

STAT 5630 Project Milestone 3

Lei Zhang, Yunsheng Lu, Yuxuan Wang

April 2023

1. Classification

We implemented several non-linear methods on the data, including non-linear Support Vector Machine (SVM) K-Nearest Neighbors (KNN), Decision Tree, and Random Forest. The corresponding training and test error are shown in table 3. For SVM, the tuning parameter is selected by 10-fold cross-validation, as shown in table 6. The kernels and the corresponding parameters to be tuned are shown in table 6 and the parameters chosen are shown in table 7. For KNN, we also used 10-fold cross-validation to find the best “k” (table 5). For Random Forest and Decision Tree, the tuning parameters are recorded in table 8 and table 9. The result suggests non-linear SVM yields the best performance for both data. The comparison with linear method (table 1) is similar to that of regression task.

2. Regression

For regression, we used non-linear SVM, KNN, neural network, and Polynomial regression. The corresponding mean squared error of the training and test set is shown in table 4. For SVM, the tuning process is the same as that of classification task, see table 6 and table 7. For KNN, we also used 10-fold cross-validation to find the best “k” (table 5). For neural network, we selected the best parameters. For Random Forest and Decision Tree, the tuning parameters are recorded in table 8 and table 9. The result suggests that in terms of regression, the regression tree yields the best performance for math data and non-linear SVM for Portuguese data. This is also true when compared with the linear methods (table 2). As the improvement is insignificant, it remains uncertain whether we should choose linear method for better interpretation and generalization.

3. Model Diagnostic & Variable Selection

We begin by checking multicollinearity in our regression model. For Portuguese, the Value Inflation Factor (VIF) output shows that all of the VIFs are below or equal to 3, which means the effect of multicollinearity is little. For math, the VIF output indicates that two variables have a VIF over 3, so we remove these variables to eliminate multicollinearity.

Next we move on to outliers identification. From the Cook’s Distance plots (figure 5), we detect observations 1, 157, 193, 199, and 266 should be removed from the math data, and observations 1, 604, 500, 524, 550 should be removed from the Portuguese data.

We then proceed to check error assumptions. We conducted the Box-Cox power transformation and based on the results (figure 6), we let $y_{trans} = \sqrt{y}$.

We further conduct a variable selection based on the summary output using regsubset function. For math, the model includes the number of past class failures, school’s extra educational support, Mother’s and Father’s job, family educational support, study time, going out with friends, and free time. For Portuguese, the model includes past class failures, student’s school, student’s will of higher education, study time, Father’s job, school’s extra educational support, number of absences, and workday alcohol consumption.

Having finished the above procedures, We refitted the models and calculated the adjusted R-squared values of the model, which are 0.30 and 0.37 for math and Portuguese.

4. Parameter Interpretation

We choose non-linear SVM to compare predictors’ influence on the training error, which we defined as the importance in these models, across tasks and datasets. As shown in table 10, table 14, table 11 and table 12. The Venn diagram shown in figure 7 indicate the common important predictors shared by different task and datasets. It is worth noting that 6 predictors are important in regression task for both data.

Appendix

0. Linear Methods Tables

	math		Português	
	training error	test error	training error	test error
baseline	0.768	0.835	0.802	0.738
SVM	0.642	0.633	0.489	0.661
Logistic regression	0.446	0.633	0.462	0.715
LR + Lasso	0.579	0.620	0.493	0.700
LR + Ridge	0.582	0.696	0.489	0.692
LDA	0.487	0.620	0.484	0.692

Table 1: Performance of different methods for classification task

	math		Português	
	training MSE	test MSE	training MSE	test MSE
baseline	0.997	0.966	0.998	0.988
SVM	0.676	0.844	0.629	0.763
OLS	0.886	0.969	0.601	0.800
Lasso	0.678	0.782	0.641	0.774
Ridge	0.671	0.803	0.617	0.772

Table 2: Performance of different methods for regression task

1. Nonlinear Methods Tables

	Math		Português	
	training error	test error	training error	test error
baseline	0.768	0.835	0.802	0.738
Non-linear SVM	0.468	0.608	0.466	0.653
KNN	0.592	0.646	0.561	0.746
Tree(unpruned)	0.437	0.722	0.410	0.677
Tree(pruned)	0.642	0.633	0.578	0.692
Random Forest	0	0.646	0	0.662

