

Online Technical Assessment Classification Based on Personalization

Lim Zheng Yu ^{#1}, Liyana Shuib ^{#2}

[#] Faculty of Computer Science and Information Technology, Universiti Malaya

50603 Kuala Lumpur, Malaysia

¹ u2102809@siswa.um.edu.my

² liyanashuib@um.edu.my

Abstract – Online technical assessments are slowly becoming the thing of the future, especially after COVID-19. Institutions are conducting more online assessments, but this is not a viable solution as traditional assessments have always been a disadvantage for most students. There is an urgency to curate a multilabel classification model to suggest suitable types of online assessments for students taking technical subjects like programming. The data is collected using an online survey form and it is targeted toward university students, asking questions about demographics, online platform preferences, learning object preferences, online instructional strategies, assessments preferences, and learning style tests. After data cleaning, only 314 entries remain. Six machine learning models were then trained, and their performances were compared on various metrics. The results show that the baseline model of gradient boosting is the best, returning a Hamming loss of 0.228, ROC AUC of 0.693, F1-Score, precision, and recall with a macro average of 0.83, 0.82, and 0.84 respectively, and an accuracy of 0.238. Another ensemble technique used is random forest, which returned similar results, which may be the core reason for the better performance.

Keywords – Online technical assessments; Multilabel classification; Gradient boosting; VAK; Learning styles; Personalization; Programming; Coding.

I. INTRODUCTION

Online assessments have become increasingly commonplace throughout the globe, with the COVID-19 pandemic confirming this trend. The reason is the pandemic has significantly hindered the progress of education, forcing institutions to shift their teaching to online platforms (Abduh, 2021), such as Microsoft Teams and Google Classroom. To further prove it, MOOC registrations are becoming more and more prevalent, with enrolments sometimes going up to three or four times higher than the previous year in the period of 2012 to 2022 (Papadakis, 2023). This in turn means that eventually, students who are in the field of Computer Science, Engineering, and other tech-related fields will have a lot of their technical subjects assessed online as well.

Here comes the problem, a lot of existing programming certification examinations, such as the Sun's Java Certification Examination and Novell's certification examinations are multiple-choice (Roberts & Verbyla, 2003), serving no significant purpose in evaluating programming students. Traditional assessments are also standardized, making them "indirect and inauthentic" (Bailey, 1997), meaning that they are effective in measuring students' capabilities at one point in time but ineffective in informing student progression (Dikli, 2003).

The VAK (Visual-Auditory-Kinaesthetic) Learning Styles Model (Barbe et al., 1979) describes that everyone is different in showcasing their skills and absorbing information. Visual learners typically learn through observation; therefore, they learn effectively by using graphs, posters, and other visual tools. Auditory learners learn by listening, and hence, they gather information through discussions, lectures, and stories. Finally, kinesthetic learners learn best by doing, and therefore, learn through physical activities and touch.

According to this model, it would be unwise to assume everyone is similar in terms of their abilities. The same can be said for students' performances in assessments. Not everyone does well in certain exam formats; some students may be excellent in memorization, while some display exceptional critical thinking skills instead. In the context of technical subjects, assessments would be futile if all programming tests were conducted in essay or multiple-choice format. Certain formats do not properly assess students, and it is certainly unfair to assume that a student who performs well in assessments that involve heavy memorization is better than a student who does badly in the same assessment. Hence, in this context, it is suggested that it would be fairer to evaluate every student differently in online technical assessments.

II. LITERATURE REVIEW

Throughout the years, many studies have been conducted about learning objects, learning styles, online students' performance, programming test performance, and so on. Yet, there is very limited research on the amalgamation of said subjects, meaning that to date, there may not be a classification model for online technical assessments based on students' personalities. Discussed below are some relevant literature.

Most online assessments are in the form of multiple-choice questionnaires (Kim et al., 2008), likely because they can be easily created and evaluated by teachers and lecturers to provide a score for a student. However, this form of assessment has been proven to be only useful in evaluating the lower levels of cognitive skills by many researchers (Nowicki & Jones, 2005). It cannot be deduced that the evaluation scores from multiple-choice assessments can sufficiently prove a student's capabilities in the subject matter as such online assessments may not indicate the achievement of higher levels of learning objectives (Zlatović et al., 2015).

It is proposed by Lee and Wu (2022) that the necessary adjustment of teaching methods to conform to students' personalities may lead to better learning effectiveness. Some factors to consider are students' personality traits, behavioral styles, and learning strategies. This study aims to find the correlation between these learning traits and learning effectiveness using association rule mining and sequential pattern mining. The data was collected from 56 students in programming courses by the Department of Information Management, National Chung Cheng University from September 2018 to January 2019. To predict the learning effectiveness, they used different datasets on decision trees and ANN and used AUC (area under the curve), precision, and recall as their metrics. The results show that across all datasets used, ANN is generally better in predicting learning effectiveness. From this prediction, it is shown that personality traits, range of progress, and problem-solving strategies were able to contribute to predicting learning effectiveness, but also shows bias whenever academic scores are introduced as a variable because every student is unique in retaining and absorbing information.

Garcia et al. (2005) realized the importance of profiling students based on their learning styles in computer-based education. He highlighted that an adaptive and personalized online education platform is highly sought-after to accommodate various types of learners. The goal of their research is to identify students' learning styles based on the Felder-Silverman model using Bayesian Networks. Eventually, they achieved 100% accuracy for the understanding dimension, 80% for the perception dimension, and 80% for the processing dimension. However, the input of their research is too few, a mere 10 computer science and engineering students, and hence, the results of the research should be taken with a grain of salt.

Ulloa-Cazarez et al. (2018) conducted an intriguing study to predict online students' performance using genetic programming. The dataset consisted of 245 students' grade records and log records. Their study was conducted to predict the final grades of online students based on their grades from the beginning stages of the course. The results of the genetic programming model returned a Mean Absolute Residuals (MAR) of 14.38 and Median Absolute Residuals (MdAR) of 9.80, while the results of the linear regression model were MAR of 17.49 and MdAR of 16.53. The genetic programming model was significantly better than the linear regression model by a

margin of 90% confidence level on the training data, and 99% confidence level on the test data.

