

STAT441-Final Project

Xiangru Tan, Ziyue Lin, Pranay Nagpal, Yiran Zhu

Introduction

Customer demographics are very useful to businesses for many purposes. Analyzing the demographics of potential customers can help enterprises to make marketing decisions. For example, older customers are more likely to purchase wheelchairs than young customers. We will be looking into Customer demographics, and how it influence purchases decision in this project.

After looking at the data, we find that 49% of the individuals were male. The age range is between 18 to 60 years old, while 50% of the individuals were between 30 and 46 years old. The salary range is between \$15,000 and \$150,000 and 50% of the range lies in \$43,000 and \$88,000. For the response, 36% of the individuals purchased.

Problem of interest

In this report, we use social network advertisement data to predict whether the product introduced has been purchased or not. We focus on three main features, age, gender, and estimated annual salary.

We split the data into training set (80%) and testing set (20%) for analyses purposes.

Methods

Logistic Regression

	(1)	(2)
(Intercept)	-12.625 *** (1.491)	-12.499 *** (1.446)
Gender	0.128 (0.341)	
Age	0.243 *** (0.029)	0.242 *** (0.029)
EstimatedSalary	0.000 *** (0.000)	0.000 *** (0.000)
N	320	320
logLik	-110.111	-110.181
AIC	228.222	226.363

*** p < 0.001; ** p < 0.01; * p < 0.05.

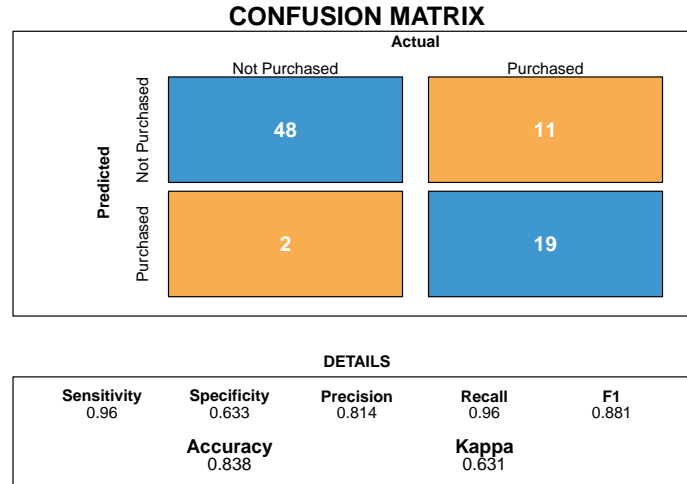
We first fit the full logistic regression model, which includes all three variables, Gender, Age and Estimated Salary.

For the full model, The log odds ratio comparing male(=1) and female(=0) is 0.1279, but we can see the p-value is 0.375 which is much greater than 0.05, so the variable Gender is not significant. For Age and

Estimated Salary, which are continuous random variables, the log odds ratio associated with one unit increase is 0.243 and $3.039 \cdot 10^{-5}$ respectively.

The value of AIC for the full model is 228.22 and the predict error rate is 16.25%.

We can see all the information in the following confusion matrix.

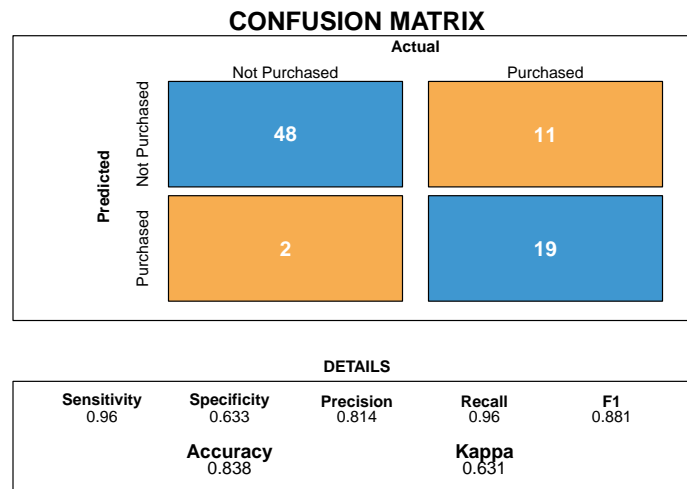


Since, we found the variable Gender being not significant, we fitted the model again excluding Gender.

