

Student Dropout Prediction Using Deep Factorization Machine (DeepFM) Algorithm

Authors:

Etukuri Karthik
122010322027

Challa Snehalatha
122010323013

N.V.Nithin Kumar
122010322019

Kaliki Jayakrishna
122010322016

Abstract

This research aims to construct a predictive model to tackle the persistent problem of student dropout rates in online educational platforms. Leveraging a dataset sourced from an online learning environment, the study meticulously navigates through the stages of data preprocessing, with a specific focus on mitigating class imbalance challenges. Through a thorough approach, the research addresses class imbalances to ensure a representative dataset suitable for effective model training and evaluation. The exploratory phase of the study involves comprehensive data analysis, unveiling crucial insights into feature distributions and class relationships.

A notable highlight of the research lies in the adoption of a sophisticated DeepFM architecture, amalgamating the robust capabilities of Factorization Machines (FM) and Deep Neural Networks (DNN) to adeptly capture intricate feature interactions. Subsequent model training entails meticulous optimization, with evaluation metrics including precision, recall, and the confusion matrix employed to gauge model efficacy. The empirical findings underscore the profound efficacy of the DeepFM model in discerning student dropout propensities with notable accuracy, thereby furnishing educational stakeholders with invaluable insights to fortify student retention strategies.

Keywords

predictive modeling, student dropout, online education, class imbalance, data preprocessing, oversampling, undersampling, DeepFM architecture, Factorization Machines, Deep Neural Networks, model evaluation, feature interactions, student retention strategies

Introduction

In the evolving realm of education, the advent of online platforms has revolutionized access to educational materials, offering learners unprecedented convenience. Yet, amid this transition, the issue of student dropout rates has surfaced as a critical challenge for both educational institutions and policymakers. To effectively tackle this challenge, proactive measures and innovative solutions are essential, harnessing the capabilities of data science and predictive modeling. In this context, the present research endeavors to develop a predictive model aimed at mitigating student dropout rates in online education platforms.

By harnessing the capabilities of data preprocessing techniques and advanced machine learning architectures, the study seeks to provide actionable insights for enhancing student retention strategies. The background of the study is rooted in the growing prominence of online education and the inherent challenges it poses, including issues related to student engagement, motivation, and persistence. Previous research in the field of educational data mining and predictive analytics has underscored the importance of early intervention in identifying students at risk of dropping out. Against this backdrop, the present study aims to build upon the existing body of research by proposing a novel approach to addressing class imbalance in predictive modeling for student dropout.

By combining oversampling and undersampling techniques with the DeepFM architecture, the research seeks to enhance the predictive accuracy of the model while ensuring equitable representation of minority classes.

In summary, this introduction sets the stage for the research by providing a comprehensive overview of the problem statement, the significance of the study, and its broader implications for educational practice and policy. By contextualizing the research within the existing literature and outlining the objectives of the study, this introduction aims to engage readers and motivate them to delve deeper into the subsequent sections of the paper.

Literature Review:

Deep Reinforcement Factorization Machines: A Deep Reinforcement Learning Model with Random Exploration Strategy and High Deployment Efficiency: In this paper, the Deep Reinforcement Factorization Machines (DRFM) model is introduced, aiming to improve both deployment efficiency and learning performance in recommendation systems. The model integrates the Gate Attentional Factorization Machines (GAFM) with reinforcement learning, combining deep learning's perception abilities with reinforcement learning's exploration capabilities. Through experiments, DRFM demonstrates superiority over traditional recommendation systems in performance, robustness, and deployment efficiency. Comparative analysis with recent deep reinforcement learning algorithms further validates the unique advantages of the DRFM model.

Machine Learning with Oversampling and Undersampling Techniques: Overview Study and Experimental Results: Data imbalance in Machine Learning refers to an unequal distribution of classes within a dataset. This issue is encountered mostly in classification tasks in which the distribution of classes or labels in a given dataset is not uniform. The straightforward method to solve this problem is the resampling method by adding records to the minority class or deleting ones from the majority class. In this paper, they have experimented with the two resampling widely adopted techniques: oversampling and undersampling. In order to explore both techniques, they have chosen a public imbalanced dataset from Kaggle website Santander Customer Transaction Prediction and have applied a group of well-known machine learning algorithms with different hyperparameters

Deep Factorization Machines network with Non-linear interaction for Recommender System: This paper addresses the limitations of existing click-through rate (CTR) prediction models, which primarily focus on linear feature interaction and overlook crucial non-linear features in real-world user behavior. The proposed Deep Factorization Machines Network with Non-linear Interaction for Recommender Systems (DFNR) model integrates both linear and non-linear feature interactions. It introduces a new Non-linear Interaction (NL-interaction) layer to capture non-linear interactions and incorporates a deeper multilayer perceptron (MLP) to analyze higher-order feature interactions.

