

Investigating the Factors Affecting Risky Levels of Alcohol Consumption among Students Using Machine Learning Approach

Sara Fariha Shanchary
Department of Electrical and
Computer Engineering , North South
University
Bangladesh
sara.shanchary@northsouth.edu

Md Naved Meraz
Department of Electrical and
Computer Engineering , North South
University
Bangladesh
naved.meraz@northsouth.edu

Ayman Ibne Hakim
Department of Electrical and
Computer Engineering , North South
University
Bangladesh
ayman.hakim@northsouth.edu

Chowdhury Nafis Faiyaz
Department of Electrical and
Computer Engineering , North South
University
Bangladesh
nafis.faiyaz@northsouth.edu

Muhammad Shafayat Oshman
Department of Electrical and
Computer Engineering , North South
University
Bangladesh
muhammad.oshman@northsouth.edu

Abstract

Addressing the pressing issue of alcohol consumption among students is crucial for the well-being of young individuals and the community. Understanding and mitigating alcohol related challenges faced by young people is essential for their healthy development. This study utilizes machine learning models to analyze data from a Portuguese school, aiming to link students' alcohol consumption levels to personal and familial factors. The primary objective is to identify key factors associated with high alcohol consumption among students. After data preprocessing on their dataset, we employed various machine learning algorithms, including hyperparameter-optimized Decision Tree, Random Forest, Boosting, and Ensemble Learning. Our findings revealed that the decision tree algorithm performed well in predicting risky alcohol consumption for our target research question. Our selected feature subset showed strong positive correlation with the target variable, achieving an accuracy of 80.9 percent on the test set and 98.43 percent on the train set. Notably, this project breaks new ground by using explainable artificial intelligence to add reason to our prediction based on students' familial relationships, expanding upon previous research that also focused on demographic factors.

CCS Concepts

• **Computing methodologies** → Artificial intelligence; Ensemble methods; Feature selection; Supervised learning.

Keywords

alcohol consumption, alcohol use disorder, machine learning, family relationships

ACM Reference Format:

Sara Fariha Shanchary, Md Naved Meraz, Ayman Ibne Hakim, Chowdhury Nafis Faiyaz, and Muhammad Shafayat Oshman. 2024. Investigating the Factors Affecting Risky Levels of Alcohol Consumption among Students Using Machine Learning Approach. In *2024 6th Asia Conference on Machine Learning and Computing (ACMLC 2024)*, July 26–28, 2024, Bangkok, Thailand. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3690771.3690777>

1 Introduction

The consumption of alcohol among students is a prevalent and well-documented phenomenon that can have grave consequences for both individuals and society as a whole. It has significant societal repercussions, including academic decline, violence, injuries, mental health issues, and fatalities [1]. Understanding the determinants of this risky behaviour is vital, especially because young individuals are often introduced to alcohol early, leading to societal challenges [2] which can continue further into adult life. This is a probable pattern for many young individuals and thus is a societal issue. Unchecked alcohol consumption can develop into risky drinking patterns that can affect work, focus, and health [3]. Recovery from alcohol use disorder (AUD) often involves medications with serious side effects [4], making it imperative to investigate the underlying reasons for AUD. We employ machine learning to predict high alcohol consumption among school students in Portugal—focusing on factors like poor academic performance and unstable family relationships. To ensure transparency and mitigate biases, we utilize LIME (Local Interpretable Model-Agnostic Explanations) and SHAP (SHapley Additive exPlanations), model-agnostic explainable artificial intelligence (XAI) techniques. Our contributions encompass:

- Utilizing multiple classifiers—Decision Tree, Random Forest, Soft voting, Hard voting, Out-of-bag error (OOB), extreme gradient boost (XGBoost), Adaptive Boosting (AdaBoost).
- Performing hyperparameter optimization for the best-performing classifiers.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ACMLC 2024, July 26–28, 2024, Bangkok, Thailand

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-1001-8/24/07

<https://doi.org/10.1145/3690771.3690777>

- Applying XAI methods (LIME and SHAP) to identify influential features.