Cetinkaya et al. (2023) did research on 600 secondary students in Turkey to recognize students who were talented in computer programming such that they could be given exposure to programming fields that are considered their strengths. Besides the spatial test that was assigned to the students, data was also retrieved from Code.org to assist with model training. Support vector machine (SVM), decision trees (DT), k-nearest neighbours (KNN), and quadratic discriminant (QD) were the machine learning algorithms used to classify students' programming skills. The results show that Cubic SVM was the best in accurately classifying their abilities, achieving 94.8% accuracy, followed by fine DT classifiers (89.0%), QD (84.3%), and medium KNN (80.8%). The best-performing algorithm of SVM also did well in other metrics: 91.5% Kappa, 93.6% precision, 94.8% recall, and 94.1% F-Score. Yet, this study was unable to solve our problem of classifying online technical assessments.

Dema (2021) invited 16 female and 24 male students enrolled in the program Bachelor of Engineering in Information Technology from the Royal University of Bhutan to conduct his study. He noticed that not only many freshmen were struggling with their C programming module, but the faculty was also ineffective in teaching it, leading to high rates of failure. The study aims to determine Neil Fleming's VARK learning styles of these first-year students and to convey the findings to the faculty such that improvements can be made in teaching. It was found the most dominant learning style is kinaesthetic, constituting 42% of the students. Unfortunately, there were many limitations to the study, including no predictive modeling being employed, no personalization methods are being used, proper sampling methods were not utilized, and hence, the findings cannot be generalized for all students in the tech-related fields. Only descriptive statistics was used to analyze the findings and the study has a lot of wasted potential.

The research of Ocepek et al. (2013) is aimed at curating an adaptive learning recommender system for 272 undergraduate students from the University of Ljubljana to suggest the appropriate multimedia type (or learning object) based on their learning styles. It incorporates a variety of learning styles: Kolb's Rancourt's, VAK, and the hemispheric to profile a student's characteristics in detail. This was achieved using a multi-target regression tree, predicting the usage of multimedia types. By combining the learning styles, they were able to determine the characteristics of students in a robust and unbiased manner. The findings from the tree were, that textual content with color content proved to be the best multimedia type, and hence all subjects should include textual content in their teachings. The findings are not definitive, however, since there may be scenarios where the educator prepares only video and audio lectures. Other findings based on the learning styles are that students who are non-visual, assimilators, left-hemispheric, and empirical are the types of people who prefer audio learning material, on the other hand, visual, rational, non-

right-hemispheric and divergers are least associated with audio materials. The limitations of this research are that the results of the tree were not evaluated and are only used to inform findings, and their classification targets are on multimedia types instead of assessment types.

Seyal et al. (2015) pioneered a study relating students' Kolb learning style and their performance in programming courses. The targeted students were first-year Bachelor of Internet Computing students. The students' first-year programming concept module results were obtained and assigned a Kolb's Learning Style Questionnaire to identify their learning style. The study suggests that the students were mostly of the convergers and assimilators learning type, their learning styles were found to have an impact on their programming test scores, and gender played a role in their programming test scores. Most importantly, the relationship between programming results and learning styles was concluded using the Chi-square test of independence, with the study showing a difference between convergers and assimilators on their programming results: 39% of convergers passed as opposed to 15% who are assimilators. It is stated that the limitations of their study include having a small sample size, and the conducted survey was only at a single point in time, eventually leading up to the problem of generalization due to the specifics of the settings.

Nandakumar et al. (2014) implied a need to improve student assessments. To achieve this, they created an adaptive grading methodology whereby the performance of students is assessed based on the duration each student used to answer questions using rough set theory. The study also includes the generation of classification rules for knowledge discovery about the students' performance based on the E-Assessment data in both predefined and adaptive grading. It was found that using the adaptive grading using Johnson's rough set classification algorithm, the accuracy went up to 97.2% as opposed to the predefined grading accuracy of 92% for 250 samples.

Buyrukoglu (2022) recognized the importance of classifying students not based only on their examination scores. Other parameters such as online learning activities, online attendance, and assignment submission time are needed to determine students' learning styles. He introduced a Student Classification Rate (SCR) formula which he deduced that may have an impact on improving students' learning. His data consisted of the e-learning activities of 80 students with features such as extra study support, family study support, weekly study time, going out with friends, and free time after school. Results show that the best mean absolute error rate is performed on the SCR + General feature set using the random forest data mining model with a value of 0.7.

III. PROBLEM STATEMENT

From the previous section, it is apparent that there is limited study on our proposed topic: classification of online technical assessments based on students' VAK learning styles. The relevant literature mostly discusses the domain of learning

styles or learning objects with classification models or online performance. They focus on performance instead of classifying the most suitable types of assessments for technical subjects and do not include personalization methods.

This study has three main objectives:

1. To develop a classification model of online technical assessments based on personalization.
2. To evaluate the performance of the classification model of online technical assessments based on personalization.
3. To develop a data product of the classification model of online technical assessments based on personalization.

IV. METHODOLOGY

The Cross Industry Standard Process for Data Mining, known as CRISP-DM, is the most widespread methodology for data-related projects, and it is used as the methodology for this data science project. It has six phases: business understanding, data understanding, data preparation, modeling, evaluation, and deployment. As shown in Figure 1, the flow is not one-way, rather, there are backward flows between phases as a lot of revisions are done in this methodology to clean the dataset to fit our requirements.

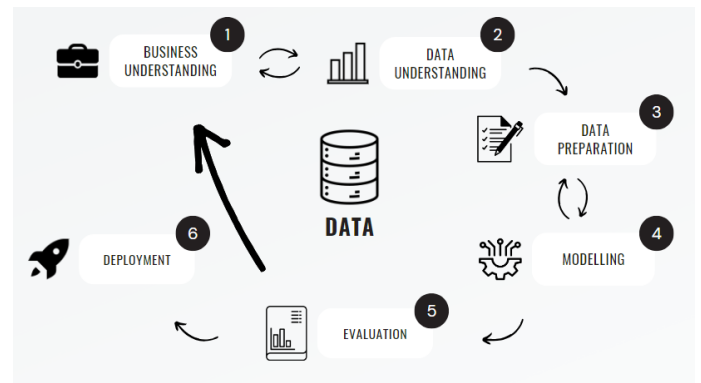


Figure 1: CRISP-DM Methodology

A. Business Understanding

The first phase of this methodology helps us identify our project objectives and requirements. It begins with collecting and obtaining data and understanding the dataset briefly to see if the dataset fits our use case. The data used for this project was collected using Google Forms from January 2021 to November 2023. The survey form contains four parts: personal details, learning object preference, learning style awareness, and learning style test.

