

1. Classification vs Regression

Your goal is to identify students who might need early intervention - which type of supervised machine learning problem is this, classification or regression? Why?
Classification. Because it predicts binary label, which are yes and no.

2. Exploring the Data

Total number of students

395

Number of students who passed

265

Number of students who failed

130

Graduation rate of the class (%)

67.09%

Number of features

30

4. Training and Evaluating Models

Choose 3 supervised learning models that are available in scikit-learn, and appropriate for this problem.

I choose Neighbours-Based, SVM, and Decision Trees models.

Neighbors-Based

What are the general applications of this model?

Neighbours-Bases classification is a instance-based learning. Which does not do learning but just store the training data. When classify it will compute k point that is near to input and vote for major label.

What are its strengths and weaknesses?

Strengths is it is very quick to learn and easy for adding new point because no need to train model.

Weakness is prediction time is slow because it have to search the nearest point and if there are more and more training data it will get slower and slower.

Given what you know about the data so far, why did you choose this model to apply?

I choose Neighbours-Based classification because data point in this project is not large(395 students).

Fit model to the training data

Training set size	F1 score on training set	F1 score on test set
100	0.785714285714	0.705882352941
200	0.839160839161	0.80303030303
300	0.880898876404	0.780141843972

SVM

What are the general applications of this model?

SVM model is a represent of data point in space that can separate class label as wide as possible. When predicting data it map to the space and predict label that belongs to side of class label.

What are its strengths and weaknesses?

The strengths are effective in high dimension spaces, memory efficient and it can use kernel function.

The weaknesses are if number of features is much greater that number of sample it is likely to give poor performance. And SVMs are expensive operation for large dataset.

Given what you know about the data so far, why did you choose this model to apply?

I choose SVM because in this dataset there are 48 features after preprocess, which is considerably high compare to number of training set (395).

Fit SVM classification model to the training data

Training set size	F1 score on training set	F1 score on test set
100	0.833333333333	0.771241830065
200	0.85342019544	0.797297297297
300	0.876068376068	0.783783783784

Decision Trees

What are the general applications of this model?

It is a graph-like model called decision tree. One node of a tree is a decision to go to one branch. And go deeper until reach leaf node which is output.

What are its strengths and weaknesses?

Strengths are simple to understand the model because model can be interpret.

Weakness are it is easy to get overfitted, and it is difficult to create optimal decision tree model.

Given what you know about the data so far, why did you choose this model to apply?

I choose decision tree model because it is quite accurate and it is white box model which can be helpful when analyse how to deal with student.

Fit this model to the training data

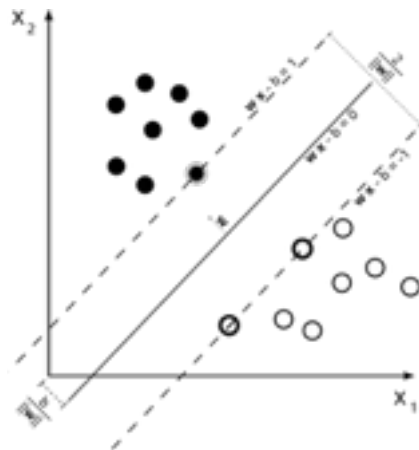
Training set size	F1 score on training set	F1 score on test set
100	1.0	0.755555555556
200	1.0	0.904347826087
300	1.0	0.627118644068

5. Choosing the Best Model

Based on the experiments, I choose SVM as the best model because, from F1-score, it performs best in both training set and test set. Neighbours-based model's F1-score is very near to SVM's but since the data has high dimension (48 features after preprocess) and training data may not be large enough to train neighbours-based model.

SVM uses more training time (0.011 s) than neighbours-based (0.001 s) and decision tree (0.004 s) and tends to increase rapidly as training set size increases. But training data is not large and I assume that to help students we should focus on correct prediction rather than time performance, so choosing SVM from others is the best model.

SVM or Support Vector Machine creates support vectors from data which are representation of the edge point in class, which is "passed" class and "failed" class in this project. The points are selected and make linear lines for each class's side. Those 2 lines are parallel and optimised to have as wide as possible gap while it still separates classes. When predicting, points that fall on the side of the centre line will be classified to be in that class.



- https://en.wikipedia.org/wiki/Support_vector_machine

If data cannot be separated by a line, we could apply a function to all points that will change it to be linearly separable in a new dimension, which is called the kernel trick.

I used the RBF kernel and fine-tuned C and gamma using gridSearchCV. The model's final F1 score is 0.808510638298.