The subsequent sections are organized as follows: Section II offers a critical literature review, followed by our research methodology in Section III. Section IV presents our experimental findings, while Section V discusses limitations. Finally, Section VI concludes by summarizing key findings and suggesting avenues for future research.

2 Literature review

In previous approaches to identify factors contributing to AUD, individual substances have fallen short in capturing complex inter-relationships, as highlighted by [5]. Hence, the authors propose a joint model to predict hazardous substance use. Trained on individual, familial, and socio-demographic data of adolescent substance users, their model surpassed traditional regression models in accuracy. It identifies risk factors like early tobacco use, reduced parental attachment, and heightened early-life stress for finding at-risk adolescents. Limitations include cross-sectional design, inability to establish causality, and a small sample size for youths regularly using alcohol and cannabis; potentially limiting generalizability. Despite these constraints, the study provides valuable insights into adolescent hazardous substance abuse.

In a study [6], machine learning and deep learning techniques were used to analyze alcohol consumption patterns among secondary school students in relation with their enrolled courses. The study explores various machine learning models with XGBoost and Random Forest exhibiting the highest accuracy levels. The study reveals a higher proportion of alcohol use among students in the Portuguese course compared to the mathematics course. The same dataset [7] was previously used in [8] to identify alcohol addiction using a Decision Tree algorithm, achieving 93% accuracy. This study identifies ‘going out with friends’ as the most influential factor in alcohol consumption, while emphasizing that it lacks a validation set and is solely focused on the Decision Tree algorithm.

As seen in [9], the study used neural networks to predict the alcohol consumption behaviors among higher secondary school students. Their study compares the prediction accuracies of two neural network algorithms: a self-tuning multi-layered perceptron classifier (AutoMLP) and the standard multi-layered perceptron (MLP), using the dataset from [7]. Both MLP and AutoMLP models constructed based on the training set achieved classification on the test set using 10-fold cross-validation. The study found that AutoMLP outperformed standard neural networks, achieving an accuracy of 64.54% compared to 61.78%.

A 2021 study, funded by Seoul St. Mary’s Hospital [10], aimed to identify hazardous drinkers and assess the severity of alcohol-related issues using deep learning. The study, based on a large South Korean population survey, employed conventional machine learning algorithms (Logistic Regression, SVM, Random Forest, and K-Nearest Neighbors) alongside deep learning, consistently finding deep learning to outperform. The study unveils that hazardous drinker had higher average intakes of energy, carbohydrates, fats, proteins, and dietary components. However, the dataset lacks crucial features like social and economic status, or family history of

alcohol abuse; limiting a comprehensive understanding of alcohol-related problems.

As seen in [11], the authors surveyed online, 4840 Brazilian medical students, with 53.03% exhibiting high-risk alcohol consumption. Utilizing a dataset of 93 predictors spanning demographic, personal, medical student domain, and mental health-related variables, machine learning models (Lasso, Elastic-Net Regularized Generalized Linear Models, artificial neural network, and Random Forest) effectively differentiated high-risk from low-risk drinkers. The average Area Under Curve (AUC) scores in cross-validation were approximately 0.72 to 0.73. The study identified the top 30 relevant features associated with high-risk alcohol consumption, including factors like tobacco or cannabis use, monthly family income, relationship status, physical activity, and sexual orientation; influencing the model’s classification of high-risk drinking.

Perceptions of peer behavior significantly influence teen drinking. In a study seen in [12], theory-driven covariate selection was combined with machine learning to reveal that adolescents often overestimate their peers’ alcohol use and abuse. Using a sample of 14,738 respondents, Generalized Linear Models, Random Forest algorithms, and XGBoost Regression were used. XGBoost Regression achieved the highest accuracy (75.83%). The findings highlighted the positive associations between both-actual peer drinking and normative misperception with adolescent alcohol consumption, underscoring the significant role of peer norms. The study suggests that interventions targeting peer norms may effectively reduce adolescent alcohol consumption.