Then, it is followed by determining our business objectives which were already identified in section III. Problem Statement, and finally to produce a project plan.

B. Data Understanding

The survey form has 103 questions, resulting in 103 variables in the initial dataset with 1052 rows and 836 missing values, a column 'dominant_VAK' is added as a new variable before proceeding with other works.

To further understand the dataset in detail, the 'ProfileReport' method from the 'ydata-profiling' package is invoked, and an HTML file describing every column is produced. Each column's description included the column name, number and percentage of distinct values, number and percentage of missing values, memory size, and a visualization showing its distribution based on the type of data. All columns in the dataset were originally read as text data. Hence, all the visualizations were word clouds.

Here is a summary of the columns of the dataset:

The first 13 columns of the dataset consisted of demographics and other information, such as Timestamp, Email Address, Gender, Level of Study, Field of Study, Household Income, Preferred Social Media Platform, and Preferred Communication Platform.

The following 13 columns are respondents' learning object preferences and they are required to answer from a scale of 'Not at All', 'Not Really', 'Undecided', 'Somewhat', and 'Very Much'.

The next 40 columns are online instructional strategies/assessment preferences based on the same scale.

The next column is targeted toward technical subjects, and respondents are required to manually input their online instructional strategies/assessment preferences. Some additional cleaning is needed on this column.

Columns 68 to 70 describe respondents' understanding of learning styles and the outputs are 'Yes' and 'No'.

Columns 71 to 100 are learning styles test questions. A scenario is provided, and respondents are required to select one response. These responses are then used to calculate the dominant_VAK (column 104) of the respondents.

Column 101 is a text field gathering feedback on the survey, and columns 102 and 103 are the faculty and department of the respondents respectively, both are text input fields.

After some cleaning in the data preparation phase, we revised the dataset and found out there were new issues. Some of the variables were greatly imbalanced as shown in Figure 2. There were also issues with the multiple selection questions in Google Forms. For example, survey respondents who chose YouTube as their preferred social media platform and survey respondents who chose both YouTube and Instagram as their selections would result in two different outputs. Eventually, 59 respondents chose only YouTube, 35 respondents chose both

Instagram and YouTube, and 41 respondents chose Facebook, Instagram, and YouTube when there is a total of $59 + 35 + 41 = 135$ respondents who chose YouTube from these 3 outputs. This resulted in many unique outputs and extra cleaning is needed. The Google survey form also accepts some text inputs from respondents, resulting in very inconsistent and dirty data for those columns.

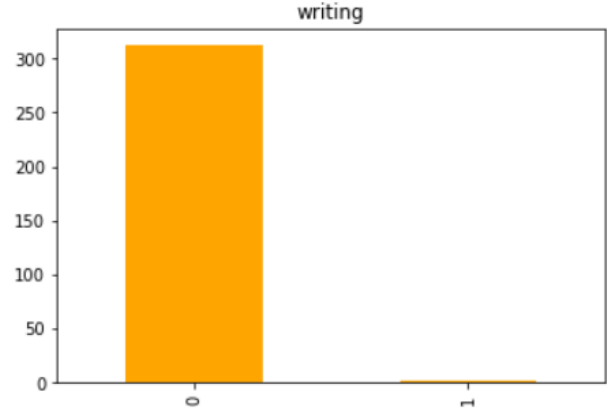


Figure 2: Imbalanced data in the 'writing' column

C. Data Preparation

This phase involves many processes, including cleaning, modifying, and transforming data into a format that would help us to fit the data into our models.

1) *Missing Values*: Some of the columns that contain missing values are 'Faculty', 'Department', 'Institutions', and 'Please share any comments or suggestions related to this issue. Thank You'. Coincidentally, they are not necessary for our analysis as they are deemed irrelevant to our study. These columns are therefore dropped, and no data imputation is needed.

2) *LabelEncoder*: As mentioned in section IV. Methodology subsection B. Data Understanding, all columns in the dataset are text data. The 'LabelEncoder' method from the scikit-learn library was used to create a function named 'LabelEncoding' which accepts the parameter 'column'. The function prints the distribution of the column using the 'value_counts' method before invoking 'LabelEncoder' and prints out the distribution after invoking. Here is an example on the 'Level of Study' column:

```
Undergraduate    900
Certificate/Diploma 134
Master           9
Postgraduate      5
PhD              3
Name: Level of Study, dtype: int64

4    900
0    134
1      9
3      5
2      3
Name: Level of Study, dtype: int64
```

Figure 3: Output of LabelEncoding(“Level of Study”)

3) *Extract Relevant Study Fields:* The ‘Field of study’ column is useful in identifying our intended target. There are also many unique values in this column. Hence, the values in this column are cross-checked with a list named ‘relevant_fields’, which contain terms such as ‘computer’, ‘engineering’, ‘tech’, ‘jurutera’, and ‘komputer’. If the values in the column matches any of the values in the ‘relevant_fields’ list, they are then kept into the data frame.

Then, the ‘Field of study’ column is passed into the ‘LabelEncoding’ function.

```
Computer Science/Information Technology    214
Engineering                               92
Medical Laboratory Technology             2
sijil sistem komputer                   2
Computer and Information Technology       1
Science and Technology Studies           1
medical lab technology                   1
Chemical engineering                     1
Name: Field of study, dtype: int64

1      214
3       92
4        2
7        2
2         1
5         1
6         1
0         1
Name: Field of study, dtype: int64
```

Figure 4: Output of LabelEncoding(“Field of study”)

We only have a subset of the original dataset to work with left, only 314 rows remain.

4) *Cleaning Multiple Choice Selections:* To solve the problem mentioned in section IV. Methodology subsection B. Data Understanding, the function ‘one_hot_encode_multiple_choice’ is created. This function produces a set which contains all unique values in the column. The set is type casted as a list, and the list is then iterated through to create new columns indicating whether that entry has that value. If the column of that entry has a substring that matches the list, the newly created column will have a value of 1, otherwise 0. This process is akin to one-hot-encoding.