Table 3: Performance of different methods for classification task

	Math		Português	
	training MSE	test MSE	training MSE	test MSE
baseline	0.997	0.966	0.998	0.988
Non-linear SVM	0.585	0.827	0.553	0.751
KNN	0.798	0.888	0.701	0.888
Neural Network	0.757	0.889	0.414	0.829
Tree(unpruned)	0.564	0.757	0.448	0.999
Tree(pruned)	0.855	0.767	0.739	0.855
Random Forest	0.160	0.798	0.139	0.763

Table 4: Performance of different methods for regression task

KNN type	tuning parameters
math regression	k=25
por regression	k=19
math classification	k=27
por classification	k=29

Table 5: Tuning parameters for the k (range from 1 to 50)

kernel	formula	tuning parameters
polynomial	$(\gamma x_1^T x_2)^{degree}$	γ , degree(2,3,4), cost
radial	$\exp(-\gamma x_1 - x_2 ^2)$	γ , cost
sigmoid	$\tanh(\gamma x_1^T x_2)$	γ , cost

Table 6: Tuning parameters for non-linear SVM (γ range from 10^{-4} to 10^{-2} , cost range from 0.1 to 5.)

task	parameter chosen
math regression	radial $\gamma = 0.01$, cost = 0.5
Português regression	radial $\gamma = 0.006$, cost = 0.8
math classification	radial, $\gamma = 0.002$, cost = 3.4
Português classification	radial, $\gamma = 0.0025$, cost = 5

Table 7: Parameters chosen for non-linear SVM

RF type	tuning parameters
math regression	ntree=950
por regression	ntree=700
math classification	ntree=700
por classification	ntree=850

Table 8: Tuning parameters for the number of trees (range from 50 to 1000, step=50.)

3. Figures

Tree type	tuning parameters
math regression	cp=0.04
por regression	cp=0.052
math classification	cp=0.05
por classification	cp=0.036

Table 9: Tuning parameters for the cost parameter (range from 50 to 1000, step=50.)

	predictors	importance
32	schoolsup_yes	0.03
8	freetime	0.03
33	famsup_yes	0.03
4	traveltime	0.02
5	studytime	0.02
10	Dalc	0.02
12	health	0.02
16	address_U	0.02
26	Fjob_teacher	0.02
29	reason_reputation	0.02

Table 10: Predictors importance of Non-linear SVM in classification task for math data(top 10)

	predictors	importance
6	failures	0.05
32	schoolsup_yes	0.04
33	famsup_yes	0.02
26	Fjob_teacher	0.02
5	studytime	0.02
37	higher_yes	0.01
15	sex_M	0.01
29	reason_reputation	0.01
36	nursery_yes	0.01
24	Fjob_other	0.01

Table 11: Predictors importance of Non-linear SVM in regression task for math data(top 10)

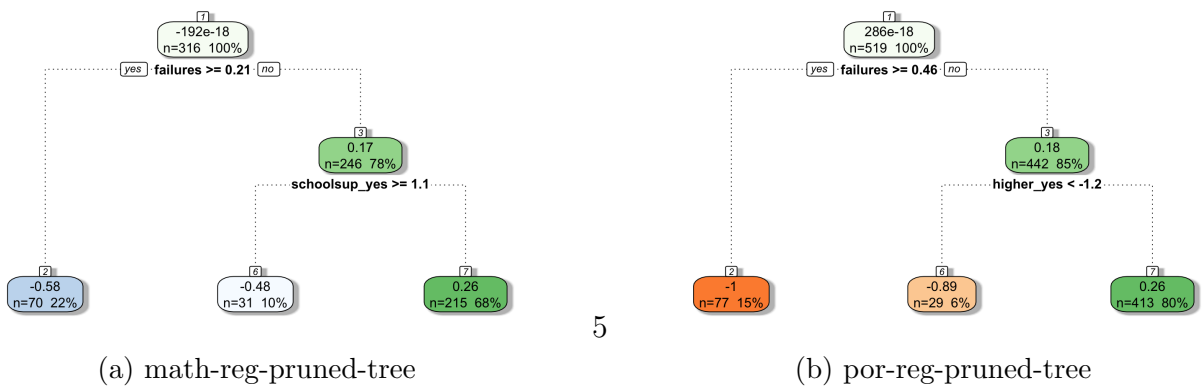


Figure 1: Pruned Trees

	predictors	importance
21	Mjob_services	0.02
14	school_MS	0.02
1	age	0.01
12	health	0.01
37	higher_yes	0.01
10	Dalc	0.01
15	sex_M	0.01
22	Mjob_teacher	0.01
39	romantic_yes	0.01
6	failures	0.01

