

Statistics 479

Project 1

Professor: Wei-Yin Loh

TA: Xiaomao Li

Name: Qiuying Li

Date: 01/12/2015

Overview:

In order to predict the INTRDVX from the Consumer Expenditure Survey, I fitted data into multiple linear regression model, randomForest model, Guide-linear split classification and Guide least square stepwise regression. Among these four models, Guide least square stepwise regression has the best performance.

Data description :

The Consumer Expenditure Survey (CE) program offers a continuous and comprehensive flow of data on the buying habits of American consumers. There are 25822 observations and .645 variables in the dataset. INTRDVX describes the the amount received in interest or dividend in the past 12 months. INTRDVX_ has four levels: A C D T. A means Valid blank, a blank field where a response is not anticipated ; C means“Don’t know,” refusal, or other nonresponse, D means Valid value, unadjusted, T means Valid value, topcoded or suppressed.

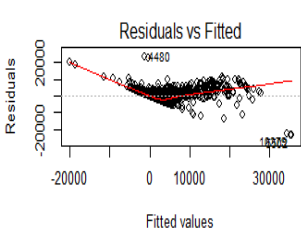
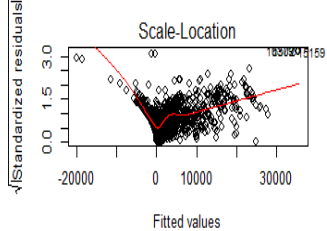
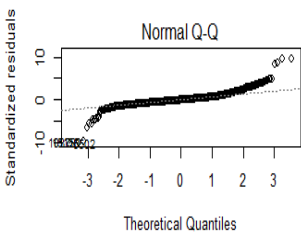
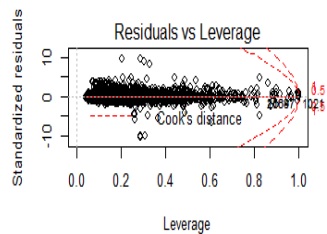
Study question:

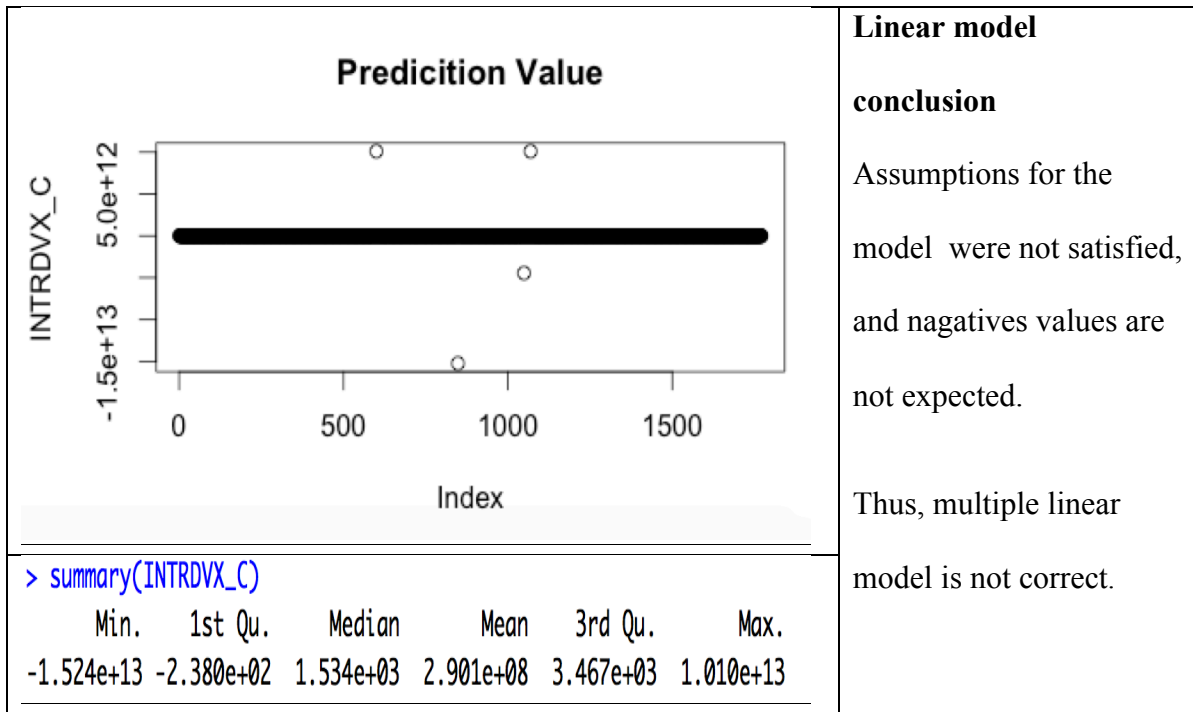
Use the CUs with INTRDVX_ = “D” or “T” to build a model to predict the values of variable INTRDVX (interest and dividends) for which INTRDVX_ = “C”