Reddit_prefsocmed	Google Classroom_prefsocmed	Tiktok_prefsocmed	google meet_prefsocmed	Twitter_prefsocmed
0	0	0	0	0
0	0	0	0	1
0	0	0	0	1
0	0	0	0	0
0	0	0	0	1

Figure 5: Example Output of one_hot_encode_multiple_choice

5) *Cleaning Manual Inputs:* Some of the columns have data which were manually input by respondents. Example outputs from the ‘Online Instructional Strategies/ Assessment Preference for Technical Subjects’ column are ‘YOUTUBE VIDEO’, ‘YouTube video’, ‘live lecture and youtube video’, ‘Youtube Video and Live Lecture’. With the help of the Natural Language Toolkit (NLTK), the functions ‘remove_stopwords’ and ‘lemmatize_text’ are created.

‘remove_stopwords’ removes stop words if it exists in the text value. Stop words are commonly used words in any language (Ganesan, 2023). Words like “how” and “to” are considered stop words in English and they do not bring out much value. By removing them, we can focus on other words which have more valuable meaning.

Lemmatization is also done to process the text data. It is the process of grouping together words in different forms. For example, ‘better’ and ‘good’ or ‘rocks’ and ‘rock’ are words that will be classified together through lemmatization (GeeksforGeeks, 2024). This process is done in Python by instantiating a ‘WordNetLemmatizer’ object and it is processed for each word in the selected column.

6) *Imbalanced Data:* From a simple look of Figure 3, it is apparent that this column has imbalanced data. The easy workaround to the problem is to compress the minority levels into one level in a way which makes sense. In this context, the minority levels of ‘Master’, ‘Postgraduate’ and ‘PhD’ are compressed or grouped together into the level or label of ‘Postgraduate’.

D. Modelling

The 11 target variables used in the modelling phase are:

- _Written assignment_
- _Case Study_
- _Real Time Online Exam_
- _Individual Project/Assignment_
- _Group Project/Assignment_
- _Online Quiz/Test – MCQ_
- _Online Quiz/Test – Essay_
- _Online Quiz/Test – Open Book_
- _Peer Review Assessment Live Presentation_
- _Recorded Presentation_
- _Portfolio_

The name of the mentioned variables begins with the prefix “6. Online Instructional Strategies/Assessment”. This makes the problem a multilabel classification problem, with 1 stating that the particular online technical assessment is suitable and 0 meaning unsuitable.

The ‘train_test_split’ method from scikit learn is used to split the final processed data frame into train and test sets, with a test_size of 0.2 and random_state of 2024. The shape of the input data is 314 rows and 115 columns, while the output data is 314 rows and 11 columns.

Six machine learning algorithms from scikit learn were used in this study, namely gradient boosting (GB), Extreme Gradient Boosting (XGBoost or XGB), random forest (RF), support vector machine (SVM), Gaussian Naïve Bayes (GNB), and logistic regression (LR). Their respective functions were created to initialize the models and their predictions at the same time, with some accepting parameters such as learning rate. All of their hyperparameters are not tuned and hence only the default values are used.

E. Evaluation

The more common evaluation metrics used are Area under the ROC Curve (AUC), F1-score (Macro average), precision, and recall. A niche metric is also adopted called Hamming loss. According to the scikit-learn documentation, the Hamming loss is the fraction of labels that are incorrectly predicted. Hence, the results of the Hamming loss are within the range of 0 to 1, and the lower the Hamming loss, the better the performance of a model. Figure 6 below shows an example of how the Hamming loss works against accuracy.

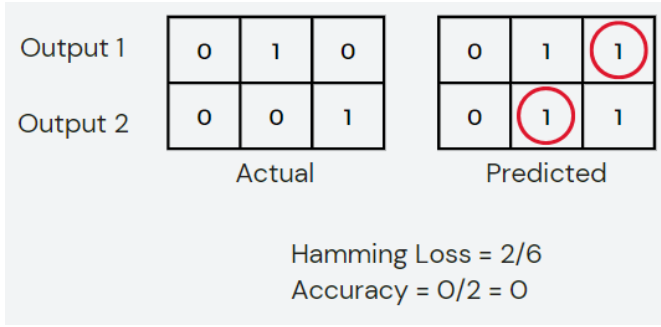


Figure 6: Hamming loss versus Accuracy

In Figure 6, we can see that the results for Hamming loss is 2/6 as there are 2 out of 6 labels which are incorrect. On the other hand, accuracy greatly punishes such inaccuracy, labelling the entire string of outputs to be incorrect even if only one of the labels are incorrect. This is an issue for this study as the outputs are as long as 11 targets, which may result in poor scores for the accuracy metric.

All the mentioned metrics are available in the scikit-learn library.

F. Deployment

The gradient boosting algorithm is found to be the best performing model out of the six, with random forest in close second. The gradient boosting model is then exported in a pickle file into a Streamlit web application and to be deployed as the data product.

Streamlit, an open-source Python web application framework for data applications is used due to its low-code solution that requires no prior experience in front-end and back-end.

A. Model Performance

The evaluation of the models using Hamming loss, ROC AUC, Macro F1-score, precision, recall, and accuracy are shown in the tables below.

Algorithm	Hamming loss	ROC AUC	Macro F1-Score
GB	0.228	0.693	0.83
XGB	0.258	0.655	0.81
RF	0.227	0.687	0.83
SVM	0.242	0.636	0.82
GNB	0.553	0.576	0.34
LR	0.276	0.669	0.79

Table 1: Model Evaluation using Hamming loss, ROC AUC, and Macro F1-Score

Algorithm	Precision	Recall	Accuracy
GB	0.82	0.84	0.238
XGB	0.79	0.83	0.254
RF	0.81	0.85	0.238
SVM	0.79	0.86	0.206
GNB	0.81	0.30	0.048
LR	0.80	0.78	0.159

Table 2: Model Evaluation using Precision, Recall, and Accuracy

Five-fold cross validation is done for all metrics to ensure that the results of the model are generalizable and valid. Below are some figures which display the results of gradient boosting from cross validation (note: Hamming loss is not an option in the 'scoring' parameter of the 'cross_validate' method).

