

# UNDER GRADUATE STUDENT DROPOUT PREDICTION USING MACHINE LEARNING

<sup>1</sup> D Mahitha, <sup>2</sup>A. Ramya Sri, <sup>3</sup>A. Rishitha, <sup>4</sup>A. Jeevani, <sup>5</sup>B. Sanya

<sup>1</sup>Assistant Professor, <sup>2,3,4,5</sup>UG Students, Dept. of CSE (AI & ML), Malla Reddy Engineering College for Women (Autonomous), Hyderabad, India. E-Mail: mahithadilli@gmail.com

## ABSTRACT

Higher education plays a vital role in national development, but promoting educational quality and addressing the issue of university dropouts pose significant challenges for governments. University dropouts not only have negative economic consequences but also impact students personally. This project focuses on finding the factors that are affecting the student dropouts from semester to semester and different classification models can be built to predict whether a student is on the verge on dropping out of college and helps in giving early warnings to the university management to help such students. By leveraging principles of data mining and machine learning, we collected student data from the admin internal sources by making sure the data has no missing or null values. The dataset comprises information from 4424 undergraduate students. Utilizing Support Vector Classifier (SVC) and Logistic Regression, we built a basic model that is capable of predicting the on the verge drop out students and the model also is capable of understanding the most affecting features to the dropout rate. Our findings highlight the top five significant features for predicting educational status: Curricular units 2nd semester (approved), Curricular units 1st semester (approved), Curricular units 2nd semester (grade), Curricular units 1st semester (grade), and Tuition fee up to date. These top five features showed us that the most affecting features are the grades in the previous semesters. In order to increase the accuracy of the predictions, we took two datasets- DBS.csv and DBS\_2000.csv. Through careful analysis, we identified key features that contribute to improving the accuracy and applicability of our predictive models.

## INTRODUCTION

The aim of this project is to address the challenges of promoting educational quality and reducing university dropouts in higher education, which are crucial for national development. The negative economic consequences and personal impact on students resulting from university dropouts necessitate a comprehensive understanding of the factors influencing student's educational status. To achieve this, the researchers utilized principles of data mining and machine learning to collect and analyse data from internal sources.

The dataset used in this project comprises information from 4424 undergraduate students, providing a robust foundation for analysis and prediction. By leveraging Support Vector Classifiers (SVC), we developed classification models to predict whether a student will be a dropout or fall into non-dropout categories.

To ensure the accuracy and applicability of the predictive models, a meticulous analysis of the dataset was conducted. We identified key features that significantly contribute to predicting educational status. These features include Curricular units 2nd sem (approved) , Curricular units 1st sem (approved), Curricular units 2nd sem (grade), Curricular units 1st sem (grade), Tuition fees up to date. By incorporating these features, the classifiers were able to make reliable predictions regarding student's educational outcomes.

The findings of this project shed light on the complex factors that influence university dropouts and provide valuable insights for implementing targeted interventions. With a deeper understanding of these factors, educational institutions and policymakers can develop strategies to enhance educational quality and support student success. By leveraging the power of data mining and machine learning, this research contributes to the

ongoing efforts to improve higher education systems, foster national development, and empower students to achieve their educational goals.

## LITERATURE SURVEY

An "ever-enrolled person" who does not finish the last level of education for which he or she enrolled and is not currently enrolled in any educational institution is considered a dropout. The implementation of a socially equal environment presents a significant problem for the Indian community in terms of student dropout. According to a recent poll by the National Statistical Office (NSO), 12.6% of Indian pupils drop out of school, 19.8% of them do so in secondary school, and 17.5% do so in upper primary. In addition to the student's bad self-perception brought on by the sense of failure and despair, university dropouts have an economic and societal consequence.

Even while the government's Right to Education Act and National Policy on Education may have encouraged everyone to receive an education, it is equally necessary to assess the system's viability and effectiveness. Dropout rates are viewed as a huge waste of the educational system's limited resources since many students quit school before gaining even the most basic abilities and because their early departure implies a significant loss of those resources.

Schools can identify students who are at risk of dropping out and concentrate on those who have difficulty performing well thanks to dropout early warning systems. Children and adolescents who drop out of college typically do so for a variety of reasons. Instead, this is a process that is shaped by numerous variables that interact in complicated, dynamic ways. These variables may relate not just to traits or situations specific to a person or family, but also to variables at the college, local, and international levels.

These can include shortcomings in educational environments and procedures, in social welfare and education systems, in general social policies for young people and employment, and in societal norms, especially gender standards that may be detrimental to education. The reason why kids and teenagers leave school is frequently due to "individual and family circumstances that structures and systems are unable to respond to or address appropriately."

## EXISTING SYSTEM

Early Warning Systems (EWS) for school children who are at risk of quitting. UNICEF supports nations and education experts in enhancing their approaches to stop school dropout. It promotes EWS as a means of fostering fair access to high-quality education and lowering dropout rates. The main objective is to spot pupils who are at risk of leaving school before they stop attending altogether.

Choosing indicators to pinpoint students at risk of dropping out is the first step in the current system for forecasting student dropouts, which normally entails a step-by-step procedure. The current system handles this step as follows:

**Formation of a school staff group:** The existing system establishes a group consisting of school management, teachers, and support staff. This group collaboratively works to establish indicators for the Early Warning System (EWS).

**Identification and prioritization of risk factors:** Based on contextual evidence, the school staff group determines and ranks the key risk factors and predictors of dropout. To comprehend dropout trends and causes, they study data on dropout from prior years that is currently available.

