# 數據科學方法期末考

總計 150 分。超過 100 分者均以 100 分計。

1. The **Students performance in exams** dataset (StudentsPerformance.csv) contains students performance in three exams (math, reading, and writing) as long as their demographic and socioeconomic information.

   (a) (15 分) Investigate the dataset by some appropriate explorative data analysis. Are math, reading, and writing scores correlated with each other?

   (b) (15 分) Develop a regression model to predict their math, reading, and writing scores by the other variables. Evaluate your regression model by 10-fold cross-validation. Notice that most of the predictors are categorical.

   (c) (10 分) Does your model underfit or overfit? Explain why.

   (d) (15 分) Is "gender" an important factor for the scores?

   (e) (15 分) Is the "parental level of education" an important factor for the scores?

2. The **Credit card customers** dataset (BankChurners.csv) consists of 10000 customers mentioning their age, salary, marital_status, credit card limit, credit card category, etc. The purpose of this dataset is to predict whether a customer is gonna get churned ("Attrition_Flag") by the other 19 features (some of them are correlated), so the bank can proactively go to the customer to provide them better services and turn customers' decisions in the opposite direction. Notice that the dataset is unbalanced: we have only 16.07% of customers who have churned.

   (a) (15 分) Build a binary classification model to predict who is going to leave their credit card services.

   (b) (25 分) Compute the accuracy, precision, recall, $F_1$-score (the harmonic mean of precision and recall), and AUC by 10-fold stratified cross-validation of your model in (a).

   (c) (25 分) Determine an appropriate decision cutoff $\delta$ so that your classifier in (a) retains at least 90% of recall while maximizes its sensitivity.

   (d) (15 分) Is education level an important factor for customer attrition?