One of the first studies to develop and validate a machine learning model for predicting adolescent alcohol use in a cross-study, cross-cultural setting [13] used data from two groups of adolescents: one in Canada and one in Australia. The study used seven machine learning algorithms, creating models for different clusters. The elastic-net algorithm demonstrated superior performance, predicting alcohol consumption levels with over 85% accuracy in both group samples. Key predictors included externalizing psychopathologies, baseline alcohol use, and sensation-seeking personality profiles. However, the study’s limitation to two countries restricts the generalizability of findings to other cultural settings.

3 Methodology

Our research’s key objective is to identify multiple models that can predict with good accuracy the noticeable traits associated with heavy drinking among students. For our research and analysis, we have used the dataset from [8], that contains records of secondary school students in Portugal in two distinct courses: Mathematics and Portuguese language courses, along with their demographic and personal information. The block diagram for our research methodology is seen in Figure 1.

3.1 Block Diagram

The preprocessing involved merging similar columns, converting categorical features into numeric using One Hot Encoder, and scaling all features with MinMaxScaler. We noticed the attribute size has increased to 45 after encoding. The dataset was split into an 80:20 train-test ratio. Synthetic Minority Oversampling Technique (SMOTE) was applied to the train set to address class imbalance.

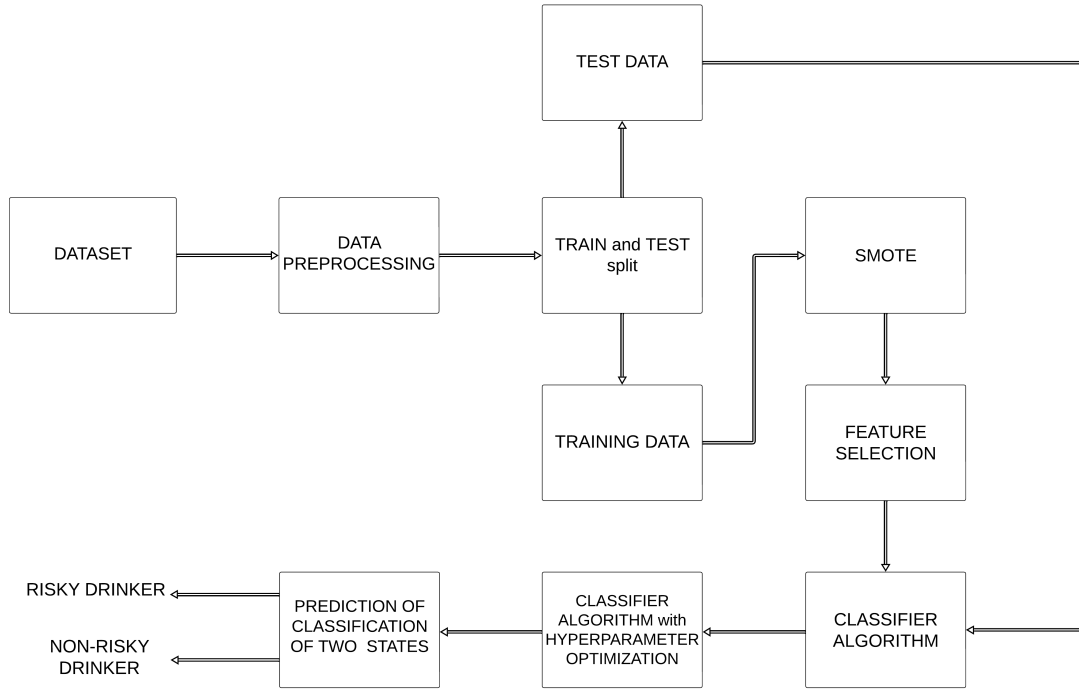


Figure 1: Block diagram of the research methodology

Table 1: Dataset Attributes

| Attribute | Description (datatype) |
|-----------|---|
| age | student's age (numeric) |
| famsize | family size (binary) |
| Pstatus | parent's cohabitation status (binary) |
| Medu | mother's education (numeric) |
| Fedu | father's education (numeric) |
| Mjob | mother's job (nominal) |
| Fjob | father's job (nominal) |
| guardian | student's guardian (nominal) |
| famsup | family educational support (binary) |
| famrel | quality of family relationships (numeric) |
| alc | weekly alcohol consumption (numeric) |

Feature selection was performed on the SMOTE dataset, and classifiers were applied to both the train and test sets. Hyperparameter optimization was conducted on the classifiers to enhance classification performance. Finally, predictions were made to categorize individuals as “risky drinker” (class 1) or “non-risky drinker” (class 0) based on classifier’s output.

