

# **PROJECT REPORT**

## **"AutoInsight.ai: No-Code ML Workflow for Smart Data Exploration & Modelling"**

Submitted in partial fulfillment for the requirement of the award of

***INTERNSHIP***

***IN***

## **ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING WITH DATA SCIENCE**



*Submitted By*

*Tunga Sai Lakshmi (Sanskriti School of Engineering)*

*Smriti Behura (Bhilai Institute of Technology)*

*Bhavye Kathuria (HMR Institute of Technology and Management)*

*Jyotiraditya Mishra (S.C.S Autonomous College)*

*Vidyadheesha M Pandurangi (Adhiyamaan College of Engineering)*

*Under the guidance of*

**Mr. Abhishek Rao**

## **ACKNOWLEDGEMENT**

Our sincere gratitude and thanks to my project report guide **Mr. Abhishek Rao**. It was only with his backing and support that we could complete the report. He provided us with all sorts of help and corrected us if ever seemed to make mistakes. We acknowledge our professor for providing me with all sorts of help and assistance when and where required. We also acknowledge my friends who supported us throughout this project. We have no such words to express our gratitude. We acknowledge our dearest parents for being such a nice source of encouragement and moral support that helped us tremendously in this respect. We also declare to the best of our knowledge and belief that Project Work has not been submitted anywhere else.

## TABLE OF CONTENT

<b>ABSTRACT.....</b>	<b>1</b>
<b>CHAPTER 1: INTERNSHIP OVERVIEW.....</b>	<b>2</b>
1.1. ORGANIZATION AND COLLABORATORS.....	2
1.2. DURATION AND MODE OF INTERNSHIP .....	2
1.3. TEAM MEMBERS .....	3
1.4. OBJECTIVES OF THE INTERNSHIP.....	3
<b>CHAPTER 2: PROJECT INTRODUCTION.....</b>	<b>4</b>
2.1. PROJECT TITLE .....	4
2.2. MOTIVATION BEHIND THE TOOL.....	4
2.3. KEY FEATURES OF AUTOINSIGHT.AI .....	5
<b>CHAPTER 3: TECHNOLOGY STACK USED .....</b>	<b>6</b>
3.1. PROGRAMMING LANGUAGES AND LIBRARIES.....	6
3.2. TOOLS AND PLATFORMS .....	7
3.3. STREAMLIT DEPLOYMENT .....	8
3.4. DEPLOYMENT HIGHLIGHTS .....	8
<b>CHAPTER 4: SYSTEM ARCHITECTURE AND WORKFLOW.....</b>	<b>9</b>
4.1. ARCHITECTURE DIAGRAM .....	9
4.2. DATA FLOW PIPELINE.....	10
4.3. MODULE INTERACTIONS.....	11
<b>CHAPTER 5: PROJECT MODULES IN DETAIL.....</b>	<b>12</b>
5.1. UPLOAD AND DATA CLEANING .....	12
5.2. EXPLORATORY DATA ANALYSIS(EDA).....	13
5.3. MODEL BUILDING WITH PYCARET .....	14
5.4. PERFORMANCE EVALUATION AND METRICS .....	15
5.5. DOWNLOAD AND EXPORT FUNCTIONALITY .....	16
<b>CHAPTER 6: DATASET HANDLING AND PREPROCESSING .....</b>	<b>17</b>
6.1. HANDLING MISSING VALUES.....	17
6.2. REMOVING DUPLICATES .....	18
6.3. ENCODING AND TRANSFORMATION .....	19

<b>CHAPTER 7: EXPLORATORY DATA ANALYSIS (EDA).....</b>	<b>21</b>
7.1. PROFILE REPORT GENERATION.....	21
7.2. INSIGHTS GAINED.....	22
<b>CHAPTER 8: AUTOMATED MACHINE LEARNING(AUTOML) MODELLING .....</b>	<b>24</b>
8.1. TARGET SELECTION .....	24
8.2. MODEL COMPARISION.....	25
8.3. BEST MODEL SELECTION .....	26
8.4. TRAINING AND SAVING THE MODEL.....	27
<b>CHAPTER 9: MODEL EVALUATION METRICS.....</b>	<b>28</b>
9.1. ACCURAACY, F1 SCORE, RECALL, PRECISION .....	28
9.2. ROC-AUC, KAPPA, MCC.....	30
9.3. CLASSIFICATION REPORT .....	31
<b>CHAPTER 10: MODEL AND REPORT EXPORT.....</b>	<b>32</b>
10.1. HDML EDA REPORT.....	32
10.2. CSV PERFORMANCE MATRICES .....	33
10.3. TRAINED MODEL (.pkl FILE) .....	34
<b>CHAPTER 11: CHALLENGES FACED AND SOLUTIONS .....</b>	<b>35</b>
<b>CHAPTER 12: KEY LEARNINGS AND SKILLS ACQUIRED .....</b>	<b>36</b>
<b>CHAPTER 13: CONCLUSION.....</b>	<b>37</b>
<b>CHAPTER 14: SOURCE CODE .....</b>	<b>38</b>
<b>CHAPTER 15: TESTING OF APPLICATION USING DATASET.....</b>	<b>44</b>
<b>CHAPTER 16: SCREENSHOTS OF THE APPLICATION.....</b>	<b>44</b>
<b>CHAPTER 17: REFERENCES.....</b>	<b>49</b>

## ABSTRACT

AutoInsight.ai is an innovative web-based tool designed to streamline the process of data analysis and machine learning model deployment, aimed at empowering users with limited technical expertise to create insightful analyses and predictive models with minimal effort. This project leverages the power of Automated Machine Learning (AutoML) to automate key stages of the data analysis pipeline, including data preprocessing, exploratory data analysis (EDA), model selection, training, and performance evaluation.

The core functionality of AutoInsight.ai is built upon a user-friendly interface using Streamlit, which allows users to interact with the tool via a simple web application. Users can upload datasets, visualize key insights through automated EDA, and choose machine learning algorithms for model building without the need for manual coding. By integrating PyCaret, an open-source low-code AutoML library, the platform automates the process of model selection, comparison, and optimization, ensuring that users can select the best-performing model based on their specific dataset.

Furthermore, the system provides comprehensive performance metrics such as accuracy, F1-score, precision, recall, and ROC-AUC, which enable users to evaluate the effectiveness of their models in a simple and understandable format. The ability to export EDA reports in HTML format, download performance metrics as CSV files, and save the trained models as pickle (.pkl) files adds a layer of versatility and ease for users to integrate these insights and models into real-world applications.

The project aims to democratize data science by providing a platform that simplifies the machine learning workflow, making it accessible to a wider audience, from beginners to intermediate users. By abstracting away the complexities of machine learning, AutoInsight.ai enhances productivity, reduces the time and effort required for model development, and offers an accessible approach to predictive analytics.

Through this project, we seek to provide an intuitive platform that bridges the gap between raw data and actionable insights, all while reducing the barrier to entry for machine learning practitioners across various domains.

# **CHAPTER 1**

## **INTERNSHIP OVERVIEW**

### **1.1 Organization and Collaborators**

The internship program was organized by EduTech Life, a premier EdTech platform focused on bridging the gap between academic learning and industry-relevant skill development. EduTech Life is known for its structured, hands-on training programs in emerging technologies such as Artificial Intelligence (AI), Machine Learning (ML), and Data Science (DS).

This internship was conducted in collaboration with IIT Mechanica 2025, a prestigious technical fest hosted by the Indian Institute of Technology (IIT). The collaboration provided the interns an opportunity to learn from top academic mentors, industry practitioners, and certified instructors in the field of AI/ML.

The program was designed to be experiential and application-driven, providing students with practical exposure to real-world tools, model deployment strategies, and project-based learning. The goal of this collaboration was not just to build foundational knowledge, but to allow students to build and deploy end-to-end AI solutions that could solve real problems, efficiently and effectively.

### **1.2 Duration and Mode of Internship**

The internship spanned over a period of 6 weeks, during which students participated in a mix of live interactive sessions, recorded modules, and hands-on project work.

**Duration:** 15 Days

**Mode:** Fully Online (Virtual Internship)

The virtual mode enabled flexibility while ensuring access to expert sessions, one-on-one mentorship, and peer collaboration through online forums. All sessions were conducted via Zoom and Google Meet platforms.

Additionally, tools like Jupyter Notebook, Anaconda Navigator, Streamlit, GitHub, and VS Code were extensively used throughout the project lifecycle.

### **1.3 TEAM MEMBERS:**

Name	College Name	Role
Tunga Sai Lakshmi	Sanskriti School of Engineering, Puttaparthi	Research and Analysis
Smriti Behura	Bhilai Institute of Technology, Chhattisgarh	Documentation and Reporting
Bhavye Kathuria	HMR Institute of Technology and Management, Delhi	Design and Development
Jyotiraditya Mishra	S.C.S Autonomous College, Puri, Odisha	Testing and Validation
Vidyadheesha M Pandurangi	Adhiyamaan College of Engineering, Hosur	Project Coordination and Integration

### **1.4 Objectives of the Internship**

The internship was structured with clear academic and industrial objectives that aimed to cultivate practical AI and ML expertise among students. The major objectives of the internship are as follows:

- To gain a comprehensive understanding of AI, Machine Learning, and Data Science concepts through structured curriculum and live mentoring.
- To apply theoretical concepts to real-world datasets and develop a full-stack AI solution that covers data loading, preprocessing, model training, evaluation, and deployment.
- To design and develop a no-code tool titled AutoInsight.ai that enables users to explore datasets, perform automatic EDA (Exploratory Data Analysis), train ML models using PyCaret, evaluate them, and download the results – all without writing a single line of code.
- To learn modern tools and frameworks like Streamlit, PyCaret, Scikit-learn, Seaborn, Matplotlib, and pandas profiling for end-to-end machine learning lifecycle development.
- To deploy the final AI tool as a web application using Streamlit Cloud at the URL: <https://autoinsightai.streamlit.app/>, making it publicly accessible.
- To gain internship experience like working in a startup/industry environment by following deadlines, documentation practices, and presenting deliverables for evaluation.

This internship was not just a learning experience but also a career-defining milestone, helping me understand how to build AI tools that democratize data science, reduce technical barriers, and enable data-driven decision-making.