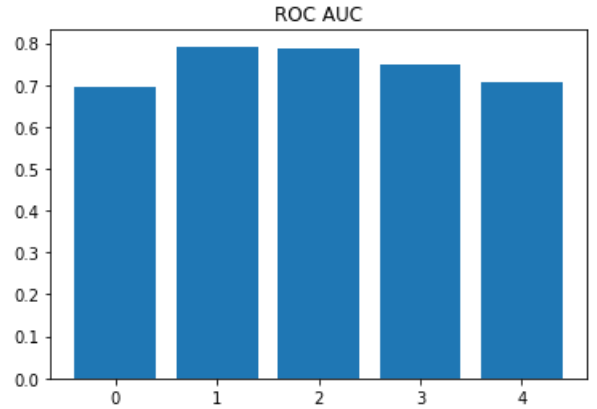


Figure 7: Plot of Five-fold Cross Validation on Gradient Boosting's ROC AUC

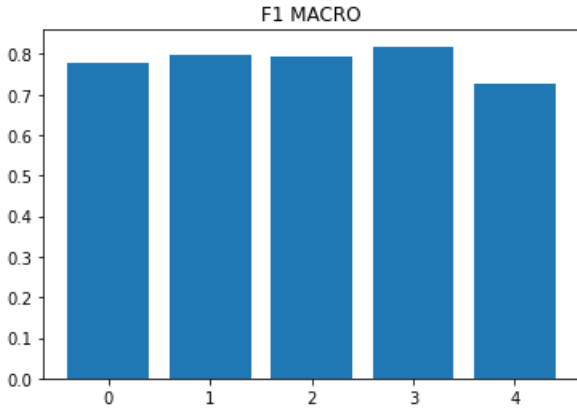


Figure 8: Plot of Five-fold Cross Validation on Gradient Boosting's Macro F1-Score

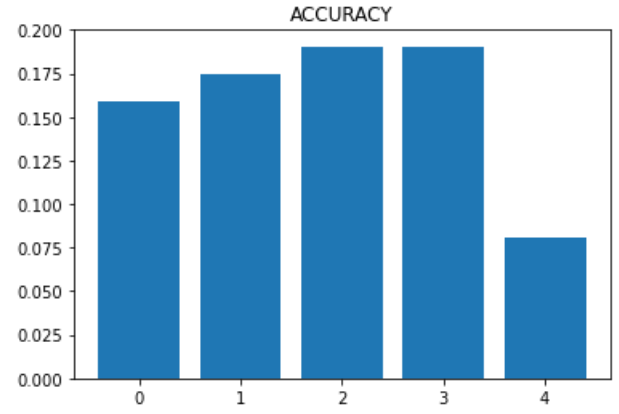


Figure 11: Plot of Five-fold Cross Validation on Gradient Boosting's Accuracy

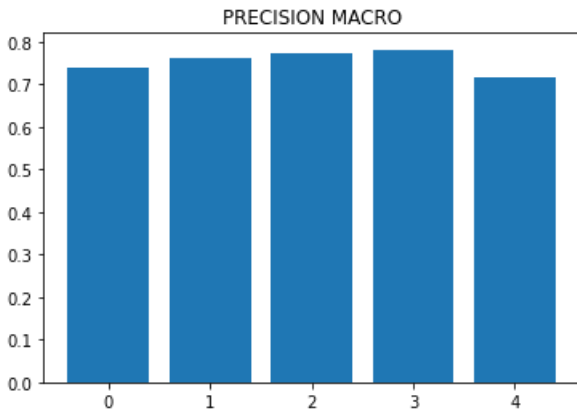


Figure 9: Plot of Five-fold Cross Validation on Gradient Boosting's Macro Precision

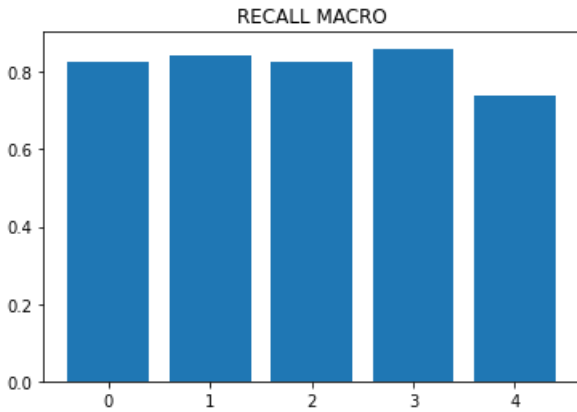


Figure 10: Plot of Five-fold Cross Validation on Gradient Boosting's Macro Recall

To further ensure the robustness of the gradient boosting model and its results reliability, the confusion matrices of all outputs are plotted. For each confusion matrix, the columns are the predicted labels, and the rows are the actual labels. Hence, from left to right, top to bottom, the quadrants each represent true negative, false positive, false negative, and true positive.

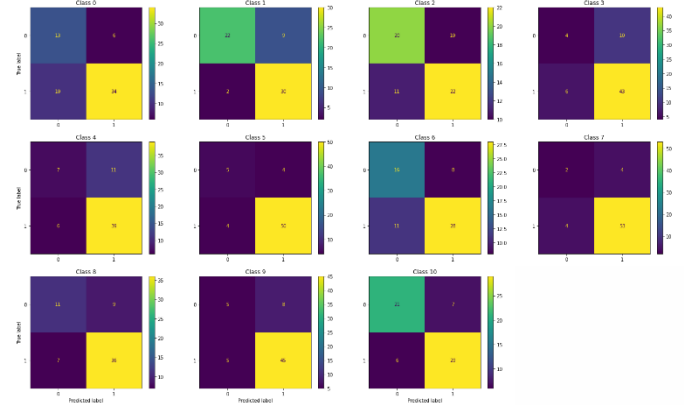


Figure 12: Confusion Matrices of All Outputs of Gradient Boosting Model

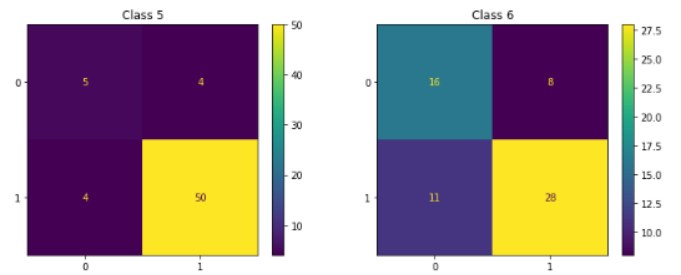


Figure 13: Confusion Matrices of Output Class 5 and 6 (Sixth and Seventh) of Gradient Boosting Model

In figure 13, the confusion matrix for class 5 displays that the counts of each quadrant are 5, 4, 4, 50 while for class 6 it is 16, 8, 11, 28. The accuracy for these classes are $55/63 = 0.873$ and $44/63 = 0.698$ respectively.