3.2 Dataset

The study began by merging data from Arithmetic and Portuguese language classes. Post-merge, we retained 11 relevant columns and discarded the rest as irrelevant. The aggregation of workday

and weekend alcohol consumption columns formed a new variable, ‘weekly alcohol consumption’. Alcohol consumption levels were measured on a scale of 1 to 5, denoting very low to very high consumption, respectively. Based on cumulative weekly intake, individuals surpassing level 4 were identified as unsafe drinkers. Key features for the analysis is seen in Table 1.

3.3 Research Question

We assessed multiple classification algorithms on our manually curated subset of features to identify the best-performing model.

Employing XAI, we sought insights into the predictions of this top-performing model. We hypothesized that alcohol consumption is influenced by a lack of family or educational support, aligning with existing research indicating the impact of parental and educational support systems on student alcohol consumption; actively involved parents correlate with lower chances of alcohol misuse in children. The target question for confirming the hypothesis is as follows:

Does the quality of family relationships, parents' cohabitation status have an impact on alcohol consumption? The classifiers used for the research question are discussed below.

3.4 Classifiers used

Decision Tree: Decision trees are a versatile supervised machine learning algorithm used for classification and regression. They create a tree-like model where internal nodes represent features and leaf nodes signify class labels or prediction values. The algorithm splits data at each internal node based on feature values, directing it to the appropriate child node. The root node is selected by maximizing information gain and minimizing entropy, as calculated using equations (1) and (2).

$$Entropy(p) = - \sum_{i=1}^c p_i(t) \log_2 p_i(t) \quad (1)$$

$$Gain = Entropy(p) - \sum_{i=1}^k \frac{n_i}{n} Entropy(i) \quad (2)$$

Here c = number of classes, k = number of partitions, n = number of parent nodes, and n_i = number of child nodes. The Gini index in equation (3) is considered for each node and hyperparameter tuning is performed for multiple max depths. In the context of this research, the Gini impurity criterion was used for splitting nodes.

$$Gini = 1 - \sum_{i=1}^c p_{i(t)}^2 \quad (3)$$

Random Forest: Random Forest, an ensemble learning approach, is used for classification, regression, and diverse tasks. It constructs multiple Decision trees during training, with each tree based on a bootstrap sample of the training data. At each node, a random subset of features is considered for splitting, mitigating overfitting and improving generalization. In this study, Random Forest was applied to predict risky alcohol consumption by using the entropy criterion for node splitting.

Soft voting: soft voting is an ensemble learning voting method that enhances prediction accuracy by combining the outputs of multiple machine learning models. It begins by averaging the predicted probabilities from each model for each class. The class with the highest average probability is then selected as the final ensemble prediction.

Hard voting: Hard voting is an ensemble learning approach that combines the outputs of multiple individual models or classifiers to make a final prediction. Each model independently predicts the class label for a given input, and the class label with the most votes is selected as the ultimate prediction.

OOB evaluation: OOB evaluation is commonly employed in ensemble learning, especially in bagging-based algorithms like random forests as it offers a performance estimate without requiring a distinct validation set. This technique predicts the target variable

for each data point in the training set using all trees in the forest, excluding the one trained on that specific data point. The OOB error rate, calculated using equation (4), represents the proportion of data points for which the predicted target variable is incorrect.

$$OOB(error) = \frac{\text{number of incorrect predictions}}{\text{total number of data points}} \quad (4)$$

XGBoost: XGBoost is a highly scalable and efficient gradient boosting implementation extensively applied in various machine learning tasks. It operates by sequentially constructing a series of weak decision trees, with each tree trained to minimize the loss of the preceding one. This iterative process continues until a pre-defined stopping criterion is satisfied. XGBoost employs gradient tree boosting, a more efficient approach to training tree models compared to traditional boosting methods. In this context, 100 parameters were utilized in XGBoost.