## CHAPTER 2

# PROJECT INTRODUCTION

### 2.1 Project Title

**Project Title:** *AutoInsight.ai: No-Code ML Workflow for Smart Data Exploration & Modeling*

This project aims to bridge the gap between non-technical users and the complex world of Machine Learning (ML) by developing a user-friendly, web-based AI tool called **AutoInsight.ai**. The tool enables users to upload datasets, perform automatic exploratory data analysis (EDA), train machine learning models, evaluate their performance, and download the results – all through a graphical user interface without writing a single line of code. Built using Streamlit, PyCaret, and other open-source libraries, this tool helps democratize data science by making intelligent data modelling accessible to everyone.

### 2.2 Motivation Behind the Tool

In today's data-driven world, organizations generate enormous amounts of data every day. However, not all professionals working with data have the technical skills to write code, build ML models, or perform statistical analysis. Traditional data analysis and machine learning workflows require a deep understanding of programming languages like Python or R, familiarity with ML algorithms, and proficiency in libraries such as pandas, scikit-learn, and matplotlib. The motivation for building AutoInsight.ai stems from three major challenges faced by learners and professionals:

- **Technical Barriers:** Many students, business analysts, and domain experts struggle to implement ML workflows due to limited programming experience.
- **Time Constraints:** Data exploration and model development can be time-consuming, especially when done manually or from scratch.
- **Lack of Automation:** Existing solutions often require multiple tools or platforms to achieve end-to-end data modelling, increasing complexity.

AutoInsight.ai was created to solve these problems by offering an end-to-end no-code solution for machine learning. The idea was to develop a platform where even a beginner or a non-programmer could explore a dataset, generate automated visualizations, train predictive models, and interpret the results easily.

The goal was clear: make data science and ML workflows intuitive, fast, and accessible to all. Whether you're a student, business analyst, or an early-stage data scientist, AutoInsight.ai aims to be your go-to tool for fast, intelligent data exploration and modeling.

## **2.3 Key Features of AutoInsight.ai**

AutoInsight.ai stands out for its simplicity, automation, and versatility. The following are the key features that define and differentiate this tool:

### **1. Intuitive Web Interface**

- Built using Streamlit, the interface is clean, responsive, and easy to navigate.
- Users can interact with the tool through simple dropdowns, file uploads, and checkboxes—no coding required.

### **2. Dataset Upload and Preview**

- Users can upload their own .csv files to begin the analysis.
- The tool provides an instant preview of the dataset (top 5 rows), column types, and dimensions (rows × columns).

### **3. Automatic Exploratory Data Analysis (EDA)**

- Using pandas profiling, users can generate a comprehensive EDA report in one click.
- The report includes missing value heatmaps, univariate and multivariate analysis, correlation matrices, and more.

### **4. Smart ML Model Selection and Training**

- With the integration of PyCaret, the tool automatically selects the best algorithms based on the data type (classification or regression).
- Users can train multiple ML models with a single click, and the tool displays a comparison of their performance metrics.

### **5. Model Evaluation and Metrics**

- Displays important metrics like Accuracy, AUC, RMSE, MAE, F1-score, and Confusion Matrix, depending on the problem type.
- Also includes interpretability features such as feature importance ranking.

### **6. Downloadable Results**

- After training and evaluation, users can download the trained model or export the results for further use.

### **7. End-to-End Automation**

- The entire process—from data input to model output—is automated.
- Users don't need to handle any manual coding, data cleaning, or hyperparameter tuning.

## **8. Deployed on Streamlit Cloud**

- The application is publicly hosted and accessible via <https://autoinsightai.streamlit.app/>
- This enables instant access from any device, anytime, without any software installation.

# **CHAPTER 3**

## **TECHNOLOGY STACK USED**

To build a robust, scalable, and user-friendly machine learning tool like **AutoInsight.ai**, a modern and efficient technology stack was selected. This stack ensured ease of development, interactive user experience, seamless ML model integration, and hassle-free cloud deployment. The following subsections outline the key components of the technology stack:

### **3.1 Programming Languages and Libraries**

**Programming Language:** Python

- The entire backend logic, data handling, and ML pipeline are implemented using Python, owing to its simplicity, readability, and vast ecosystem of data science libraries.
- Python's dynamic typing, strong community support, and suitability for rapid prototyping made it the ideal choice for this no-code ML platform.

**Key Python Libraries Used:**

#### **1. Streamlit**

- Used to build the interactive web UI.
- Allows creation of intuitive, fast, and responsive data apps with minimal lines of code.
- Provides built-in widgets like file uploaders, dropdowns, buttons, and more for user interaction.

#### **2. Pandas**

- Handles data loading, manipulation, and analysis.
- Allows efficient reading of CSV files and data cleaning operations.
- Used to generate dataset previews, summaries, and column-wise statistics.

#### **3. Pandas Profiling**

- Enables automatic generation of detailed exploratory data analysis (EDA) reports.
- Includes correlation matrices, missing value charts, and variable distributions.

#### **4. PyCaret**

- Serves as the low-code machine learning framework.
- Simplifies the entire ML workflow including data preprocessing, model training, comparison, evaluation, and deployment.
- Supports both classification and regression problem types.

#### **5. Scikit-learn**

- Works as the underlying engine behind PyCaret for many ML algorithms.
- Used indirectly for model training, validation, and performance measurement.

#### **6. Matplotlib / Seaborn**

- Visualization libraries for plotting data distributions and performance metrics.
- Used for enhancing the visual presentation of data and results.

### **3.2 Tools and Platforms**

A combination of open-source tools and cloud-based platforms was used to support development, version control, and deployment.

#### **Development Environment:**

- **Visual Studio Code (VS Code):**

- Main IDE used for writing and debugging code.
- Integrated with Python extensions for better code intelligence and linting.
- Offers a user-friendly interface for rapid development.

#### **Version Control:**

- **Git & GitHub:**

- Git was used for version control, ensuring proper code management and collaboration.
- GitHub hosted the project repository, enabling easy deployment to external platforms like Streamlit Cloud.
- Repository link (example): <https://github.com/yourusername/AutoInsightAI>

#### **Cloud Platforms:**

- **Streamlit Community Cloud:**

- Used to deploy the application publicly.
- Enables any user with a browser to access the tool without needing any installations.

- Offers continuous deployment from GitHub, making updates seamless and automatic.

### **3.3 Streamlit Deployment**

Deploying **AutoInsight.ai** to the cloud using Streamlit Community Cloud was a critical step in ensuring accessibility, scalability, and usability of the tool.

#### **3.3.1 Benefits of Streamlit Cloud Deployment:**

- Serverless deployment.
- Integrated GitHub workflows for CI/CD.
- Free tier allows for educational and demo purposes.
- Eliminates the need for complex backend setup or DevOps pipelines.

Here is a breakdown of the deployment process:

#### **Deployment Steps:**

##### **1. Code Preparation:**

- Ensured all dependencies were listed in requirements.txt.
- Maintained a proper main.py (entry point) script containing the Streamlit app.

##### **2. GitHub Integration:**

- Uploaded the complete project to a GitHub repository.
- Repository included all scripts, models, and config files.

##### **3. Streamlit Cloud Configuration:**

- Signed in to <https://streamlit.io/cloud>.
- Linked the GitHub repository to Streamlit.
- Selected the main script (main.py) and configured environment settings.

##### **4. Automatic Deployment:**

- Streamlit Cloud automatically installed dependencies.
- Deployed the app and provided a public URL.

#### **3.4 Deployment Highlights:**

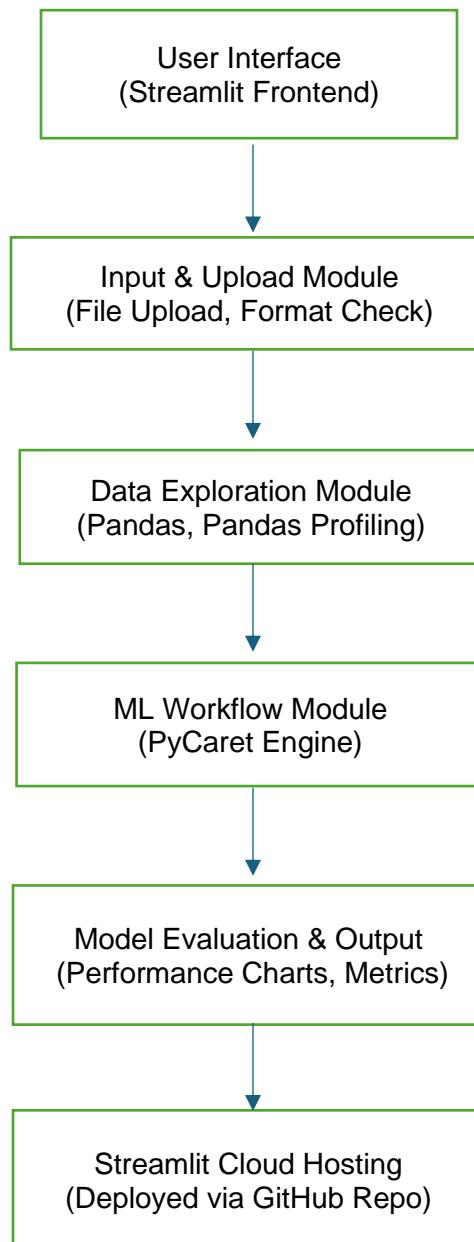
- **URL:** <https://autoinsightai.streamlit.app/>
- The deployed tool is accessible to users globally. No installation is required; only an internet connection and web browser are needed.

## CHAPTER 4

# SYSTEM ARCHITECTURE AND WORKFLOW

The system architecture of **AutoInsight.ai** is designed to streamline the process of loading a dataset, performing exploratory data analysis (EDA), building machine learning models, and evaluating them—all without requiring the user to write a single line of code. It encapsulates modular interaction, seamless data flow, and a user-friendly interface.

### 4.1 Architecture Diagram



## **4.2 Data Flow Pipeline**

The data flow in AutoInsight.ai follows a linear yet interactive pipeline designed to process the input dataset and yield valuable insights and predictive models. The following steps summarize the flow of data through the system:

### **Step-by-Step Pipeline:**

#### **1. Data Ingestion:**

- The user uploads a .csv file via the Streamlit file uploader widget.
- The system reads and validates the file using pandas.

#### **2. Data Preview & Summary:**

- The application shows the first few rows, basic statistics, and data types.
- This step helps the user understand the structure and nature of their dataset.

#### **3. Automated EDA:**

- Using pandas\_profiling, an in-depth profiling report is generated.
- The report includes:
  - Missing values
  - Correlation heatmaps
  - Categorical vs numerical analysis
  - Skewness and distribution plots

#### **4. Problem Selection:**

- The user selects the target variable and problem type (classification or regression).
- The system dynamically adapts based on the selection.

