# STUDENTS PERFORMANCE IN EXAMS

*using classification Algorithms*

<u>Developed by:</u> *Kola Venkata Rohith*

## 1.INTRODUCTION

The concept of machine learning is something born out ofthis environment. Computers can analyze digital data to find patterns and laws in ways thatis too complex for a human to do. The basic idea of machine learning is that a computer can automatically learn from experience (Mitchell, 1997).Although machine learning applications vary, its general function is similarthroughoutits application .The computer analyzes a large amount of data, and finds patterns and rules hidden in the data. The creation ofrules from data is an automatic process, and itis something that continuously improves with newly presented data.

To understand the influence of various factors like economic, personal and social on the students performance. Every educational system consists of an examination system through which the qualities and abilities ofthe students are assessed by giving them grades and positions.

"Examination tests the efficiency of the education provides ,we shall need to test what it is, studentscan do, rather what he knows."It means, the ultimate objectives of the examination is to measure the performance level of the students and without this, we cannot know what the students attain from their educational system. So examination are

doing the job offinal appraisal of student achievement.

"Examination are conducted to test the ability of the student and find out if he or she has reached a certain standard of academic learning and knowledge. They scrutinize and measure the student 's capabilities against skill in answering a question under the condition imposed by the examiner." The aim of the examination is to evaluate the ability of candidates. Examination supplies a tangible proof of fitness of a student for high class or a particular professional course. According to dictionary of education(1998)Examination is defined as under: "It is test of a person's knowledge or proficiency in which he or she is required to answer questions or perform tasks." "There is an usual degree unanimity as to the serious limitation of our examination system." Teachers,students and general public are all a like of the view that it is neither provide an accurate test of the scholastic attainment of the students, nor it is designed to assess the intellectual development." "The inappropriate structure of subjective marks and individual difference in evaluating the answers, dishonest invigilating staff, wrong marking of scripts etc are the main factors which affect student's performance in examination." Above factors are creating obstacles to measure the real performance of the students. As a result many students fail in the examination. It is fact that failure of student is not their fate rather as there are some problems, which becomes hurdle in their successes. Deserving students are deprived to get actual performance in spite oftheir good. Extensive efforts have been made in order to predict student performance for different aims, like: detecting at risk students, assurance of student retention, course and re-source allocations, and many others. This research aims to predict student perform-ance to engage distinct students in researches and innovative projects that could improve universities reputation and ranking nationally and internationally. Additionally, in most researches that were aimed to classify or predict, researchers used to spend much efforts just to extract the important indicators that could be more useful in constructing

reasonable accurate predictive models. They will either use features ranking algorithms or will look at the selected features while training the data-set on different machine learning algorithms . The economic success of any country highly depends on making higher education more affordable and that considers one of the main concerns for any government.One ofthe factors that contributes to the educational expenses is the studying time spent by students in order to graduate. For example, the loan debt of the American students has been increased due to the failure of many students in getting graduated on time . Higher education is provided for free to the students in Iraq by the government. Yet, failing of graduating on time costs the government extra expenses. To avoid these expenses, the government has to ensure that the student graduate on time. Machine learning techniques can be used to forecast the performance of the students and identifying the at risk students as early as possible so appropriate actions can be taken to enhance their performance. One of the most important steps when using these techniques is choosing the attributes or thedescriptive features which used as input to the machine learning algorithm. The attributes can be categorized into GPA and grades, demographics, psychological profile, cultural, academic progress, and educational background . This research introduces two new attributes thatfocus on to the effect of using the internet as a learning resource and the effect ofthe time spent by students on social networks on the students' performance. Four machine learning techniques, fully connected feed forward Artificial Neural Network, Naïve Bayes, Decision Tree, and Logistic Regression, have been used to build the machine learning model. The dataset used to build the models is collected from the students at the College Of Humanities during 2015 and 2016 academic years using a survey and the student's grade book. The dataset has the information of students. The activities of this research include feature engineering to create the students dataset, data collecting, data preprocessing, creating and evaluating four machine learning models, and finding the best model and analyzing the results

## 1.1 Overview

There are many risks related to students performance in exams, for the students and institutions like parent level of education. The analysis of risk in performance of students is need to understand what is the meaning of risk . In addition, "Examination are conducted to test the ability of the student and find out if he has reached a certain standard of academic learning and knowledge. "The inappropriate structure of subjective marks and individual difference in evaluating the answers, dishonest invigilating staff, wrong marking of scripts etc are the main factors which affect student's performance in examination."

As a result many students fail in the examination.. One of the factors that contributes to the educational expenses is the studying time spent by students in order to graduate. For example, the loan debt of the students has been increased due to the failure of many students in getting graduated on time,Failing of graduating on time costs the government extra expenses.

