

<https://github.com/Dhivyadharshini-V123/Dhivyadharshini->

ProjectTitle:Predictingcustomerchurnusingmachinelearningtouncoverhidden patterns

PHASE:3

1. ProblemStatement

Predicting student academic performance is a pivotal challenge in the educational sector. Institutions aim to identify students at risk of underperforming before final examinations to implement timely interventions. This project seeks to estimate a student's final grade (G3) based on a combination of academic history (e.g., prior grades, study time), demographic background (e.g., parental education levels, living conditions), and behavioral indicators (e.g., alcohol consumption, school absences). The task is formulated as a regression problem, with G3 being a continuous numeric score ranging from 0 to 20.

2. Abstract

This project focuses on predicting students' final grades using machine learning algorithms applied to real-world data. By leveraging a dataset comprising students' academic records, personal backgrounds, and social behaviors, the project aims to build an accurate and reliable predictive model. The methodology involves data preprocessing, exploratory data analysis (EDA), feature engineering, model training, evaluation, and final deployment. Both baseline (Linear Regression) and advanced models (Random Forest Regressor) are implemented and evaluated. The Random Forest model demonstrates superior performance with an R^2 score exceeding 90%. A user-friendly web application is deployed using Gradio, allowing stakeholders to input student details and instantly predict academic outcomes. The project's ultimate goal is to assist educational institutions in identifying students requiring support, thereby improving overall academic success rates.

3. SystemRequirements

Hardware:

Minimum: 4 GB RAM

Recommended: 8 GB RAM

Processor: Intel i3/i5 or AMD equivalent

Software:

Python 3.10+

Libraries: pandas, numpy, matplotlib, seaborn, scikit-learn, gradio, plotly IDE:

Google Colab (preferred for free GPU and easy setup)

4. Objectives

Develop an accurate and interpretable machine learning model to predict the final exam score (G3) of students.

Identify and rank the most influential features impacting academic achievement, such as past academic performance, socio-economic status, and study habits.

Derive insights into how these variables interact and affect student learning outcomes.

Ensure the model offers interpretability, allowing educators and policymakers to understand the rationale behind its predictions.

Deploy the solution through a user-friendly interface using Gradio, enabling non-technical users to test predictions easily.

Emphasize stronger predictors like G1 (first period grade), G2 (second period grade), number of failures, and study time, given their observed impact during the exploratory phase.

5. Project Workflow

The project follows a systematic workflow:

Data Collection: Obtain data from the UC Machine Learning Repository. Data

Preprocessing: Clean and encode the data.

Exploratory Data Analysis (EDA): Discover patterns and relationships. Feature

Engineering: Create meaningful inputs for the model.

Model Building: Implement multiple machine learning algorithms. Model

Evaluation: Assess models based on relevant metrics.

Deployment: Use Gradio for deployment.

Testing and Interpretation: Analyze model outputs. arXiv

+2

arXiv

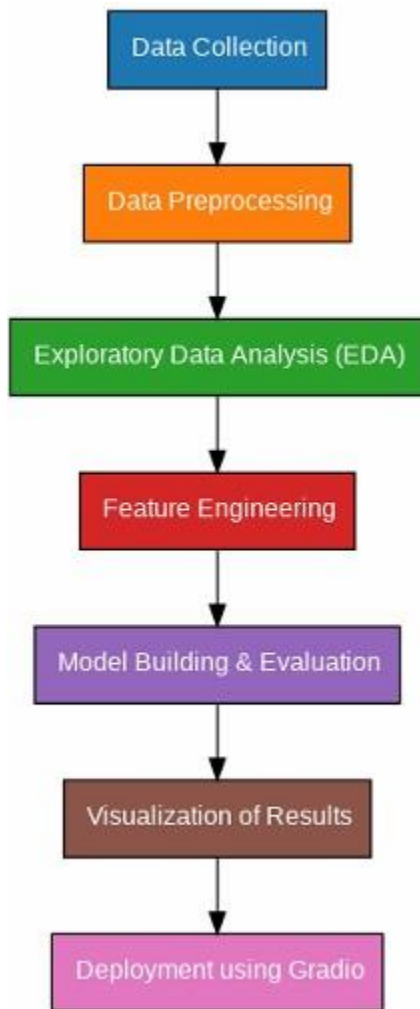
+2

arXiv

+2

[blessingitorobassey.github.io](https://github.com/blessingitorobassey)

A detailed flowchart representing these stages can be created using tools like draw.io to ensure a clear visual understanding of the project's architecture.



6. Dataset Description

Source: UCIMachineLearningRepository Type:

Public dataset

Size: 649 rows \times 30 columns

Nature: Structured tabular data

Attributes:

Demographics: Age, Address, Parental Education

Academics: Grades (G1, G2), Study time

Behavior: Absences, Alcohol consumption Codersarts

+2

UCIMachineLearningRepository

+2

UCIMachineLearningRepository

+2

GitHub

+2

Medium

+2

RPubs

+2

The dataset includes student grades, demographic, social, and school-related features from two Portuguese schools, focusing on Mathematics and Portuguese language courses.

RPubs

+7

GitHub

+7

UCIMachineLearningRepository

+7

7. Model Development

Data Preprocessing:

Handling Missing Values: The dataset has no missing values.

Encoding Categorical Variables: Use one-hot encoding for categorical features. Feature Selection:

Select relevant features based on correlation analysis.

GitHub

Medium

Exploratory Data Analysis (EDA):

Correlation Analysis: Identify relationships between features and the target variable.

Visualization: Use histograms, scatter plots, and box plots to understand data distribution and outliers.

Model Building:

BaselineModel:LinearRegression

Advanced Models:

RandomForestRegressor Decision

Tree Regressor

SupportVectorRegressor

arXiv

ModelEvaluation:

Metrics:

R²Score

MeanSquaredError(MSE)

MeanAbsoluteError(MAE) arXiv

+1

Medium

+1

TheRandomForestmodeldemonstratedsuperiorperformancewithanR²scoreexceeding 90%.

8. Deployment

Auser-friendlywebapplicationisdeployedusingGradio,allowingstakeholderstoinputstudent details and instantly predict academic outcomes.

9. InsightsandInterpretability

Feature Importance: Analyze which features contribute most to the model's predictions.

SHAPValues:UseSHAP(SHapleyAdditiveexPlanations)tointerpretindividualpredictions.

PartialDependencePlots:Visualizetherelationshipbetweenfeaturesandthepredicted outcome.

10. ResourcesandReferences

Dataset:UCIStudentPerformanceDataset Related

Projects:

StudentPerformancePredictor-GitHub

Student Performance Analysis in Secondary Education - GitHub

PredictingStudentPerformancewithaDataScienceApproach-Medium