Analysis methods:

1. Multiple linear regression ;
2. RandomForest
3. Guide-linear split classification;
4. Guide-least square stepwise regression

1. Multiple Linear Regression

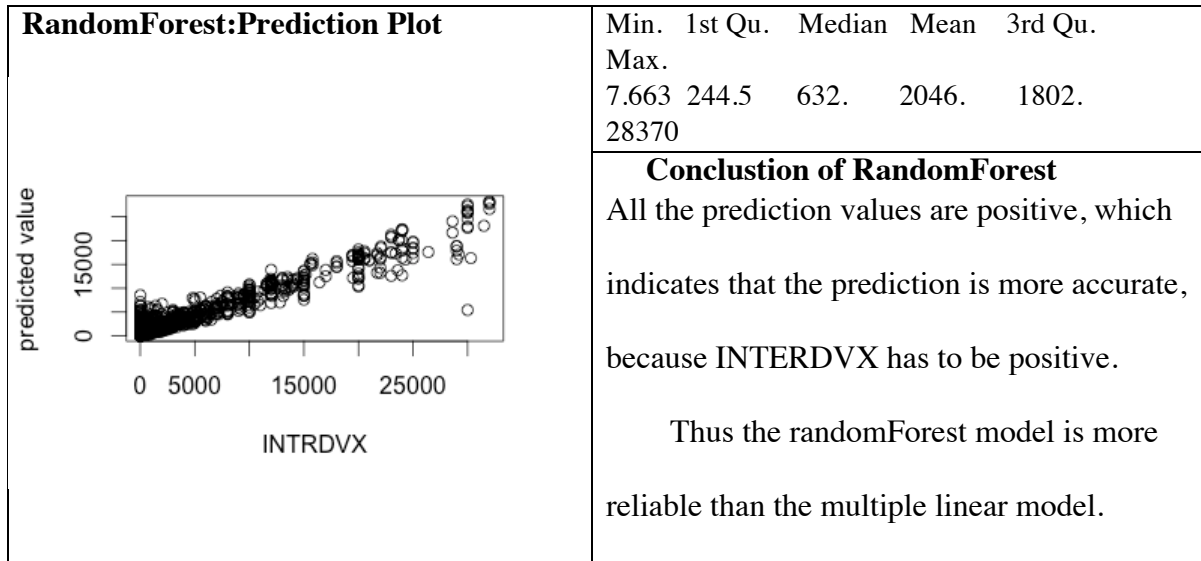
<div>INTRDVX ~ DIRACC + DIRACC_ + AGE_REF + AS_COMP2 + AS_COMP3 + AS_COMP4 + BUILDING + CUTENURE + EARNCOMP + EDUC_REF + EDUCA2 + EDUCA2_ + FAM_SIZE + FAM_TYPE + FAMTFEDX + FAMT_EDX + +FEDRFNDX + FEDR_NDX + FEDTAXX + FEDTAXX_ + FGOVRETX + FGOV_ETX + FINCATAX + FINCAT_X + FINCBTAX + FINCBT_X + FINDRETX + FIND_ETX + FINLWT21 + FJSSDEDX + FJSS_EDX + FPRIPENX <div>(656 predictors in total)</div></div>	<div>Residuals:</div> <table><tr><td>Min</td><td>1Q</td><td>Median</td><td>3Q</td><td>Max</td></tr><tr><td>-23756.1</td><td>-1192.3</td><td>-137.7</td><td>898.2</td><td>23409.0</td></tr></table> <div>Coefficients:</div> <table><tr><td>Estimate</td><td>Std. Error</td><td>t value</td><td>Pr(> t)</td></tr><tr><td>(Intercept)</td><td>-1.064e+04</td><td>1.648e+05</td><td>-0.0650.948520</td></tr><tr><td>BEDR_OMQC</td><td>NA</td><td>NA</td><td>NA</td></tr><tr><td>NA</td><td></td><td></td><td></td></tr><tr><td>.....</td><td></td><td></td><td></td></tr><tr><td>BUILDING5</td><td>3.118e+02</td><td>7.045e+02</td><td>0.4430.658147</td></tr></table> <div>Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1</div> <div>Residual standard error: 2734 on 2142 degrees of freedom</div> <div>Multiple R-squared: 0.7707, Adjusted R-squared: 0.6964</div> <div>F-statistic: 10.38 on 694 and 2142 DF, p-value: < 2.2e-16</div>	Min	1Q	Median	3Q	Max	-23756.1	-1192.3	-137.7	898.2	23409.0	Estimate	Std. Error	t value	Pr(> t)	(Intercept)	-1.064e+04	1.648e+05	-0.0650.948520	BEDR_OMQC	NA	NA	NA	NA							BUILDING5	3.118e+02	7.045e+02	0.4430.658147
Min	1Q	Median	3Q	Max																															
-23756.1	-1192.3	-137.7	898.2	23409.0																															
Estimate	Std. Error	t value	Pr(> t)																																
(Intercept)	-1.064e+04	1.648e+05	-0.0650.948520																																
BEDR_OMQC	NA	NA	NA																																
NA																																			
.....																																			
BUILDING5	3.118e+02	7.045e+02	0.4430.658147																																
<div><div><div>Residuals vs Fitted</div></div><div><div>Scale-Location</div></div></div> <div><div><div>Normal Q-Q</div></div><div><div>Residuals vs Leverage</div></div></div>	<div>Output Description</div> <div><div>1. Large R^2 and small p-value might due to the hundreds of predictors.</div><div>2. QQ plot shows a heavy tails on the two side</div><div>3. Independence and normality test are failed</div></div>																																		