To avoid these expenses, the government has to ensure that the student graduate on time. Machine learning techniques can be used to forecast the performance of the students and identifying the at risk, students as early as possible so appropriate actions can be taken to enhance their performance.

The model has been built using data form to predict performance of students .Algorithms havebeen used to build the proposed model:RandomForest, Decision Tree, Navie bayes,KNN,SVM By using the algorithm a Flask model has been implemented and tested.,The results also has been discussed .

**Inferences would be :**

1. How to imporve the students performance in each test ?

2. What are the majorfactors influencing the test scores ?

3. Effectiveness oftest preparation course?

4.Otherinferences

## 1.2 Purpose

Our aim from the projectis to make use of pandas, seaborn libraries from python to extractthe libraries for machine learning forthe prediction on students performance

secondly, to learn how to analyse the scores impacted based on different variables which include gender,race ,lunch,test preparation course etc…. Each column is picked and has been analysed how they affect the scores. For easy understanding I have used graphs and plots.After all visualization is the best way to understand from multiple machine learning algorithms and withdrawing the conclusions.

## 2.LITERATURE SURVEY

Much research has been done in the area of educational data mining where a predictive model is built to forecast the performance of students to identify the at risk students. This problem can be considered a hard problem because the performance depends on many characteristics related to the students. These characteristics can be categorized into student's Grade demographics, psychological profile, culture, academic progress, and educational background .

The student's Grade is the most important attribute used to predict the performance. The writing score can represent the real value forthe future educational and career possibilities and progression. In addition, the academic potentials can be evaluated by the student writing score.

This research introduces two new attributes that focus on using descriptive features related to the internet and social network usage and their effect on the performance. On the other hand, many machine learning and data mining techniques have been used to predict the students' performance such as: K-Nearest Neighbor (KNN); Decision Tree (DT); Random Forest (RF); setc..Classification is the most commonly applied data mining technique, This approach frequently employs Decision tree Classification Algorithm. In Classifier, a training set is used to build the model as the Classifier which can classify the data items into its appropriate classes. Atest setis used to validate the model.

## 2.1 Existing Problem

Student retention is an important issue in education. While intervention programs can improve retention rates, such programs need prior knowledge of students performance That is where performance prediction becomes important. The usage of machine learning to predict either the student performance or the student dropout is a commonly found subject in academic literature. Dropout prediction in virtual learning, or e-learning is a particularly common focus in such studies, due to both high dropout rates and easily available data . Areas outside of virtual learning are also common contexts where dropout or performance predictions are used for research. The purpose ofthe research ofthese studies varies.

In some of them, the aim is to find the best method for prediction .However, predicting student performance instead of student dropouts is more related with this thesis, Study was made by researches in india and analyzed the usage of machine learning,The previous models have high time complexity and space complexity whereas this model is constrained with the lot of advantages and with a higher accuracy than any other model already proposed. In this model we used Machinelearning algorithm named LogisticRegression,RandomForest,Decision Tree,Navie Bayes,KNN,SVM which give an accuracy more then 80%ofthe previously predicted problem The pattern is similar in most of these studies. First, different algorithms are applied to a data set to build prediction models. Then, predictions made by these models are compared using common evaluation criteria, such as accuracy, precision, and recall. Feature selection is also a commonly compared criteria. However, what these studies are missing is a more comprehensive comparison between distinct approaches such as method selection and feature engineering. This is the part where this thesis can introduce a new approach. By comparing the effectiveness of different processes used in machine learning, this thesis can provide insightinto the more efficient ways to improve predictions in student performance.

## 2.2 Proposed Solution

**Machine Learning(SVM)**

svm which is support vector classifier algorithm in machine learning methods which efficiently performs both Regression and classification tasks.A SVM is an ensemble technique capable of performing classification .In this project it predicts the grade of a student as an output.