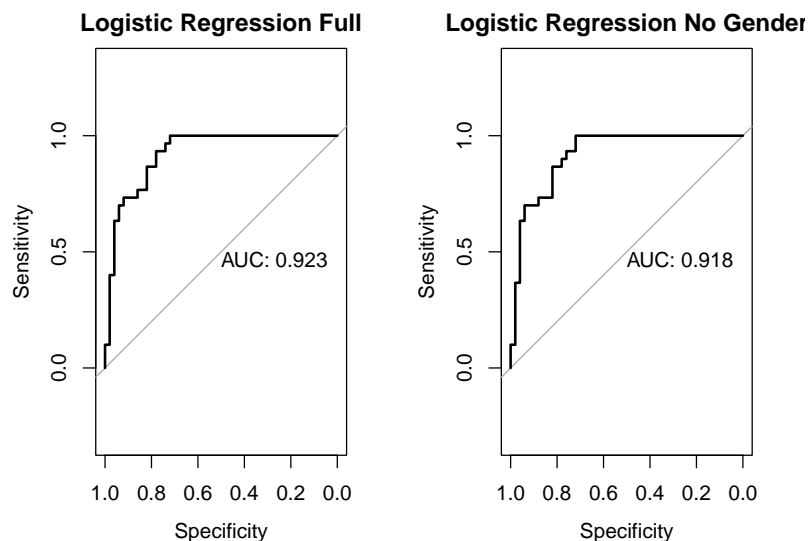
For the new model, the coefficient is similar and all the estimates are significant.

AIC improves to 226.36 and the prediction error rate is also 16.25%.

We resulted having the following confusion matrix.



We also plotted the ROC curve for both the full model and the model without the Gender predictor.



We can see that the AUC for the model without Gender predictor decreased, but not by a lot, so it is a reasonable decision.

Simple Classifier

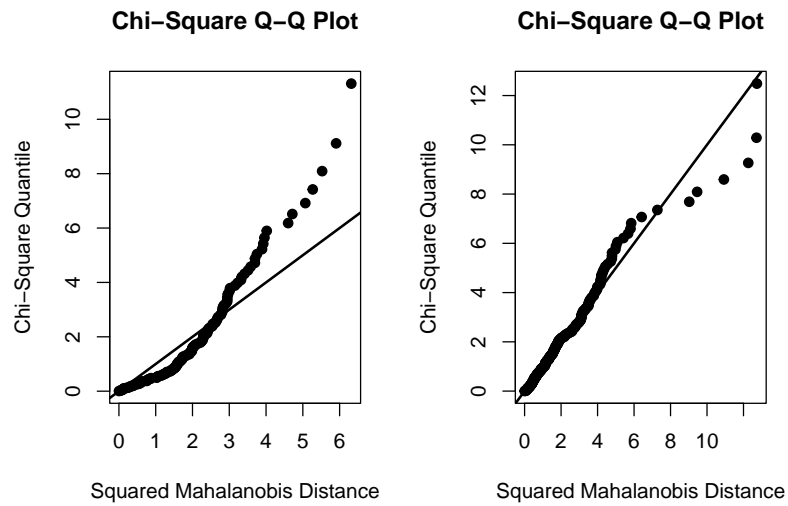
Now let's look at the performance of simple classifiers. We first look at LDA and QDA. Previously, we have seen that Gender is not significant at all, so we will rule that out when analyzing so it's easier to visualize.

In order for LDA and QDA to work properly, the key assumptions we made is that the predictors(X) given Y follows a normal distribution. In our case, it means that the Age and Estimated Salary for people who purchased and people who didn't purchased follows a normal distribution. Lets check that assumptions first.



We can see that only the Estimated Salary for people who purchased doesn't follow normality.

We can also plot the ordered Mahalanobis distances versus estimated quantiles(percentiles), the plot should resemble a straight-line.

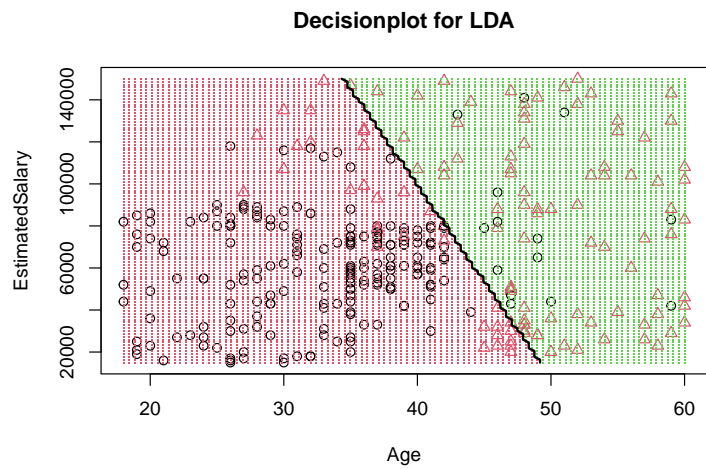


It somewhat did, and together with the normal qq-plot we just drew, we can conclude the normality assumption holds.