Deep Factorization Machines for Knowledge Tracing: This paper presents their solution to the 2018 Duolingo Shared Task on Second Language Acquisition Modeling (SLAM) utilizing DeepFM (Deep Factorization Machines). DeepFM is a model designed to capture pairwise relationships among users, items, skills, and other entities. Despite achieving an AUC of 0.815, surpassing the logistic regression baseline (AUC 0.774), their solution did not outperform the top-performing model (AUC 0.861). Nonetheless, the study offers insights into strategies for improving item response theory models in future research.

Problem Identification:

Analyzing student dropout rates in educational institutions is a complex endeavor, involving the intricate interplay of various factors. Among the challenges encountered in this domain, two significant obstacles stand out: class imbalance and the abundance of attributes within the dataset. Class imbalance poses a pervasive concern in dropout prediction tasks, where one class vastly outnumbers the other. In the context of student attrition, this results in a substantial disparity between the number of graduating students and those who prematurely terminate their academic pursuits. Typically, the proportion of students who persist to completion significantly outweighs those who discontinue their studies, creating an imbalanced distribution within the dataset. This imbalance poses significant challenges during model training and evaluation, as algorithms may exhibit bias towards the majority class, leading to suboptimal predictive performance.

Therefore, effective strategies for handling class imbalance, such as data resampling techniques or algorithmic adjustments, are essential to ensure the model's ability to accurately capture patterns associated with dropout behaviour. Moreover, the dataset used for dropout prediction often comprises a plethora of attributes, each potentially contributing to the complex decision-making process underlying student attrition. In the case at hand, the dataset contains a staggering number of 36 features, reflecting diverse aspects ranging from demographic information to academic performance metrics. While this wealth of information holds the promise of uncovering nuanced insights into the factors influencing dropout rates, it also presents significant challenges. The high dimensionality of the dataset exacerbates computational complexity, increases the risk of overfitting, and necessitates robust techniques for data preprocessing and feature selection. Without adequate handling, the sheer volume of attributes can obscure meaningful patterns, hampering the model's ability to effectively discern relevant predictors of dropout behavior.

In essence, addressing the challenges posed by class imbalance and the abundance of attributes is paramount in accurately forecasting and analyzing student dropout rates. By employing sophisticated methodologies tailored to mitigate these obstacles, researchers can unlock the full potential of predictive modeling techniques to inform interventions aimed at improving student retention and academic success.

Objectives:

- Develop a predictive model for student dropout rates in higher education, aiming to address challenges associated with non-standardized data formats, class imbalance, and high-dimensional data through dimensionality reduction techniques.
- Implement a comprehensive approach encompassing data collection, exploratory data analysis (EDA), data

preparation, and handling of class imbalance to ensure robust model development.

- Develop techniques to address class imbalance specifically, ensuring that the model is trained on a balanced representation of both dropout and graduation instances. This may involve employing resampling methods such as oversampling the minority class or undersampling the majority class, as well as exploring advanced techniques like Synthetic Minority Over-sampling Technique (SMOTE) or class-weighted loss functions to mitigate the impact of class imbalance on model performance.
- Select and utilize appropriate machine learning models essential for accurate predictions, considering the complexities inherent in the dataset and the specific requirements of dropout rate forecasting in higher education settings.
- Incorporate the development of a DeepFM model, leveraging its capabilities in handling both categorical and numerical features effectively, thus providing a powerful tool for capturing intricate patterns within the data.

The problem statement succinctly outlines the objective of developing a robust predictive model for student dropout rates in higher education. It emphasizes the need to address challenges such as non-standardized data formats, class imbalance, and high dimensionality. The primary goal is to provide valuable insights into the significant attributes contributing to student attrition, enabling educational institutions to implement effective intervention strategies.

