



Subject Stream Prediction: A Machine learning Approach to Select the Suitable Subject Stream for Senior Secondary Students in Sri Lanka

K.G.Kaushalya Abeywardhane
(Reg. No.: MS21902420)
M.Sc. in IT

Supervisor: Ms.Anjalie Gamage

September 2022

Department of Graduate Studies and Research
Sri Lanka Institute of Information Technology

I certify that I have read this thesis and that in my opinion, on it is fully adequate, in scope and quality, as a thesis for the degree of Master of Science.

.....
Ms.Ajalie Gamage (Supervisor)

Approved for MSc. Research Project:

.....
Head/<Department >

Approved for MSc:

.....
Head – Graduate Studies

DECLARATION

I declare that this is my research idea, research paper and thesis do not incorporate without acknowledgment any material previously published or submitted for a degree or Diploma in any other university or institute of higher learning, and to the best of my knowledge and belief, it does not contain any material previously published or written by another person except where the acknowledgment is made in the text.

Sign: 

K.G.Kaushalya Abeywardhane

Date: 29.03.2023

ABSTRACT

Education is an important factor that measures the nation's wealth and directly affects the country's future development. According to the Sri Lankan government, free education provides to students at all levels up to the university level. General Certificate of Examination (Ordinary Level) – G.C.E.(O/L) and General Certificate of Examination (Advanced Level) – G.C.E.(A/L) are essential exams that complete senior secondary education. G.C.E.(A/L) is the examination that causes one to enter a university for higher education. According to the Sri Lankan education schemas, students happen to select one subject stream and related subjects relevant to that subject stream to continue their senior secondary education key stage 2. That selection is caused to the students' whole lives because students happen to face G.C.E.(A/L) from that subject stream. Most of the students have taken this decision according to the force of someone or comparing it with their own. I think it may be caused to break the senior secondary education key stage (2) in the middle or change the subject stream in the middle. These kinds of reasons affect to keep away the students from their target careers. From my point of view, students should pay attention to O/L results and their inborn talents, skills, and relevant working field that they hope for their job when selecting the subject stream for continuing their senior secondary education. I have developed a machine learning model to suggest the best subject stream based on the above features. The implemented model which is called the SubjectStreamPredict system predicts the best subject stream for students. As well as the implemented model suggest another suitable ten solutions including an appropriate career path according to the user's input values.

To implement the model, I have trained and tested four machine learning algorithms: K-Nearest Neighbors, Decision Tree, Random Forest, and Support Vector Machine Algorithm for the same data set. The Random Forest algorithm outperformed other algorithms and gave high accuracy (0.70). According to the analysis results I implemented my model using Random Forest Classifier algorithm and I improve the output generated from Random Forest by predicting more than one feature.

Keywords –Machine Learning Algorithm, Subject Stream, Prediction System

ACKNOWLEDGEMENT

My sincere thanks go to the authority of the Sri Lanka Institute of Information Technology (SLIIT) for pointing me to the correct path to conduct and complete this research project in my final year degree program M.Sc. in Information Technology.

I would pay my gratitude to the supervisor Ms. Anjalie Gamage for her valuable guidance and assistance in this work. Your feedback helped me narrow down my research topic, select one particular area, and develop it as the major part. It pushed me to think deeply about one major area, make a clear path to reach my research objectives, and complete the research successfully.

Always my parents were a big strength to me to finish my research up to this level. They were a giant shadow for me and helped me overcome the challenges that came my way. Also, my sister, who was my closest friend, helped me to successfully overcome those difficulties by sharing her knowledge with me when I faced difficulties.

Thank You very much my friends who are following the same course for helping me every time I am stuck on some issues. Your feedback also has been very important to me to succeed in this work.

