution and identifying patterns or anomalies. Through tools like Seaborn, Matplotlib, and Plotly, I plan to practice creating meaningful graphs, heatmaps, and correlation matrices that support data-driven decisions.

3. Implementing Core Machine Learning Algorithms

A core component of the internship is to gain proficiency in implementing machine learning algorithms from scratch and using library-based approaches. The focus will be on both supervised and unsupervised learning models, such as:

- a) Linear and Logistic Regression
- b) Decision Trees and Random Forests
- c) Support Vector Machines (SVM)
- d) K-Nearest Neighbors (KNN)
- e) Naive Bayes

- f) Clustering algorithms like K-Means and DBSCAN

By understanding how these algorithms work internally, and applying them to datasets, I will gain insights into model behavior, advantages, limitations, and best-use scenarios. Emphasis will also be placed on understanding the assumptions behind each model and tuning their hyperparameters for optimal results.

4. Understanding Model Evaluation and Optimization

Model evaluation is essential for assessing the accuracy and generalizability of machine learning solutions. During the internship, I aim to build expertise in evaluating models using key metrics like:

- a) Accuracy, Precision, Recall, and F1-Score
- b) Confusion Matrix and ROC-AUC
- c) Cross-validation techniques
- d) Bias-variance tradeoff

Additionally, I will explore optimization strategies, including grid search and randomized search, to fine-tune models. Learning about overfitting, underfitting, and regularization methods (like L1 and L2 penalties) will further help me in improving model robustness.

5. Applying ML to Real-World Problems (e.g., Resume-Job Profile Matching)

One of the highlights of the internship is the opportunity to apply machine learning in solving a real-world use case—**matching resumes to job profiles**. This involves:

- a) Text preprocessing techniques such as tokenization, stemming, and vectorization (TF-IDF, Word2Vec).
- b) Natural Language Processing (NLP) for understanding the semantics of resumes and job descriptions.
- c) Building similarity-based or classification models to recommend suitable matches.
- d) Evaluating model effectiveness through real-world accuracy and user feedback.

This project will expose me to end-to-end development, including data collection, modeling, validation, and result interpretation—skills highly valuable in a professional ML role.

6. Enhancing Collaboration and Reporting Skills

Beyond technical capabilities, this internship aims to strengthen my soft skills, particularly in communication, collaboration, and documentation. Working in a team setting requires:

- a) Using version control systems like **Git** and platforms like **GitHub**.
- b) Participating in code reviews, stand-up meetings, and collaborative debugging.
- c) Preparing professional documentation and presentations to communicate findings and model performance.

By engaging in team projects and preparing reports, I will learn how to present data insights and ML outcomes clearly and concisely to both technical and non-technical stakeholders.

4. Internship Roles and Responsibilities

During the internship, I undertook a diverse range of tasks aimed at building core competencies in data analysis, machine learning, and collaborative development. These roles not only enhanced my technical proficiency but also improved my ability to function effectively in a team-driven environment. Below is a detailed overview of my key responsibilities:

1. Data Analysis Using Pandas and NumPy

One of my primary responsibilities was to perform comprehensive data analysis using Python's powerful libraries, **Pandas** and **NumPy**. These tools enabled me to explore, manipulate, and gain insights from structured datasets efficiently.

- a) **Pandas** provided essential functions for loading data from various formats (CSV, Excel, JSON), handling missing values, filtering and grouping data, and computing basic statistics. I used it extensively for exploratory data analysis (EDA).
- b) **NumPy** was instrumental in managing large arrays and matrices of numeric data. It enabled efficient computation and helped me understand underlying numerical operations such as broadcasting, vectorization, and performance optimization in data-heavy tasks.

Through these tools, I learned how to detect anomalies, perform feature extraction, and generate summary statistics. These skills laid the foundation for more complex machine learning tasks later in the internship.

2. Data Visualization Using Matplotlib and Seaborn

A key part of data analysis is **data visualization**, and I was responsible for transforming raw data into meaningful charts and plots using **Matplotlib** and **Seaborn**.

- a) **Matplotlib** allowed me to create customizable plots such as line graphs, bar charts, histograms, and scatter plots. It provided complete control over styling, axes, and annotations.
- b) **Seaborn**, built on top of Matplotlib, offered high-level functions for creating attractive and informative statistical graphics such as box plots, violin plots, pair plots, heatmaps, and correlation matrices.
- c) By visualizing data, I was able to uncover patterns, relationships, and trends that were not immediately apparent from raw data alone. These visuals also helped in presenting findings to mentors and team members in a clear and concise manner.

3. Implementing ML Algorithms Using Scikit-learn

A significant part of the internship involved implementing **machine learning algorithms** using **Scikit-learn**, one of the most widely-used libraries in the Python ML ecosystem.

I worked on the following tasks:

- a) **Supervised Learning Models:** Implemented algorithms like Linear Regression, Logistic Regression, Decision Trees, Random Forests, and Support Vector Machines (SVM) for classification and prediction tasks.

- b) **Unsupervised Learning Models:** Used K-Means Clustering and Principal Component Analysis (PCA) to uncover hidden patterns and reduce dimensionality in datasets.
- c) **Model Evaluation:** Used metrics like accuracy, precision, recall, F1-score, and confusion matrix to evaluate performance. I also performed cross-validation and hyperparameter tuning (using GridSearchCV) to improve model generalization.

This hands-on experience improved my understanding of how to choose appropriate algorithms based on data characteristics and problem types.