#### **5. Model Training (PyCaret Integration):**

- PyCaret sets up the data pipeline internally (encoding, scaling, splitting).
- Trains multiple models and compares them based on performance metrics.
- Displays top-performing models and allows the user to choose one.

#### **6. Model Evaluation & Results Display:**

- Performance metrics such as Accuracy, RMSE, Precision, Recall, or R<sup>2</sup> are displayed.
- Plots such as confusion matrix, ROC curve, and feature importance are shown interactively.

#### **7. Deployment Preview:** Allows model export or prediction interface.

### **4.3 Module Interactions**

Each module in AutoInsight.ai interacts with the others in a tightly coupled yet logically modular fashion. Here's how they collaborate to offer a seamless no-code ML experience:

#### **1. User Interface Module (Streamlit):**

- Acts as the main control and display layer.
- Interacts with:
  - Input Module: for uploading and selecting files.
  - EDA Module: to display data summaries and visualizations.
  - ML Module: to select the target variable and display models/results.

#### **2. Input & Upload Module:**

- Checks file type and reads data using pandas.
- Sends valid data to:
  - EDA Module for analysis.
  - ML Workflow Module for modeling.

#### **3. Data Exploration Module:**

- Analyzes the data using pandas\_profiling.
- Outputs:
  - HTML/interactive reports rendered via Streamlit.
- Prepares data insights to support ML decisions.

#### **4. ML Workflow Module (PyCaret):**

- Receives clean data and target selection from the interface.
- Internally manages:
  - Preprocessing
  - Feature engineering
  - Model training and comparison
- Sends model summaries and plots to the results display module.

#### **5. Results Display Module:**

- Shows model metrics interactively using Streamlit.
- Enables users to interpret the best model's performance easily.

#### **6. Streamlit Hosting Layer:**

- Entire workflow is hosted on Streamlit Cloud.

- Ensures real-time responsiveness and public accessibility via web.

## CHAPTER 5

### PROJECT MODULES IN DETAIL

The **AutoInsight.ai** platform consists of several key modules that provide users with a streamlined, no-code process for data exploration, model building, and performance evaluation. Each module is designed to serve a specific purpose in the workflow, from uploading data to generating insights and downloading the results.

#### **5.1 Upload and Data Cleaning**

**Purpose:** The Upload and Data Cleaning module serves as the entry point for users to upload their datasets and prepares the data for further processing. This module ensures that data is correctly formatted, and any necessary preprocessing is handled automatically.

#### **Functionality:**

- **File Upload:**
  - The user uploads a .csv file through the Streamlit interface using the st.file\_uploader widget.
  - Supports common file formats, primarily CSV, but can be expanded for other formats like Excel (.xlsx).
- **Data Integrity Check:**
  - Once the file is uploaded, the data is validated for compatibility, ensuring that it is a tabular dataset with rows and columns.
  - Basic validation includes checking for missing values, correct data types (numeric, categorical), and file integrity.
- **Data Preview:**
  - Displays the first few rows of the dataset to give the user an initial view of the data and its structure.
  - Column names, data types, and sample values are displayed.
- **Handling Missing Data:**
  - Identifies and reports missing values within the dataset.
  - Offers automatic data imputation strategies (mean, median, or mode for numerical columns, mode for categorical columns).

- Alternatively, users can manually remove rows or columns with excessive missing data.
- **Data Type Conversion:**
  - Automatically detects the types of data (numeric, categorical) and suggests necessary conversions.
  - For example, it can convert a numeric column into a categorical one based on specific conditions.
- **Data Cleaning Summary:**
  - Generates a summary report detailing:
    - Number of missing values per column.
    - Columns with constant values that might be redundant.
    - Identified outliers and duplicates.

## 5.2 Exploratory Data Analysis (EDA)

**Purpose:** The Exploratory Data Analysis (EDA) module automates the process of gaining insights into the dataset. EDA is vital for understanding the data's structure, detecting patterns, and identifying any issues that might affect model training.

### Functionality:

- **Pandas Profiling Report:**
  - Automatically generates an interactive profiling report using the pandas\_profiling library.
  - The report includes a wide range of analyses, such as:
    - **Descriptive statistics:** Mean, median, standard deviation, etc.
    - **Correlation analysis:** Identifies relationships between numerical features.
    - **Distribution plots:** Visualizes the distribution of numerical variables (e.g., histograms).
    - **Categorical analysis:** Frequency distribution for categorical variables.
    - **Missing value analysis:** Detailed visualizations to identify patterns of missing data.
    - **Data quality checks:** Checks for constant, highly skewed, or duplicated columns.

- **Visualizations:**
  - Interactive visualizations, such as:
    - **Box plots** to identify outliers in numeric columns.
    - **Heatmaps** for correlation analysis.
    - **Bar charts** for categorical feature distributions.
    - **Pair plots** for pairwise comparison between numerical features.
- **Identifying Key Insights:**
  - The system highlights key insights from the report, such as highly correlated features, distribution abnormalities, or imbalances in categorical classes.
- **User Customization:**
  - Users can filter the types of EDA analysis they want to perform, for instance, focusing only on numerical features or only on correlations.

### **5.3 Model Building with PyCaret**

**Purpose:** The Model Building module is the core functionality of AutoInsight.ai, where users can create machine learning models with minimal effort. PyCaret, an open-source machine learning library, powers the model building process, allowing automatic handling of preprocessing, feature engineering, model selection, and training.

#### **Functionality:**

- **Automatic Data Preprocessing:**
  - PyCaret automatically handles tasks like:
    - Encoding categorical variables.
    - Feature scaling (normalization/standardization).
    - Splitting the data into training and test sets.
    - Feature selection.
- **Model Initialization:**
  - Users specify the target variable (dependent variable) and the problem type (classification or regression).
  - Based on this input, PyCaret initializes the appropriate pipeline and algorithm selection.
- **Model Comparison:**

- PyCaret automatically compares a wide range of machine learning models (e.g., Random Forest, Logistic Regression, XGBoost, LightGBM, etc.).
  - Each model is evaluated based on predefined metrics, such as accuracy for classification or RMSE for regression.
- **Best Model Selection:**
  - After comparing multiple models, the system displays a table of performance scores, allowing users to select the best-performing model.
  - The model with the best accuracy or least error (depending on the task type) is highlighted.
- **Hyperparameter Tuning:**
  - PyCaret supports automatic hyperparameter tuning to optimize model performance.
  - Users can choose to perform fine-tuning on the best model with grid search or random search techniques.

## 5.4 Performance Evaluation & Metrics

**Purpose:** The Performance Evaluation & Metrics module evaluates the performance of the trained models and provides various metrics to help users assess the effectiveness of their model.

### Functionality:

- **Evaluation Metrics (for Classification and Regression):**
  - **Classification Models:**
    - Accuracy
    - Precision, Recall, F1-Score
    - Confusion Matrix
    - ROC Curve and AUC Score
  - **Regression Models:**
    - R-squared
    - Mean Absolute Error (MAE)
    - Mean Squared Error (MSE)
    - Root Mean Squared Error (RMSE)
- **Confusion Matrix:**
  - A confusion matrix is shown for classification models to visualize the performance in terms of true positives, false positives, true negatives, and false negatives.

- **ROC Curve:**
  - For binary classification problems, the Receiver Operating Characteristic (ROC) curve is plotted to evaluate the model's ability to distinguish between the classes.
- **Model Comparison:**
  - The evaluation metrics are presented in a table, allowing users to compare different models' performance.
  - Interactive visualizations like **model performance plots** allow users to easily identify the best model.
- **Feature Importance:**
  - The system generates a feature importance plot (for tree-based models), showing which features contribute most to the model's predictions.

## 5.5 Download and Export Functionality

**Purpose:** The Download and Export module allows users to download their results, including trained models, performance metrics, and prediction outputs. This functionality makes AutoInsight.ai a flexible platform for users to export their insights for further use or deployment.

### Functionality:

- **Export Model:**
  - The trained model can be exported as a Pickle file (.pkl) for later use.
  - Users can download the model and integrate it into production systems or share it for further analysis.
- **Export Predictions:**
  - Users can download a CSV file containing predictions made by the model on the test dataset.
  - This can be useful for generating reports or sharing results with stakeholders.
- **Export EDA Report:**
  - The interactive EDA report generated by pandas\_profiling can be downloaded as a HTML file.
  - This report can be shared for further review or archiving.
- **Metrics Summary:**
  - A summary of model performance metrics can be exported as a **CSV** file.
    - This allows the user to have a formalized version of their model evaluation.

## CHAPTER 6

### DATASET HANDLING AND PREPROCESSING

Dataset preprocessing is a critical step in the machine learning pipeline as it ensures the data is clean, structured, and ready for analysis and model training. **AutoInsight.ai** handles several key preprocessing steps automatically, helping users efficiently prepare their data for modeling without needing to write any code. This section discusses how **AutoInsight.ai** handles missing values, removes duplicates, and performs encoding and transformation.

#### 6.1 Handling Missing Values

**Purpose:** Missing values can severely impact the quality of a model's predictions. Handling them appropriately is essential for ensuring accurate model performance. The Handling Missing Values module in **AutoInsight.ai** is designed to automatically identify, report, and handle missing data within a dataset.

##### Functionality:

- **Identification of Missing Data:**
  - **AutoInsight.ai** uses pandas functionality to identify missing values (NaN or None) in the dataset.
  - The platform provides a detailed summary of missing values per column, which helps the user visualize the extent of missing data.
  - For each column with missing data, the percentage of missing values is displayed.
- **Imputation Strategies:**
  - **Automatic Imputation:**
    - **Numerical Columns:**
      - **Mean Imputation:** Replaces missing values in numerical columns with the mean of that column (standard practice for columns with normally distributed data).
      - **Median Imputation:** For columns with skewed distributions, missing values are filled with the median to avoid biasing the dataset with outliers.
    - **Categorical Columns:**
      - **Mode Imputation:** Missing values in categorical columns are filled with the most frequent (mode) value in the column.

- **User-Defined Imputation:**
  - Users can customize the imputation strategy if they prefer a different approach, such as using a constant value or filling missing values based on external data.
- **Removal of Rows/Columns:**
  - If the missing data is excessive (e.g., a column with more than 50% missing values), the system suggests removing the column entirely.
  - For rows with missing values, users can choose to drop them or replace the missing values with a default imputed value, depending on the proportion of missing data.
- **Visualization of Missing Data:**
  - The system generates visualizations (e.g., heatmaps or bar charts) to illustrate missing data patterns, helping users understand if the missing values are random or follow a specific pattern.