Let do LDA now.

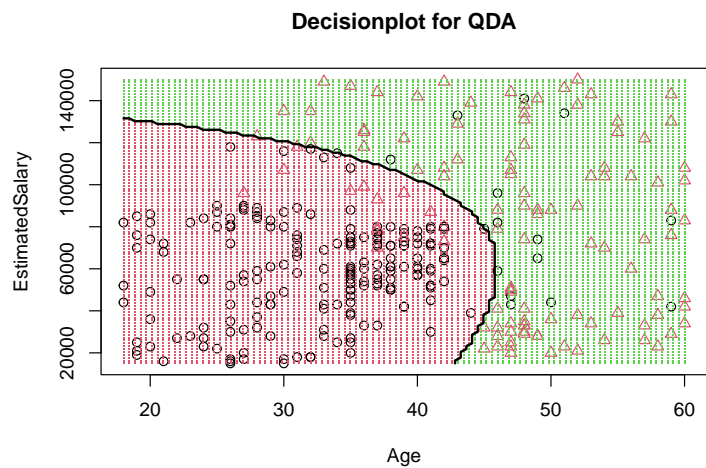
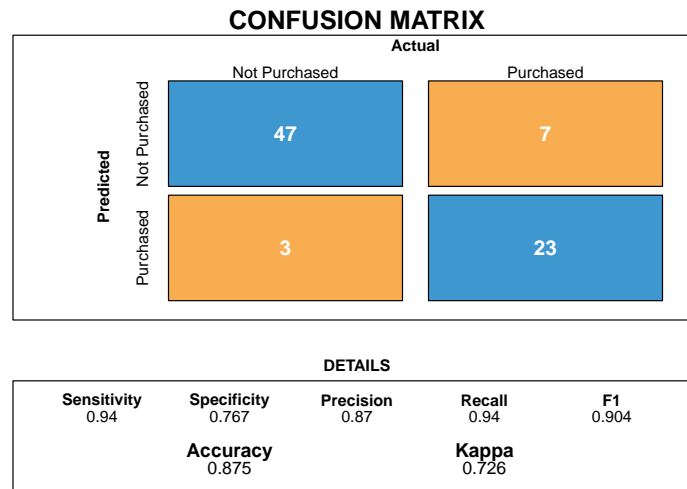
CONFUSION MATRIX		
	Actual	
	Not Purchased	Purchased
Predicted Not Purchased	48	13
Predicted Purchased	2	17

DETAILS				
Sensitivity 0.96	Specificity 0.567	Precision 0.787	Recall 0.96	F1 0.865
Accuracy 0.812		Kappa 0.568		



We can see that the test error of LDA is 18.75%.

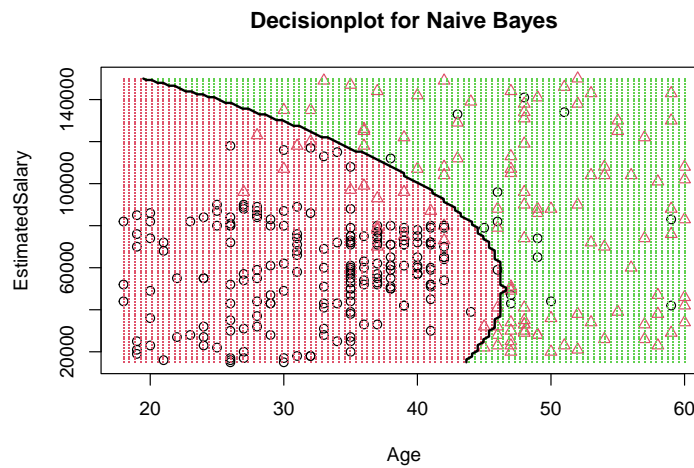
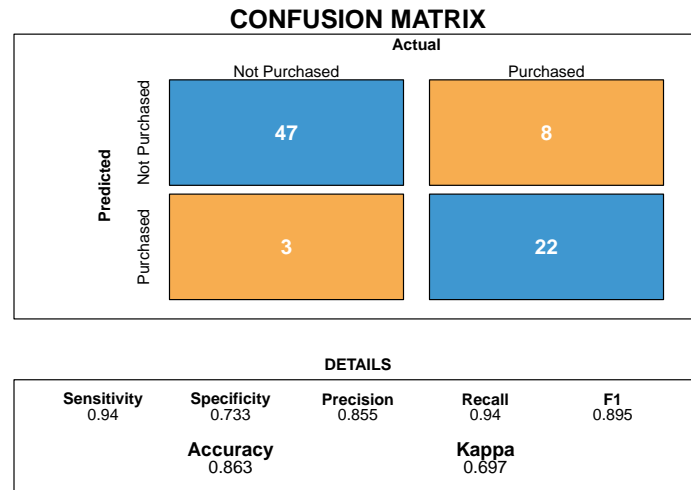
Let's use QDA.