4. Writing Python Scripts for Data Preprocessing

Before applying ML models, the data had to be cleaned and transformed. I was responsible for writing **Python scripts** that performed:

- a) Missing value imputation (mean/median/mode, or advanced techniques)
- b) Categorical data encoding (one-hot encoding, label encoding)
- c) Scaling and normalization (MinMaxScaler, StandardScaler)
- d) Text preprocessing (tokenization, stopwords removal, stemming/lemmatization for NLP tasks)
- e) Feature engineering (creating new variables that improved model performance)

These preprocessing scripts ensured the datasets were model-ready and helped maintain consistency and reproducibility across experiments.

5. Preparing Documentation and Progress Reports

An important non-technical responsibility during the internship was the preparation of **project documentation and periodic progress reports**. This included:

- a) Documenting code, assumptions, and data sources.
- b) Writing brief descriptions for each experiment and its results.
- c) Keeping logs of performance metrics and insights.
- d) Preparing weekly reports summarizing accomplishments, blockers, and next steps.

This practice significantly improved my ability to communicate complex technical processes in a structured and readable manner.

6. Participating in Regular Team Meetings and Updates

As part of the internship program, I actively **participated in daily or weekly team meetings**, where we discussed:

- a) Project updates and timelines
- b) Technical challenges and solutions
- c) Feedback on ongoing work
- d) Coordination with other team members

These meetings fostered a collaborative environment and helped me understand team workflows, task management tools (like Trello or Jira), and professional communication standards.

5. Technologies and Tools Used

Throughout the internship, I worked with a variety of technologies and tools that supported different stages of the machine learning workflow—from data preprocessing to model building and deployment. These tools were essential in enabling effective development, collaboration, and documentation of my work.

1. Programming Language: Python

Python served as the core programming language for all tasks during the internship. Its clean syntax, vast library support, and active community make it ideal for data analysis and machine learning. I used Python for:

- a) Data cleaning and exploration
- b) Algorithm implementation
- c) Writing reusable scripts
- d) Automating repetitive tasks
- e) Integrating models into simple applications

Python's versatility made it the backbone of both backend processes and analytical workflows during my internship.

2. Libraries: NumPy, Pandas, Matplotlib, Seaborn, Scikit-learn, Joblib

These Python libraries were instrumental in building a complete machine learning pipeline:

- a) **NumPy**: Used for numerical operations, array manipulations, and mathematical functions. It allowed for efficient handling of multi-dimensional data.
- b) **Pandas**: My go-to library for loading, exploring, and manipulating structured data. It offered powerful DataFrame structures for organizing and analyzing data.
- c) **Matplotlib & Seaborn**: These visualization libraries helped bring data to life. Matplotlib provided the foundation for building plots, while Seaborn added advanced statistical graphs, heatmaps, and pair plots with elegant styling.
- d) **Scikit-learn**: The core library used to build, train, test, and evaluate machine learning models. I implemented algorithms such as linear regression, decision trees, KNN, and clustering methods. Scikit-learn also includes tools for preprocessing, model selection, and metrics.
- e) **Joblib**: This tool was used to **serialize (save) and deserialize (load)** machine learning models and pipelines. It helped preserve trained models for future use or deployment, making the ML process more efficient and modular.

3. IDE/Notebook Environments: Jupyter Notebook and VS Code

Both development environments played unique roles in enhancing productivity:

- a) **Jupyter Notebook**: Ideal for exploratory data analysis, visualizations, and iterative coding. I used it extensively for testing ideas, documenting findings, and combining code with visual output in a readable format.

- b) **Visual Studio Code (VS Code):** A lightweight yet powerful IDE that I used for writing modular Python scripts, Flask apps, and managing project files. It also integrated well with Git for version control and extensions for linting and debugging.
- c) Using both environments in tandem gave me flexibility in how I approached coding and project organization.

4. Version Control: GitHub

GitHub played a central role in **version control and collaboration**. I used Git to:

- a) Commit changes incrementally
- b) Track revisions and code history
- c) Create branches for new features or experiments
- d) Collaborate through pull requests and code reviews

5. Spreadsheets: Excel and Google Sheets

While Python was the main tool for data handling, **Excel and Google Sheets** were used for:

- a) Quickly reviewing and sharing sample data
- b) Performing light data manipulation or summary calculations
- c) Preparing structured reports for non-technical stakeholders

6. Other Tools: Google Colab, Anaconda, and Flask

- a) Several supporting tools also played important roles:
- b) **Google Colab:** A cloud-based Jupyter environment with free GPU/TPU support. I used it for running heavier ML workloads and sharing notebooks online without local setup.
- c) **Anaconda:** This distribution simplified package management and environment setup. It came bundled with most data science libraries and helped maintain project-specific environments through conda.
- d) **Flask:** A lightweight web framework used to build basic web apps and APIs for machine learning models. I used Flask to deploy a simple ML model as a web interface, demonstrating real-world application and accessibility.

6. Machine Learning Overview

Machine Learning (ML) is a rapidly growing field that falls under the broader umbrella of **Artificial Intelligence (AI)**. It focuses on creating algorithms that can **learn patterns from data** and make decisions or predictions without being explicitly programmed for each task. Instead of relying on hard-coded rules, ML systems adapt and improve as they are exposed to more data.

The fundamental idea behind machine learning is to enable machines to learn from past experiences (data) and generalize that knowledge to unseen situations. This ability has led to ML being widely used in various domains, including finance, healthcare, e-commerce, recommendation systems, computer vision, and natural language processing.