## 6.2 Removing Duplicates

**Purpose:** Duplicate records in a dataset can distort statistical analyses and model training, leading to overfitting or biased predictions. The Removing Duplicates module ensures that all redundant entries are identified and removed from the dataset to maintain the integrity of the analysis.

### Functionality:

- **Identification of Duplicates:**
  - **AutoInsight.ai** automatically detects duplicate rows in the dataset using pandas' `duplicated()` function.
  - It checks for exact duplicates (rows with the same values across all columns) and provides a count of how many duplicate rows exist.
- **Removal of Duplicates:**
  - The system provides the option to remove all duplicate rows in one step. This is especially useful when dealing with large datasets.
  - **AutoInsight.ai** also gives the option to keep the first or last occurrence of the duplicated rows while removing the others.
- **User Customization:**
  - Users can choose to remove duplicates based on selected columns instead of the entire row. This is particularly useful when certain columns (e.g., ID numbers) should

not be duplicated, even if other data in the row might differ.

- **Post-Removal Summary:**

- After removing duplicates, the system provides a summary report detailing the number of duplicates removed and the total number of rows remaining.
- Users can view the updated dataset to ensure that no necessary data was lost in the process.

### 6.3 Encoding and Transformation

**Purpose:** Machine learning models require numerical data to process input. Therefore, categorical data (strings or labels) must be encoded into numerical formats. The Encoding and Transformation module in **AutoInsight.ai** handles categorical encoding and feature transformations automatically to ensure compatibility with machine learning algorithms.

#### Functionality:

- **Categorical Encoding:**

- **Label Encoding:**

- For columns with ordinal data (where there is a natural order, like "low," "medium," "high"), **AutoInsight.ai** uses Label Encoding to convert categories into numerical values. Each category is assigned a unique integer value (e.g., 0, 1, 2).

- **One-Hot Encoding:**

- For nominal data (where categories do not have an inherent order), **AutoInsight.ai** applies One-Hot Encoding.
- This process creates a new binary column for each category. For example, a column with categories ["red," "blue," "green"] will be converted into three new columns (red, blue, green), with 1 indicating the presence of that category in a row and 0 indicating its absence.

- **Target Encoding (Optional):**

- **AutoInsight.ai** also supports Target Encoding, where categorical features are encoded based on their relationship with the target variable. This technique can help improve the model's predictive power by capturing the correlation between categories and the target.
- This is particularly useful for categorical variables with many unique values.

- **Feature Scaling (Normalization/Standardization):**
  - For numerical features, **AutoInsight.ai** provides the option to perform scaling:
    - **Normalization** (Min-Max Scaling): Rescales numerical features to a range between 0 and 1.
    - **Standardization** (Z-score Scaling): Centers the data around 0 and scales it according to the standard deviation (useful for data with a normal distribution).
- **Data Transformation:**
  - The system offers transformation techniques like:
    - **Log Transformation:** Applied to reduce the skewness in the data and make the distribution more normal.
    - **Polynomial Features:** Users can choose to create interaction terms or higher-order features to capture nonlinear relationships between variables.
- **User Customization:**
  - Users can manually choose the encoding method for each categorical column (Label Encoding, One-Hot Encoding, or Target Encoding).
  - The system also allows the user to select whether they want to apply scaling or leave the data in its original form.
- **Post-Transformation Summary:**
  - After encoding and transforming the dataset, a summary is provided to the user showing the updated dataset's structure. This includes:
    - The number of columns that were transformed.
    - Details on which columns were encoded or scaled.
    - A preview of the newly transformed dataset.

## CHAPTER 7

### EXPLORATORY DATA ANALYSIS (EDA)

Exploratory Data Analysis (EDA) is a critical process in understanding the underlying patterns in data, detecting anomalies, checking assumptions, and visualizing distributions before applying machine learning algorithms. **AutoInsight.ai** automates key steps in EDA, enabling users to explore their data intuitively and gain valuable insights without needing to write any code. This section delves into the Profile Report Generation and the Insights Gained from the EDA process.

## 7.1 Profile Report Generation

**Purpose:** The Profile Report Generation is a comprehensive report that automatically generates detailed statistics and visualizations for every feature in the dataset. This report helps users quickly understand the characteristics and quality of the data, which is essential for guiding the next steps in the data science pipeline.

### Functionality:

- **Data Overview:**
  - **AutoInsight.ai** automatically generates an overview of the dataset, displaying key information like:
    - Number of rows and columns in the dataset.
    - Data types for each column (e.g., numerical, categorical, datetime).
    - Non-null counts for each column (indicating how many values are missing).
- **Descriptive Statistics:**
  - For each numerical column, the profile report includes essential statistics such as:
    - **Mean:** The average of the column.
    - **Median:** The middle value when the data is ordered.
    - **Standard deviation:** A measure of how much the values deviate from the mean.
    - **Minimum and Maximum values:** To understand the range of data.
    - **Quantiles (25th, 50th, 75th percentiles):** To show the distribution of the data.
  - For categorical columns, the profile report includes:
    - **Unique values:** The number of unique categories in each column.
    - **Top categories:** The most frequent categories.
    - **Frequency count:** How often each category appears in the dataset.
- **Correlation Matrix:**
  - **AutoInsight.ai** computes the correlation between numerical features to identify relationships. The system displays this as a correlation heatmap, which visually represents correlations (positive or negative) between features, allowing users to spot multicollinearity or redundant features early in the analysis.
- **Missing Data Visualization:**

- The profile report provides a visualization of missing data, often displayed as a heatmap or bar chart, showing which columns have missing values and their proportion.
- **Data Distribution Visualization:**
  - **AutoInsight.ai** provides visual summaries of the data distribution using histograms, boxplots, and density plots for numerical columns. This allows users to understand:
    - The shape of the distribution (whether it's normal, skewed, etc.).
    - Outliers in the data (visible in boxplots).
    - Range and spread of the data.
- **Data Type Summary:**
  - A summary table is included, showing how each feature is categorized (e.g., numerical, categorical, boolean) and its distinct characteristics. This helps users understand what preprocessing steps might be needed, such as encoding categorical variables or scaling numerical ones.

#### **Benefits:**

- The Profile Report is generated in an automated manner, allowing users to quickly comprehend the structure and health of the dataset.
- Users gain insights into missing data patterns, correlations, distributions, and basic statistics, enabling them to make informed decisions about the next steps in data cleaning.

## **7.2 Insights Gained**

**Purpose:** The Insights Gained section focuses on providing valuable, actionable conclusions based on the EDA process. It includes understanding the distribution, relationships, and quality of the data, which directly influences model selection, feature engineering, and the overall approach to solving the business problem.

#### **Functionality:**

- **Understanding Data Quality:**
  - **Missing Data Insights:**
    - **AutoInsight.ai** highlights the columns with the most missing values and suggests potential imputation strategies. If a significant portion of data is missing, the system might suggest removing the feature altogether or using imputation methods to fill in missing values.

- For example, if a column has over 30% missing values, it might be recommended to remove it from further analysis.
- **Duplicate Data Insights:**
  - If duplicates exist in the dataset, the platform flags them and suggests removing the duplicates to ensure data integrity.
- **Feature Distribution Insights:**
  - **Numerical Features:**
    - By analysing the distribution of numerical columns (via histograms, boxplots, etc.), the system helps identify:
      - **Skewed distributions:** If a numerical column has a skewed distribution, it may require log transformation or other techniques to normalize the data.
      - **Outliers:** Features that exhibit extreme outliers may need special handling such as clipping or removal.
    - Users gain insights into which features have a normal distribution and which ones might need scaling or transformation for optimal model performance.
- **Correlation Insights:**
  - **Multicollinearity Detection:**
    - By computing correlations between numerical features, **AutoInsight.ai** helps identify highly correlated features. Strong correlations (close to +1 or -1) can be problematic because they introduce multicollinearity, which may lead to overfitting or instability in models.
  - **Feature Interactions:**
    - Understanding feature interactions allows the system to highlight features that may need to be combined into new derived features (e.g., adding interaction terms like X \* Y or creating polynomial features).
- **Categorical Feature Insights:**
  - The system provides insights into categorical features, such as:
    - The **number of unique values** in each column: If a categorical feature has too many categories (e.g., a column with hundreds of unique values), it might need encoding or a dimensionality reduction technique.

- The **distribution of categories**: Identifying if certain categories dominate and may need to be handled as outliers or grouped into a smaller number of categories.
- **Outliers & Anomalies:**
  - The system identifies potential outliers using boxplots or Z-scores for numerical columns. The user is notified of any extreme outliers, which can distort model training and performance. The system suggests appropriate methods (e.g., clipping or transformation) to handle outliers.
- **Feature Engineering Suggestions:**
  - Based on the EDA insights, **AutoInsight.ai** may suggest creating new features by combining existing ones (e.g., ratio of two numerical features) or encoding categorical features in a particular way.
  - The system also suggests possible interactions between features that could improve the model's predictive power.

**Benefits:**

- The Insights Gained section automates the process of interpreting data and providing actionable feedback.
- It gives the user a clear understanding of how the features relate to one another, potential issues with the data (such as outliers or missing values).

## CHAPTER 8

### AUTOMATED MACHINE LEARNING (**AutoML**) MODELING

Automated Machine Learning (AutoML) refers to the process of automating the end-to-end process of applying machine learning to real-world problems. In **AutoInsight.ai**, the AutoML pipeline is designed to provide users with a simplified and intuitive approach to model building, evaluation, and selection. Users without in-depth machine learning knowledge can leverage this tool to quickly create, compare, and select the best-performing model. This section explains the key steps in AutoML Modeling:

#### 8.1 Target Selection

**Purpose:** Target Selection refers to identifying the dependent or target variable (the one to be predicted) in the dataset. This step is essential as the selected target determines the type of machine

learning task (classification or regression) that will be performed.

#### **Functionality:**

- **AutoInsight.ai** automatically analyzes the dataset to detect possible target variables based on data types and user input.
- **User Input for Target:**
  - Users are prompted to select the target variable manually. The system presents a list of all potential target variables, which can be categorical (for classification tasks) or numerical (for regression tasks).
  - The user can choose the column that represents the outcome they want to predict (e.g., predicting a customer's buying behavior or the price of a product).
- **Target Variable Analysis:**
  - Once the target variable is selected, the system assesses its data type, its relationship with other features, and its suitability for the chosen machine learning task.
  - If the selected target variable is categorical (e.g., class labels), **AutoInsight.ai** will set up a classification pipeline.
  - If the target is numerical (e.g., price, age), the system will configure a regression pipeline.
- **Auto Insight:** If no explicit target is provided, **AutoInsight.ai** can recommend the most likely target variable based on feature correlation, class distribution, and business context.

