

2. Background

2.1 Logistic Regression

Logistic regression is a supervised machine learning model and a classification variant of the linear regression model in which the outcome variable is dichotomous [1]. This is done by finding the boundary line of the classification between two classes [2] and giving the probability of a data point belonging to one class. This is represented by the conditional probability $P(Y|X)$ where X is a vector of input features and the random variable Y takes the value of 1 or 0 [3]. The input variable X can be described as a set of features $(x_0, x_1, x_2, \dots, x_n)$ which are multiplied by a set of associated weights $(w_0, w_1, w_2, \dots, w_n)$, and summed together to produce the dot product [2]. Therefore:

$$z = w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$
$$z = w \cdot x = w^T x$$

This is then passed to a function which converts z to a value between 0 and 1. The Sigmoid function is the most widely used function for this [2], and the distribution is shown in Figure 1. Since it is asymptotic, it ensures that the output will always be between 0 and 1. The Sigmoid function is given by:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

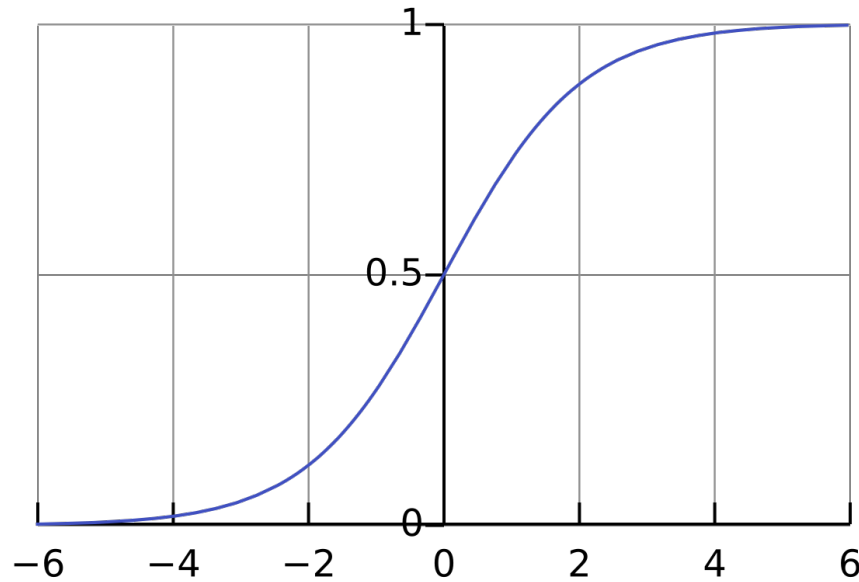


Figure 1: Sigmoid function.

After being passed through the Sigmoid function, the data point is then given the prediction of belonging to the class 0 or 1. Therefore, if $\sigma(w^T x) \geq 0.5$ then the outcome is class 1, and if $\sigma(w^T x) < 0.5$ then the outcome is class 0. The point 0.5 is where the point $z = 0$.

The standard cost function for logistic regression is the cross-entropy loss, derived from the Maximum Likelihood Estimation (MLE) principle. [3]. **Insert Standard Cost function?**

$$\begin{cases} -\log(\hat{y}), & \text{if } y = 1 \\ -\log(1 - \hat{y}), & \text{if } y = 0 \end{cases}$$

The parameter weights (w) are optimised during training the model, usually by using the gradient descent method [2] [3].

While standard logistic regression gives a binary outcome variable, it can also be used for multi-class problems. A way to do this is to decompose the problem into multiple logistic regression models, known as the One-vs-Rest (OvR) approach [4]. The OvR model trains separate binary classifiers for each class, with each classifier distinguishing one class from the rest. The class with the highest confidence score is chosen for the final prediction. **Another common approach is using Multinomial Logistic Regression, a unified model that computes the probabilities for multiple class simultaneously. How? Softmax regression. Find source.**

3. Data Preparation

The initial dataset consisted of 4424 records with 36 independent variables, with no missing values. The dataset was processed using the *pandas* library in Python.