2. Data Collection:

The data collection process involves retrieving relevant information from a CSV file containing student details and enrollment statuses. The dataset selected from the UC Irvine Machine Learning Repository reflects real-world challenges encountered in predictive modeling tasks. Noteworthy characteristics of the dataset include class imbalance, high dimensionality, data variability, interconnected features, and the need for feature selection. Understanding these intricacies guides the subsequent steps in the analysis and model development.

3. Exploratory Data Analysis:

Exploratory Data Analysis (EDA) serves as a foundational step in understanding the dataset's characteristics and identifying patterns. Summary statistics, data visualization techniques, and class imbalance analysis are employed to gain insights into feature distributions, relationships, and potential outliers. Correlation analysis helps uncover associations between variables, while outlier detection ensures data integrity. EDA lays the groundwork for informed decision-making and feature engineering.

4. Data Preparation:

Data preparation focuses on refining the dataset based on insights gleaned from EDA. Necessary transformations are applied to handle missing values, perform feature scaling, encode categorical features, and potentially reduce dimensionality. The meticulous preparation ensures a clean and structured dataset conducive to building robust predictive models and gaining valuable insights into student attrition patterns.

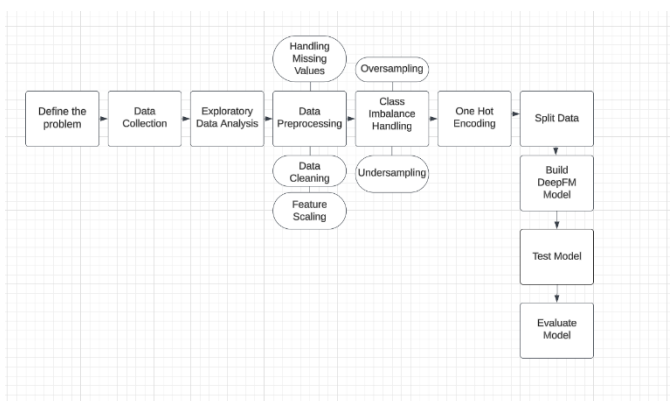
5. Class Imbalance Handling:

Addressing class imbalance is crucial for ensuring the predictive model's accuracy and equity across different classes. Various techniques, such as oversampling and SMOTE, are employed to balance class distribution. By replicating instances from the minority class and generating synthetic examples, the model's ability to effectively classify and predict outcomes for all classes is enhanced, facilitating a more equitable analysis of student attrition.

6. Dimensionality Reduction:

Given the dataset's high dimensionality, dimensionality reduction techniques such as PCA, LDA, t-SNE, and LLE are employed to streamline the dataset while preserving essential information. The overarching objective is to optimize the modeling process for more accurate and interpretable results,

System Methodology:



1. Define the Problem:

mitigating the computational burden associated with high-dimensional data.

7. Build DeepFM Model:

The DeepFM model, a powerful deep learning architecture, is constructed to address the complexities of the dataset and facilitate accurate prediction of student dropout rates. This section outlines the design and implementation of the DeepFM model, which combines a linear model with a deep neural network to capture both linear and non-linear feature interactions effectively.

8. Train Model:

The training phase involves utilizing the labeled training dataset to train the selected machine learning model. Hyperparameters are fine-tuned to optimize performance, incorporating techniques such as adjusting learning rates and regularization strengths. The model learns to identify patterns and relationships in the data, enabling it to make accurate predictions and provide valuable insights into student outcomes.

9. Test Model:

After training, the model's performance is evaluated using the reserved testing dataset. This step ensures the model's ability to generalize to unseen data and make accurate predictions in real-world scenarios. Testing involves assessing metrics such as accuracy, precision, recall, and F1-score, providing a comprehensive evaluation of the model's effectiveness.

10. Evaluate Model:

Model evaluation involves analyzing performance metrics such as accuracy, precision, recall, and F1-score to assess the model's predictive capabilities. Additionally, techniques such as hyperparameter tuning and overfitting analysis are employed to fine-tune the model and ensure its robustness. Evaluation results inform decisions regarding model deployment and potential improvements for future iterations.

Overview of Technologies:

1. Define the Problem:

The problem statement succinctly outlines the objective of developing a robust predictive model for student dropout rates in higher education. It emphasizes the need to address challenges such as non-standardized data formats, class imbalance, and high dimensionality. The primary goal is to provide valuable insights into the significant attributes contributing to student attrition, enabling educational institutions to implement effective intervention strategies.