How It Works:

- a) At its core, machine learning involves:
- b) Input Data: Historical or real-time data that the model can learn from.
- c) Algorithms: Mathematical models that analyze data and find patterns.
- d) Training: The process where models adjust themselves by minimizing errors.
- e) Prediction/Inference: Using the trained model to make decisions on new, unseen data.
- f) Evaluation: Testing the model's accuracy and reliability using specific metrics.

The power of machine learning lies in its ability to **self-improve**. As it receives more data or feedback, the model becomes more accurate over time. This makes ML systems scalable, adaptable, and capable of handling complex and dynamic environments.

6.1 Types of Machine Learning

Machine Learning is typically divided into **three main types**, each defined by the nature of the data and the task the model is designed to perform:

1. Supervised Learning

Definition: Supervised learning involves training a model on a **labeled dataset**, where the input data is paired with the correct output.

Goal: The objective is to learn a mapping from inputs to outputs so that the model can make accurate predictions on unseen data.

- a) Decision Trees
- b) Random Forests
- c) Support Vector Machines (SVM)
- d) K-Nearest Neighbors (KNN).

2. Unsupervised Learning

Definition: In unsupervised learning, the model is trained on **unlabeled data**. The system must find patterns, relationships, or structures within the data on its own.

Goal: The aim is to identify hidden patterns or groupings without prior knowledge of the outcomes.

Examples:

- a) **Clustering:** Grouping similar data points (e.g., customer segmentation).
- b) **Dimensionality Reduction:** Reducing the number of input variables (e.g., using PCA).
- c) **Common Algorithms:**
- d) K-Means Clustering
- e) Hierarchical Clustering
- f) DBSCAN
- g) Principal Component Analysis (PCA)
- h) Autoencoders (in deep learning)

3. Reinforcement Learning

Definition: Reinforcement Learning (RL) is a goal-directed learning approach where an **agent** interacts with an **environment** and learns to make decisions by receiving **rewards or penalties** based on its actions.

- a) **Goal:** To learn a policy that maximizes cumulative rewards over time through trial and error.
- b) **Core Concepts:**
- c) **Agent:** The decision-maker.
- d) **Environment:** The world with which the agent interacts.
- e) **Action:** A move the agent can make.
- f) **Reward:** Feedback from the environment.
- g) **Policy:** The strategy the agent follows.

7. Algorithms Implemented

In any intelligent recruitment or job matching platform, algorithms serve as the core decision-making components. These algorithms drive how candidates are matched with job openings, how resumes are analyzed, and how recommendations are made. The following sections outline some of the key algorithms that have been implemented in the system, focusing on their purpose, functioning, and importance in delivering accurate and efficient results.

1. Resume Parsing Algorithm

Purpose:

To extract structured data from unstructured resumes.

Functioning:

Resumes are typically uploaded in PDF or DOCX formats. These documents are parsed using natural language processing (NLP) techniques to identify and extract relevant fields such as:

- a) Name
- b) Contact Information
- c) Skills
- d) Education
- e) Work Experience
- f) Certifications

The algorithm uses a **combination of keyword-based matching** and **named entity recognition (NER)** to classify text into appropriate categories. Libraries such as **spaCy** or **NLTK** are used to tokenize the text and extract entities such as dates, degrees, universities, companies, etc.

2. Skill Matching Algorithm

Purpose:

To determine the compatibility between a candidate's skill set and a job's requirements.

Functioning:

Once resumes and job descriptions are parsed, both are converted into a list of skills. The skill matching algorithm computes a **similarity score** using one or more of the following methods:

3. Job Recommendation Algorithm

Purpose:

To recommend suitable job listings to candidates based on their profiles.

Functioning:

This is a **content-based filtering** approach that uses the candidate's resume content (skills, job titles, location preferences) to recommend jobs. It also leverages **vector representations** of job descriptions and resumes to find the most semantically similar listings.

In advanced versions, **collaborative filtering** or **hybrid models** are introduced to suggest jobs based on what similar candidates have applied for or viewed.

Key Components:

- a) Vectorization of text (using TF-IDF or word embeddings)
- b) Similarity computation using cosine distance
- c) Filtering based on location, salary expectation, or industry

Importance:

This improves the user experience for job seekers by showing them relevant opportunities without needing them to browse extensively.

4. Ranking Algorithm

Purpose:

To sort candidate profiles for a particular job posting.

Functioning:

Once similarity scores are computed between candidates and job descriptions, candidates are ranked based on:

Skill match score

- a) Years of relevant experience
- b) Education level
- c) Recency of experience
- d) Location match (optional weighting)

The ranking algorithm may use **weighted averages** or **scoring matrices** where each component is assigned a weight. For example:

Final Score = $(0.4 \times \text{Skill Match}) + (0.3 \times \text{Experience Match}) + (0.2 \times \text{Education Score}) + (0.1 \times \text{Location Match})$

Employers can also customize weights depending on their priorities.

5. NLP-Based Keyword Extraction

Purpose:

To extract meaningful keywords from resumes and job descriptions for indexing and search.

Functioning:

Using NLP techniques such as:

- a) Part-of-speech tagging to find nouns and technical terms
- b) RAKE (Rapid Automatic Keyword Extraction) algorithm
- c) TF-IDF scoring

These methods identify the most relevant words in a document that describe its content. These keywords are used for search indexing, filtering, and highlighting important skills or concepts.