**Formulation of indicators:** Indicators are formulated based on the identified predictors of dropout. Each predictor is translated into an indicator that can be measured or observed. For example, indicators may include absenteeism, academic performance, behaviour issues, socio-economic status, special alerts for specific circumstances, and a sense of belonging to the school.

**Data collection methodology:** The existing system determines how data will be collected for each indicator. It

considers existing data collection systems, such as electronic or paper-based systems, to ensure the accessibility and reliability of the necessary data. Special attention is given to reliable absenteeism data, including justified and non-justified absences, tardiness, and reasons for these.

**Weighting and thresholds:** The school staff group may give various indications varying weights in order to increase the EWS's sensitivity to dropout concerns. They concur on cutoff points that specify when pupils will be classified as "at risk" or "at high risk" of failing to graduate. This aids in classifying students according to their degree of danger.

Based on data obtained on early leaver profiles and the most significant risk factors, the current system continuously adjusts the indicator sets, weighting, and scoring systems. Through this iterative approach, the system's capacity to forecast student dropouts can be improved continuously.

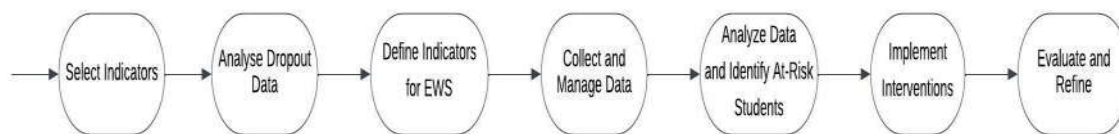


Fig.1. Existing System flowchart

## PROPOSED SYSTEM

There is a need to identify multiple dimensions that influence dropout, including individual, family, community, and school factors. By considering a wide range of indicators, EWS can effectively capture the complexity of dropout risks. The existing system being a manual procedure takes a lot of time and effort into finding out the drop outs leaving none for the rectification of problems causing the dropouts.

Machine learning (ML) approaches can significantly aid Early Warning Systems (EWS) by automating and enhancing the process of predicting student dropouts. This way, the efforts of the school/university committees can be used where it is actually needed. Here's how ML can assist EWS and reduce the need for manual work:

1. **Feature selection:** ML algorithms can automatically identify the most relevant features or indicators for predicting student dropout. Instead of relying on manual selection and prioritization of indicators, ML algorithms can analyse large amounts of data to identify patterns and relationships that are predictive of dropout.
2. **Data preprocessing:** ML techniques can handle data preprocessing tasks such as handling missing values, normalizing data, and encoding categorical variables. This automation streamlines the data preparation phase, saving time and effort compared to manual data cleaning and transformation.
3. **Predictive modelling:** ML algorithms, such as logistic regression, decision trees, random forests, or neural networks, can be trained on historical data to create predictive models. These models can learn complex relationships between predictors and dropout outcomes, providing accurate predictions for individual students.
4. **Continuous learning:** ML models can continuously learn and adapt based on new data. As more data becomes available, the models can be updated to incorporate the latest information, ensuring that the predictions remain relevant and accurate over time. This dynamic learning capability reduces the need for manual model adjustments.
5. **Identification of complex patterns:** ML algorithms can uncover hidden patterns and interactions among predictors that may not be evident through manual analysis. This ability to detect complex

relationships can lead to more accurate predictions and identification of at-risk students who may not fit conventional risk profiles.

6. **Scalability:** ML approaches allow for the efficient analysis of large-scale datasets. They can process vast amounts of student data quickly, enabling EWS to handle data from multiple schools or districts. This scalability is especially beneficial when dealing with extensive student populations.
7. **Performance evaluation:** ML techniques provide automated evaluation metrics to assess the performance of the predictive models. These metrics, such as accuracy, precision, recall, and F1 score, help quantify the effectiveness of the ML-based EWS and guide further improvements.

It's important to note that while ML approaches offer significant advantages, they still require careful consideration and domain expertise. Interpretability of ML models and ethical considerations should be taken into account to ensure transparency, fairness, and accountability in the decision-making process of EWS.

By leveraging ML techniques, EWS can streamline the process of predicting student dropouts, reduce manual efforts in feature selection and indicator formulation, and improve the accuracy and efficiency of dropout risk identification. Here are the steps:

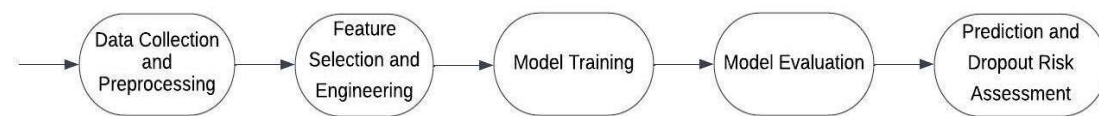


Fig.2. Proposed System Flowchart

#### Step 1: Data Collection and Preprocessing

- Gather student data (demographics, academic records, attendance, behavior, etc.)
- Clean and preprocess the data (handle missing values, normalize features, etc.)
- Split the data into training and testing sets
- Label encoding and ordinal encoding of all the values
- Data scaling into the interval -1

#### Step 2: Feature Selection and Engineering

- Identify relevant features for dropout prediction
- Perform feature selection techniques (e.g., correlation analysis, feature importance)