- **Note: Should I include means of variance tables despite not all variables being in the final MI dataset?**

Each student entry has one of 3 Target variables: Graduated, Dropout, or Enrolled (where the student took another three years to complete the course). The distribution of the Target variable is shown in Figure 2, which shows that it is imbalanced, where the Graduated Target makes up 50% of the data records. This was addressed using the Synthetic Minority Over-sampling Technique (SMOTE), an oversampling algorithm which generates synthetic data using k nearest neighbours [5]. The SMOTE implementation from *imbalanced-learn* in the *scikit-learn* library was used.

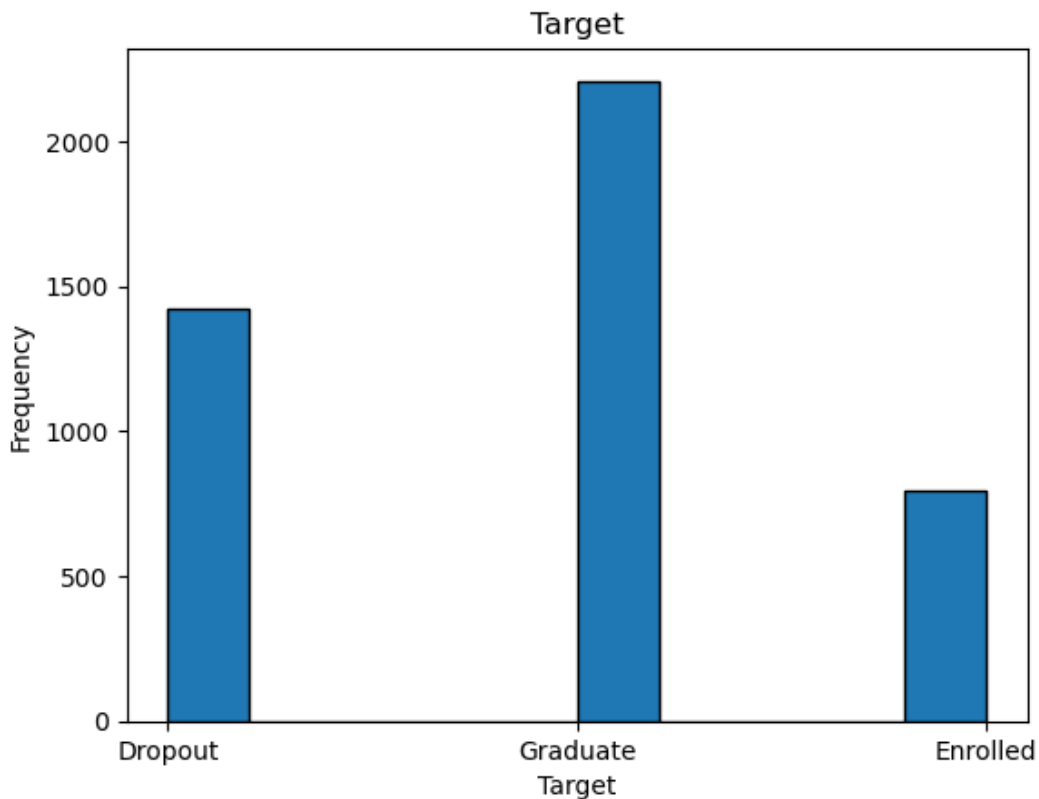


Figure 2: Count of Target variable.

Mutual Information (MI) is a filter-based feature selection method which measures how mutually dependent two variables are [6]. This method was used to filter out redundant variables, and was opted for since MI can be used for both discrete and continuous variables. The MI score with target was computed for each variable using *scikit-learn*. Each feature is assigned a scoring value, with the resulting features being organised in descending order based on the scores and are assigned rankings for the features [6]. A limitation of this method is that while it measures a feature's importance by its correlation with the target, it assumes independence of features from each other [7]. Therefore, choosing the highest scoring variables from MI can lead to redundant features can be selected.