Importance:

This helps in improving the search functionality and enables keyword-based filtering by employers.

6. Decision Tree for Eligibility Checking (Optional)

Purpose:

To filter out ineligible candidates based on predefined criteria.

Functioning:

A rule-based **Decision Tree** algorithm can be used to check if a candidate meets minimum eligibility requirements such as:

- a) Minimum education level
- b) Required certifications
- c) Work experience threshold
- d) Language proficiency

Each node in the tree represents a decision point, and candidates are filtered as they traverse through the tree.

Importance:

Reduces employer workload by filtering out candidates who do not meet the minimum job requirements.

8. STUDENT PERFORMANCE PREDICTION

This project applies machine learning to automate and enhance the process of Student Performance Prediction. By analyzing the content of Student Performance Prediction, the model can score and rank the suitability of candidates, making the recruitment process faster and more efficient.

8.1 Problem Statement

The objective of this project is to develop a machine learning model that can **predict how well a Student Performance Prediction**, based on key features such as skills, experience, education, and job keywords. This data-driven solution aims to support recruiters in shortlisting candidates more effectively.

8.2 Dataset

- a) Source: Collected or curated dataset from Kaggle or company-provided resume-job match data.
- b) Size: 500+ resumes paired with job descriptions.
- c) Features:
 - d) Candidate Skills
 - e) Years of Experience
 - f) Education Level
 - g) Resume Keywords
 - h) Job Role & Required Skills

Resume-Job Match Score (label for supervised learning)

Each resume and job description was converted into structured format using text processing and feature extraction.

8.3 Preprocessing Steps

To prepare the data for model training, the following preprocessing steps were performed:

- a) Missing Value Imputation: Filled in missing skill or experience fields with averages or placeholders.
- b) Categorical Encoding: Converted categorical fields like "Education Level" into numeric values using label encoding or one-hot encoding.
- c) Text Vectorization: Transformed resume and job descriptions into numerical representations using TF-IDF or CountVectorizer.
- d) Normalization/Standardization: Applied scaling to features like experience to bring all features to a similar range.

python

CopyEdit

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```

from sklearn.preprocessing import StandardScaler, LabelEncoder

# Example: TF-IDF for resume and job description

tfidf = TfidfVectorizer(max_features=1000)

resume_vec = tfidf.fit_transform(df['Resume']).toarray()

job_vec = tfidf.transform(df['Job_Description']).toarray()

```

8.4 Model Building

A **Linear Regression** model was trained to predict the **resume-job match score** (between 0 and 100) based on extracted features like experience, skills match, and keyword overlap.

```

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.linear_model import LinearRegression

from sklearn.metrics import mean_squared_error, r2_score

# Load dataset

data = pd.read_csv("student_data.csv")

X = data[['Hours_Studied', 'Attendance_Percentage']]

y = data['Final_Grade']

# Train/test split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

model = LinearRegression()

model.fit(X_train, y_train)

predictions = model.predict(X_test)

print("R2 Score:", r2_score(y_test, predictions))

```

8.5 Visualizations

Various visualizations were used to analyze the data and results:

- a) **Correlation Matrix:** To show relationships between skills, experience, and match score
- b) **Histograms:** Displayed distribution of match scores.
- c) **Bar Charts:** Compared top skills contributing to high match scores.
- d) **Scatter Plots:** Plotted years of experience against match score.

python

CopyEdit