B. Web Application

This simple yet functional data product consists of two web pages.

1) *Introduction Page*: This page provides a brief description of the VAK Learning Styles. It explains the characteristics of each type of learner.



Figure 14: Introduction Page

2) *Prediction Page*: The prediction page is essentially an input form, with the top section describing the form's structure. The following sections are parts of the form, which enables users to input their details.

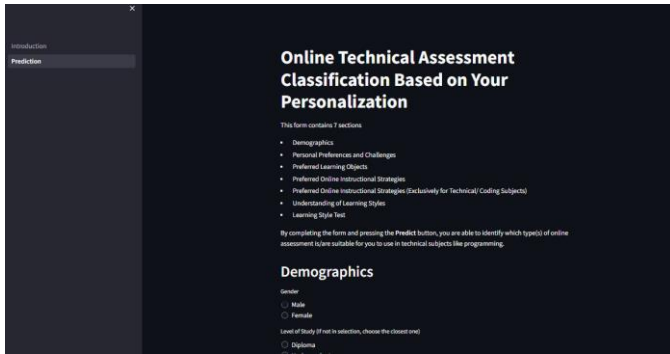


Figure 15: Prediction Page

In the same page, users can see the output by clicking the “Predict” button after completing the form. The output will suggest the suitable types of online technical assessment for the user based on their input information.

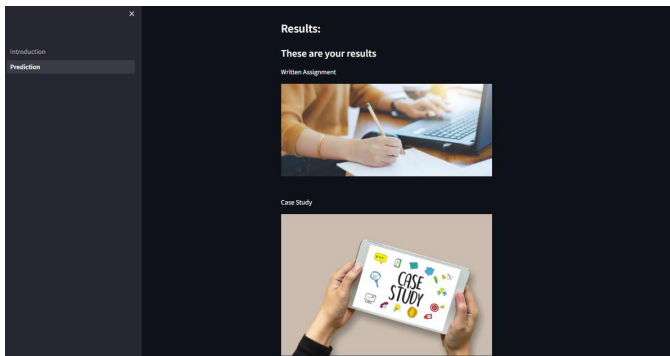


Figure 16: Output from Prediction Page

VI. DISCUSSION

The gradient boosting and the random forest models are very similar in performance, and it is decided that GB is better with the ROC AUC having a more significant margin. The performance is then ranked as follows: XGB, SVM, LR, and the worst performing by a wide margin GNB.

The reason the gradient boosting and random forest models perform better is because they are both ensemble methods, which have a good reputation of improving performance in the research world and in other practical applications (Bentéjac et al., 2021; Natekin & Knoll, 2013). Based on the amount of features the dataset has, it is highly possible that the complex relationships between features are difficult to handle and ensemble-type algorithms do quite well in managing this. It is worthy to note that random forest algorithms also utilise feature importance (Couronné et al., 2018), which may be imperative to this use case given the sheer amount of variables to work with.

For this study, the metrics Hamming loss, ROC AUC, and F1-Score are given more importance. Hamming loss is emphasized based on the explanation in section IV. Methodology subsection E. Evaluation whereby accuracy may be unsuitable in assessing the performance of the models. ROC AUC and F1-Score are used as they are more robust evaluation metrics than precision and recall.

The plots of the five-fold cross-validation in Figures 7 to 11 show that GB performs well on unseen data for all metrics except accuracy. This goes to show that GB is a robust model, and its performance scores are reliable. It also further concretizes the stance that accuracy should not be used in this study.

Further observation on Figure 12 shows that the distribution of most output classes is imbalanced, with more extreme cases in classes 5, 7, and 9 inclining towards an output of 1. Regardless, the confusion matrices show that the dataset is not overfitted into the model as there are still predictions in the false negative and false positive quadrants, leading to a more reliable accuracy score for each output label.

It should be mentioned that there are flaws in this study. The Google Forms survey was distributed mainly around Universiti Malaya, and some other public universities in Malaysia. Hence, the generalizability of the results is not certain. It would be a major improvement if similar future research may be conducted under different demographics to further concrete the findings of this study.

VII. CONCLUSION

A. Outcome

In conclusion, gradient boosting and random forest algorithms do well on this dataset to classify suitable online technical assessments for students. Their results rounded up to 2 decimal places are Hamming loss of 0.23, ROC AUC of 0.69, and F1-Score of 0.83.

B. Limitations and Improvements

This study although revised many times, is not perfect as there are still other optimization works that can be done. The models did not have their hyperparameters tuned. It is suggested to perform GridSearchCV or genetic algorithms to search for the optimal hyperparameters.

There was an attempt to resolve the imbalanced data issue, yet it is still insufficient. Some technique suggestions are resampling techniques through oversampling and undersampling, and data augmentation. Dimensionality reduction and feature selection should also be explored due to the large number of variables that were at disposal. Finally, other algorithms such as ensemble methods can also be explored in more detail to improve the results of this study.

C. Future Work

Future researchers on the topic may delve into how other learning styles models will affect the classifications of the online technical assessments. They may also propose other relevant target variables based on the trend of their time. Hopefully soon, we may be able to see such recommender systems being implemented in universities to help improve the learning experience of university students.

ACKNOWLEDGMENT

I would like to express my utmost gratitude to my project supervisor, Associate Prof. Dr. Nor Liyana bt Mohd Shuib for her dedicated and invaluable guidance throughout the development of this data science project. Her expertise and encouragement have been instrumental in enhancing the quality of this work.

I would also like to extend my many thanks to my coursemates Janice Chong See Wai and Jasmine Chong See Yan, whom I have worked with under the same lecturer, for their unconditional help from time to time.