#### Step 3: Model Training

- Apply SVM algorithm for training the model
- Choose appropriate SVM kernel (linear, polynomial, radial basis function, etc.)
- Optimize hyper parameters through techniques like cross-validation

#### Step 4: Model Evaluation

- Evaluate the trained model using the testing dataset
- Evaluating the model using a confusion matrix
- Calculate performance metrics (accuracy, precision, recall, F1-score, etc.)
- Analyse the model's performance and identify any shortcomings

#### Step 5: Prediction and Dropout Risk Assessment

- Use the trained SVM model to predict dropout probabilities for new student data
- Classify students into risk levels (low, medium, high) based on predicted probabilities
- Identify students at high risk of dropout for intervention

#### Next steps that can be followed

##### Step 1: Intervention and Support

- Implement targeted interventions and support strategies for at-risk students

- Provide personalized counselling, academic assistance, or mentoring
- Collaborate with relevant stakeholders to address identified risk factors

**Step 2: Monitoring and Refinement**

- Monitor the effectiveness of interventions in reducing dropout rates
- Collect feedback from students, teachers, and other stakeholders
- Continuously refine the ML model and update the training data

**SYSTEM ARCHITECTURE**

The process begins at the time of admissions. When a student joins into the student, he/she submits details like the salary of parents, education history, etc. For this model, along with this information, it needs more info like the marital status of the student, the relationship status of the parents, etc. The admin office can give a form to be filled to every student explaining the need.

Thus, all the details of the student are right in the hands of the student which are needed to execute this model. The university can now appoint an EWS analyzer who needs to make a separate dataset with all the student details required for the model. The next task for the analyzer is to keep updating the dataset with course grade with every semester.

Now, the analyzer has to keep training and executing the model every semester in order to predict the students with a high chance of dropping out in the next semester. the analyzer and the college management can work out the most effective features which are determining the dropout rate with statistical analysis techniques like correlation matrix, hypothesis testing to reduce the rate of drop outs.

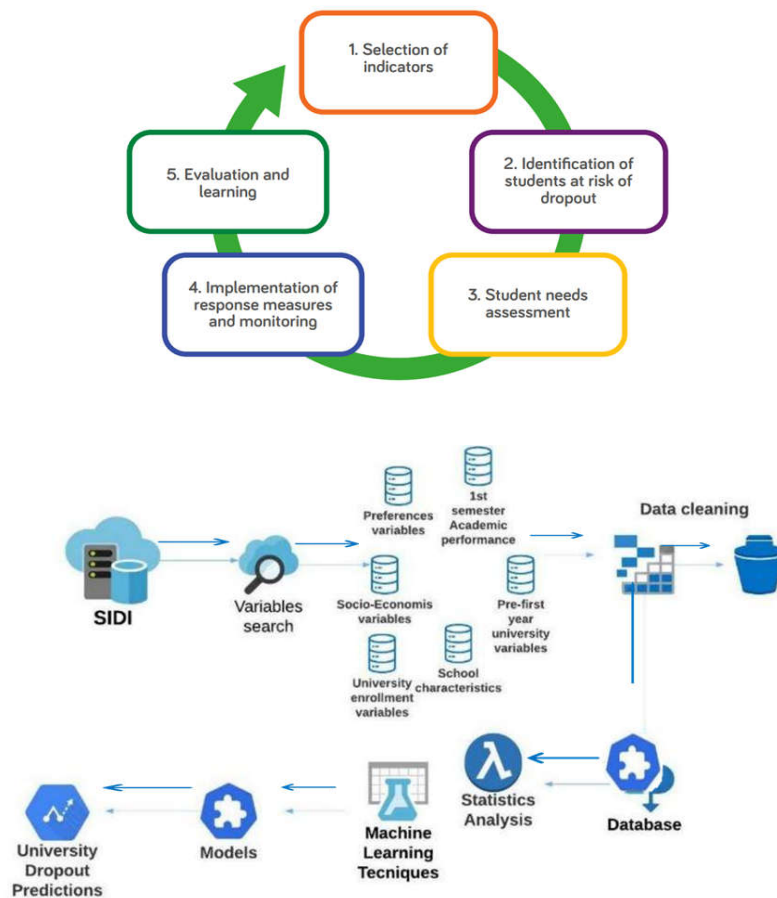


Fig.3. System Architecture

**DEFINE MODULES & FUNCTIONALITIES:****From the library Scikit-learn (sklearn):**

- a. **MinMaxScaler:** A module for data normalization or standardization.
- b. **train\_test\_split:** A function for splitting the data into training and testing sets.
- c. **LabelEncoder:** A class for encoding categorical variables as integers.
- d. **OrdinalEncoder:** A class for encoding categorical variables as integers with an ordinal relationship.
- e. **svm:** A module for Support Vector Machine algorithms.
- f. **datasets:** A module providing various datasets for testing and practicing machine learning algorithms.
- g. **GridSearchCV:** A class for performing hyperparameter tuning using grid search and cross-validation.
- h. **csv:** The csv module provides functionality for reading and writing CSV (Comma-Separated Values) files. It is a built-in module in Python.
- i. **Axes3D:** Axes3D is a class in the `mpl_toolkits.mplot3d` module of Matplotlib that allows plotting in 3D.
- j. **Pyplot:** Pyplot is a module within the Matplotlib library. Matplotlib is a plotting library in Python that provides a wide range of functions for creating static, animated, and interactive visualizations.