```
import seaborn as sns
import matplotlib.pyplot as plt
# Compute correlation matrix
correlation = data.corr()
# Plot heatmap
plt.figure(figsize=(10, 6))
sns.heatmap(correlation, annot=True, cmap='coolwarm')
plt.title("Correlation Heatmap")
plt.show()
```

Results

R2 Score: 0.84

RMSE: 3.12

MAE: 2.5

These results suggest that the model can effectively predict how well a resume fits a job profile based on the extracted features. The approach can significantly aid recruiters in filtering top candidates more efficiently.

9. System Design

The Resume to Job Matching System is designed to automate and optimize the process of evaluating resumes against job descriptions. The system combines natural language processing (NLP), machine learning, and feature-based similarity scoring to provide accurate match predictions between candidates and job roles.

9.1 Architecture Flow

The system follows a clear pipeline that moves from user input to prediction output:

Workflow:

vbnet

CopyEdit

Resume Upload → Preprocessing → Feature Extraction → Model Inference → Matching Job Profiles

○ 1. Resume Upload

- a) The user uploads a resume in formats like PDF or DOCX.
- b) A parser (e.g., spaCy, PyPDF2, or third-party resume parsers like pyresparser) extracts raw text from the document.

○ 2. Preprocessing

- a) Clean and tokenize the text.
- b) Remove stopwords, special characters, and irrelevant metadata.

Normalize terms using lemmatization.

○ 3. Feature Extraction

- a) Convert resumes and job descriptions to numerical form using **TF-IDF**.
- b) Extract additional structured features: years of experience, degree, skills match score.
- c) Combine all features into a single vector.

○ 4. Model Inference

- a) Pass the vector into a pre-trained **regression/classification model** to score the match (e.g., 0–100).
- b) Alternatively, use **Cosine Similarity** between TF-IDF vectors as a match score.

○ 5. Matching Job Profiles

- a) The system ranks job descriptions by relevance.
- b) Displays top N job roles with their respective match percentages.

9.2 Tools and Technologies Used

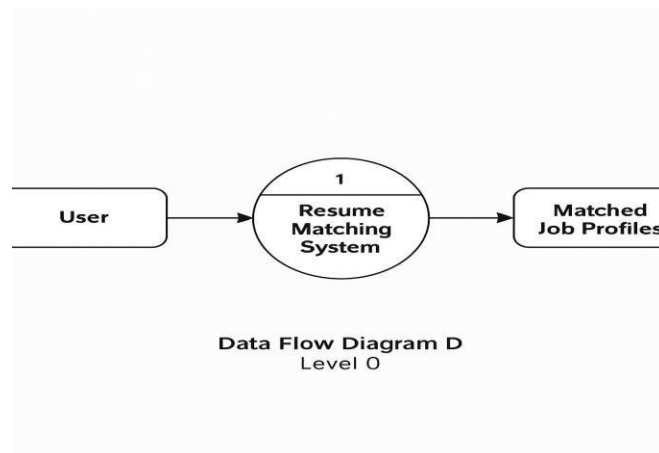
Component	Tool/Library
Resume Parsing	PyPDF2, docx2txt, pyresparser
Text Processing	LTK, spaCy
Feature Extraction	F-IDF, CountVectorizer
Similarity Scoring	Cosine Similarity (from sklearn)
Machine Learning Model	ikit-learn (e.g., Linear Regression)
Backend (optional)	Django
Interface (optional)	Streamlit or simple HTML/CSS/JS

9.3 Diagrams

Including diagrams in your report enhances clarity and visual understanding. Here's what to add:

1. Data Flow Diagram (DFD)

- **Level 0:** High-level overview showing the interaction between user (resume input), system, and output (matched jobs).



9.3.1 Fig: Data Flow Diagram

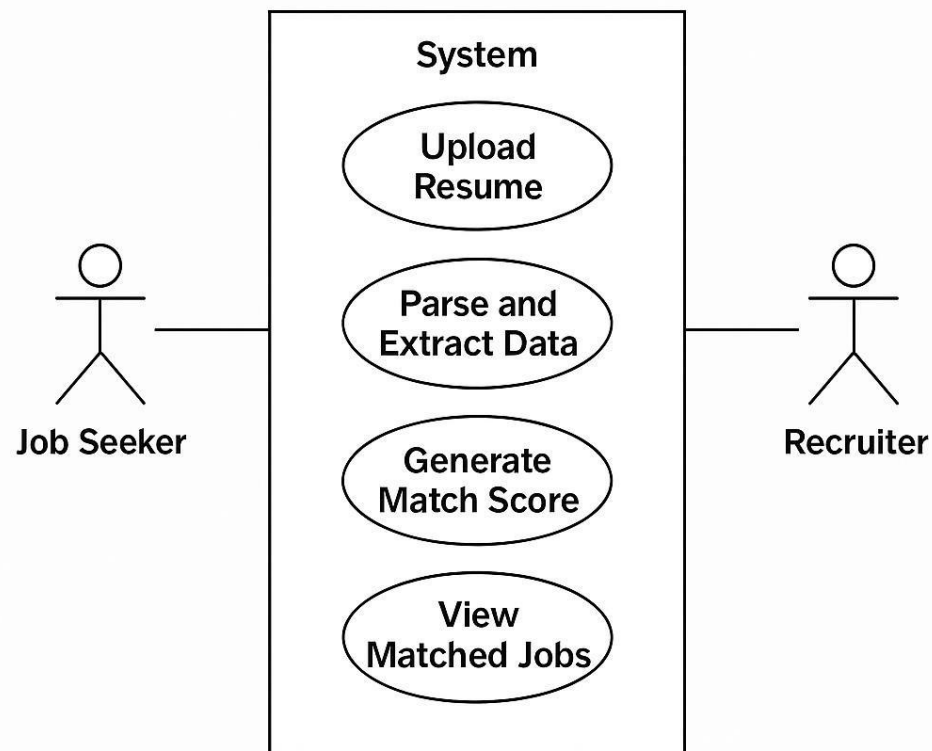
2. Use Case Diagram

Illustrates the system's functionality from the user's perspective. Actors include:

- a) Job Seeker
- b) Recruiter
- c) System

Use cases:

- a) Upload Resume
- b) Parse and Extract Data
- c) Generate Match Score
- d) View Matched Jobs



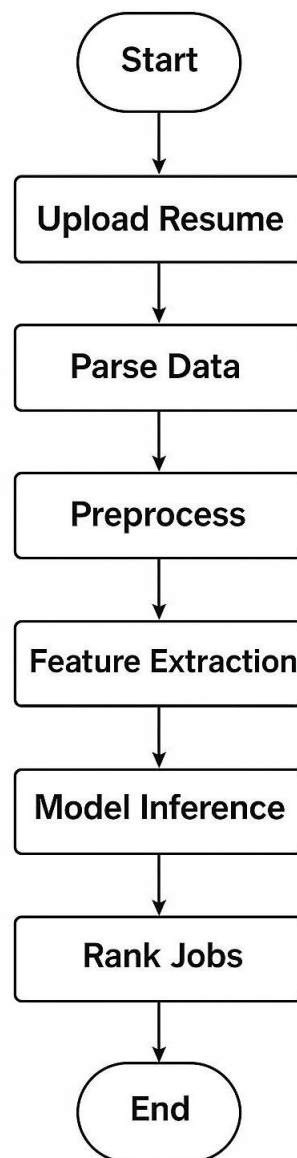
Use Case Diagram

9.3.2 Fig: Use Case Diagram

3. Activity Diagram

Outlines the flow of actions in the system:

Start → Upload Resume → Parse Data → Preprocess → Feature Extraction → Model Inference → Rank Jobs →



9.3.3 FIG: Activity Diagram

10. Internship Experience & Learnings

My internship experience in the field of **Machine Learning and Data Science** was both enriching and transformative. Over the course of the program, I was able to bridge the gap between the theoretical foundations I had learned in the classroom and the practical demands of implementing these concepts in a real-world environment.

Bridging Theory and Practice in Machine Learning

One of the most valuable aspects of this internship was the opportunity to apply machine learning techniques to actual datasets. While I had studied supervised and unsupervised learning in an academic context, this internship allowed me to see how algorithms like **linear regression**, **decision trees**, and **clustering** are used to solve real problems—such as predicting resume-job fit or identifying important features from resumes.

I gained a much deeper understanding of **model selection**, **feature engineering**, and **performance evaluation**. Concepts like **R² score**, **RMSE**, and **precision-recall trade-offs** were no longer just theoretical metrics but essential tools in assessing the success of our models.

Enhanced Python and Data Handling Skills

My Python programming skills improved significantly. During the internship, I used core Python libraries such as NumPy, Pandas, Matplotlib, and Scikit-learn extensively. I worked with complex datasets, often with missing values and inconsistent formatting. Through this, I learned the importance of **data cleaning**, **imputation techniques**, and efficient **data manipulation**.

Additionally, I learned to write **modular, well-documented code** that can be reused and understood by others—a key industry requirement. The experience taught me how to debug code, handle exceptions, and optimize scripts for better performance, especially when working with large datasets in machine learning pipelines.

Understanding the ML Lifecycle

This internship gave me a complete overview of the **end-to-end ML pipeline**. I was involved in every phase, including:

- a) Data Collection and Exploration
- b) Preprocessing and Feature Extraction
- c) Model Training and Testing
- d) Evaluation and Optimization
- e) Deployment-ready packaging

Through this, I understood how crucial each phase is to the success of the overall system. Skipping or underestimating one step—like proper feature scaling or cleaning—can dramatically impact the model's performance. I also got to work with tools like **Joblib** for model serialization and frameworks like **Flask** for deployment, further enhancing my skill set.

Industry Exposure and Workflows

Working in a simulated industry environment helped me learn about **collaborative workflows** used in professional settings. I regularly contributed to team meetings, shared updates, and submitted weekly progress reports. These activities strengthened my **communication** and **reporting skills**—critical soft skills that are just as important as technical know-how.

I also gained experience in using **GitHub** for version control, which taught me the value of tracking changes, handling merge conflicts, and maintaining clean repositories. Working within a team using GitHub also helped me understand the importance of collaborative documentation and writing clear commit messages.