2. Data Collection:

The data collection process involves retrieving relevant information from a CSV file containing student details and enrollment statuses. The dataset selected from the UC Irvine Machine Learning Repository reflects real-world challenges encountered in predictive modeling tasks. Noteworthy characteristics of the dataset include class imbalance, high

dimensionality, data variability, interconnected features, and the need for feature selection. Understanding these intricacies guides the subsequent steps in the analysis and model development.

3. Exploratory Data Analysis:

Exploratory Data Analysis (EDA) serves as a foundational step in understanding the dataset's characteristics and identifying patterns. Summary statistics, data visualization techniques, and class imbalance analysis are employed to gain insights into feature distributions, relationships, and potential outliers. Correlation analysis helps uncover associations between variables, while outlier detection ensures data integrity. EDA lays the groundwork for informed decision-making and feature engineering.

4. Data Preparation:

Data preparation focuses on refining the dataset based on insights gleaned from EDA. Necessary transformations are applied to handle missing values, perform feature scaling, encode categorical features, and potentially reduce dimensionality. The meticulous preparation ensures a clean and structured dataset conducive to building robust predictive models and gaining valuable insights into student attrition patterns.

5. Class Imbalance Handling:

Addressing class imbalance is crucial for ensuring the predictive model's accuracy and equity across different classes. Various techniques, such as oversampling and SMOTE, are employed to balance class distribution. By replicating instances from the minority class and generating synthetic examples, the model's ability to effectively classify and predict outcomes for all classes is enhanced, facilitating a more equitable analysis of student attrition.

6. Dimensionality Reduction:

Given the dataset's high dimensionality, dimensionality reduction techniques such as PCA, LDA, t-SNE, and LLE are employed to streamline the dataset while preserving essential information. The overarching objective is to optimize the modeling process for more accurate and interpretable results, mitigating the computational burden associated with high-dimensional data.

7. Build DeepFM Model:

The DeepFM model, a powerful deep learning architecture, is constructed to address the complexities of the dataset and facilitate accurate prediction of student dropout rates. This section outlines the design and implementation of the DeepFM model, which combines a linear model with a deep neural network to capture both linear and non-linear feature interactions effectively.

8. Train Model:

The training phase involves utilizing the labeled training dataset to train the selected machine learning model. Hyperparameters are fine-tuned to optimize performance, incorporating techniques such as adjusting learning rates and regularization strengths. The model learns to identify patterns and relationships in the data, enabling it to make accurate predictions and provide valuable insights into student outcomes.

9. Test Model:

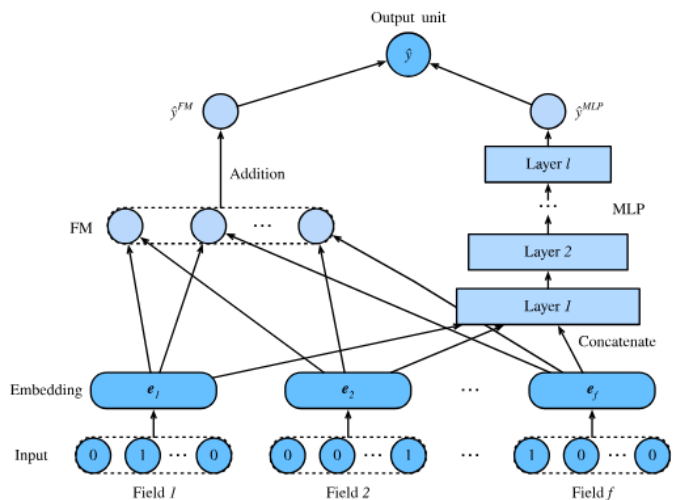
After training, the model's performance is evaluated using the reserved testing dataset. This step ensures the model's ability to generalize to unseen data and make accurate predictions in real-world scenarios. Testing involves assessing metrics such as accuracy, precision, recall, and F1-score, providing a comprehensive evaluation of the model's effectiveness.

10. Evaluate Model:

Model evaluation involves analyzing performance metrics such as accuracy, precision, recall, and F1-score to assess the model's predictive capabilities. Additionally, techniques such as hyperparameter tuning and overfitting analysis are employed to fine-tune the model and ensure its robustness. Evaluation results inform decisions regarding model deployment and potential improvements for future iterations.