And also we have created an UI using the flask for the student performance prediction,this UI will allow the users to predict the student performance very easily and the User interface is User friendly not at least one complication in using the interface,and it can be used just by entering some necessary details into the UI in real time it will give the predicted value, like if the student written exam then it will predict the scores based on student performance in test and it also includes the educational background and student economical status for their success or else for failure in graduating ,Then this model predicts the valuewith student details in realtime.
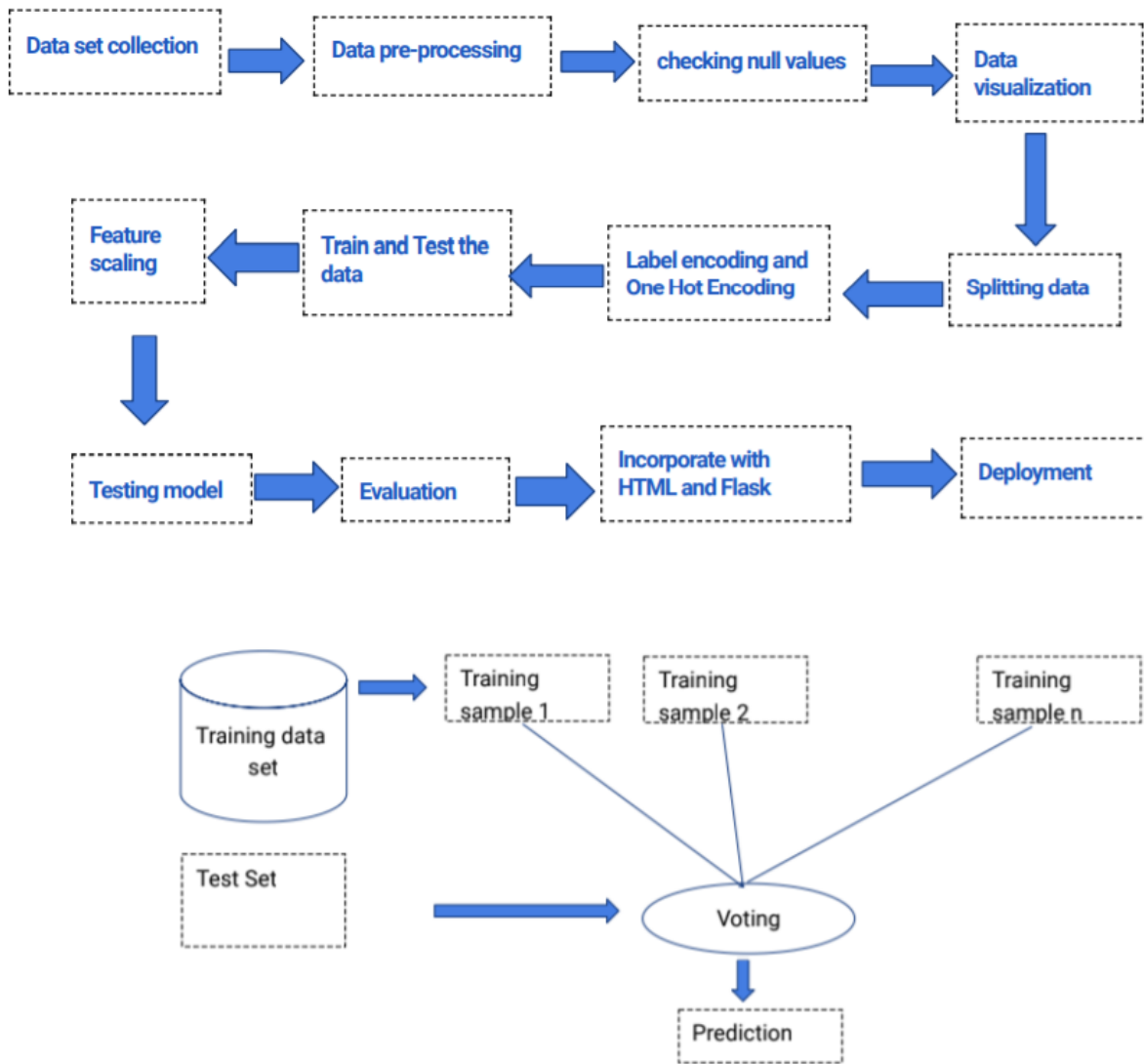
## 3.THEORETICAL ANALYSIS

While selecting the algorithm that gives an accurate prediction we gone through lot of algorithms which gives the results abruptly accurate and from them we selected only one algorithm forthe prediction problem thatis SVM and also LogisticRegression, it assumes that the presence of a particularfeature in a class is unrelated to the presence of any otherfeature. thats how the prediction work great with the SVM and LogisticAlgorithms.The peculiarity ofthis problem is collecting the student details in realtime and working with the prediction atthe same time, so we developed an userinterface forthe people who'll be accesssing forthe loan status prediction.Accuracy is defined as the ratio ofthe number of samples correctly classified by the classifierto the total number of samples for a given test data set. The formula is as follows

**Accuracy=TP+TN/TP+TN+FT+FN**

 Atfirst we gotlike lot of worst accuracies because we tried lot of algorithms forthe best accurate algorithm ,finally after all ofthat we tried the best suitable algorithm which gives the prediction accurately is Classification algorthims and developed itto use as a realtime.