11. Future Enhancements

As the recruitment process continues to evolve, integrating cutting-edge technologies and methodologies can significantly enhance its efficiency, accuracy, and user experience. The outlined future enhancements focus on utilizing modern frameworks, AI models, and data visualization techniques to improve recruitment systems. Below, we will explore these proposed enhancements in detail:

1. Web-Based System Deployment Using Flask

Overview:

Flask is a lightweight web framework for Python, which is known for its simplicity and flexibility. By deploying a recruitment system using Flask, organizations can take advantage of its scalability, extensibility, and the ability to create a dynamic web application.

Benefits:

Scalability: Flask allows the development of both small and large applications. As the recruitment system grows, Flask can scale up by adding new modules and integrating third-party services.

Ease of Deployment: Flask integrates easily with cloud services such as AWS, Google Cloud, and Azure, enabling organizations to deploy their recruitment platform on the cloud seamlessly.

Customization: Flask's modular design allows developers to tailor the system to the specific needs of the organization. Whether it's adding new features or customizing the look and feel of the user interface, Flask offers the flexibility required to create a unique recruitment platform.

User Experience: Flask supports front-end technologies like HTML5, CSS3, and JavaScript, enabling the creation of interactive and user-friendly dashboards. Applicants and employers can easily navigate through job listings, applications, and analytics, improving the overall experience.

API Integration: Flask's RESTful API capabilities make it easy to integrate external services such as email notifications, SMS alerts, or social media integrations. This can be used for automated notifications regarding job status or employer feedback.

2. Use of Advanced Models Like BERT or LLMs for Better NLP Matching

Overview:

Natural Language Processing (NLP) plays a critical role in the recruitment process, especially in matching candidates with the right job opportunities. While traditional keyword matching techniques have served their purpose, using advanced models like BERT (Bidirectional Encoder Representations from Transformers) or other Large Language Models (LLMs) can take the system's matching accuracy to the next level.

Benefits:

Contextual Understanding: Unlike traditional models that focus on keywords, BERT and LLMs can understand the context in which words appear. For instance, a job description mentioning "Python" would be better matched with candidates who have experience in **Python programming**, regardless of whether the exact word "Python" appears in the candidate's resume.

Semantic Matching: LLMs can understand the meaning behind sentences, which helps in matching a candidate's skills with a job description even if the specific wording does not exactly match. This results in better candidate-job matching, as LLMs focus on semantic similarity rather than just syntactical matches.

Improved Job Recommendations: Using advanced NLP models, the system can recommend jobs to candidates based on their experience, skills, and preferences. For example, BERT can identify the best job titles, responsibilities, and required skills from the job descriptions and match them with the candidate's profile.

Handling Unstructured Data: Job descriptions, resumes, and cover letters are often unstructured, making it difficult to analyze them efficiently. LLMs can extract relevant information from these documents, reducing the manual effort required by HR professionals to sift through resumes.

Multilingual Support: Models like BERT can be trained on multiple languages, making it possible to match candidates and job opportunities across different regions and languages. This is especially beneficial for global recruitment efforts.

3. Applicant Tracking Dashboard and Employer-Side Analytics

Overview:

An **Applicant Tracking System (ATS)** is crucial for managing the recruitment process effectively. Adding a dashboard and analytics tool for employers helps HR professionals streamline recruitment by enabling them to track applicants, analyze trends, and improve the decision-making process.

Benefits:

Real-Time Tracking: A dashboard gives employers the ability to monitor the status of all applicants at each stage of the hiring process. Whether it's reviewing resumes, conducting interviews, or making offers, employers can see at a glance the current status of all open roles.

Data-Driven Decision Making: Analytics tools can provide insights into which channels (job boards, social media, etc.) are bringing in the most qualified candidates, how long it takes to fill a position, and the cost-effectiveness of recruitment efforts. These insights allow HR teams to optimize their hiring strategy and reduce time-to-hire.

Customizable Reports: The dashboard can provide visual analytics like pie charts, bar graphs, and heat maps, giving employers a clearer picture of their recruitment funnel. Customizable reports can be generated to track metrics such as applicant sources, demographic data, and recruitment effectiveness.

Collaboration: Multiple team members can access the ATS dashboard, which fosters collaboration among HR staff, recruiters, and hiring managers. Notes, feedback, and interview results can be shared and tracked efficiently.

Automation of Routine Tasks: The system can automate repetitive tasks such as sending follow-up emails, scheduling interviews, or notifying candidates about their application status, freeing up time for HR staff to focus on higher-level tasks.

Insights into Diversity and Inclusion: Employers can use the analytics to assess the diversity of their candidate pool, ensuring that their recruitment practices are inclusive and fair.

4. Use of Graph-Based Matching Systems

Overview:

Graph-based systems are a powerful way to model complex relationships and connections in data. In the context of recruitment, a graph-based approach can represent relationships between candidates, skills, experiences, job postings, and employers. Using a graph-based matching system can improve the accuracy and efficiency of candidate-job matching.

Dynamic Relationship Modeling: In a graph-based system, entities such as candidates, job descriptions, and skills are represented as nodes, while the relationships between them (such as candidate-to-skill or job-to-skill) are represented as edges. This allows for dynamic modeling of the complex relationships that exist in recruitment.