REFERENCES

- Abduh, M. Y. M. (2021). Full-time online assessment during COVID-19 lockdown: EFL teachers' perceptions. *Asian EFL Journal*, 28(1.1), 26-46.
- Bailey, K. M. (1997). *Learning About Language Assessment: Dilemmas, Decisions, and Directions*. Heinle & Heinle.
- Barbe, W. B., Swassing, R. H., & Milone, M. N. (1979). Teaching Through Modality Strengths: Concepts and Practices.
- Bentéjac, C., Csörgő, A., & Martínez-Muñoz, G. (2021). A comparative analysis of gradient boosting algorithms. *Artificial Intelligence Review*, 54(3), 1937-1967. <https://doi.org/10.1007/s10462-020-09896-5>
- Buyrukoğlu, S. (2023). A novel color labeled student modeling approach using e-learning activities for data mining. *Universal Access in the Information Society*, 22(2), 569-579. <https://doi.org/10.1007/s10209-022-00894-8>
- Çetinkaya, A., Baykan, Ö. K., & Kırız, H. (2023). Analysis of Machine Learning Classification Approaches for Predicting Students' Programming Aptitude. *Sustainability*, 15(17), 12917. <https://www.mdpi.com/2071-1050/15/17/12917>
- Couronné, R., Probst, P., & Boulesteix, A.-L. (2018). Random forest versus logistic regression: a large-scale benchmark experiment. *BMC bioinformatics*, 19, 1-14.
- Dema, K. (2021). Understanding Students' C Programming Language Learning Styles: A Case Study in College of Science and Technology. *Journal of Information Engineering and Applications*, 11(1), 7-14. <https://doi.org/https://doi.org/10.7176/JIEA/11-1-02>
- Dikli, S. (2003). Assessment at a distance: Traditional vs. alternative assessments. *Turkish Online Journal of Educational Technology-TOJET*, 2(3), 13-19.
- Ganesan, K. (2023, March 16). *What are Stop Words?* Kavita Ganesan, PhD. <https://kavita-ganesan.com/what-are-stop-words/>
- García, P., Amandi, A., Schiaffino, S., & Campo, M. (2005). Using Bayesian networks to detect students' learning styles in a web-based education system. *Proc of ASAI, Rosario*, 115, 126.
- GeeksforGeeks. (2024, January 2). *Python Lemmatization with NLTK*. <https://www.geeksforgeeks.org/python-lemmatization-with-nltk/>
- Kim, N., Smith, M. J., & Maeng, K. (2008). Assessment in Online Distance Education: A Comparison of Three Online Programs at a University. *Online Journal of Distance Learning Administration*, 11(1).
- Lee, P.-J., & Wu, T.-Y. (2022). Mining relations between personality traits and learning styles. *Information Processing & Management*, 59(5), 103045. <https://doi.org/https://doi.org/10.1016/j.ipm.2022.103045>
- Nandakumar, S. B., Geetha, V., & Thangasamy, S. (2014). A Rough Set Based Classification Model for Grading in Adaptive E-Assessment. *International Review on Computers and Software*, vol 9, 1169-1177.
- Natekin, A., & Knoll, A. (2013). Gradient boosting machines, a tutorial [Methods]. *Frontiers in Neurobotics*, 7. <https://doi.org/10.3389/fnbot.2013.00021>
- Nowicki, C., & Jones, E. L. (2005). A framework for extending online assessment. 2005 6th International Conference on Information Technology Based Higher Education and Training.
- Ocepek, U., Bosnić, Z., Šerbec, I. N., & Rugelj, J. (2013). Exploring the relation between learning style models and preferred multimedia types. *Computers & Education*, 69, 343-355.
- Papadakis, S. (2023). MOOCs 2012-2022: An overview. *Advances in Mobile Learning Educational Research*, 3, 682-693. <https://doi.org/10.25082/AMLER.2023.01.017>

- Roberts, G.H.B. and Verbyla, J.L.M. (2003). An Online Programming Assessment Tool. In Proc. Fifth Australasian Computing Education Conference (ACE2003), Adelaide, Australia. Conferences in Research and Practice in Information Technology, 20. Greening, T. and Lister, R., Eds., ACS. pp. 69-75.
- Seyal, A. H., Mey, Y. S., Matusin, M. H., Siau, N. H., & Rahman, A. A. (2015). Understanding students learning style and their performance in computer programming course: Evidence from Bruneian technical institution of higher learning. *International Journal of Computer Theory and Engineering*, 7(3), 241.
- Ulloa-Cazarez, R. L., López-Martín, C., Abran, A., & Yáñez-Márquez, C. (2018). Prediction of online students performance by means of genetic programming. *Applied Artificial Intelligence*, 32(9-10), 858-881.
- Zlatović, M., Balaban, I., & Kermek, D. (2015). Using online assessments to stimulate learning strategies and achievement of learning goals. *Computers & Education*, 91, 32-45. <https://doi.org/https://doi.org/10.1016/j.compedu.2015.09.012>