# 3.1 Block Diagram

```
Data set collection → Data pre-processing → checking null values → Data visualization
                                                                         ↓
Feature scaling ← Train and Test the data ← Label encoding and One Hot Encoding ← Splitting data
      ↓
Testing model → Evaluation → Incorporate with HTML and Flask → Deployment
```

```
Training data set → Training sample 1    Training sample 2    Training sample n
Test Set →                              Voting
                                          ↓
                                      Prediction
```

# 3.2 Software Desining

● jupiter Notebook Environment

● SpyderIDE

● Machine LearningAlgorithms

● Python(pandas,numpy,sklearn,matplotlib,seaborn)

● HTML

● Flask

We developed this student performance prediction in exams by using the python language which is interpreted and high level programming language and using Machine learning algorithms. For coding we used jupyter notebook environment of Anaconda distributions and the Spyder,itis an integrated scientific programming in the python language.

For creating an user interface for the prediction we used the Flask.It is a micro web framework written in python.It does not required particular libraries.where pre-existing third-party libraries provide common functions,and a scripting language to create awebpage is HTML by creating the templates to use in the functions ofthe Flask and HTML.

## 4.EXPERIMENTAL INVESTIGATION

In this project the data set we used is derived from https://www.kaggle.com/spscientist/student-performance-in-exams/data It contains 1000 original performance data of students with 15 attributes including Pass_status attributes.

After that checked for any null values i.e 'zero',if have any null values are filled in by means of mode interpolation, and the duplicate or meaningless attributes are deleted, Datapreprocessing ,LabelEncoding is done to convert the textual data to integer values and the feature scaling is done to ignore the outliers,Declaring passmark and based on students percentage assigning suitable grades then after training and testing the

model using best model of regression algorthims and SVM .Attributes were shown below in the screenshot of the data set we used.