**IMPLEMENTATION****LIBRARIES:**

1. **SEABORN:**  
Seaborn is a data visualization library built on top of Matplotlib. It provides a high-level interface for creating informative and attractive statistical graphics. We used this library to construct a heatmap
2. **PANDAS:**  
Pandas is a powerful library for data manipulation and analysis. It provides data structures like DataFrames that allows us to efficiently work with structured data.
3. **MATPLOTLIB:**  
Matplotlib is a plotting library that provides a flexible and comprehensive set of plotting functions. It is widely used for creating static, animated, and interactive visualizations in Python.
4. **SCIKIT-LEARN (SKLEARN):**  
Scikit-learn is a popular machine learning library in Python. It provides a wide range of algorithms for classification, regression, clustering, dimensionality reduction, and more. Our project implementation is mostly dependent on the modules of this library.
5. **NUMPY:**  
NumPy is a fundamental library for scientific computing in Python. It provides support for large, multi-dimensional arrays and matrices, along with a collection of mathematical functions.

**ALGORITHMS****Multi-dimensional multi-class Support Vector Classifier:**

The Multi-dimensional multi-class Support Vector Classifier is a supervised machine learning algorithm used for classification problems. It is an extension of the Support Vector Machine (SVM) algorithm. SVM aims to find an optimal boundary between different output classes by performing complex data transformations using kernel functions. These transformations maximize the separation between data points based on the defined classes or labels.

In its basic form, SVM seeks to find a line that maximizes the separation between two classes in a two-dimensional space. However, in the case of multi-dimensional data, the objective is to find a hyperplane that maximizes the separation of data points into their respective classes in an n-dimensional space. The data points that are closest to the hyperplane are known as Support Vectors. In the case of multi-class classification with

three target classes (dropout, graduate, and enrolled), binary classification cannot be directly applied. However, the Support Vector Classifier (SVC) algorithm provides two approaches for multi-class classification:

### One-to-One approach

This approach breaks down the multi-class problem into multiple binary classification problems. For each pair of classes, a binary classifier is trained to separate the instances of one class from the instances of the other class. In this approach, multiple hyperplanes are created to separate each pair of classes, ignoring the points belonging to the third class.

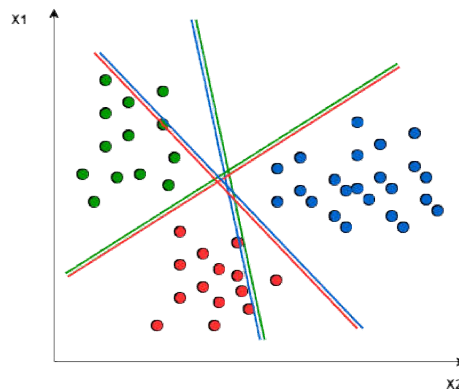


Fig.4. One-to-One approach:

### One-to-Rest approach

In this approach, a hyperplane is created to separate one class from the rest of the classes all at once. The separation takes into account all the points and divides them into two groups: one group for the points of the class under consideration and another group for all the remaining points representing the other classes.

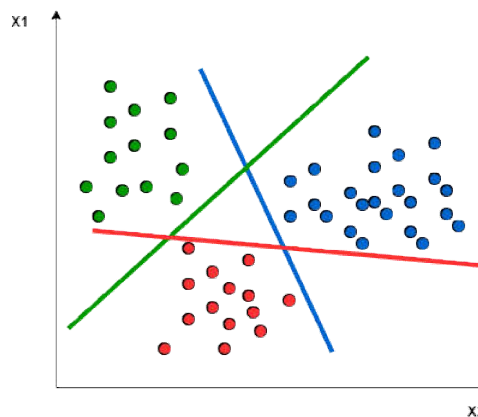


Fig.5. One-to-Rest approach:

To implement these approaches, you can use the SVC function provided by the sklearn library. This function allows you to specify the desired kernel functions for the SVM algorithm, such as linear, polynomial, or radial basis function (RBF), among others.

### Logistic Regression

Based on a given dataset of independent variables, logistic regression calculates the likelihood that an event will occur, such as voting or not voting. A supervised ML model is what it is. Given that the result is a probability, the dependent variable's range is 0 to 1. In logistic regression, the odds—that is, the probability of success divided by the probability of failure—are transformed using the logit formula.

In contrast to linear regression, logistic regression does not simply fit a straight line to the data. Instead, we used the Sigmoid curve to fit an S-shaped curve to our observations. The following formulas are used to represent this logistic function, which is also referred to as the log odds or the natural logarithm of odds:

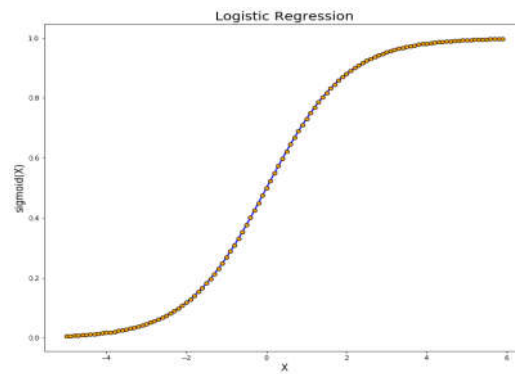


Fig.6. Logistic regression

$$\ln(\pi/(1-\pi)) = \text{Beta}_0 + \text{Beta}_1 * X_1 + \dots + B_k * K_k$$

Logit( $\pi$ ) is the dependent or response variable in this logistic regression equation while  $x$  is the independent variable. Most frequently, maximum likelihood estimation (MLE) is used to estimate the beta parameter, or coefficient, in this model. In order to find the best fit for the log odds, this approach iteratively evaluates various beta values. The log likelihood function is created after each of these iterations, and logistic regression aims to maximise this function to get the most accurate parameter estimate.