APPENDIX

USER MANUAL

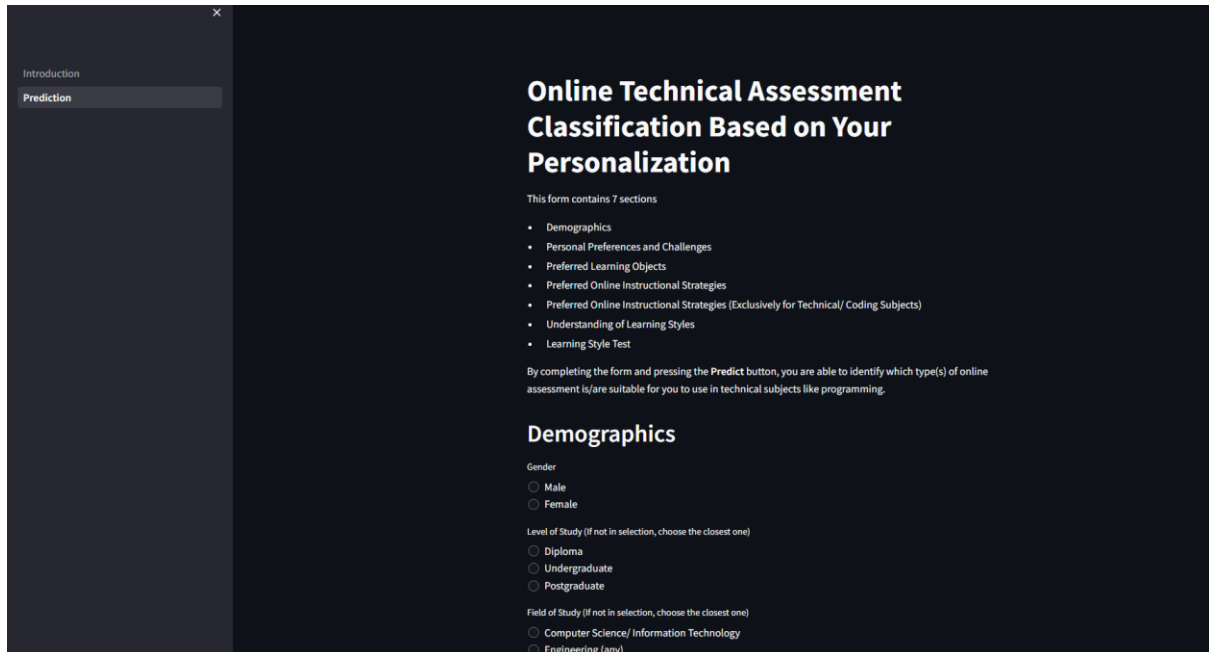
Running the Application

1. To deploy the “Online Technical Assessment Classification Based on Personalization” web application, you may download the source code from the link:
https://drive.google.com/drive/folders/1aSSnSUCw_3ucd2s29FogUpUjEfMbMndk?usp=sharing
2. Include the file into a directory.
3. In that directory, execute a command prompt.
4. Enter the command “streamlit run Introduction.py” and press Enter.
5. When the command executes, the home page of the application will appear from your browser



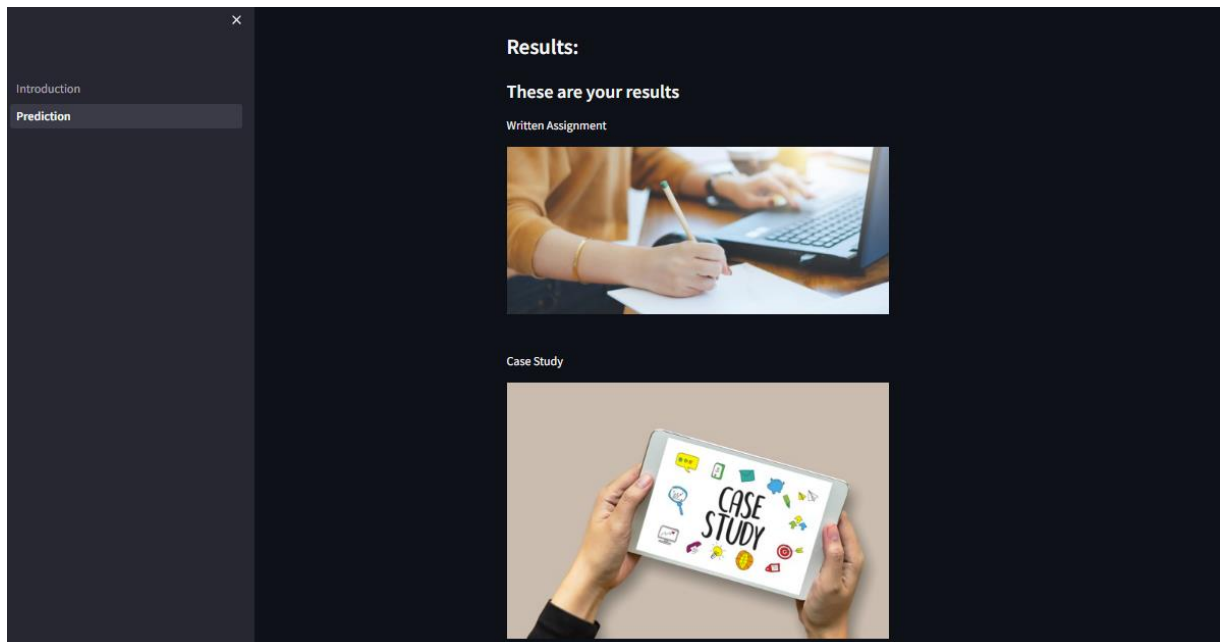
Using the Application

1. Once the application is up and running, select the “Prediction” tab on the left panel of the website.
2. The page will display questions about demographics, personal preferences and challenges, preferred learning objects, preferred online instructional strategies, understanding of learning styles, and learning style test. The top section will also provide a brief description of the layout of the form.



The screenshot shows a web application interface. On the left, a dark sidebar contains two tabs: 'Introduction' and 'Prediction', with 'Prediction' being the active tab. The main content area has a dark background. At the top, the title 'Online Technical Assessment Classification Based on Your Personalization' is displayed in white. Below the title, a list of 7 sections is provided: Demographics, Personal Preferences and Challenges, Preferred Learning Objects, Preferred Online Instructional Strategies, Preferred Online Instructional Strategies (Exclusively for Technical/ Coding Subjects), Understanding of Learning Styles, and Learning Style Test. A paragraph explains that by completing the form and pressing the 'Predict' button, users can identify suitable online assessment types for technical subjects like programming. The 'Demographics' section is expanded, showing radio button options for Gender (Male, Female), Level of Study (Diploma, Undergraduate, Postgraduate), and Field of Study (Computer Science/ Information Technology, Engineering (any)).

3. Fill in the form accordingly.
4. Once completing the form, click on the “Predict” button.



The screenshot shows the 'Results' page of the application. The left sidebar remains the same. The main content area displays the title 'Results:' followed by 'These are your results'. Under the heading 'Written Assignment', there is an image of a person's hands writing on a notepad next to a laptop. Below this, under the heading 'Case Study', there is an image of a person's hands holding a tablet that displays the words 'CASE STUDY' surrounded by various colorful icons.

5. Based on your input, the application will display and recommend what types of online technical assessment(s) suits you best.

6. To do another prediction, you may change the input in the forms, and click “Predict” again to see the results change.