AdaBoost: AdaBoost is an ensemble machine learning algorithm that aggregates multiple weak learners to form a robust learner. It functions through iterative training of weak learners on weighted versions of the training data. The data point weights are adjusted after each iteration to emphasize misclassified points from the current weak learner. This cycle persists until a stopping condition is fulfilled, like reaching a maximum iteration limit or achieving a specific accuracy threshold. In this study, the decision tree was employed as the weak learner, with a learning rate set at 0.4.

4 Results and discussions

Training and Test dataset contains 835 and 209 instances respectively and 11 attributes from original dataset were chosen for our research question.

The accuracies of the top-performing classification algorithms following grid search on both training and test data is seen in Table 2.

In addition to accuracy, the study assessed various performance metrics such as TPR (True Positive Rate), FPR (False Positive Rate), Precision (Positive Predictive Value), Recall (Sensitivity), and F-Measure (evaluation metric that measures a model's accuracy) to enable a comprehensive comparison of the utilized algorithms.

Decision Tree and Random Forest algorithm performing the best among all other algorithms on the train set is seen in Table 3.

Again, Decision Tree and Random Forest algorithm performing the best among all other algorithms on the test set is seen in Table 4.

A receiver operating characteristic (ROC) curve is a graphical representation of the performance of a classification model. It plots the sensitivity (TPR) against the specificity ($1 - \text{FPR}$) for different cut-off values. The ROC curve for our model is seen in Figure 2. The AUC is 0.89, which indicates that the model is performing well. The predictions were further analyzed using XAI which will be discussed below.

4.1 Explainable artificial intelligence

With machine learning models becoming "black boxes", understanding their predictions has become challenging. To address this, there is a growing emphasis on machine learning interpretability and explainability, achieved through tools like XAI [14]. While many notable works predicted alcohol consumption using this dataset, we have employed two XAI tools—SHAP and LIME—to analyze

Table 2: Accuracy for Train and Test set

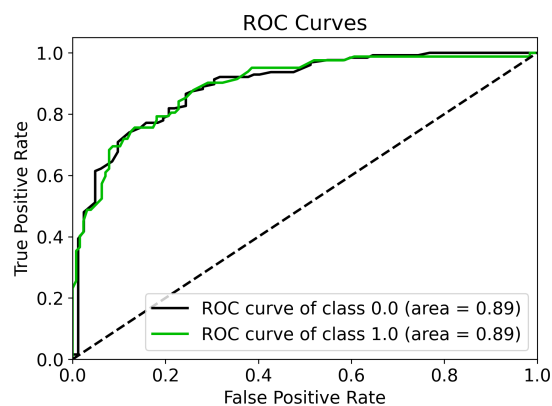
| Algorithms | Train Accuracy | Test Accuracy |
|----------------|----------------|---------------|
| Decision Tree | 98.43% | 80.90% |
| Random Forest | 98.43% | 77.50% |
| OOB evaluation | 98.43% | 77.50% |
| Adaboost | 98.43% | 78.95% |
| XGBoost | 96.64% | 78.47% |
| Soft Voting | 96.86% | 79.43% |
| Hard Voting | 96.86% | 77.99% |

Table 3: Detailed evaluation on train set

| Algorithms | TP rate | FP rate | Precision | Recall | F1 score |
|----------------|---------|---------|-----------|--------|----------|
| Decision Tree | 0.97 | 0.004 | 1.00 | 0.97 | 0.98 |
| Random Forest | 0.98 | 0.01 | 0.99 | 0.98 | 0.98 |
| OOB evaluation | 0.98 | 0.01 | 0.99 | 0.98 | 0.98 |
| Adaboost | 0.98 | 0.006 | 0.99 | 0.98 | 0.98 |
| XGBoost | 0.95 | 0.02 | 0.98 | 0.95 | 0.97 |
| Soft Voting | 0.94 | 0.004 | 1.00 | 0.94 | 0.97 |
| Hard Voting | 0.94 | 0.004 | 1.00 | 0.94 | 0.97 |