To address this, a pairwise MI matrix was computed. The values were normalised to a range between 0 and 1, presented in Figure 3. Redundant variables were removed by checking each MI pairwise score according to a set threshold (0.99) and removing the variable with the lower MI score with the target. The final selected predictors resulted in 17 independent variables.

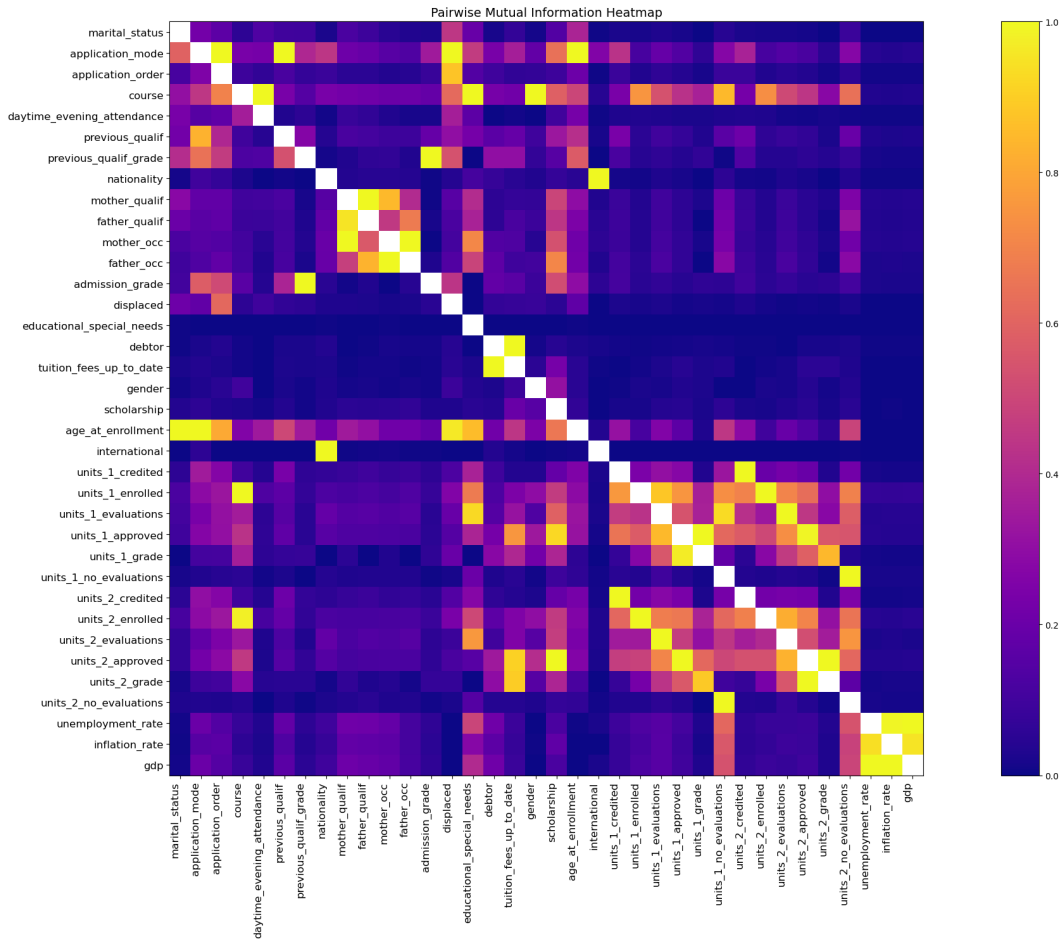


Figure 3: Heatmap of pairwise Mutual Information scores.

Principal Component Analysis (PCA) was then performed on the dataset for dimensionality reduction [8]. **Explain PCA process briefly.** The final resulting dataset consisted of 8 independent variables.