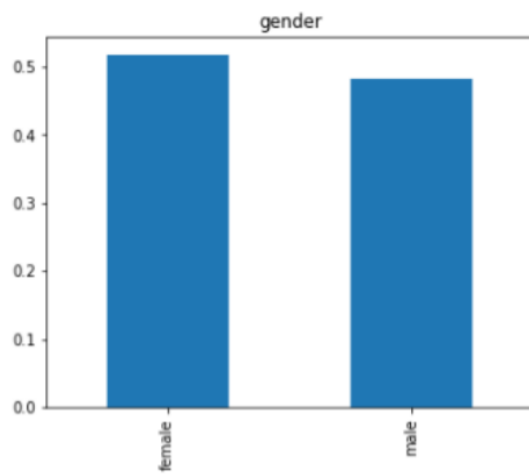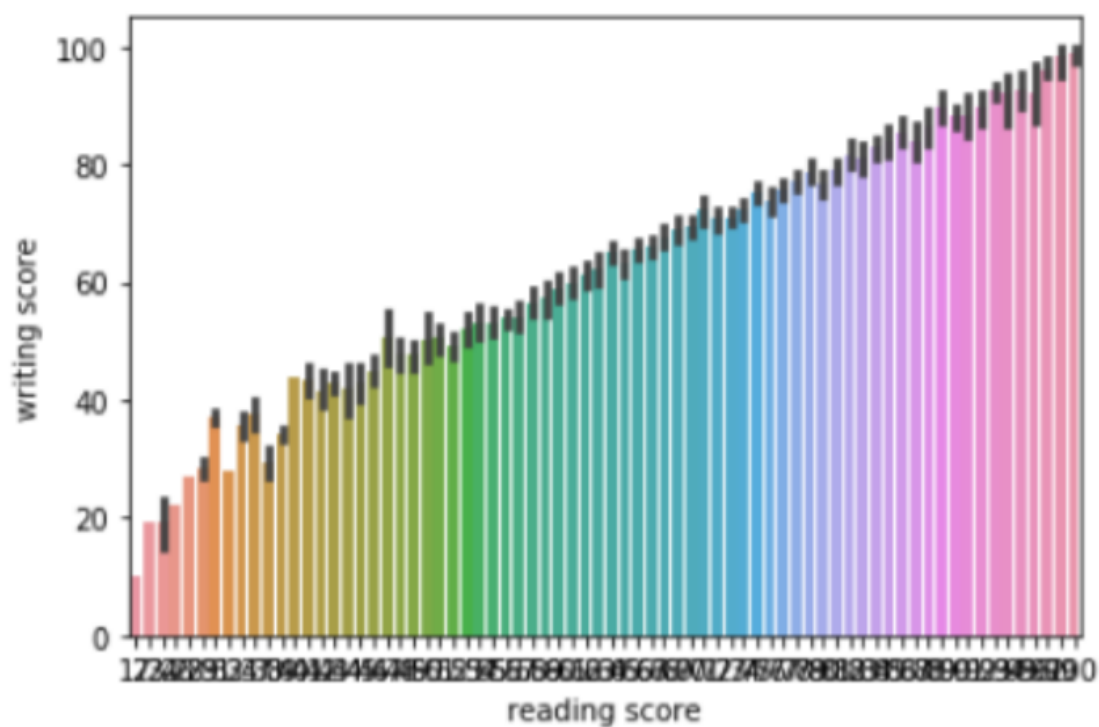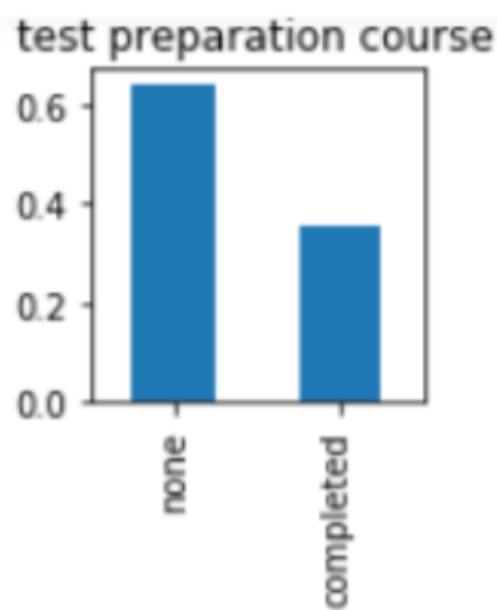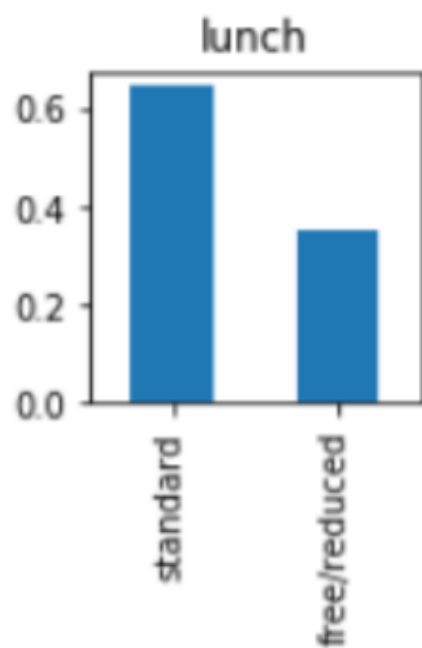| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | gender | race/ethnicity | parental level of education | lunch | test preparation course | math score | reading score | writing score |
| 2 | female | group B | bachelor's degree | standard | none | 72 | 72 | 74 |
| 3 | female | group C | some college | standard | completed | 69 | 90 | 88 |
| 4 | female | group B | master's degree | standard | none | 90 | 95 | 93 |
| 5 | male | group A | associate's degree | free/reduced | none | 47 | 57 | 44 |
| 6 | male | group C | some college | standard | none | 76 | 78 | 75 |
| 7 | female | group B | associate's degree | standard | none | 71 | 83 | 78 |
| 8 | female | group B | some college | standard | completed | 88 | 95 | 92 |
| 9 | male | group B | some college | free/reduced | none | 40 | 43 | 39 |
| 10 | male | group D | high school | free/reduced | completed | 64 | 64 | 67 |
| 11 | female | group B | high school | free/reduced | none | 38 | 60 | 50 |
| 12 | male | group C | associate's degree | standard | none | 58 | 54 | 52 |
| 13 | male | group D | associate's degree | standard | none | 40 | 52 | 43 |
| 14 | female | group B | high school | standard | none | 65 | 81 | 73 |
| 15 | male | group A | some college | standard | completed | 78 | 72 | 70 |
| 16 | female | group A | master's degree | standard | none | 50 | 53 | 58 |
| 17 | female | group C | some high school | standard | none | 69 | 75 | 78 |
| 18 | male | group C | high school | standard | none | 88 | 89 | 86 |
| 19 | female | group B | some high school | free/reduced | none | 18 | 32 | 28 |
| 20 | male | group C | master's degree | free/reduced | completed | 46 | 42 | 46 |
| 21 | female | group C | associate's degree | free/reduced | none | 54 | 58 | 61 |
| 22 | male | group D | high school | standard | none | 66 | 69 | 63 |
| 23 | female | group B | some college | free/reduced | completed | 65 | 75 | 70 |
| 24 | male | group D | some college | standard | none | 44 | 54 | 53 |
| 25 | female | group C | some high school | standard | none | 69 | 73 | 73 |