Table 4: Detailed evaluation on test set

| Algorithms | TP rate | FP rate | Precision | Recall | F1 score |
|----------------|---------|---------|-----------|--------|----------|
| Decision Tree | 0.79 | 0.18 | 0.74 | 0.79 | 0.76 |
| Random Forest | 0.76 | 0.21 | 0.70 | 0.76 | 0.73 |
| OOB evaluation | 0.76 | 0.21 | 0.70 | 0.76 | 0.73 |
| Adaboost | 0.78 | 0.20 | 0.71 | 0.78 | 0.78 |
| XGBoost | 0.76 | 0.20 | 0.71 | 0.76 | 0.73 |
| Soft Voting | 0.80 | 0.21 | 0.71 | 0.80 | 0.75 |
| Hard Voting | 0.72 | 0.18 | 0.72 | 0.72 | 0.72 |

**Figure 2: ROC curve on test set for Random Forest**

the best grid of the random forest classifier was used for model fitting in both LIME and SHAP explanations.

4.1.1 LIME. LIME is a model-agnostic machine learning tool that helps us interpret our ML models. LIME helps interpret reason for a specific prediction.

The results of one specific prediction from test set obtained show that the attribute, Fjob health (father's job in health), holds the highest feature importance with a score of 14% is seen in Figure 3. This is followed by Pstatus_A (parental status living apart) at 14% and famrel (quality of family relationships) at 7%. Notably, when the father's job is not health-related (Fjob health value 0.00 in the chart), it significantly contributes to the high label. Features highlighted in orange contribute to class 1, risky drinker, while those in blue impact class 0, non-risky drinker. Consequently, the cumulative feature importance scores lead to this prediction of high alcohol consumption. In summary, both parents working in service jobs, low family relationship quality, and parents living apart collectively led the model to predict this student as a high-risk drinker.

prediction reasons. Feature values help identify these reasons, and

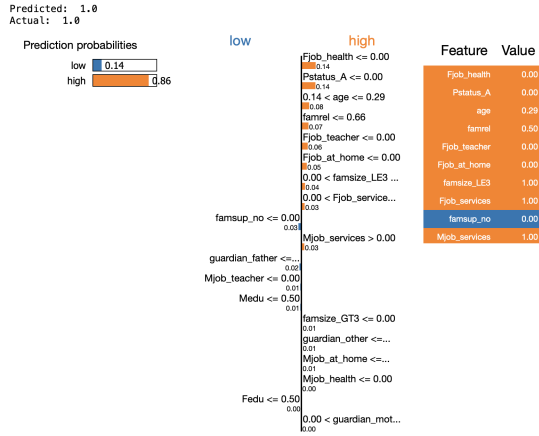


Figure 3: LIME Visualization for Prediction Reason

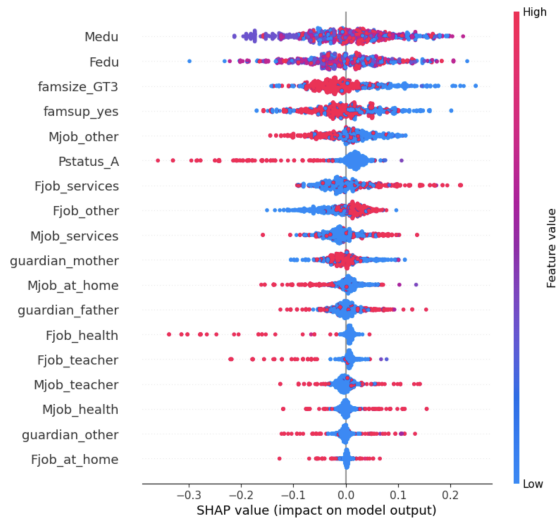


Figure 4: SHAP beeswarm plot

4.1.2 SHAP. SHAP is a mathematical method to explain the predictions of machine learning models. SHAP helps visualize across all samples.

To explain the prediction rationale, we used the SHAP explainer method. The beeswarm plot displays the contribution of the most critical features to the model, showing SHAP values across all samples is seen in Figure 4.

The color coding (red: high, blue: low) indicates feature importance. The plot effectively illustrates the distribution of each feature's impact on the model output [15]. A positive impact predicts risky drinking, while a negative impact suggests non-risky drinking.