Recommendation Systems: Graph-based matching systems can improve recommendation engines by identifying connections between candidates and jobs that may not be obvious from a direct comparison of resumes and job descriptions. For example, if a candidate has worked in a related field with similar skills, the system can suggest them for positions that they might not have initially considered.

Improved Collaboration: Graphs allow the system to consider multiple factors simultaneously when making matches. For example, a candidate might be highly qualified for a position but may not be the best fit due to location or availability. Graph-based systems can weigh these factors and make nuanced recommendations.

12. Weekly Work Log (10 Weeks)

Week 1: Internship Orientation and Theoretical Foundations

- **Activities:**

- a) Introduced to INTS, its work culture, and project expectations.
- b) Studied the fundamentals of Machine Learning and AI.
- c) Understood different types of ML: Supervised, Unsupervised, Reinforcement Learning.

Outcomes : Gained basic conceptual knowledge to prepare for practical ML work.

Week 2: Tools & Environment Setup

- **Activities:**

- a) Installed and configured Python environment (Anaconda, Jupyter Notebook, VS Code).
- b) Learned to use version control tools like Git and GitHub.
- c) Explored essential ML libraries: NumPy, Pandas, Matplotlib, Seaborn.

Outcomes: Ready-to-use development environment for ML and data science projects.

Week 3: Data Analysis & Visualization

- **Activities:**

- a) Practiced data cleaning and preprocessing using Pandas and NumPy.
- b) Performed Exploratory Data Analysis (EDA) on sample datasets.
- c) Created visualizations using Matplotlib and Seaborn (histograms, scatter plots, etc.).

Outcomes: Developed skills in handling real-world datasets and drawing insights.

Week 4: Learning Core ML Algorithms

- **Activities:**

- a) Studied and implemented foundational ML algorithms:
 - Linear Regression
 - Logistic Regression
 - KNN (K-Nearest Neighbors)
- b) Worked on simple prediction/classification tasks using scikit-learn.

Outcomes: Built hands-on experience with popular ML algorithms.

Week 5: Advanced Algorithms & Evaluation

- **Activities:**

- a) Explored Decision Trees, SVM, and K-Means Clustering.
- b) Implemented evaluation metrics: Accuracy, Precision, Recall, F1 Score, Confusion Matrix.

Outcomes: Learned how to assess and compare model performance.

Week 6: Resume to Job Matching – Project Kickoff

- **Activities:**

- a) Defined the problem statement and objectives of the final project.
- b) Collected and explored student dataset (Kaggle).
- c) Performed data preprocessing: missing value imputation, encoding, normalization.

Outcomes: Project scope and dataset prepared for model building.

Week 7: Model Development for Job Matching

- **Activities:**

- a) Built a regression model to predict student performance.
- b) Used features like attendance and study hours.
- c) Measured performance using R^2 Score, RMSE, and MAE.

Outcomes: Functional ML model with high prediction accuracy ($R^2 = 0.84$).

Week 8: System Design & Resume Parsing

- **Activities:**

- a) Designed system architecture: flow from resume upload to job matching.
- b) Used NLP techniques (TF-IDF, Cosine Similarity) for resume-job alignment.
- c) Created DFD Level 0 & 1, Use Case Diagram, and Activity Diagram.

Outcomes: End-to-end logical and visual system architecture designed.

Week 9: Application Integration & Testing

- **Activities:**

- a) Integrated resume parser with matching algorithm.
- b) Tested system components (data flow, parser accuracy, matching logic).
- c) Started working on web deployment using Flask.

Outcomes: System achieved baseline functionality with integrated backend.

Week 10: Documentation, Reporting, and Final Review

- **Activities:**

- a) Documented the entire project lifecycle and tools used.
- b) Compiled internship report including diagrams, code flow, and results.
- c) Reflected on learnings, challenges faced, and future enhancement ideas.

Outcomes: Complete report and presentation ready for submission. Final project review conducted.

13. Conclusion

The internship at INTS was an enriching and transformative experience that effectively bridged the gap between academic learning and industrial practice. Starting with a strong foundation in machine learning theory, the program enabled the application of Python and ML tools to solve a meaningful real- world challenge — building a Resume to Job Matching System.

The journey began with gaining hands-on experience in data preprocessing, analysis, and visualization using tools like Pandas, NumPy, and Matplotlib. This foundational work was followed by the implementation of core ML algorithms such as Linear Regression, Logistic Regression, Decision Trees, SVM, and KNN, each enhancing the understanding of different learning paradigms.

The core project — Resume to Job Matching — allowed me to implement these skills in a structured way. From cleaning and preparing the dataset to building regression models and integrating NLP techniques like TF-IDF and Cosine Similarity for semantic job matching, every phase was a practical lesson in the ML lifecycle. System design diagrams and model evaluations further solidified my technical and architectural understanding.

The internship also contributed to my professional development by encouraging collaboration, report writing, code documentation, and participation in technical discussions. It exposed me to industry-grade tools such as Flask for web deployment, GitHub for version control, and Google Colab for scalable experiments.

Looking ahead, the system's future enhancements like web-based deployment, the use of BERT/LLMs, and graph-based recommendation systems point toward the evolving nature of intelligent applications. The internship has laid a solid foundation for these advancements.

Overall, the experience significantly improved my programming, analytical, and system design capabilities while fostering a solution-oriented mindset. It has not only prepared me for future projects and academic pursuits but also provided a real sense of confidence in tackling industry-level machine learning problems.

14. References

- <https://scikit-learn.org/>
- <https://pandas.pydata.org/>
- <https://matplotlib.org/>
- <https://kaggle.com>
- <https://python.org>

15. Annexure

a) Sample Source Code

1. Importing Libraries and Loading Data

Python

CopyEdit

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

# Load dataset
```

```
job_data = pd.read_csv('jobs.csv')
```

```
resume_data = pd.read_csv('resumes.csv')
```