The conditional probabilities for each observation can be calculated, logged, and added together to produce a forecast probability once the best coefficient (or coefficients, if there are multiple independent variables) has been identified. If the categorization is binary, a probability of less than .5 predicts 0 and a probability of more than 0 predicts 1. It is recommended to assess the model's goodness of fit, or how well it predicts the dependent variable, once the model has been computed.

**RESULT**

```
[ ] Mounting the drive.
from google.colab import drive
drive.mount('/content/drive/')
Mounted at /content/drive/

+ Loading the dataset.

[ ] school_data = pd.read_csv('/content/drive/myColab/Notebooks/dataset1.csv', dtype = str)
school_data.head()

Marital  Application  Application  Course  Daytime/evening  Previous  Previous  Nationality  Mother's  Father's  Curricular  Curricular  Curricular  Curricular  Curricular  Curricular  Unemployment  Inflation  GDP  Ta
status    Application  order      attendance  qualification  qualification  (grade)  qualification  qualification  ...  units 2nd  units 2nd  units 2nd  units 2nd  units 2nd  units 2nd  rate      rate      (GDP  Ta
[0] 1 17 5 171 1 1 122 1 19 12 0 0 0 0 0 0 10.8 1.4 1.74 Di
[1] 1 15 1 1954 1 1 188 1 1 3 0 6 6 6 13.66555556 0 13.9 -3 0.79 Gi
[2] 1 1 5 1678 1 1 122 1 37 37 0 0 0 0 0 0 10.8 1.4 1.74 Di