4. Experiments

4.7 Logistic Regression

The multi-class logistic regression model was implemented using the *scikit-learn* library in Python. The dataset was split into 80% training and 20% testing sets. A validation set was not used since logistic regression does not use this for refining parameters.

Model evaluation was performed on the test set using the precision, recall, and F1-score metrics to assess classification performance. The *scikit-learn* `classification_report` function was used for this.

The *scikit* `LogisticRegression` class was utilised. The model was run twice. The first run used the multinomial approach, which is now default for the `LogisticRegression` class. A second model was run using the One-vs-Rest approach using the `OneVsRestClassifier` class. The results are presented in Tables 1 and 2 respectively.

Results show that, for the Multinomial approach, the overall model F1-Score accuracy is 68%. The results show that all three classes produce an F1-Score of 72 to 73%. In comparison, the OvR approach had an overall F1-score accuracy is 67%. However, the model performed significantly worse on the Dropout class, with a 56% F1-score, in comparison to the Multinomial version. It should be noted that the precision and recall for both versions had very similar results. For comparison with the other models, the Multinomial version was chosen.

Target	Precision	Recall	F1-Score
Graduate	0.75	0.72	0.73
Dropout	0.60	0.54	0.72
Enrolled	0.68	0.77	0.72
Accuracy			0.68
Weighted Avg	0.67	0.68	0.67

Table 1: Classification report for Logistic Regression model using Multinomial Logistic Regression.

Target	Precision	Recall	F1-Score
Graduate	0.75	0.71	0.73
Dropout	0.60	0.53	0.56
Enrolled	0.67	0.79	0.72
Accuracy			0.67
Weighted Avg	0.67	0.67	0.67

Table 2: Classification report for Logistic Regression model using One-vs-Rest.

References

- [1] D. W. Hosmer Jr., S. Lemeshow, and R. X. Sturdivant, *Applied Logistic Regression* (Wiley Series in Probability and Statistics). Wiley, 2013, ISBN: 9780470582473.

- [2] X. Zou, Y. Hu, Z. Tian, and K. Shen, "Logistic regression model optimization and case analysis," in *2019 IEEE 7th International Conference on Computer Science and Network Technology (ICCSNT)*, 2019, pp. 135–139. DOI: 10.1109/ICCSNT47585.2019.8962457.
- [3] H. Li, "Logistic regression and maximum entropy model," in *Machine Learning Methods*. Singapore: Springer Nature Singapore, 2024, pp. 103–125, ISBN: 978-981-99-3917-6. DOI: 10.1007/978-981-99-3917-6_6. [Online]. Available: https://doi.org/10.1007/978-981-99-3917-6_6.
- [4] M. Aly, "Survey on multiclass classification methods," California Institute of Technology, Technical Report, 2005, Technical Report. [Online]. Available: <https://authors.caltech.edu/xyz>.
- [5] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: Synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, Jun. 2002, ISSN: 1076-9757. DOI: 10.1613/jair.953. [Online]. Available: <http://dx.doi.org/10.1613/jair.953>.
- [6] B. Venkatesh and J. Anuradha, "A review of feature selection and its methods," *Cybernetics and Information Technologies*, vol. 19, no. 1, pp. 3–26, 2019. DOI: 10.2478/cait-2019-0001. [Online]. Available: <https://doi.org/10.2478/cait-2019-0001>.
- [7] J. Li, K. Cheng, S. Wang, *et al.*, "Feature selection: A data perspective," vol. 50, no. 6, Dec. 2017, ISSN: 0360-0300. DOI: 10.1145/3136625. [Online]. Available: <https://doi.org/10.1145/3136625>.
- [8] F. L. Gewers, G. R. Ferreira, H. F. D. Arruda, *et al.*, "Principal component analysis: A natural approach to data exploration," *ACM Comput. Surv.*, vol. 54, no. 4, May 2021, ISSN: 0360-0300. DOI: 10.1145/3447755. [Online]. Available: <https://doi-org.ejournals.um.edu.mt/10.1145/3447755>.