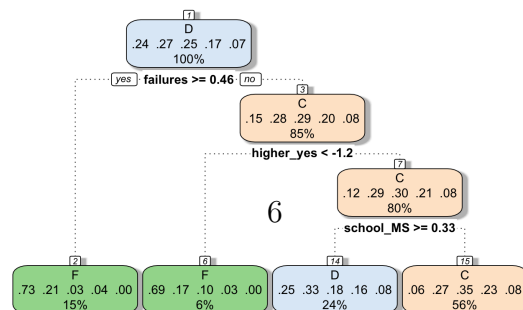
Table 12: Predictors importance of Non-linear SVM in classification for Português data(top 10)

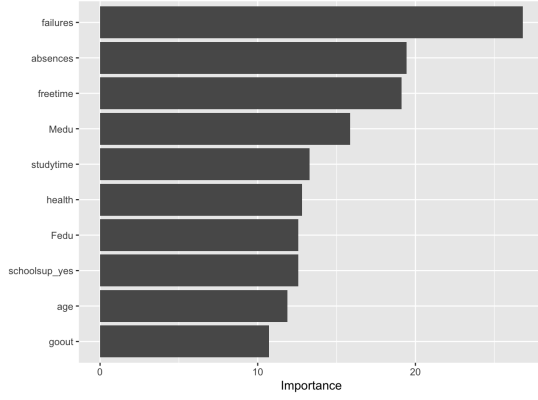
	predictors	importance
6	failures	0.04
14	school_MS	0.03
37	higher_yes	0.03
32	schoolsup_yes	0.01
5	studytime	0.01
13	absences	0.01
26	Fjob_teacher	0.01
31	guardian_other	0.01
38	internet_yes	0.01
15	sex_M	0.01

Table 13: Predictors importance of Non-linear SVM in regression for Português data(top 10)

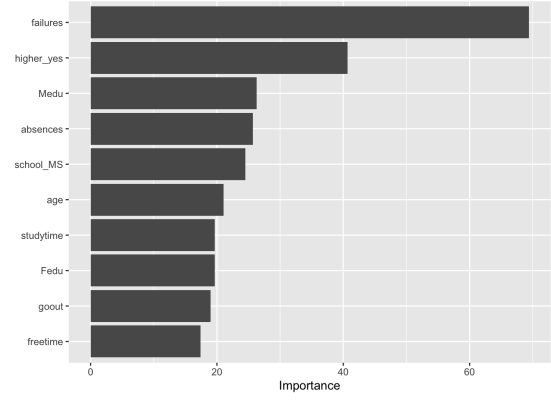
Number of Predictors	Math	Português
1	failures	failures
2	schoolsup_yes	school_MS
3	Mjob_other	higher_yes
4	Fjob_teacher	studytime
5	famsup_yes	Fjob_teacher
6	studytime	schoolsup_yes
7	goout	absences
8	freetime	Dalc

Table 14: best models for each number of variables selected by regsubset()