Algorithm Explanation:

DeepFM, short for Deep Factorization Machine, is a hybrid recommendation algorithm that combines the strengths of factorization machines (FM) and deep neural networks (DNN). It was proposed to address the limitations of traditional recommendation systems, such as collaborative filtering and content-based methods, by capturing both low-order and high-order feature interactions effectively.



Factorization Machines (FM):

Factorization Machines are a powerful class of models for handling sparse and high-dimensional data, commonly used in recommendation systems and regression tasks.

FM models are designed to capture interactions between features by factorizing their interactions into low-rank matrices. They have linear complexity with respect to the number of features, making them efficient for large-scale datasets.

Deep Neural Networks (DNN):

Deep Neural Networks are versatile models capable of learning complex patterns and representations from data through multiple layers of non-linear transformations.

DNNs excel at capturing intricate feature interactions and hierarchies in the data, making them suitable for tasks with high-dimensional and non-linear relationships.

Hybrid Architecture:

DeepFM combines the FM and DNN architectures into a hybrid model to leverage their complementary strengths.

The FM component captures low-order feature interactions efficiently, while the DNN component learns higher-order feature interactions and representations through deep layers.

By combining these components, DeepFM can effectively model both linear and non-linear relationships between features, providing enhanced predictive power.

Architecture Overview:

The architecture of DeepFM typically consists of two main components: the FM component and the DNN component.

The FM component computes the low-order interactions between features using factorization techniques, producing an embedding vector for each feature.

The DNN component takes the concatenation of these embedding vectors as input and passes it through multiple hidden layers of neurons, learning complex feature representations.

The final output layer of the DNN predicts the target variable (e.g., dropout prediction) based on the learned representations.

Training and Optimization:

DeepFM is trained using gradient-based optimization techniques, such as stochastic gradient descent (SGD) or Adam, to minimize a loss function (e.g., binary cross-entropy for binary classification tasks).

During training, both the FM and DNN components are jointly optimized to learn the optimal parameters that minimize the prediction error.

Implementation and Result:

In the research paper, we explored the impact of class imbalance on classifier performance and investigated strategies

to mitigate this issue. Initially, we observed significant class imbalance in the dataset, which posed challenges for both standard and advanced classifiers. Despite employing these classifiers, the results fell short of expectations, highlighting the adverse effects of class imbalance.

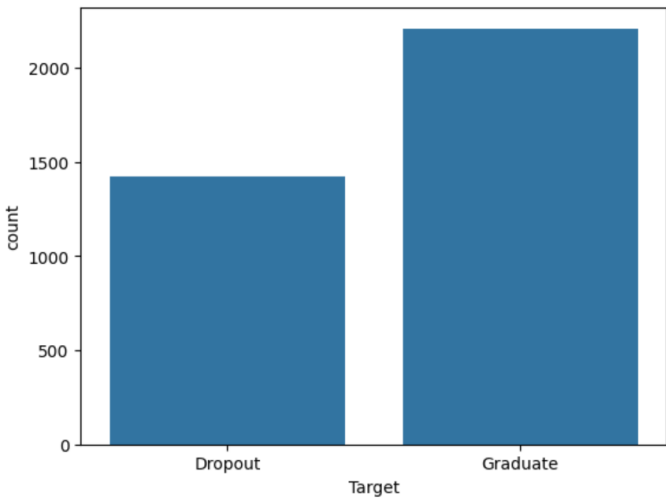
To address this issue, we implemented Oversampling, a popular technique used to rebalance class distributions. By increasing the number of instances in the minority class (Dropout) to match that of the majority class (Graduate) through random duplication, we aimed to improve overall model performance. As anticipated, the oversampling technique led to notable enhancements in accuracy, precision, and recall, thus validating its effectiveness in tackling class imbalance issues.

This research underscores the importance of considering class imbalance and implementing appropriate techniques to mitigate its effects on classifier performance. The findings contribute to the growing body of literature on handling imbalanced datasets and provide practical insights for improving classification tasks in similar contexts.