**5.FLOWCHART**

## 6.RESULT

In this, the SVM algorithm is used to predict its performance, and compared with anothe five machine learning methods namely the logistic regression random forest ,decision tree,and KNN, Navie baye.The obtained results are displayed in Table below. The results show that, the performance of logistic regression and SVM have performance than that of KNN regression, decision tree,Navie bayes. but the SVM still performs the best,with an accuracy of 93%, higher than the logistic regression with an accuracy of 90% and Decisioon Tree with an accuracy of 89%.The prediction model based on SVM indicating that the model has strong ability of generalization. In this paper,the SVM is used to predict performance,and compared with three machine learning algorithms namely Decisiontree,Random forest,KNN etc.The obtained results are displayed in table below

| S.NO | Algorithms used | Accuracy(%) |
|---|---|---|
| 1. | Logistic regression | 90 |
| 2. | Random Forest | 84 |
| 3. | Decision Tree | 89 |
| 4. | KNN | 74 |
| 5. | Navie Bayes | 87 |
| 6. | SVM | 93 |

Accuracy ofthis model is best out ofthe two algorithms at 90%and 89%. The below plots are showing us from dataset is 50% students in the dataset are female. Around 300% of the students in the dataset are group c. Around 200% students in the dataset are high school. Around 60% students have lunch standard . Around 60% of the students are not with Test preparatin course and also shows barplot between writing score and reading score from that we can considered thatreading score is in huge increase than writing score.

## gender

female
male

## race/ethnicity

group C
group D
group B
group E
group A

## parental level of education

some college
associate's degree
high school
some high school
bachelor's degree
master's degree

## lunch

standard · free/reduced

## test preparation course

none · completed

writing score vs reading score

# 7.ADVANTAGES AND DISADVANTAGES

## Advantages:

● SVM's are very good when we have no idea on the data.
● Works well with even unstructured and semi structured data like text, Images and trees.
● The kernel trick is real strength of SVM. With an appropriate kernel function, we can solve any complex problem.
● Unlike in neural networks, SVM is not solved for local optima.
● It scales relatively well to high dimensional data.
● SVM models have generalization in practice, the risk of over-fitting is less in SVM.
● SVM is always compared with ANN. When compared to ANN models, SVMs give better result.

## Disadvantages:

● Choosing a "good" kernel function is not easy.
 ● Long training time for large datasets.
● Difficult to understand and interpret the final model, variable weights and individual impact.
 ● Since the final model is not so easy to see, we can not do small calibrations to the model hence its tough to incorporate our business logic.
 ● The SVM hyper parameters are Cost -C and gamma. It is not that easy to fine-tune these hyper-parameters. It is hard to visualize their impact.

## 8.APPLICATIONS

1. Applications of machine learning cover a wide range of areas.
 2. Search engines use machine learning to better constructrelations between search phrases and web pages.
 3. By analyzing the content ofthe websites, search engines can define which words and phrases are the mostimportantin defining a certain web page, and they can use this information to return the mostrelevantresults for a given search phrase.
 4. In addition, machine learning can be used to understand the kind of products a customer might be interested in.
5.The computer processes data and learns to identify this data, and then uses this knowledge to make decisions aboutfuture data.
6.The increase in data has made these applications more effective, and thus more common in use.

## 9.CONCLUSION

The success of machine learning in predicting student performance relies on the good use ofthe data and machine learning algorithms. Selecting the right machine learning method forthe right problem is necessary to achieve the bestresults. However,the algorithm alone can not provide the best prediction results.

Feature engineering,the process of modifying data for machine learning, is also an importantfactorin getting the best prediction results. In this project,the SVM algorithm is adopted to build a UI modelfor predicting the grades.And the results are compared with other algorithm thatis logistic,decision tree classifier,knn others

- Many software tools are available for SVM implementation.
- SVMs are really good for text classification.
- SVMs are good at finding the best linear separator. The kernel trick makes SVMs non-linear learning algorithms.
- Choosing an appropriate kernel is the key for good SVM and choosing the right kernel function is not easy.
- We need to be patient while building SVMs on large datasets. They take a lot of time for training.

This research study explored the relationships of student grades from the examinations such as math score,reading score, writing score.The results indicate that predictive model is high degree of unexplained variable.

## 10.FUTURE SCOPE

This research has certain limitations that must be noted. There was not an access to a dedicated student data set, and the study relies on public data sources. In addition, bothdata sets were small, having less than thousand records. A research that has access tomore comprehensive data may offer more conclusive results.

Another area that future research can improve is the variety of the machine learning methods. This research used linear regression, decision trees, and the naïve Bayes classification. Other methods, such as clustering and artificial neural networks can be used to have a better understanding ofthe importance of method selection.

Final area that can be improved is the process of feature creation. Since the data is limited,the amount of feature modification that can be made is also limited. Both data sources used in this research consists of a single table, and custom variables were created using

variables from the same table. With a more comprehensive data set that spans multiple tables, there will be more potential to create new custom variables,while keeping in mind thatthe more a custom variable is,the more difficultitis to interpretthe relation between it and the dependent variable.