#### **Benefits:**

- Ensures the user selects the correct target variable for the task, avoiding issues related to misalignment in model training.
- Helps automate the selection process for users with limited ML experience.

## **8.2 Model Comparison**

**Purpose:** Model comparison refers to evaluating multiple machine learning models to find the one that performs best on the dataset. **AutoInsight.ai** automates this process, training and comparing a wide range of algorithms for both classification and regression tasks.

#### **Functionality:**

- **Model Selection and Comparison:**
  - **AutoInsight.ai** uses an extensive library of machine learning models and

automatically trains them on the dataset. This includes popular algorithms such as:

- For **Classification**: Logistic Regression, Decision Trees, Random Forest, Support Vector Machines (SVM), Gradient Boosting, k-Nearest Neighbors (KNN), etc.
- For **Regression**: Linear Regression, Decision Trees, Random Forest, Support Vector Regression (SVR), etc.
- **Auto Selection of Hyperparameters:**
  - The system automatically tunes hyperparameters for each model using techniques like grid search or random search, which increases the likelihood of discovering the best model configuration.
  - Users are presented with a table of different models, including their hyperparameters, training time, and performance metrics.
- **Model Evaluation Metrics:**
  - For classification models, performance metrics such as accuracy, precision, recall, F1-score, and AUC-ROC curve are computed.
  - For regression models, metrics like R-squared ( $R^2$ ), Mean Squared Error (MSE), Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE) are computed.

#### **Benefits:**

- Simplifies the process of trying and comparing multiple machine learning models, helping users identify the best model for their problem.
- Removes the need for deep expertise in different algorithms and their configurations.

### **8.3 Best Model Selection**

**Purpose:** Best Model Selection refers to automatically selecting the model that achieves the best performance based on the evaluation metrics. This step ensures that the final model used for predictions is the one that provides the most accurate, reliable, and optimized results.

#### **Functionality:**

- **AutoInsight.ai** uses the performance metrics (such as accuracy for classification or  $R^2$  squared for regression) to automatically rank the trained models.
- The system evaluates which model performs best based on the highest accuracy (or other relevant metric), lowest error, or the most consistent performance across validation folds.
- **Selection Criteria:**

- The best model is selected based on:
  - **Validation performance:** The model's accuracy or error on the validation data (not just the training set).
  - **Complexity:** Some models, like deep learning models, may have higher accuracy but require more time and resources to train. **AutoInsight.ai** takes both performance and computational efficiency into account.
  - **Business context:** The system can also weigh different metrics depending on the business context, such as prioritizing precision over recall for fraud detection.
- **User Notification:**
  - The system presents the top-performing model and allows the user to inspect the model's details, including the hyperparameters and its corresponding performance metrics.

#### **Benefits:**

- Automates the selection of the best model, which helps ensure that the most effective model is used without requiring manual intervention.
- Allows users to trust that the model chosen is optimal based on various performance metrics.

#### **8.4 Training and Saving the Model**

**Purpose:** Once the best model is selected, the next step is to train the model on the full dataset and save it for future use. This step ensures the model is ready for deployment or further predictions.

#### **Functionality:**

- **Model Training on Full Dataset:**
  - After selecting the best model, **AutoInsight.ai** retrains the model using the entire dataset (including the training and validation data) to maximize its performance.
  - This ensures the model is trained with all available data, leading to better generalization and accuracy when deployed.
- **Hyperparameter Optimization:**
  - The model is fine-tuned using the optimal hyperparameters identified during the model comparison phase.
  - The system may employ techniques such as cross-validation to further improve the robustness and performance of the model.

- **Model Saving:**
  - Once the model is trained, **AutoInsight.ai** saves the model in a serialized format (e.g., using Pickle or Joblib for Python models). This allows the model to be easily deployed or used for prediction at any time without retraining.
  - The saved model includes the trained parameters and the preprocessing pipeline (if applicable), so it can be used to make predictions on new, unseen data.
- **Model Export:**
  - Users can download the trained model for further use or deploy it on their own platforms. The system offers model export options in Pickle format.

**Benefits:**

- Automates the retraining and saving of the best model, making it ready for deployment or further usage.

## CHAPTER 9

### MODEL EVALUATION METRICS

In machine learning, evaluating the performance of a trained model is crucial to determine how well it generalizes to unseen data. In the context of **AutoInsight.ai**, we utilize various performance metrics to assess the quality of classification models. These metrics help in understanding the strengths and weaknesses of the model, guiding the user in making informed decisions.

#### 9.1 Accuracy, F1 Score, Recall, Precision

**9.1.1 Accuracy:** **Accuracy** is one of the most commonly used metrics for classification models. It is the proportion of correct predictions made by the model to the total number of predictions. It is useful when the dataset is balanced (i.e., the classes are equally represented).

**Formula:**

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Number of Predictions}}$$

**Interpretation:**

- High accuracy means the model is making predictions that closely match the actual outcomes.
- **Limitation:** Accuracy can be misleading in cases of imbalanced datasets, where one class

dominates.

**9.1.2 F1 Score:** The F1 Score is the harmonic mean of precision and recall, providing a balance between the two. It is particularly useful when the classes are imbalanced or when the cost of false positives and false negatives is not equal.

**Formula:**

$$\text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

**Interpretation:**

- The F1 score helps in assessing a model's overall performance in terms of its ability to balance false positives and false negatives.
- Higher F1 score means a better balance between precision and recall, making it ideal for applications where both false positives and false negatives are costly.

**9.1.3 Recall (Sensitivity):** Recall (also called sensitivity) is the proportion of actual positive instances that were correctly identified by the model. It measures the ability of a model to identify all relevant instances within the dataset.

**Formula:**

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

**Interpretation:**

- High recall means fewer false negatives, making it important for situations where detecting the positive class is critical (e.g., medical diagnosis, fraud detection).
- Limitation: It may lead to many false positives, which can reduce precision.

**9.1.4 Precision:** Precision is the proportion of true positive predictions among all predicted positives. It focuses on how accurate the positive predictions are and helps understand how many of the predicted positive cases are relevant.

**Formula:**

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

### **Interpretation:**

- **High precision** means fewer false positives, making it essential in applications where incorrectly labeling a negative instance as positive has a significant cost (e.g., spam email classification).
- **Limitation:** Precision may be high at the cost of recall, especially in imbalanced datasets.

## **9.2 ROC-AUC, Kappa, MCC**

**9.2.1 ROC-AUC (Receiver Operating Characteristic - Area Under Curve):** The ROC Curve is a graphical representation of a model's performance, plotting the True Positive Rate (Recall) against the False Positive Rate (1 - Specificity). The AUC (Area Under the Curve) is a scalar value between 0 and 1 that summarizes the performance of the model across all classification thresholds.

### **Interpretation:**

- An AUC value of 1 means perfect performance, while an AUC of 0.5 indicates that the model is no better than random guessing.
- Higher AUC signifies better model performance, especially when dealing with imbalanced datasets.

**9.2.2 Kappa (Cohen's Kappa):** Cohen's Kappa is a metric that measures the agreement between the predicted and actual labels while accounting for the agreement that could happen by chance.

### **Formula:**

$$\text{Kappa} = \frac{P_o - P_e}{1 - P_e}$$

- $P_{oPo}$  is the observed agreement.
- $P_{ePe}$  is the expected agreement by chance.

### **Interpretation:**

- **Kappa = 1** means perfect agreement, while **Kappa = 0** means the agreement is no better than random chance.
- **Negative Kappa** values indicate worse-than-chance agreement, which is uncommon but can

happen in certain datasets.

**9.2.3 MCC (Matthews Correlation Coefficient):** MCC is a balanced metric that considers all four quadrants of the confusion matrix (True Positives, True Negatives, False Positives, False Negatives). It ranges from -1 to 1, with **1** indicating perfect classification, 0 indicating no better than random prediction, and **-1** indicating complete disagreement between prediction and actual values.

**Formula:**

$$MCC = \frac{(True\ Positives) * (True\ Negatives) - (False\ Positives) * (False\ Negatives)}{\sqrt{(True\ Positives + False\ Positives) \times (True\ Positives + False\ Negatives) \times (True\ Negatives + False\ Positives) \times (True\ Negatives + False\ Negatives)}}$$

**Interpretation:**

- **MCC close to 1:** Model is making predictions close to actual values.
- **MCC close to 0:** Model has poor predictive power, close to random guessing.
- **MCC close to -1:** Model's predictions are worse than random guessing.

### 9.3 Classification Report

The Classification Report is a comprehensive summary that combines various evaluation metrics, including precision, recall, f1-score, and support for each class. This report allows users to understand the model's performance on each class separately, especially useful in multiclass classification problems.

**Components of the Classification Report:**

- **Precision:** The percentage of true positives among predicted positives.
- **Recall:** The percentage of actual positives correctly identified by the model.
- **F1-Score:** The harmonic mean of precision and recall.
- **Support:** The number of actual occurrences of the class in the dataset.

**Interpretation:**

- The classification report gives a complete view of the model's performance, highlighting the strengths and weaknesses in predicting different classes.
- It also includes a bar graph which compares the Accuracy, F1 Score and AUC Curve Score.

## CHAPTER 10

### MODEL AND REPORT EXPORT

Once the machine learning model has been trained and evaluated, it's important to provide users with the capability to export the results, metrics, and the trained model itself. This allows the user to save their work, share it with others, or deploy the model into production environments. In **AutoInsight.ai**, the system provides three major export functionalities:

- **HTML EDA Report:** Allows users to generate a detailed report based on the exploratory data analysis performed on the uploaded dataset.
- **CSV Performance Metrics:** Enables users to download a CSV file containing the evaluation metrics of the trained model.
- **Trained Model (.pkl File):** Provides the ability to download the trained machine learning model in a .pkl (Pickle) file format, which can later be used for predictions or deployment.

#### 10.1 HTML EDA Report

The **HTML EDA Report** is an automated, easy-to-understand document that summarizes the Exploratory Data Analysis (EDA) performed on the uploaded dataset. This report helps users gain insights into their dataset, such as distributions, correlations, missing values, and potential data imbalances. The **AutoInsight.ai** platform allows users to generate this report at the click of a button.