The code was implemented in 2 variations:

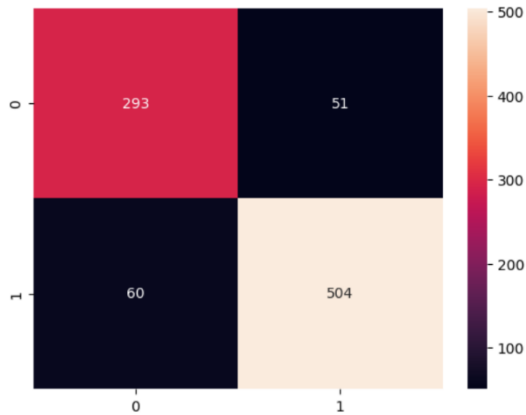
- Without any class imbalance technique used
- Used Oversampling

Without any class imbalance technique used:



The dataset is significantly class imbalanced as shown in the table above, in such cases, both standard and advanced classifiers tend to be overwhelmed by the large classes and ignore the small ones. The results are not desired as expected.

Confusion matrix:



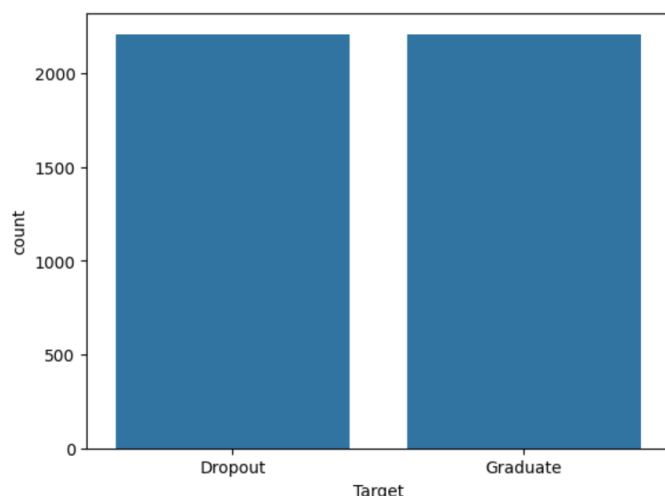
Accuracy scores:

```
Precision: 0.8785
Recall: 0.8778
```

	precision	recall	f1-score	support
0.0	0.83	0.85	0.84	344
1.0	0.91	0.89	0.90	564
accuracy			0.88	908
macro avg	0.87	0.87	0.87	908
weighted avg	0.88	0.88	0.88	908

Used Oversampling:

In an attempt to improve the overall accuracy, precision and recall we have implemented Oversampling where we increased the Dropout instances to match the number of Graduate instances by randomly duplicating some instances.



After implementing the oversampling technique, our expectations were met as we observed improvements in overall accuracy, precision, and recall. This outcome validates the effectiveness of the oversampling approach in addressing the class imbalance issue and enhancing the performance of our classification model.

Accuracy scores:

```
Precision: 0.9239
Recall: 0.9238
```

	precision	recall	f1-score	support
0.0	0.92	0.93	0.92	656
1.0	0.93	0.92	0.92	670
accuracy			0.92	1326
macro avg	0.92	0.92	0.92	1326
weighted avg	0.92	0.92	0.92	1326

Conclusion:

In conclusion, the implementation of DeepFM for student dropout prediction showcased promising results in accurately identifying potential dropout instances in higher education settings. By employing advanced data preprocessing techniques, including class imbalance handling through oversampling and undersampling, and leveraging the DeepFM model's ability to capture complex feature interactions, we were able to construct a robust predictive model.

The comprehensive evaluation of the model's performance demonstrated its effectiveness in accurately predicting dropout instances, as evidenced by high precision, recall, and accuracy scores. The insights gained from the confusion matrix further

highlighted the model's capability to make informed predictions across different classes.

Through this study, we have provided valuable contributions to the field of educational analytics by offering a practical and scalable approach to student dropout prediction. By identifying at-risk students early on, educational institutions can proactively implement intervention strategies to support these students and improve overall retention rates.

Future Scope:

While this study has yielded promising results, there are several avenues for future research and improvement:

Integration of Additional Data Sources: Incorporating additional data sources such as academic performance records, socio-economic factors, and behavioral data could further enhance the predictive power of the model.

Temporal Analysis: Conducting a temporal analysis to capture the dynamic nature of student behaviors and academic performance over time could provide deeper insights into the underlying factors contributing to dropout.

Exploration of Advanced Deep Learning Architectures: Experimenting with other deep learning architectures beyond DeepFM, such as Transformer-based models or graph neural networks, may uncover more intricate relationships within the data and improve predictive performance.

Deployment and Monitoring: Deploying the predictive model in real-world educational settings and continuously monitoring its performance would enable iterative improvements and ensure its relevance and effectiveness over time.