We can see the test error for QDA is 12.5%. Which is 6.25% larger than that of LDA, so QDA obviously does a better job.

Since the normality may not hold, we may need other methods that do not rely on normality, such as Naive Bayes and KNN.

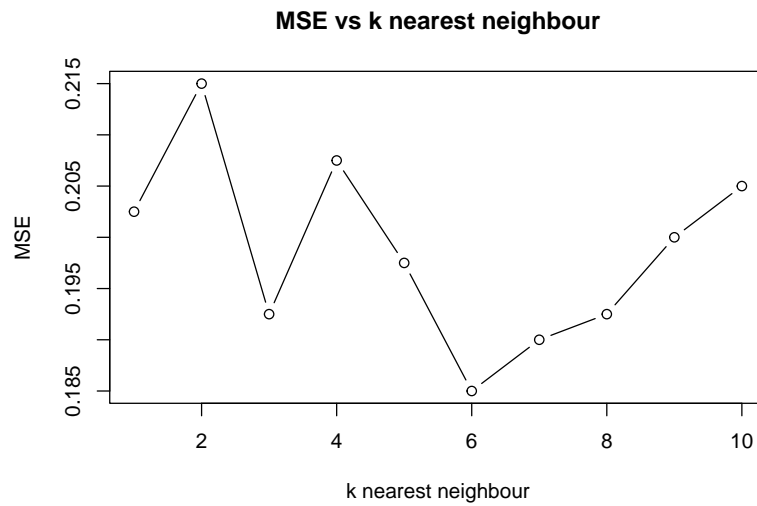
Naive Bayes.



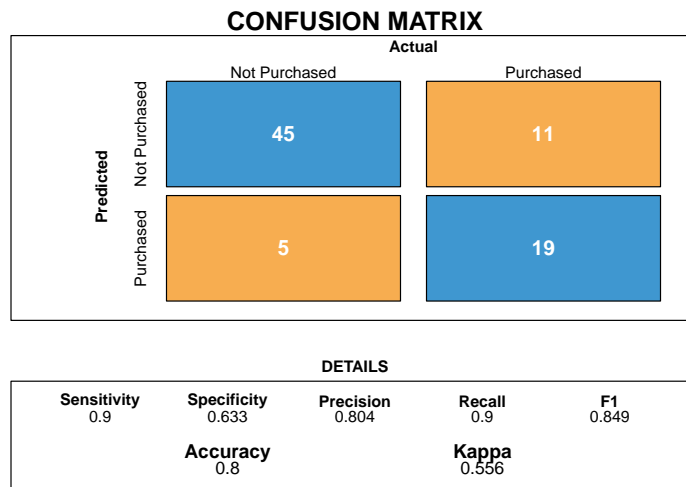
Naive Bayes has a test error of 13.75%, which performs similar to QDA, and better than LDA.

Finally, let's look at KNN.

We first need to do a cross validation to tune the number of neighbors we should use.