[3] 1 17 2 1973 1 1 122 1 38 37 0 6 10 5 12.4 0 9.4 -8 -3.12 Gi
[4] 2 35 1 8814 0 1 188 1 37 38 0 6 6 6 6 13 0 13.9 -3 0.79 Gi
5 rows * 37 columns

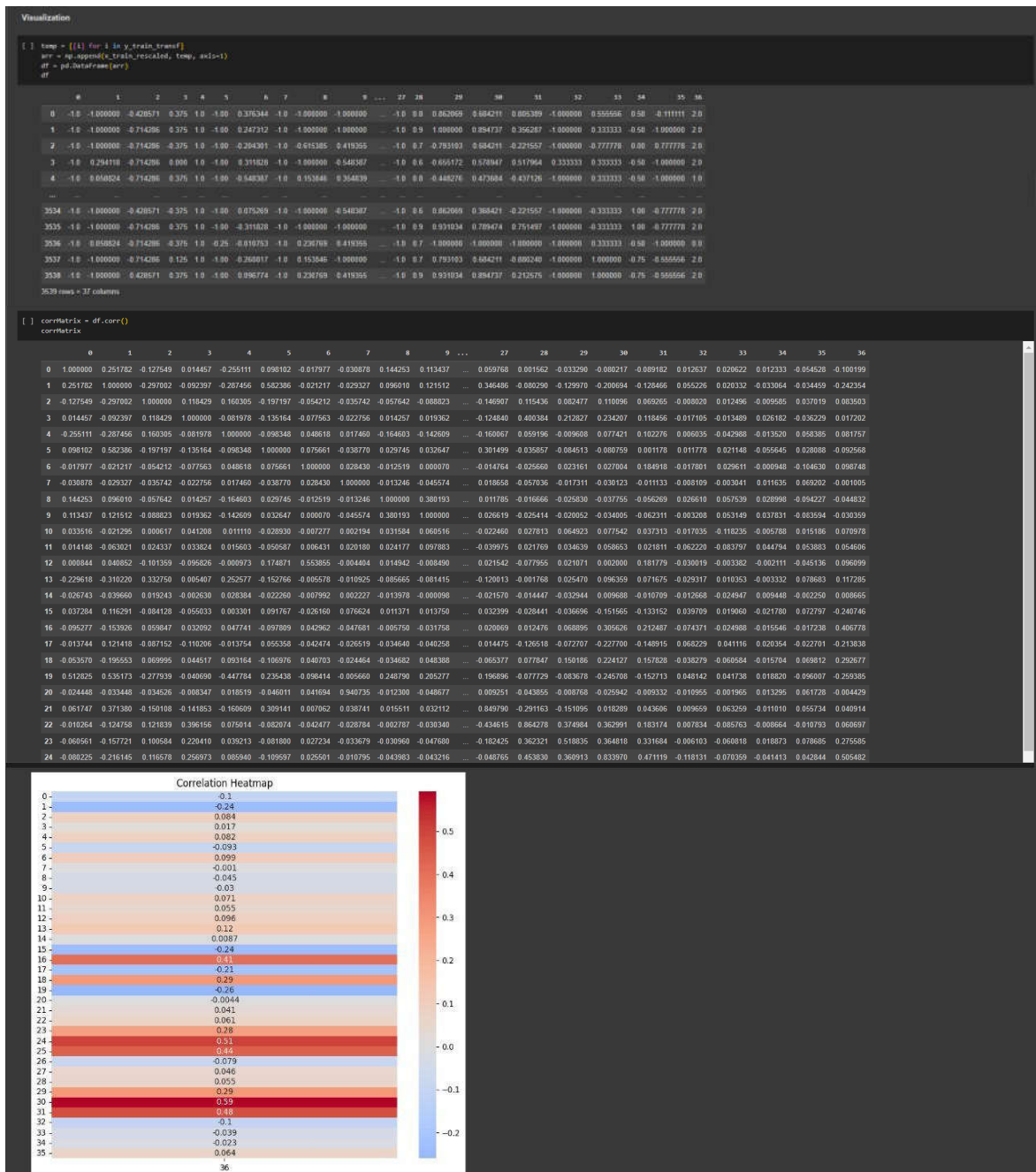
[ ] #Size of the training dataset.
y_train.value_counts()
x_train

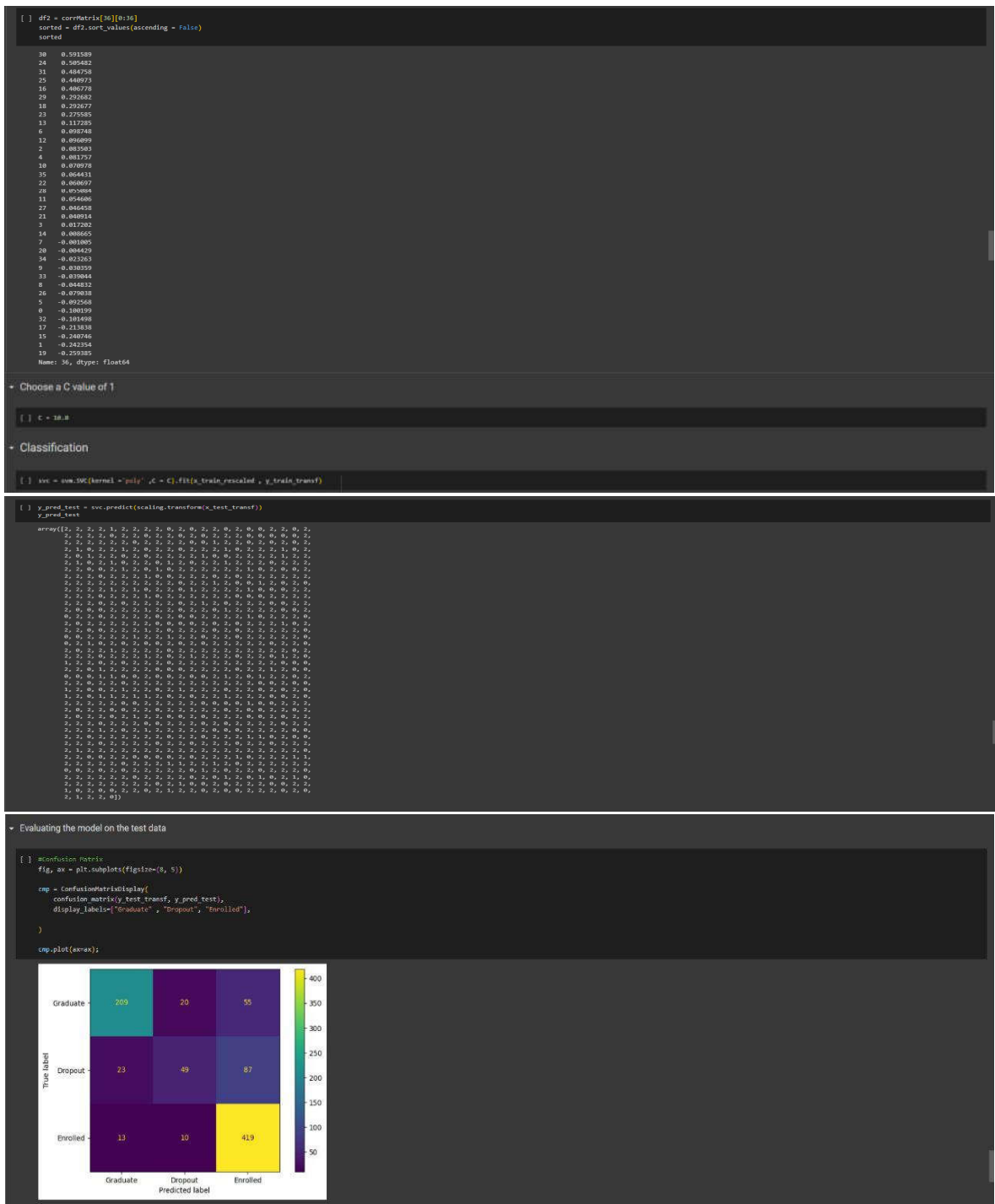
Marital  Application  Application  Course  Daytime/evening  Previous  Previous  Nationality  Mother's  Father's  Curricular  Curricular  Curricular  Curricular  Curricular  Curricular  Unemployment  Inf
status    Application  order      attendance  qualification  qualification  (grade)  qualification  qualification  ...  units 1st  units 2nd  units 2nd  units 2nd  units 2nd  units 2nd  rate  rate
3987 1 1 2 9500 1 1 155 1 1 0 0 7 7 6 15.25 0 7.6
322 1 1 1 9500 1 1 150 1 1 0 0 8 9 8 14.05555556 0 16.2
3762 1 1 1 9085 1 1 133 1 1 19 38 0 0 6 11 6 13 0 11.1
3949 1 43 1 9147 1 1 153 1 1 19 0 5 13 5 14.4 5 16.2
2445 1 39 1 9500 1 1 120 1 37 37 0 0 7 16 4 12.6 0 16.2
...
3656 1 1 2 9085 1 1 144 1 1 19 0 5 7 3 13 0 12.7
1371 1 1 1 9500 1 1 129 1 1 0 0 8 8 7 15.02857143 0 12.7
2420 1 39 1 9085 1 2 140 1 38 38 0 0 6 0 0 0 0 16.2
4332 1 1 1 9238 1 1 131 1 37 1 0 0 6 6 6 11.33333333 0 9.4
1322 1 1 5 9500 1 1 145 1 38 38 0 0 8 8 8 13.76375 0 9.4
3539 rows * 36 columns

[ ] #Size of the test dataset.
y_test.value_counts()

Target
Graduate 442
Dropout 284
Enrolled 159
Name: count, dtype: int64
```