Parents' education did not exhibit a clear pattern for students' alcohol consumption. Students with positive family relationships were less prone to risky consumption of alcohol, while those from families having fewer than 3 members and lacking family support were more likely to engage in risky drinking. Living apart from

parents correlated with a reduced likelihood of risky alcohol behavior. Students with fathers in healthcare and education sectors tended to maintain safer alcohol intake, while those with fathers in services and other jobs did not ensure safe intake. Mothers' occupation at home indicated a lower likelihood of dangerous alcohol consumption, while other categories of mother's job showed mixed trends in students' alcohol consumption.

It is to be noted that the targeted demographic was students, and not a specific location. Technically, the dataset is a set of numerical values and the factors that affect the prediction are justified using the algorithms, regardless of the location. Any contextual dataset should give us similar results. Therefore, we can acknowledge that the main target was to investigate factors because those have a direct correlation to the prediction of various models used.

5 Limitations

Several limitations of this study should be considered.

- Due to the predominantly categorical nature of the dataset, the encoding process resulted in a significant expansion of features. This complexity presented challenges in the manual selection of optimal features and the exploration of various feature subsets, impacting our ability to identify the best-performing model.
- It is important to note that the dataset is obtained from the population of Portugal. Consequently, the generalizability of our model to individuals from other geographical regions may be limited, as demographic and cultural differences could affect its applicability in different contexts.
- One notable limitation of our research is the absence of a dedicated validation set for assessing the performance of the machine learning algorithms on an unseen dataset. This omission restricts our ability to evaluate the model's generalization capability beyond the training data.

6 Conclusion and Future Work

In conclusion, this study introduces a machine learning framework tailored to predict risky alcohol consumption patterns among students. Through rigorous analysis, we have demonstrated the efficacy of machine learning in identifying students prone to hazardous drinking behaviors. Notably, our investigation highlights the superior performance of Random Forest and Decision Tree algorithms for our specific inquiry. Moreover, familial dynamics, particularly the quality of family relationships and parental cohabitation status, emerged as significant determinants of youths' alcohol consumption habits. We anticipate that this research will substantially contribute to mitigating risky alcohol consumption among students.

Looking ahead, we plan to expand the dataset to enhance model robustness and accuracy. Additionally, collaboration with domain experts, including psychologists and addiction specialists, will be pivotal in refining our approach and ensuring its relevance to real-world scenarios.

Acknowledgments

The authors extend their sincerest appreciation to their project and research supervisor, Muhammad Shafayat Oshman, Lecturer in

the Department of Electrical and Computer Engineering at North South University, Bangladesh. His unwavering support, precise guidance, and insightful advice have been instrumental throughout the experimental, research, and theoretical phases of this study, as well as in the preparation of this paper. Additionally, the authors gratefully acknowledge the Department of Electrical and Computer Engineering at North South University, Bangladesh, for their full funding provided for the publication of this research.