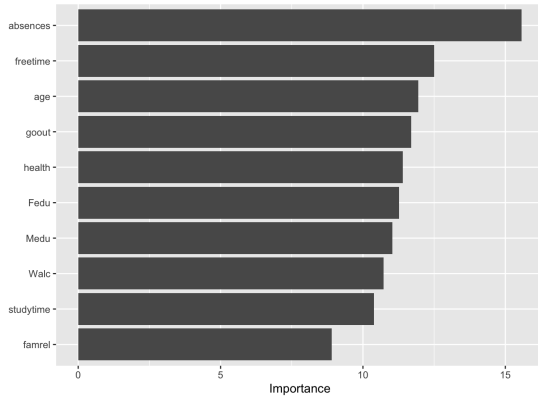


(a) math-reg-vp

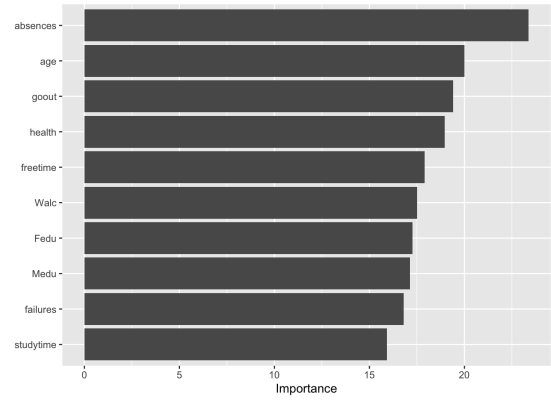


(b) por-reg-vp

Figure 3: Variable Importance Plots for Regression

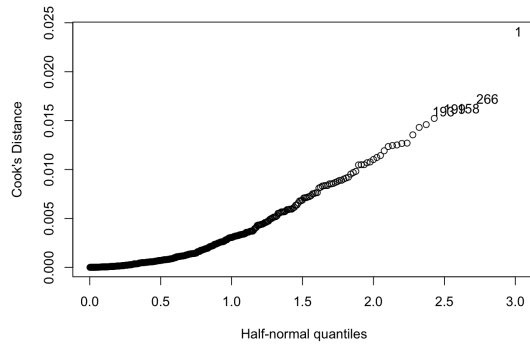


(a) math-class-vp

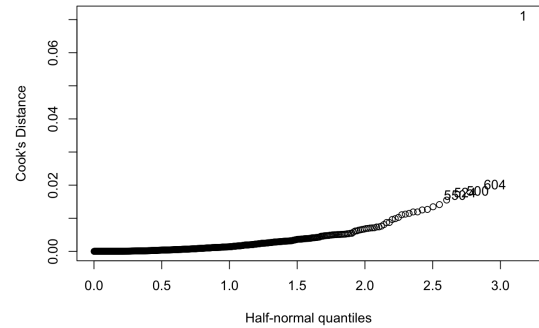


(b) por-class-vp

Figure 4: Variable Importance Plots for Classification

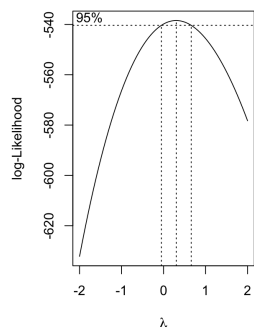


(a) math-outliers

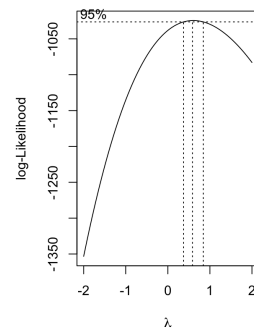


(b) por-outliers

Figure 5: Identification of Outliers using Cook's Distance



(a) math lambda



(b) por lambda

Figure 6: Box-Cox Power Transformation

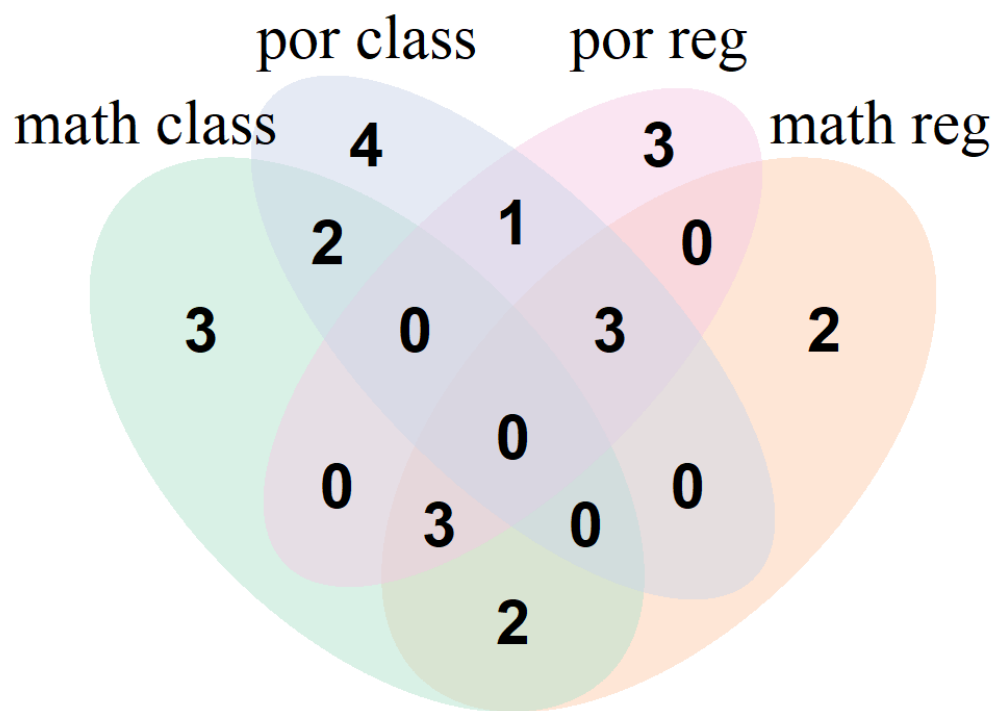


Figure 7: Venn Diagram of important predictors