**Using Logistic Regression:**

```
[3] from google.colab import drive
drive.mount('/content/drive')

Mounted at /content/drive

[5] # importing and cleaning the data
train_data = pd.read_csv('/content/drive/MyDrive/colab Notebooks/DBS.csv', sep=';')
test_data = pd.read_csv('/content/drive/MyDrive/colab Notebooks/DBS_2020.csv', sep=';')
train_data.head()
```

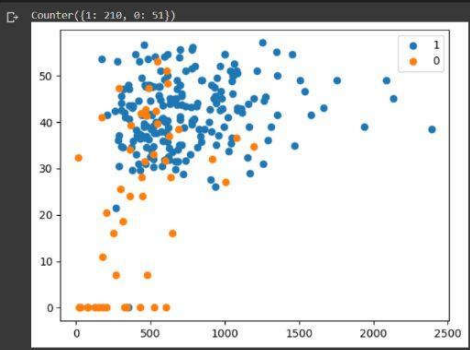
	access	tests	tests_grade	exam	project	project_grade	assignments	result_points	result_grade	graduate	year	acad_year
0	1256	57.00	A	19	91.54	A	40.0	189.82	A	1	2019	2019/2020
1	985	42.87	B	19	75.96	A	13.7	189.43	A	1	2017	2017/2018
2	1455	54.50	A	16	86.79	A	40.0	188.91	A	1	2019	2019/2020
3	998	54.50	A	16	93.36	A	40.0	186.85	A	1	2019	2019/2020
4	1347	55.00	A	16	92.86	A	39.0	186.38	A	1	2019	2019/2020

```
# Data normalization with sklearn
from sklearn import preprocessing
standardized_X = preprocessing.scale(X_test)
df = pd.DataFrame(standardized_X)
df.describe()
```

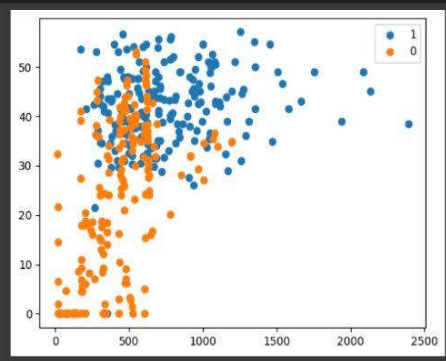
	0	1	2
count	6.000000e+01	6.000000e+01	6.000000e+01
mean	-1.850372e-17	4.218847e-16	1.591320e-16
std	1.008439e+00	1.008439e+00	1.008439e+00
min	-2.201088e+00	-2.870371e+00	-2.841506e+00
25%	-5.506751e-01	-6.288193e-01	-8.207221e-01
50%	-1.676396e-01	3.510590e-01	2.428466e-01
75%	5.836475e-01	7.641450e-01	8.809911e-01
max	2.820037e+00	1.491048e+00	1.306419e+00