##### Features of HTML EDA Report:

1. **Dataset Summary:**
  - Displays basic information about the dataset, including the number of rows and columns, data types, and the presence of any missing values.
  - Provides an overview of the statistical summary (mean, median, standard deviation, etc.) of numeric features.
2. **Missing Values:**
  - A section showing the count and percentage of missing values in each feature, along with visualizations (such as heatmaps) to highlight the missing data.
3. **Data Distribution:**
  - Visualizes the distribution of each feature using histograms for numerical data and **bar charts** for categorical data. This helps users understand the spread and frequency of the values.

#### **4. Correlation Analysis:**

- Displays a correlation heatmap to indicate how different numerical features are related to each other. This helps identify potential collinearity issues.

#### **5. Outlier Detection:**

- Includes box plots to visualize and identify potential outliers in the dataset.

#### **6. Class Distribution:**

- Visualizes the distribution of the target variable (in the case of classification) or continuous target (in the case of regression). This is especially useful for detecting class imbalance.

#### **7. Additional Visualizations:**

- Users can view scatter plots, pair plots, or other visualizations that can help in understanding relationships between features and the target variable.

#### **Benefits:**

- **Comprehensive Understanding:** Provides an easy-to-read summary of the dataset, allowing the user to make better-informed decisions about preprocessing and model selection.
- **Reusable:** The HTML report can be saved and shared with others for collaboration or documentation purposes.
- **Interactive:** Users can explore visualizations interactively, which helps in uncovering more insights.

#### **Export Functionality:**

- The generated EDA report can be exported as a downloadable HTML file with a simple click. This file can be opened in any browser to view the complete report.

## **10.2 CSV Performance Metrics**

Once the model has been trained and evaluated, **AutoInsight.ai** generates a set of performance metrics, including accuracy, precision, recall, F1 score, AUC, confusion matrix, and others. These metrics give the user a thorough understanding of how well their model is performing.

#### **Features of CSV Performance Metrics:**

##### **1. Metric Names:**

- The CSV file includes column headers corresponding to different performance metrics (accuracy, precision, recall, F1 score, etc.) for the model evaluated.

## 2. Model Selection:

- The best-performing model (based on the highest F1 score, accuracy, or AUC) can be highlighted in the CSV file for easy identification.

## 3. Export Functionality:

- The user can download the CSV file by clicking an "Export Metrics" button after evaluating the models. The CSV file can be opened in Excel, Google Sheets, or any other data analysis tool for further examination or reporting.

### Benefits:

- **Easy Analysis:** The CSV format makes it easy for users to analyze the results further using spreadsheets or statistical tools.
- **Flexible Use:** The CSV file can be shared with team members or included in project reports for further evaluation.

## 10.3 Trained Model (.pkl File)

Once the model has been trained and evaluated, **AutoInsight.ai** allows users to export the trained model as a Pickle (.pkl) file. This file contains the serialized version of the trained machine learning model, which can be easily loaded into other environments for prediction or deployment.

### Features of Trained Model Export:

#### 1. Pickle Serialization:

- The trained model is serialized using Python's Pickle library, which allows it to be saved in a format that can be loaded later without needing to retrain the model.

#### 2. Exported Model Content:

- The Pickle file contains the entire model object, including the trained weights, feature importance, hyperparameters, and other model parameters, making it easy to reload the exact same model.

#### 3. Model Deployment:

- The Pickle file can be used for real-time predictions by loading the model into a Python environment or web server.
- It can be deployed into production systems where the model can provide predictions on new incoming data without requiring retraining.

#### 4. Export Functionality:

- The trained model can be exported with a single click via the "Export Model" button, allowing the user to download the .pkl file to their local machine or cloud storage.

### **Benefits:**

- **Convenience:** The Pickle file simplifies the process of saving and reusing trained models, making it easier to integrate machine learning into production workflows.
- **Portability:** The .pkl file can be easily shared and loaded into any Python environment, allowing for deployment or further evaluation without retraining the model.
- **Consistency:** By exporting the trained model, users can ensure the exact same model is used in different environments (e.g., development, testing, and production).

## **CHAPTER 11**

### **CHALLENGES FACED AND SOLUTIONS**

Throughout the development of the **AutoInsight.ai** platform, various technical and non-technical challenges arose. However, each challenge was addressed through targeted problem-solving approaches to ensure a seamless experience for the end users.

#### **11.1 Data Preprocessing Challenges**

- **Issue:** One of the most common challenges encountered during the initial stages was dealing with inconsistent or messy data. Datasets frequently contained missing values, duplicates, and irrelevant features that needed to be handled before applying machine learning models.
- **Solution:** The solution involved creating automated data preprocessing pipelines that could handle common issues such as filling in missing values, removing duplicates, and encoding categorical variables. Additionally, we incorporated the PyCaret library to automate much of the data cleaning and transformation process, ensuring that users didn't have to write any code.

#### **11.2 Model Evaluation and Comparison**

- **Issue:** While comparing multiple machine learning models, it was difficult to make meaningful comparisons between models, especially when evaluating models with imbalanced datasets or specific performance metrics that differed across models.
- **Solution:** The platform integrated several evaluation metrics (such as accuracy, F1 score, ROC-AUC, etc.) and provided a clear visualization of the results. We also implemented cross-validation techniques to ensure a fair comparison of the models.

### **11.3 Scalability and Performance**

- **Issue:** As the dataset size grew, the platform began to slow down, especially during the model training phase. This became a bottleneck for users uploading large datasets for analysis.
- **Solution:** To address this, we optimized the backend processing using multi-threading and **parallel processing** techniques. Additionally, by using cloud-based **services** for model training and evaluation, we ensured the platform could handle larger datasets efficiently.

### **11.4 User Experience (UX)**

- **Issue:** Initially, the platform was too complex for non-technical users, which led to difficulty in navigating the user interface, particularly during the model selection and report generation stages.
- **Solution:** The user interface was simplified with a clean, intuitive layout. Interactive tooltips and step-by-step guidance were added to assist users, especially beginners, in understanding the processes involved in uploading datasets, selecting models, and interpreting results.

## **CHAPTER 12**

### **KEY LEARNINGS AND SKILLS ACQUIRED**

The development of **AutoInsight.ai** provided an immense learning opportunity in both technical and soft skills. Below are the key learnings and skills acquired during the internship:

#### **12.1 Technical Skills**

- **Machine Learning Concepts:**
  - Gained a deeper understanding of machine learning algorithms, including regression, classification, clustering, and model evaluation techniques.
  - Learned how to integrate the PyCaret library for automated machine learning, simplifying model selection, training, and evaluation.
- **Data Preprocessing:**
  - Developed strong knowledge in data preprocessing techniques, including handling missing values, encoding categorical variables, and scaling features for optimal model performance.
- **Model Evaluation and Metrics:**
  - Improved my understanding of different evaluation metrics such as F1 score, ROC-AUC, confusion matrix, and cross-validation methods, and how they can be applied

to select the best model.

- **Deployment Using Streamlit:**
  - Acquired experience in deploying machine learning models and data applications using Streamlit, making it easy to create interactive, user-friendly interfaces for non-technical users.
- **Cloud Integration:**
  - Gained experience in integrating cloud platforms for scalable model training and data processing, allowing for the handling of large datasets effectively.

## 12.2 Soft Skills

- **Problem-Solving:**
  - The project enhanced my problem-solving abilities, particularly when dealing with unstructured data, implementing model evaluation, and finding solutions for user-related challenges.
- **Collaboration and Communication:**
  - Worked closely with mentors, colleagues, and other team members to address issues, plan milestones, and meet deadlines. Improved my communication skills, especially when explaining technical details to non-technical stakeholders.
- **Project Management:**
  - Enhanced my project management skills by contributing to planning, execution, and testing stages, ensuring that the platform met both functional and non-functional requirements.

## CHAPTER 13 CONCLUSION

In conclusion, the development of **AutoInsight.ai** has provided a comprehensive and hands-on learning experience in democratizing machine learning through no-code solutions. By focusing on automating data preprocessing, model selection, and performance evaluation, the platform makes machine learning accessible to users with minimal technical expertise. The project has highlighted the importance of user-friendly design, ensuring that both technical and non-technical users can seamlessly interact with the platform. Moreover, the integration of cloud-based services and

deployment through Streamlit has further strengthened the platform's scalability and usability. Moving forward, the platform holds great potential for enhancements, such as hyperparameter tuning, NLP integration, and personalized user experiences, making it a versatile tool for data exploration and predictive modeling. The knowledge gained throughout this internship not only strengthened my technical skills in machine learning but also enhanced my ability to develop intuitive, user-centric applications that bridge the gap between advanced technologies and end-users.

## CHAPTER 14

### SOURCE CODE

```
import streamlit as st
import pandas as pd
import os
import base64
from fpdf import FPDF
import ydata_profiling
from streamlit_pandas_profiling import st_profile_report
from pycaret.classification import setup, compare_models, pull, save_model,predict_model
from sklearn.metrics import accuracy_score, f1_score,
roc_auc_score,classification_report,recall_score, precision_score, cohen_kappa_score,
matthews_corrcoef
import matplotlib.pyplot as plt
import seaborn as sns

# Initialize empty dataframe
df = None

# Load dataset if it exists
if os.path.exists("dataset.csv"):
    df = pd.read_csv("dataset.csv", index_col=None)

# Sidebar
with st.sidebar:
    st.image("Logo.jpg")
    st.title("AutoInsight.ai")
    st.info("AUTO TRAIN YOUR DATASET WITHIN A FEW CLICKS.")
    choice = st.radio("Select the task", ["Upload", "Profiling", "Modeling", "Download"])
    st.info("Analyze --> Train --> View Insights --> Download")
```

```

# Upload section
if choice == "Upload":
    st.title("Upload Your Dataset")
    file = st.file_uploader("Upload Your Dataset", type=['csv'])
    if file:
        df = pd.read_csv(file, index_col=None)
        df.to_csv("dataset.csv", index=None)
        st.success("File uploaded and saved successfully!")
        st.dataframe(df)
    if df is not None:
        st.subheader("🔧 Data Preprocessing & Cleaning")
        original_shape = df.shape
        duplicates = df.duplicated().sum()
        null_columns = df.columns[df.isnull().all()]
        null_values = df.isnull().sum().sum()

        # Remove duplicates
        df.drop_duplicates(inplace=True)

        # Drop fully null columns
        df.dropna(axis=1, how='all', inplace=True)

        # Fill missing values (basic strategy)
        for col in df.select_dtypes(include='number').columns:
            df[col].fillna(df[col].median(), inplace=True)

        for col in df.select_dtypes(include='object').columns:
            df[col].fillna(df[col].mode()[0], inplace=True)

        cleaned_shape = df.shape

        # Save the cleaned data
        df.to_csv("dataset.csv", index=False)
        if 'Unnamed: 0' in df.columns:
            df.drop(columns=['Unnamed: 0'], inplace=True)

        # Summary of cleaning
        st.success("✅ Data cleaned successfully!")
        st.markdown(f"**Original Shape:** {original_shape} → **Cleaned Shape:** {cleaned_shape}")