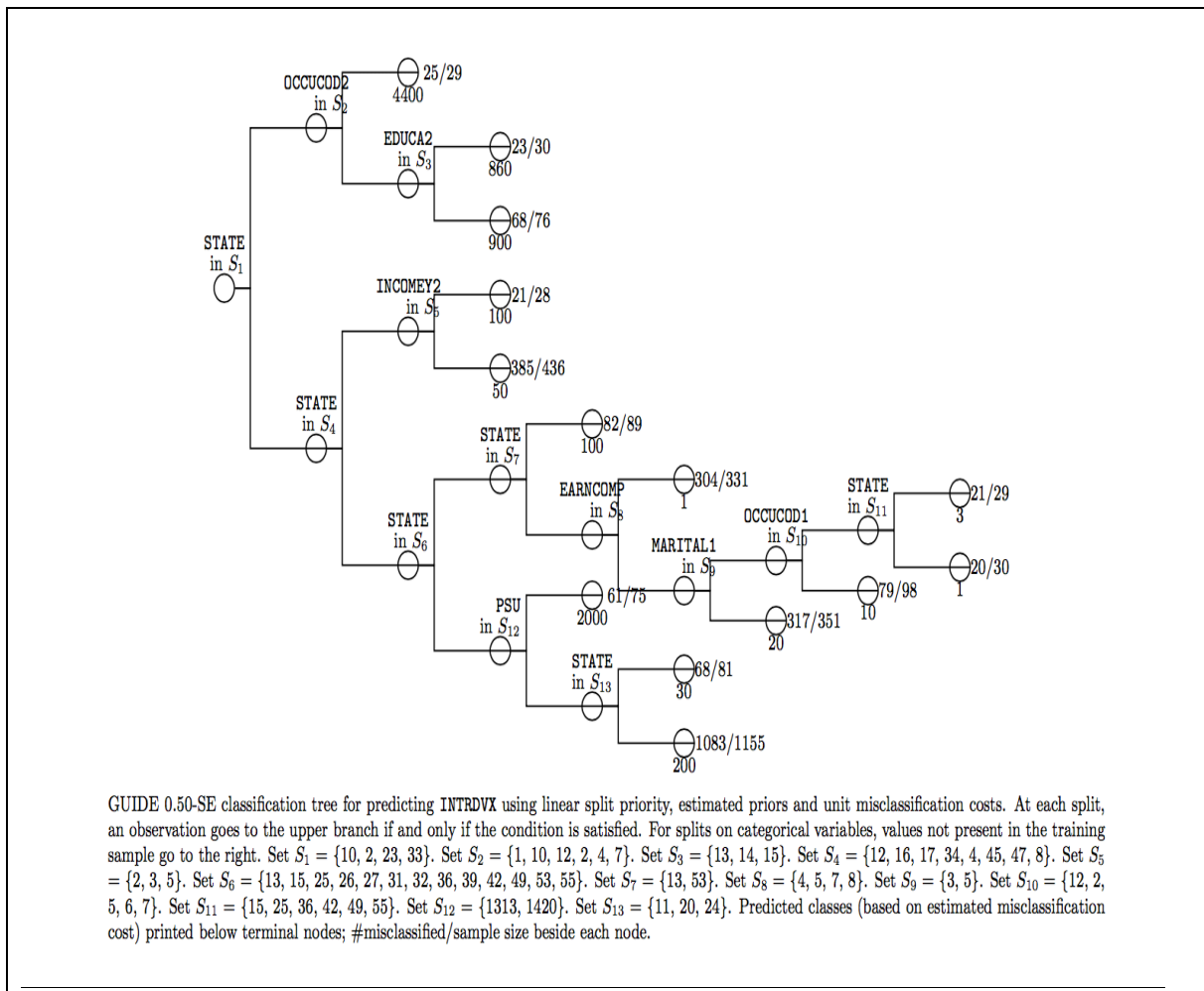
2. RandomForest model:

- Considering that linear regression model is not accurate, use random Foreset based on the linear model.
- However, RandomForest model can only handle dataset contains 53 categories. So I decided to select 38 important variables based on the ranking to be the predictors.
- Below are the 38 selected variables

AS_COMP1	AS_COMP2	AS_COMP4	FINDRETX	FPRI_ENXT	FRRETIRX	FSALARYX	FSSIX
9	10	12	69	75	76	77	81
INC_HRS2	INC_RANK	INCOMEY12	MISC_AXXT	LUMPSUMX	LUMP_UMXD	LUMP_UMXT	OCCUCOD12
86	87	99	111	112	114	115	128
OCCUCOD214	OCCUCOD23	OTHRINCX	OTHR_NCXD	ALCBEVPQ	FEEADMPQ	POV_CY_D	FSMPFRMX
140	143	150	152	215	304	368	562
NETRENTX	NETR_NTXD	OTHREGBX	OTHREGX	OTHREGX_D	RETSURVX	RETS_RVXD	RETSRVBX
599	600	621	624	625	636	637	639



3. Guide-Classification : Linear split

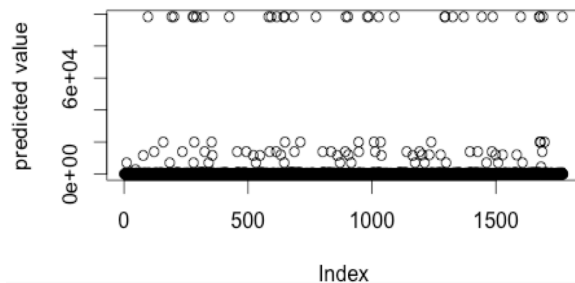


Number of cases used for tree construction = 2838
 Number misclassified = 2557
 Resubstitution est. of mean misclassification cost = 0.9009866102889392

The misclassified number is fairly larger, more importantly, the resubstitution estimation of mean misclassification cost is close to 1, which suggests the model is not accurate.

Prediction of Guide classification

Classification model prediction



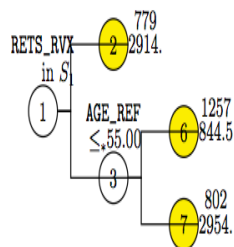
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1	20	50	2340	200	98340

Conclusion : Guide classification

There is a huge variation in the predicted values.

Moreover, there is a huge misclassification cost. So the Guide-Classification is not a good candidate model.

4. Guide-Regression: Linear Square Stepwise



GUIDE 0.50-SE piecewise linear least-squares regression tree with stepwise variable selection for predicting INTRDVX. At each split, an observation goes to the upper branch if and only if the condition is satisfied. The symbol ' \leq ' stands for ' \leq or missing'. Set $S_1 = \{C, D, T\}$. Number in italics beside terminal node is sample mean of INTRDVX. Number above terminal node is sample size.