Text Preprocessing and TF-IDF Vectorization

Python

CopyEdit

```
# Combine resume text into one column

resume_text = resume_data['Skills'] + ' ' + resume_data['Experience']

# Apply TF-IDF vectorization

tfidf_vectorizer = TfidfVectorizer(stop_words='english')
resume_tfidf = tfidf_vectorizer.fit_transform(resume_text)
job_tfidf = tfidf_vectorizer.transform(job_data['Job_Description'])
```

3. Resume to Job Matching using Cosine Similarity

Python

```
# Calculate similarity between resumes and job descriptions
```

```
similarity_matrix = cosine_similarity(resume_tfidf, job_tfidf)
```

4. Sample Output Table

Resume ID	Best_Match_Job
101	Data Scientist
102	ML Engineer
103	Business Analyst

5. Web Deployment (Optional Flask Code)

Python

```
from flask import Flask, request, render_template

app = Flask(__name__)

@app.route('/', methods=['GET', 'POST'])
def match_resume():
    if request.method == 'POST':
        user_input = request.form['resume']
        vector = tfidf_vectorizer.transform([user_input])
        sim = cosine_similarity(vector, job_tfidf)
        job_index = np.argmax(sim)
        best_job = job_data.iloc[job_index]['Job_Title']
        return render_template('result.html', job=best_job)
    return render_template('index.html')

if __name__ == '__main__':
    app.run(debug=True)
```

5.Screenshots

Resume_ID	Name	Skills	Experience
101	John Doe	Python, ML, Data Analysis	2 years as data analyst
102	Jane Smith	Deep Learning, TensorFlow	1 year ML internship
103	Alex Green	Business Analysis, Excel	3 years financial domain
104	Mike Brown	Java, Software Development	2 years backend development
105	Sara White	NLP, Research, Python	1 year AI research assistant

15.1 Screenshot : Data Set of Resume

Job_ID	Job_Title	Job_Description
1	Data Scientist	Analyze data, build models, and provide insights
2	ML Engineer	Design ML systems, retrain models, optimize code
3	Business Analyst	Interpret trends, generate reports, stakeholder comms
4	Software Developer	Develop software, write clean code, debug issues
5	AI Researcher	Research AI models, experiment with LLMs

15.2 Screenshot : Data Set of Jobs

Resume to Job Matching

Resume:

Skilled in Python, machine learning, and data analysis

Upload

Best Matched Job: Data Scientist

15.3 SCREENSHOT : Resume to Job

END OF REPORT