## 11.BIBILOGRAPHY

Ethem Alpaydin. 2004. Introduction to Machine Learning. Cambridge, MA

ElafAbuAmrieh, Thair Hamtini, and Ibrahim Aljarah. 2016. Mining educational data to predict student's academic performance using ensemble methods. International Journal of Database Theory andApplication 9(8), 119-136.

S. K. Card and J. Mackinlay. 1997. The structure ofthe information visualization design space. In: Proceedings ofthe 1997 IEEE Symposium on Information Visualization.

Paulo Cortez andAlice MariaGonçalves Silva. 2008. Using data mining to predict secondary school student performance. In: Proceedings of 5thAnnual Future Business Technology Conference, Porto, 5-12.

G. Dekker, M. Pechenizkiy, and J. Vleeshouwers. 2009. Predicting students drop out:A case study. In: Educational Data Mining 2009, 41-50.

Pedro Domingos. 2012.Afew usefulthings to know about machine learning. Communications oftheACM 55(10), 78-87.

Wayne W. Eckerson. 2007. Predictive analytics. Extending the Value of Your Data Warehousing Investment. TDWI Best Practices Report 1, 1-36

IsabelleGuyon andAndré Elisseeff. 2003.An introduction to variable and

feature selection. Journal of Machine Learning Research 3, 1157-1182

M. J. Islam,Q. J. Wu, M.Ahmadi, and M.A. Sid-Ahmed. 2007. Investigating the performance of naive-Bayes classifiers and k-nearest neighbor classifiers. In: International Conference on Convergence Information Technology. IEEE, 1541-1546.

Tom M. Mitchell. 1997. Machine Learning. McGraw-Hil.