Plot description: In order to visualize the fitted regression function and the data at the same time, fitting a piecewise simple linear model in the guide, where the best single regressor is selected to fit a straight line in each node, as shown in the picture.

Proportion of variance (R-squared) explained by tree model = .6189

Number of terminal nodes of final tree: 3

Total number of nodes of final tree: 5

Regression tree:

Node 1: RETS_RVX = "C", "D", "T"

Node 2: INTRDVX-mean = 2.91405E+03

Node 1: RETS_RVX /= "C", "D", "T"

Node 3: AGE_REF <= 55.00000 or NA

Node 6: INTRDVX-mean = 8.44452E+02

Node 3: AGE_REF > 55.00000 and not NA

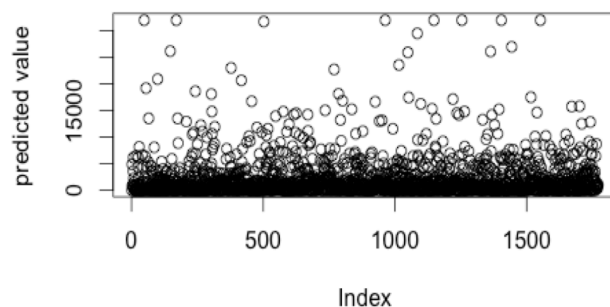
Node 7: INTRDVX-mean = 2.95427E+03

Conclusion : Guide Regression

1. $R^2 = 0.6189$, which is relatively large.
2. Predicted values are all positive with a small variance

Thus, Guide regression model is a good candidate model.

Prediction:GUIDE-Least Square Stepwise



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.0	1.0	955.5	2335.0	2705.0	32000.0

5. Models Comparisons

Models	Advantages	Disadvantages
Linear regression	<ul style="list-style-type: none">● Simple and easy to use	<ul style="list-style-type: none">● Hard to deal with the missing values● Inaccurate when containing too many variables
RandomForest	<ul style="list-style-type: none">● More accurate than linear model● Faster than linear model	<ul style="list-style-type: none">● RandomForest can only handle less than 53 categories
Guide (Linear split Classofocatpn & Least Square Regression)	<ul style="list-style-type: none">● Guide is able to deal with many unseen cases● High accuracy for the prediction● Guide is powerful to the missing values● Running Guide is faster than linear models.● Powerful to huge dataset	<ul style="list-style-type: none">● Guide-classification for this dataset has relatively large variation.● Guide-classification has huge misclassification cost for this dataset.

Conclusion:

1. Based on the discussion above, I decided to drop linear regression model for the inaccuracy.
2. Guide-least square stepwise regression model has the best performance.

According to the summary information of the predicted for 3 different models below, Guide-least square stepwise has the smallest standard deviation, so the this model has the smallest variation to the mean. In addition, over 60% variation in the response can be explained by the model.

Models	Summary of the predicted INTRDVX						sd
RandomForest	Min. 7.663	1st Qu. 244.5	Median 632.	Mean 2046.	3rd Qu. 1802.	Max. 28370	5690. 23
Classification	Min. 1	1st Qu. 20	Median 50	Mean 2340	3rd Qu. 2003.0	Max. 98340	13259.28
Guide-Regression	Min. 1.0	1st Qu. 1.0	Median 955.5	Mean 2335.0	3rd Qu. 2705.0	Max. 32000.0	4024.159

Below are the 95% confidence intervals for the INTRDVX.

Models	95% Confidence Interval
RandomForest	(-22982.82, 13246.2)
Guide-Classification	(-24178.42, 28858.7)
Guide-Regression	(-5552.352, 10222.35)

3. INTRDVX describes the the amount received in interest or dividend in the past 12 months.

INTRDVX_ = “D” or “T” menas people answered the questions, and there are 2922 observations of it. INTRDVX_ = “C” means peope who did not answer the questions. In order to predict the the INTRDVX of people who did not answer the question, we fitted the valid data into the guide least squire regression model, and the we have 95% confidence that the true INTRDVX , which is the amount of the interest or dividend received of the people who did not answer the questions are between (-5552.352, 10222.350) .