TABLE OF CONTENTS

DECLARATION	ii
ABSTRACT.....	iii
ACKNOWLEDGEMENT	iv
TABLE OF CONTENTS.....	v
List of Figures	vii
List of Tables	viii
Chapter 1 Introduction	1
1.1 Background of The Studies.....	1
1.2 Problem Identification	1
1.3 Research Questions and Research Objectives.....	4
1.3.1 Research Questions	4
1.3.2 Research Objectives.....	4
Chapter 2 Literature Review	5
2.1 Related Works and Research Gap.....	5
Chapter 3 Methodology	8
3.1 Flow Diagram of the Proposed Model.....	8
3.2 Data Gathering	8
3.3 Data Preprocessing.....	12
3.4 Feature Engineering and Feature Scaling	21
3.5 Building the Model	23
3.5.1 Machine Learning	23
<i>Support Vector Machine Algorithm (SVM)</i>	24
3.5.2 <i>Tools and Techniques</i>	24
<i>Jupyter Notebook</i>	24
<i>Special Functions</i>	26
3.6 Experiments and Results.....	26
3.6.1 Analyzing the Dataset	26
3.6.2 Predictive Analysis	34
<i>Confusion Matrix</i>	34
3.6.3 <i>Selecting Best Independent Variables</i>	35
3.7 Main User Interfaces.....	37
Interfaces for Predicted Output.....	41
3.7 Improve the Predicted Outcomes.....	42
3.8 Testing and Evaluation	43
3.8.1 Cross Validation.....	43
3.8.2 Execution Time for the Final Model.....	45

3.8.3 Evaluation of the User Friendliness of the Model	45
Chapter 4 Conclusion.....	48
Chapter 5 References	49
Appendix.....	51
Appendix 1: Circulars 2008/12.....	51
Appendix 2: Circulars 2016-13s	54
Appendix 3: Circulars 2016-20e.....	75
Appendix 4: Google Form	77
78	
Appendix 5: Evaluation Results of Random Grid	83

List of Figures

Figure 1.1 Performance of Candidate (G.C.E.(O/L) and G.C.E.(A/L) 2019,2020).....	2
Figure 1.2 Pass and Failed Percentages of Students year wise.....	3
Figure 3.1:Flow Diagram of Proposed Model	8
Figure 3.2:Used Data Gathering Methods	9
Figure 3.3:Percentage of Subject Stream	12
Figure 3.4:Percentage of Satisfaction/Unsatisfaction	12
Figure 3.5:Feature Categorization in Excel	13
Figure 3.6:Subject wised Grade Count	19
Figure 3.7:Percentage Responses of Extra Curriculum Activity -Category 1	19
Figure 3.8:Percentage Responses of Extra Curriculum Activity - Category 2	20
Figure 3.9:The Percentage Responses of Extra Curriculum Activity - Category 3	20
Figure 3.10:The Percentage Responses of Extra Curriculum Activity - Category 4	20
Figure 3.11:The Percentage Responses of Inborn Talents.....	21
Figure 3.12:Percentage of Different Kinds of Job Fields	21
Figure 3.13: Dataset that not applying the Label Encoding.....	22
Figure 3.14:Applying the Label Encoding.....	22
Figure 3.15:After Applying the Label Encoding	23
Figure 3.16:Subject wised Total Count for Extra Curriculum Activity 1	31
Figure 3.17:Confusion Matrix Generated Results	35
Figure 3.18:Accuracy Score of Four ML Algorithms.....	35
Figure 3.19:Applying Only Main Subjects in O/L	36
Figure 3.20:Applying Only Extra Curriculum Activities & Talents.....	36
Figure 3.21:Applying Both Results and Skills.....	37
Figure 3.22 The Accuracy of the Independent Features	37
Figure 3.23:Menu Bar of the Web App	38
Figure 3.24:User Input Getting for Main Subject Grades.....	38
Figure 3.25:User Input Getting for Skills and Talents.....	39
Figure 3.26:Analysing User Inputs part 1	40
Figure 3.27:Analysing User Inputs part 2.....	40
Figure 3.28:Predicted Output 1.....	41
Figure 3.29:Predicted Output 2.....	42

List of Tables

Table 3.1: Independent Feature Categorization	12
Table 3.2: Unique Features of each Streams.....	13
Table 3.3: Main Subjects for O/L	15
Table 3.4: Optional Subjects for O/L.....	15
Table 3.5: Marks Range and Relevant Grades.....	18
Table 3.6: Total Calculations for Independent Features	26
Table 3.7: Extra Curriculum Activity Features.....	29
Table 3.8: Subject wised Total Count for Extra Curriculum Activity 1	30
Table 3.9: Subject wised Total Count for Extra Curriculum Activity 2	31
Table 3.10: Subject wised Total Count for Extra Curriculum Activity 4	31
Table 3.11: Subject wised Total Count for In born Talents.....	32
Table 3.12: Each Grade Count for O/L Main Subjects.....	33
Table 3.13: Confusion Matrix Results	34
Table 3.14: User Feedback Summarizing	46