## APPENDIX

## HTML:

```html
<html>
<!--From https://codepen.io/frytyler/pen/EGdtg-->
<head>
    <meta charset="UTF-8">
    <title>ML API</title>
    <link href='https://fonts.googleapis.com/css?family=Pacifico' rel='stylesheet' type='text/css'>
    <link href='https://fonts.googleapis.com/css?family=Arimo' rel='stylesheet' type='text/css'>
    <link href='https://fonts.googleapis.com/css?family=Hind:300' rel='stylesheet' type='text/css'>
    <link href='https://fonts.googleapis.com/css?family=Open+Sans+Condensed:300' rel='stylesheet' type='text/css'>
    <link rel="stylesheet" href="../static/style.css">
    <style>

        .login {
            top: 20%;
```

```html
        }
    </style>
</head>
<body style ="background-image:
url('https://images.unsplash.com/photo-1496469888073-80de7e952517?ix
lib=rb-1.2.1&ixid=eyJhcHBfaWQiOjEyMDd9&auto=format&fit=crop&w=500&
q=60'); background-size: 100%
100%;background-repeat:no-repeat;background-size: cover;">
<h1>Students Performance in Exams</h1>
<form action ="/login" method ="post">
<p><label for ="gender"> choose a gender:</label>
<select name ="g">
<option  value = "female">female</option>
<option  value ="male">male</option>
</select></p>
<p><label for ="race/ethnicity"> choose a race/ethnicity:</label>
<select name ="re">
<option value = "group A">group A</option>
<option value = "group B">group B</option>
<option value = "group c">group C</option>
<option value = "group D">group D</option>
<option value = "group E">group E</option>
</select></p>
<p><label for ="parental level of education"> choose a parental level of
education:</label>
<select name ="edu">
<option value = "bachelor's degree">bachelor's degree</option>
<option value = "some college">some college</option>
<option value = "master's degree">master's degree</option>
<option value = "associate's degree">associate's degree </option>
<option value = "high school">high school</option>
<option value = "some high school">some high school</option>
```

```html
</select></p>
<p><label for ="lunch"> choose a lunch:</label>
<select  name ="d">
<option value = "standard">standard</option>
<option value = "free/reduced">free/reduced</option>
</select></p>
<p><label for ="test preparation course"> choose a test preparation
course:</label>
<select  name ="tpc">
<option value ="none">none</option>
<option value ="completed">completed</option>
</select></p>
<p>math score</p>
<p><input type = "text" name = "ms"/></p>
<p>reading score</p>
<p><input type = "text" name = "rs"/></p>
<p>writing score</p>
<p><input type = "text" name = "ws"/></p>
<p><input type = "submit" value = "click"/></p>
<p><b>{{y}}</b></p>
</form>
</body>
</html>
```
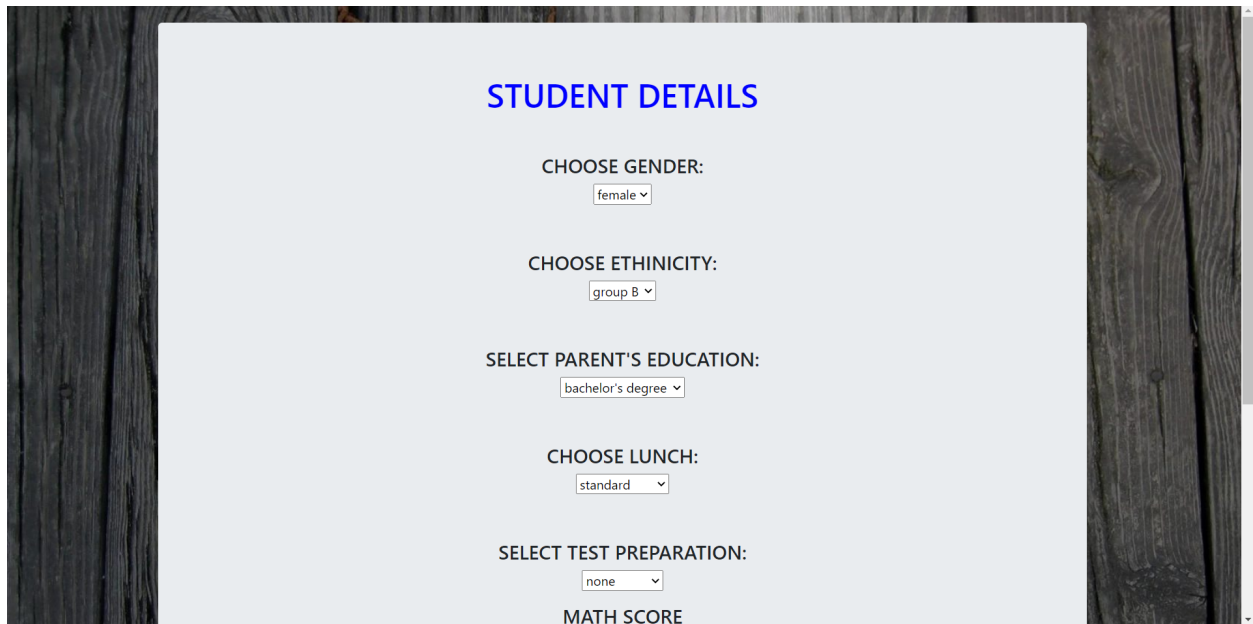
```
App.py
from flask import Flask,render_template,request
import pickle from joblib importload
from itertools import chain
app = Flask(__name__)
model =pickle.load(open('Grade.pkl','rb'))
sc=load('trans.save')
@app.route('/')
def hello_world():
    return render_template("index.html")
@app.route('/login',methods=["POST"])
deffunc2():
    x_test = [[str(x)for x in request.form.values()]]
    g = request.form['g']
    if(g=="male"):
        a1=1
    if(g=="female"):
        a1=0
    re = request.form['re']
    if(re=="groupA"):
        b1=0
    elif(re=="group B"):
        b1=1
    elif(re=="group C"):
        b1=2
    elif(re=="group D"):
        b1=3
    elif(re=="group E"):
        b1=4
    edu = request.form['edu']
    if(edu=="associate's degree"):
        c1=0
```

```python
elif(edu=="bachelor's degree"):
    c1=1
elif(edu=="high school"):
    c1=2
elif(edu=="master's degree"):
    c1=3
elif(edu=="some college"):
    c1=4
elif(edu=="some high school"):
    c1=5 d= request.form['d']
if(d=="standard"):
    d=0
if(d=="free/reduced"):
    d=1 tpc=request.form['tpc']
if(tpc=="none"):
    e1=1
if(tpc=="completed"):
    e1=0
rs=request.form['rs']
ms = request.form['ms']
ws = request.form['ws']
l = [[a1,b1,c1,d,e1,rs,ms,ws]]
prediction = model.predict(sc.transform(l))
if prediction == 0:
    my_prediction='A'
elif prediction ==1:
    my_prediction='B'
elif prediction ==2:
    my_prediction='C'
elif prediction ==3:
    my_prediction='D'
elif prediction ==4:
```

```python
        my_prediction='E'
    else:
        my_prediction = 'Re appear'
    return render_template("index.html",y=my_prediction)
if __name__ == '__main__':
    app.run(debug = True)
```

OUTPUT:

**SELECT PARENT'S EDUCATION:**

bachelor's degree ▾

**CHOOSE LUNCH:**

standard ▾

**SELECT TEST PREPARATION:**

none ▾

**MATH SCORE**

72

**READING SCORE**

72

**WRITING SCORE**

74

Submit

**Grade A**