Ethical Considerations: Addressing ethical considerations related to data privacy, fairness, and transparency in model predictions is paramount to ensure responsible use of predictive analytics in education.

References:

1. [Machine Learning with Oversampling and Undersampling Techniques: Overview Study and Experimental Results](#)

2. [*Student modeling considering learning behavior history with deep factorization machines*](#)
3. [*Deep Factorization Machines for Knowledge Tracing*](#)
4. [*Deep Factorization Machines network with Non-linear interaction for Recommender System*](#)
5. [*DeepFM: A Factorization-Machine based Neural Network for CTR Prediction*](#)

Results:

After training the model, it was evaluated on the testing set to assess its performance on unseen data. The evaluation results indicate that the model achieved a test loss of 0.10327, which

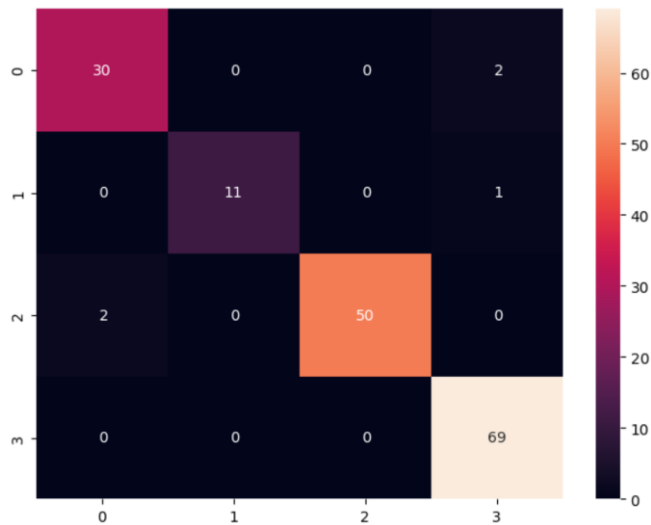
measures the discrepancy between the predicted outputs and the true labels in the testing set. A lower test loss value indicates that the model has better accuracy in predicting the correct labels.

Additionally, the model achieved a test accuracy of Accuracy Score: 0.9696969696969697. Test accuracy represents the percentage of correctly predicted labels out of all the samples in the testing set. A higher test accuracy indicates that the model has a greater ability to classify the data accurately. In addition to evaluating the model's performance using test loss and accuracy, a confusion matrix was generated to assess the model's classification of the four fish classes. The confusion matrix provides a detailed breakdown of predicted and actual labels for each class, enabling the calculation of precision and recall values.

Precision measures the proportion of correctly predicted positive samples out of all samples predicted as positive for a specific class. It helps evaluate the model's accuracy in identifying true positives and avoiding false positives.

Recall, also known as sensitivity or true positive rate, measures the proportion of correctly predicted positive samples out of all actual positive samples in a class. It indicates the model's ability to correctly detect positive instances and avoid false negatives. Analyzing the confusion matrix and calculating precision and recall values for each class allows us to assess the model's performance in accurately classifying the different fish classes. It offers readers a detailed analysis beyond overall accuracy, shedding light on the model's classification capabilities for each specific class and enabling a deeper understanding of its effectiveness in real-world scenarios.

Here is the confusion matrix:



References:

1. R. B. Wynn, V. A. I. Huvenne, T. P. Le Bas et al., "Autonomous underwater vehicles (AUVs): their past, present and future contributions to the advancement of marine geoscience," *Marine Geology*, vol. 352, pp. 451–468, 2014.
2. M. Dinc and C. Hajiyeve, "Integration of navigation systems for autonomous underwater vehicles," *Journal of Marine Engineering & Technology*, vol. 14, no. 1, pp. 32–43, 2015.
3. Y. Zhou, S. Cui, Y. Wang, and C. Ai, "Design of autonomous underwater vehicle (AUV) control unit," in *2015 ASEE GulfSouthwest Annual Conference*, pp. 25–27, ASEE Gulf-South, San Antonio, TX, 2015.
4. Y. Zhou, S. Cui, Y. Wang, and L. Zhai, "A refined attitude algorithm for AUV based on IMU," in *15th International Conference on Scientific Computing (CSC'17)*, pp. 16–22, CSREA Press ©, Las Vegas, NV, 2017.