```

```

st.markdown(f"- Removed **{duplicates}** duplicate rows")
st.markdown(f"- Dropped **{len(null_columns)}** columns with all nulls")
st.markdown(f"- Filled **{null_values}** missing values")

# Profiling section
if choice == "Profiling":
    if df is not None:
        st.title("Exploratory Data Analysis")
        profile = ydata_profiling.ProfileReport(df, explorative=True)
        profile.to_file("eda_report.html")
        st_profile_report(profile)
    else:
        st.warning("Please upload a dataset first!")

# Modeling section
elif choice == "Modeling":
    if df is not None:
        st.title("🔧 Model Building")

        # Ensure session_state keys exist
        if "model_trained" not in st.session_state:
            st.session_state.model_trained = False
        if "model_metrics" not in st.session_state:
            st.session_state.model_metrics = None

    chosen_target = st.selectbox("🎯 Select the Target Column", df.columns)

    if st.button("🚀 Run Modeling"):
        # Step 1: Check for rare classes
        value_counts = df[chosen_target].value_counts()
        rare_classes = value_counts[value_counts < 2]
        if not rare_classes.empty:
            st.error("❗ The following classes have less than 2 samples, which will cause an error during model training:")
            st.write(rare_classes)
            st.warning("Please choose a different target column or ensure each class has at least 2 samples.")
        else:
            # Step 2: PyCaret setup
            setup(df, target=chosen_target, verbose=False, session_id=123)

```

```

# Step 3: Show setup config
setup_df = pull()
st.subheader("⚙️ Setup Configurations")
st.dataframe(setup_df)

# Step 4: Compare models
best_model = compare_models()
compare_df = pull()
st.subheader("📊 Model Comparison")
st.dataframe(compare_df)

# Step 5: Save best model
save_model(best_model, "trained_model")

# Step 6: Show best model
st.subheader("🏆 Best Performing Model")
st.write(type(best_model).__name__)
params_df = pd.DataFrame(best_model.get_params().items(), columns=["Parameter",
"Value"])
st.markdown("### 🔪 Model Parameters")
st.dataframe(params_df)

# Step 7: Predict and Evaluate
predicted_df = predict_model(best_model)
y_true = predicted_df[chosen_target]
y_pred = predicted_df["prediction_label"]

# Step 8: Calculate metrics
accuracy = accuracy_score(y_true, y_pred)
f1 = f1_score(y_true, y_pred, average="weighted")
recall = recall_score(y_true, y_pred, average="weighted")
precision = precision_score(y_true, y_pred, average="weighted")
kappa = cohen_kappa_score(y_true, y_pred)
mcc = matthews_corrcoef(y_true, y_pred)
try:
    auc = roc_auc_score(pd.get_dummies(y_true), pd.get_dummies(y_pred),
average="macro", multi_class="ovr")
except:
    auc = None

```

```

# Save metrics to session_state
metrics_df = pd.DataFrame({
    "Metric": ["Accuracy", "F1 Score", "AUC Score", "Recall", "Precision", "Kappa",
    "MCC"],
    "Score": [accuracy, f1, auc if auc is not None else "N/A", recall, precision, kappa,
mcc]
})

st.session_state.model_trained = True
st.session_state.model_metrics = metrics_df

# Step 9: Classification report
report_df = pd.DataFrame(classification_report(y_true, y_pred,
output_dict=True)).transpose()
report_df.to_csv("model_metrics.csv")

# Step 10: Graphical visualization
st.subheader("📈 Model Performance Metrics")
st.dataframe(metrics_df)

graph_df = metrics_df[metrics_df["Metric"].isin(["Accuracy", "F1 Score"])]
if auc is not None:
    graph_df = pd.concat([
        graph_df,
        pd.DataFrame([{"Metric": "AUC Score", "Score": auc}])
    ], ignore_index=True)

st.bar_chart(graph_df.set_index("Metric"))

st.success("✅ Model trained and metrics saved!")

else:
    st.warning("Please upload and clean the dataset first.")

# Download section - Replaced functionality
if choice == "Download":
    st.title("Download Reports")

    # Download EDA Report as PDF
    if os.path.exists("eda_report.html"):
        st.subheader("EDA Report (HTML)")

```

```

# Convert HTML report to PDF using external tools (example assumes it's saved externally)
# For demo, show the HTML option with download button
with open("eda_report.html", "r", encoding='utf-8') as f:
    html_content = f.read()
    b64 = base64.b64encode(html_content.encode()).decode()
    href = f'<a href="data:file/html;base64,{b64}" download="EDA_Report.html">Download EDA Report as HTML</a>'
    st.markdown(href, unsafe_allow_html=True)
else:
    st.warning("EDA report not found. Please run profiling first.")

# Download Model Metrics
if os.path.exists("model_metrics.csv"):
    st.subheader("")
    if os.path.exists("model_metrics.csv"):
        with open("model_metrics.csv", "rb") as f:
            st.download_button(
                label=" Download Model Metrics (CSV)",
                data=f,
                file_name="model_metrics.csv",
                mime="text/csv" )
    else:
        st.warning("No model metrics to download. Run modeling first.")

#  Trained model download
st.subheader("Download Trained Model")
if os.path.exists("trained_model.pkl"):
    with open("trained_model.pkl", "rb") as f:
        st.download_button("📦 Download Trained Model (.pkl)", f,
file_name="trained_model.pkl")
else:
    st.warning("❌ No trained model found. Please run modeling first.")

```

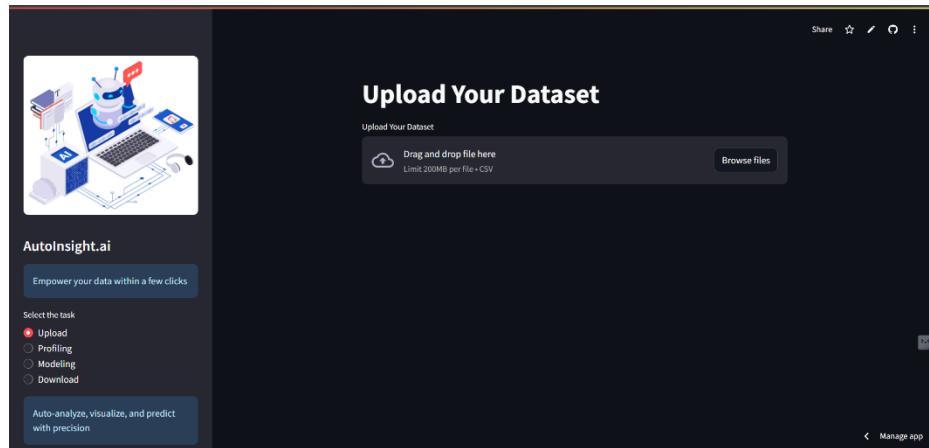
## CHAPTER 15

### TESTING OF APPLICATION USING DATA SET

To evaluate and validate the functionality of the developed application, an employee attrition dataset was utilized during the testing phase. This dataset contains structured records of employee information including variables such as age, job role, department, marital status, job satisfaction, years at company, overtime status, and the target variable 'Attrition', which indicates whether an employee has left the organization. The dataset provides a realistic and diverse mix of categorical and numerical features, enabling thorough testing of the application's modules—ranging from data cleaning and preprocessing to model training and evaluation. It effectively tested the robustness of AutoInsight.ai in automating exploratory data analysis (EDA), generating insights, and predicting employee attrition using AutoML, ensuring that all functionalities operated as intended.

## CHAPTER 16

### SCREENSHOTS OF THE APPLICATION



**Upload Your Dataset**

Dataset.csv (222 KB)

File uploaded and saved successfully!

Age	Attrition	BusinessTravel	Department	DistanceFromHome	Educat
1	41	No	Travel_Single	100	1
1	41	No	Travel_Frequently	379	5
2	31	No	Travel_Single	1313	3
3	21	No	Travel_Frequently	1302	3
4	21	No	Travel_Single	161	7
5	30	No	Travel_Frequently	1099	2
6	50	No	Travel_Single	1724	3
7	30	No	Travel_Single	1308	24
8	30	No	Travel_Frequently	216	19
9	30	No	Travel_Single	1299	27

**Data Preprocessing & Cleaning**

Data cleaned successfully!

Original Shape: (1470, 35) → Cleaned Shape: (1470, 35)

- Removed 0 duplicate rows
- Dropped 0 columns with all nulls
- Filled 0 missing values

**Exploratory Data Analysis**

YData Profiling Report

**Overview**

Brought to you by YData

Dataset statistics		Variable types	
Number of variables	23	Numeric	15
Number of observations	1470	Boolean	8
Missing cells	0	Categorical	17
Missing cells (%)	0.0%		
Duplicate rows	0		
Duplicate rows (%)	0.0%		
Total size in memory	1.1 MB		
Average record size in memory	796.8 B		