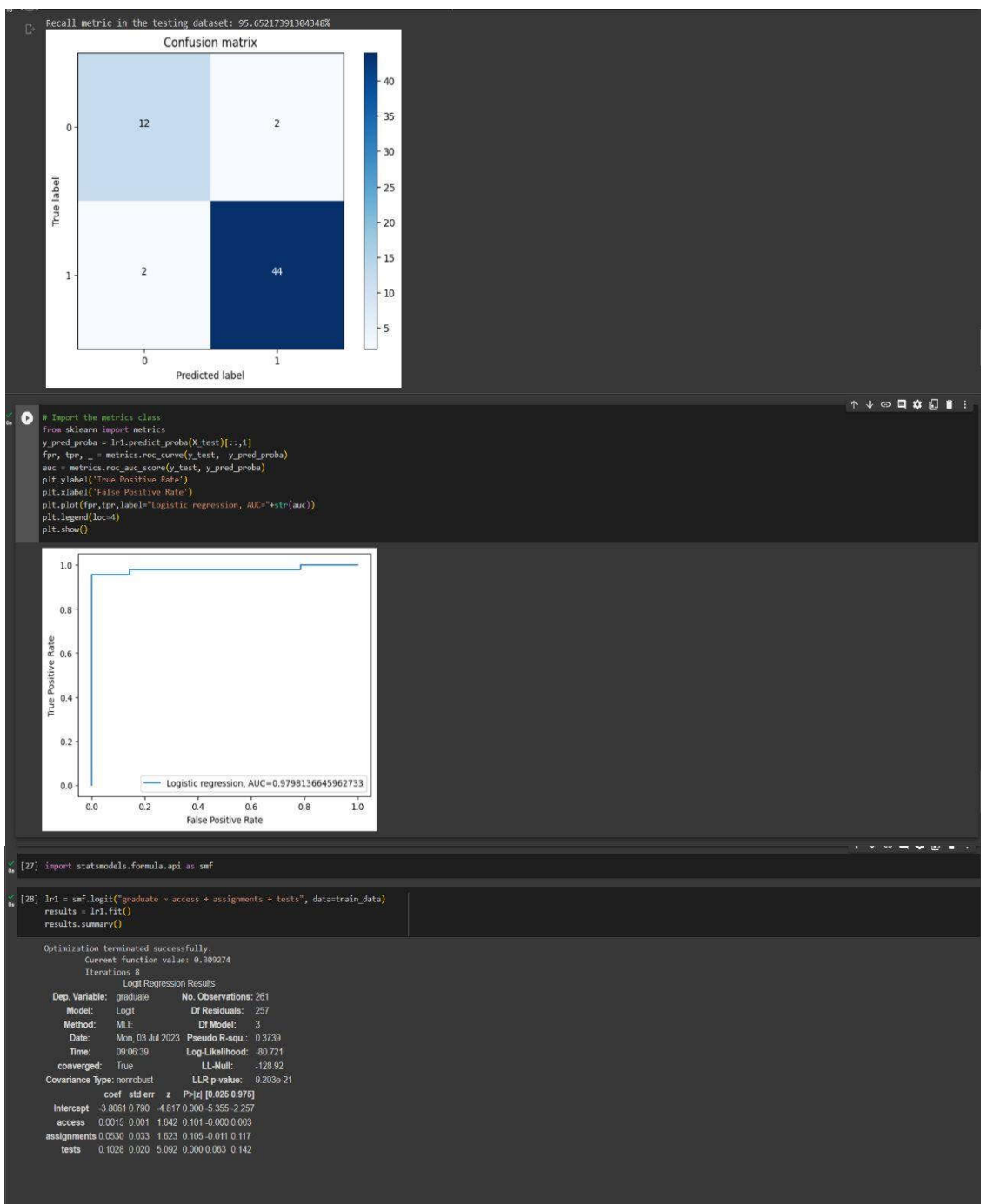
```
# Summarize class distribution
counter = Counter(y_train)
print(counter)

for label, _ in counter.items():
    row_ix = where(y_train == label)[0]
    plt.scatter(X_train[row_ix, 0], X_train[row_ix, 1], label=str(label))
plt.legend()
plt.show()
```



```
[17] for label, _ in counter.items():
    row_ix = where(y_train_res == label)[0]
    plt.scatter(X_train_res[row_ix, 0], X_train_res[row_ix, 1], label=str(label))
plt.legend()
plt.show()
```





### CONCLUSION

The whole idea of this project is to use Machine learning approaches in understanding the problems students face to keep their education at pace. If a university is able to help a student in even a single way, it will help make a big change in a student’s life. This model is a tool for the universities to make this change possible. The model makes the analysis part easy for EWS (Early Warning System) analyzers. Using this simple model we executed, the analyzer will be able to pose questions about different features and find the answers using the model.

Using this SVM model at the simplest level:

1. The Analyzer can get details of the students on the verge of dropping out.

2. Analysis of the top features affecting the dropout rate can be understood

## FUTURE SCOPE

### **Predictive Models for Different Student Groups:**

Customize predictive models for specific student populations, such as first-generation college students, international students, or students from underrepresented backgrounds. These tailored models can account for unique challenges and provide more accurate predictions for targeted groups.

### **Cross-Institution Collaboration:**

Foster collaboration among universities and educational institutions to share data, best practices, and research findings related to dropout prediction and prevention. Collaborative efforts can lead to a broader understanding of dropout factors and the development of standardized frameworks.

### **Comparative Analysis:**

Conduct comparative studies across different universities, courses, and student demographics to identify common trends and factors contributing to dropout rates. This analysis can help institutions develop targeted interventions and policies to reduce attrition.

### **Policy Formulation:**

Dropout prediction models can provide valuable insights into the factors contributing to student attrition. Government agencies responsible for education can use this information to formulate policies that target the identified risk factors. For example, if the data shows that financial constraints are a major reason for dropout, the government can introduce scholarships or financial aid programs to support students in need.

### **Addition of Psychological Features:**

Include psychological assessments or surveys to collect data on student's psychological characteristics such as personality traits, self-efficacy, motivation, resilience, and mental health indicators. These features can be administered periodically to gather information throughout student's academic journey.

## REFERENCES

- [1]. (baeldung, Multiclass Classification Using Support Vector Machines, 2022) <https://www.baeldung.com/cs/svm-multiclass-classification>
- [2]. (Valentim Realinho, Predict students' dropout and academic success, 2021) <https://archive-beta.ics.uci.edu/dataset/697/predict+students+dropout+and+academic+success>
- [3]. (Samuel-Njoroge, Machine Learning with Support Vector Classifier., 2022) [https://github.com/Samuel-Njoroge/School\\_Drop\\_out\\_Prediction\\_with\\_SVMs/blob/main/School\\_Dropout\\_.ipynb](https://github.com/Samuel-Njoroge/School_Drop_out_Prediction_with_SVMs/blob/main/School_Dropout_.ipynb)
- [4]. (Wagavkar, Introduction to the Correlation Matrix, 2023) <https://builtin.com/data-science/correlation-matrix>
- [5]. (Vinayak Hegde, Higher education student dropout prediction and analysis through educational data mining, 2018) <https://ieeexplore.ieee.org/document/8398887>
- [6]. (pypi,pandas.DataFrame,2023)<https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.html>