References

- [1] Ralph Hingson, Wenxing Zha, and Daniel Smyth. 2017. Magnitude and trends in heavy episodic drinking, alcohol-impaired driving, and alcohol-related mortality and overdose hospitalizations among emerging adults of college ages 18–24 in the United States, 1998–2014. In *Journal of studies on alcohol and drugs* 78, 4 (2017), New Jersey, 540–548. <https://doi.org/10.15288/jsad.2017.78.540>.
- [2] Christopher M Jones, Heather B. Clayton, Nicholas P. Deputy, Douglas R. Roehler, Jean Y. Ko, Marissa B. Esser, Kathryn A. Brookmeyer, Marci Feldman Hertz. Prescription opioid misuse and use of alcohol and other substances among high school students—Youth Risk Behavior Survey, United States, 2019. *MMWR supplements* 69, 1 (2020), 38. <http://dx.doi.org/10.15585/mmwr.su6901a5>.
- [3] Ali Ebrahimi, Uffe Kock Wiil, Thomas Schmidt, Amin Naemi, Anette Søgaard Nielsen, Ghulam Mujtaba Shaikh, and Marjan Mansourvar. 2021. Predicting the risk of alcohol use disorder using machine learning: a systematic literature review. In *IEEE Access* 9 (2021), 151697–151712. doi: 10.1109/ACCESS.2021.3126777.
- [4] K Witkiewitz, RZ Litten, and L Leggio. Advances in the science and treatment of alcohol use disorder. *Science advances* 5, 9 (2019). <https://doi.org/10.1126/sciadv.aax4043>.
- [5] Ruberu, Thanthirige Lakshika Maduwanthi and Kenyon, Emily A and Hudson, Karen A and Filbey, Francesca and Ewing, Sarah W Feldstein and Biswas, Swati and Choudhary, Pankaj K. 2022. Joint risk prediction for hazardous use of alcohol, cannabis, and tobacco among adolescents: A preliminary study using statistical and machine learning. In *Preventive Medicine Reports* 25 (2022), 101674. <https://doi.org/10.1016/j.pmedr.2021.101674>.
- [6] Singh, Advait, Vinayak Singh, Mahendra Kumar Gourisaria, and Ashish Sharma. 2022. Alcohol Consumption Rate Prediction using Machine Learning Algorithms. In *2022 OITS International Conference on Information Technology (OCIT)*. Bhubaneswar, India, 2022. 85–90. doi: 10.1109/OCIT56763.2022.00026.
- [7] Fabio Pagnotta and HM Amran. 2016. Using data mining to predict secondary school student alcohol consumption. Department of Computer Science, University of Camerino (2016), 1–9. doi:10.13140/RG.2.1.1465.8328.
- [8] Rijad Sarić, Dejan Jokić, and Edhem Čustović. 2020. Identification of alcohol addicts among high school students using decision tree based algorithm. In *CMBEBIH 2019: Proceedings of the International Conference on Medical and Biological Engineering*, May 16–18, 2019, Banja Luka, Bosnia and Herzegovina. Springer. 2020. 459–467. https://doi.org/10.1007/978-3-030-17971-7_69.
- [9] Shamala Palaniappan, Norhamreeza A. Hameed, Aida Mustapha, and Noor Azah Samsudin. 2017. Classification of alcohol consumption among secondary school students. *JOIV: International Journal on Informatics Visualization* 1, 4–2, (2017), 224–226. <http://dx.doi.org/10.30630/joiv.1.4-2.64>.
- [10] Kim, Suk-Young, Taesung Park, Kwonyoung Kim, Jihoon Oh, Yoonjae Park, and Dai-Jin Kim. A deep learning algorithm to predict hazardous drinkers and the severity of alcohol-related problems using K-NHANES. In *Frontiers in psychiatry* 12 (2021), 684406. <https://doi.org/10.3389/fpsyt.2021.684406>.
- [11] Marcon, Grasiela, Flávia de Ávila Pereira, Aline Zimmerman, Bruno Castro da Silva, Lisia von Diemen, Ives Cavalcante Passos, and Mariana Recamonde-Mendoza. Patterns of high-risk drinking among medical students: A web-based survey with machine learning. In *Computers in Biology and Medicine* 136 (2021), 104747. <https://doi.org/10.1016/j.compbiomed.2021.104747>.
- [12] Aliaksandr Amialchuk, Onur Sapci, and Jon D Elhai. Applying machine learning methods to model social interactions in alcohol consumption among adolescents. *Addiction Research & Theory* 29, 5 (2021), 436–443. <https://doi.org/10.1080/16066359.2021.1887147>.
- [13] Mohammad H Afzali, Matthew Sunderland, Sherry Stewart, Benoit Masse, Jean Seguin, Nicola Newton, Maree Teesson, Patricia Conrod. 2019. Machine-learning prediction of adolescent alcohol use: A cross-study, cross-cultural validation. *Addiction* 114, 4 (April. 2019), 662–671. <https://doi.org/10.1111/add.14504>.
- [14] Abid Ali Awan. 2023. An Introduction to SHAP Values and Machine Learning Interpretability. Retrieved November, 2023 from <https://www.datacamp.com/tutorial/introduction-to-shap-values-machine-learning-interpretability>.
- [15] Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 4768–4777. oi: 10.5555/3295222.3295230