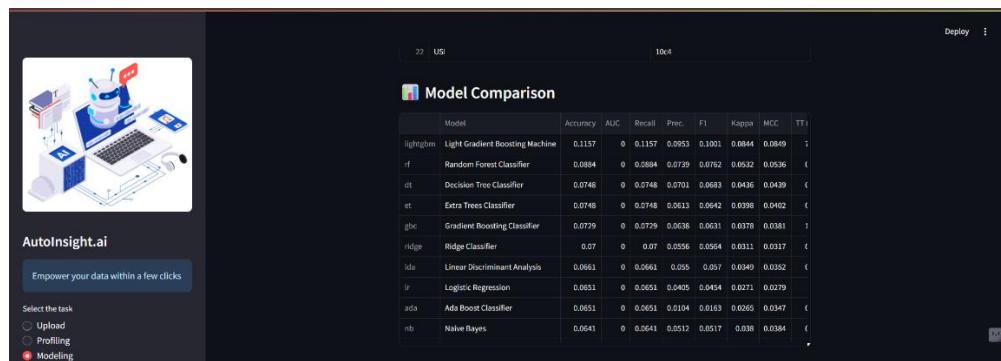
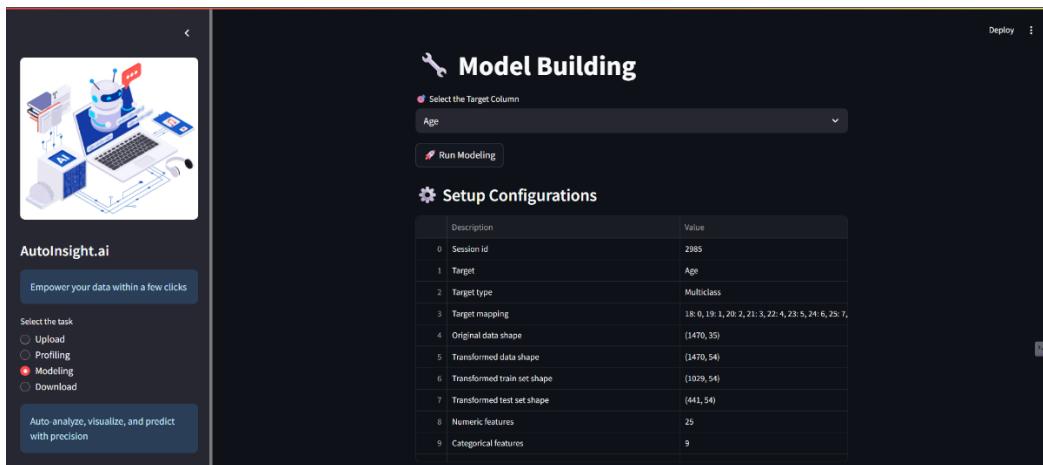
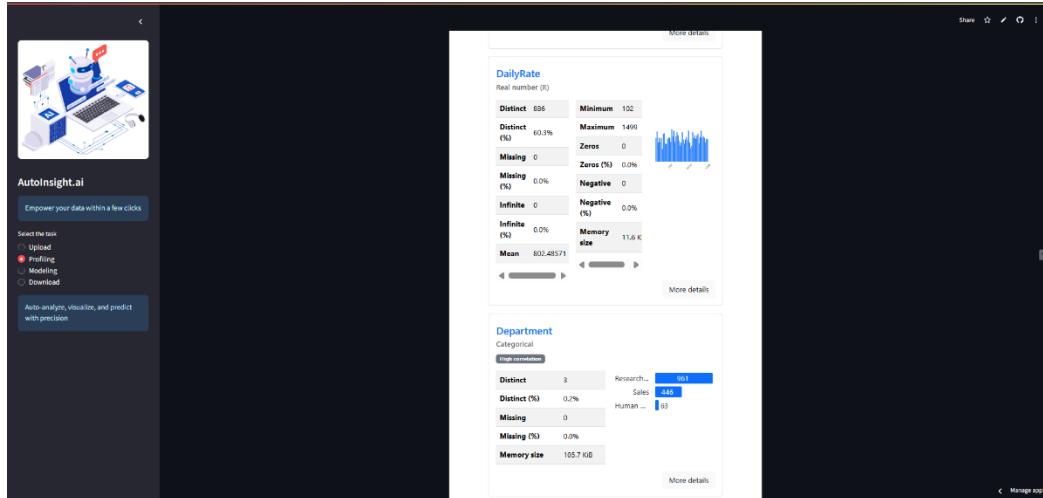
**Select Columns**

**Age**  
Real number (R)  
High correlation

Distinct	29%	Minimum	18
Distinct (%)	2.9%	Maximum	60
Missing	0	Zeros	0
Missing (%)	0.0%	Zeros (%)	0.0%
Infinites	0	Negative	0
Infinites (%)	0.0%	Negative (%)	0.0%
Mean	36.92881	Memory size	11.6 K

**Attrition**  
Boolean

Distinct	2	False	12/33
Distinct (%)	0.1%	True	1/33
Missing	0		
Missing (%)	0.0%		
Memory size	1.6 KB		



**Best Performing Model**  
LGBMClassifier

**Model Parameters**

Parameter	Value
0 boosting_type	gbdt
1 class_weight	None
2 colsample_bytree	1.0
3 importance_type	split
4 learning_rate	0.1
5 max_depth	-1
6 min_child_samples	20
7 min_child_weight	0.001
8 min_split_gain	0.0
9 n_estimators	100

Deploy ⚙️

**Prediction Results**

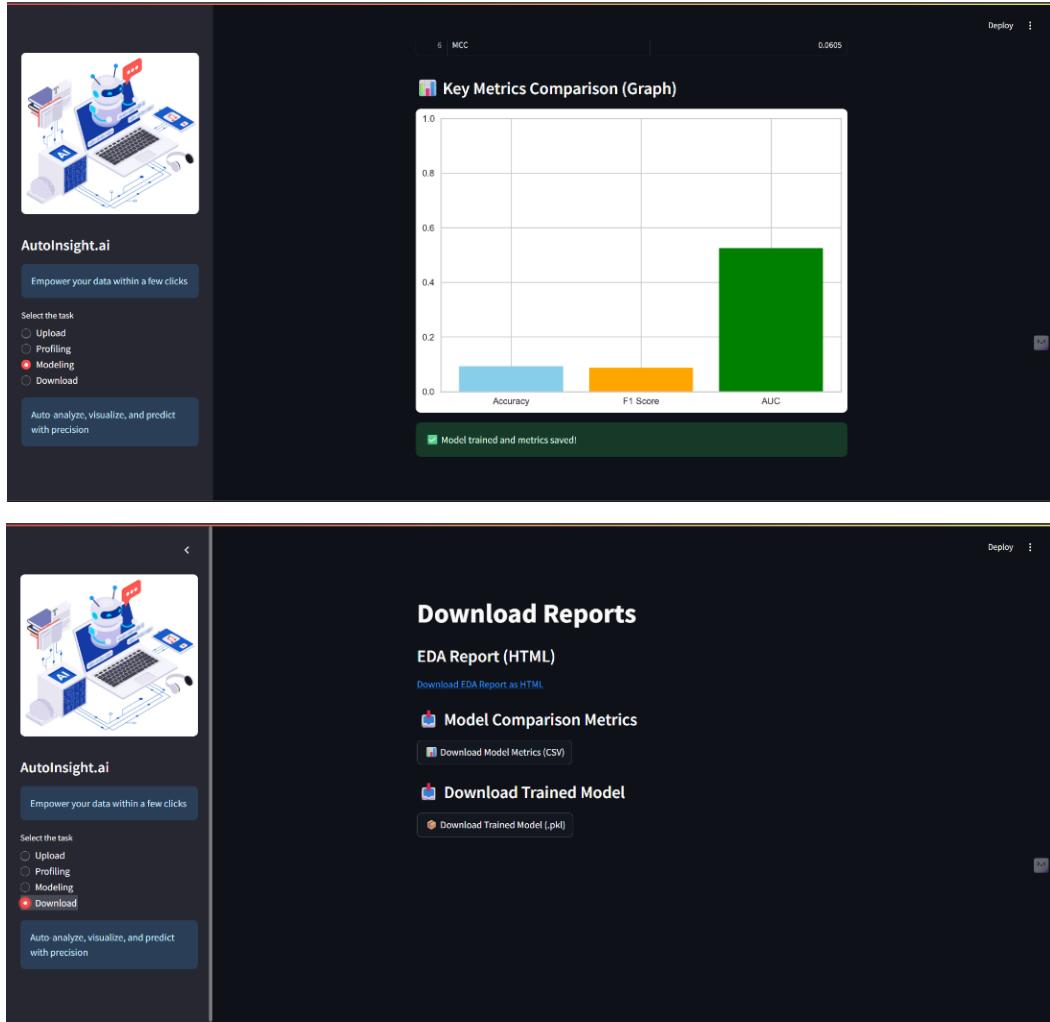
Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EmployeeCount
285 No	Travel_Rarely	1372	Research & Development	1	3	L
595 Yes	Travel_Rarely	286	Research & Development	2	4	L
1447 No	Non-Travel	301	Sales	15	4	M
328 No	Travel_Frequently	504	Sales	10	3	M
1332 Yes	Travel_Frequently	459	Research & Development	24	2	L
1426 No	Travel_Rarely	267	Research & Development	29	4	L
593 No	Travel_Rarely	676	Research & Development	1	3	O
154 No	Travel_Frequently	967	Sales	8	3	M
1106 Yes	Travel_Rarely	740	Sales	1	3	L
541 No	Non-Travel	427	Research & Development	8	3	L

Deploy ⚙️

**Model Performance Metrics**

Metric	Score
0 Accuracy	0.0748
1 F1 Score	0.0718
2 AUC Score	0.5066
3 Recall	0.0748
4 Precision	0.0733
5 Kappa	0.0618
6 MCC	0.0419

Deploy ⚙️



**APPLICATION URL:** <https://autoinsightai.streamlit.app/>

To check the results and documentations related to the tested dataset on **Employee Attrition**  
[\(click here\)](#)

**GitHub Repository:** <https://github.com/Vidyadheesha-M-Pandurangi/AutoInsight.ai.git>

## CHAPTER 17

### REFERENCES

#### LINKS

<https://docs.streamlit.io/>  
<https://pandas.pydata.org/docs/>  
<https://pycaret.org/>  
<https://github.com/ydataai/ydata-profiling/>  
<https://scikit-learn.org/stable/documentation.html>  
<https://matplotlib.org/stable/users/index.html>  
<https://seaborn.pydata.org/>  
<https://pyfpdf.readthedocs.io/en/latest/>  
<https://docs.streamlit.io/streamlit-community-cloud/deploy>  
<https://scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics>

#### RESEARCH PAPERS

1. Moez Ali et al., “*PyCaret: An Open Source, Low-Code Machine Learning Library in Python*”
2. F. Pedregosa et al., “Scikit-learn: Machine Learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
3. J. D. Hunter, “Matplotlib: A 2D Graphics Environment,” *Computing in Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.
4. M. Waskom, “Seaborn: Statistical Data Visualization,” *Journal of Open Source Software*, vol. 6, no. 60, pp. 3021, 2021.
5. D. Powers, “Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness and Correlation,” *Journal of Machine Learning Technologies*, vol. 2, no. 1, pp. 37–63, 2011.
6. Maloney and R. Barga, “The Rise of No-Code Machine Learning,” *IEEE Internet Computing*, vol. 25, no. 1, pp. 4–9, Jan.–Feb. 2021.
7. "AutoML: Methods, Tools, and Applications" by Hutter et al.
8. "The State of AutoML" from Towards Data Science  
(Available Online: <https://towardsdatascience.com/>)