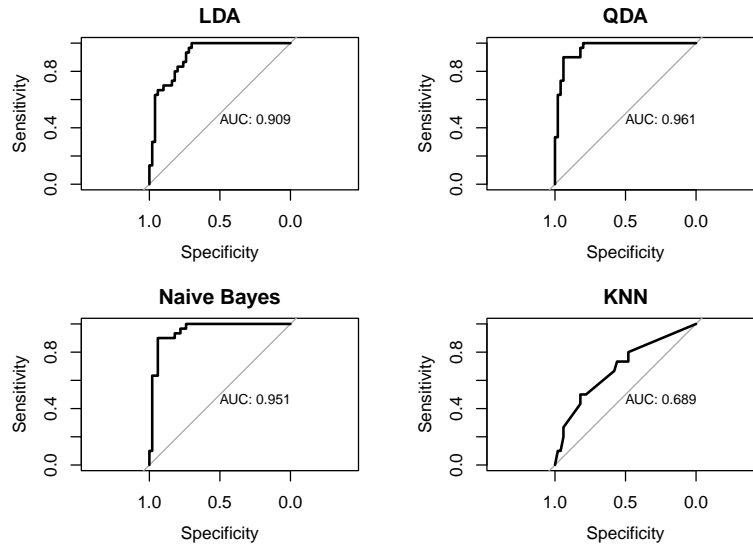


According to cross validation, the optimum number of neighbors we should use is 6.



The test error is 20%, which is even worse than LDA.

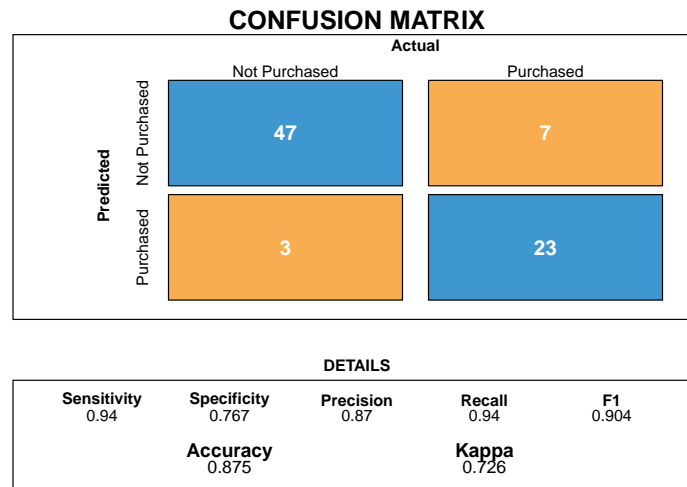
Finally, let's look at the AUC of all the models we used.



We can see that QDA has the largest AUC, which is consistent with the result we had when examining test error. So of all the models in Simple Classifiers, we should use QDA.

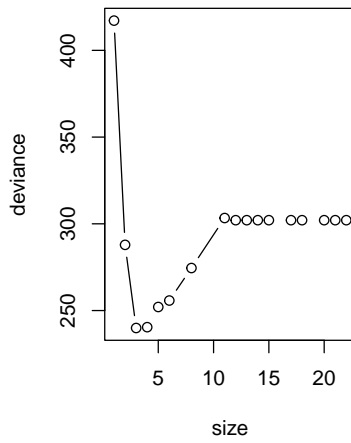
Tree

We now try tree classifier for this problem. We use CART to partition the space and making binary splits by choosing minimum Gini index at each split, the test error is 12.5%. Here we over-build the tree to its full extent and prune the tree later.

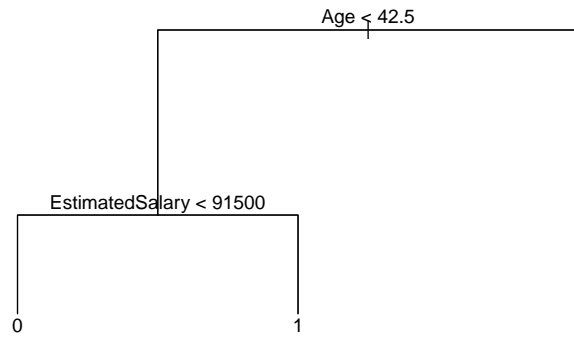
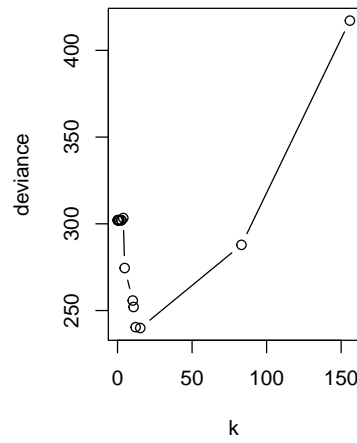


We prune the tree and find the best sized tree based on the deviance with 5-fold cross validation. We plot the error rate as a function of both size and k. The optimal size for the tree is 3.

Plot of deviance vs size of tree



Plot of deviance vs k of tree



CONFUSION MATRIX

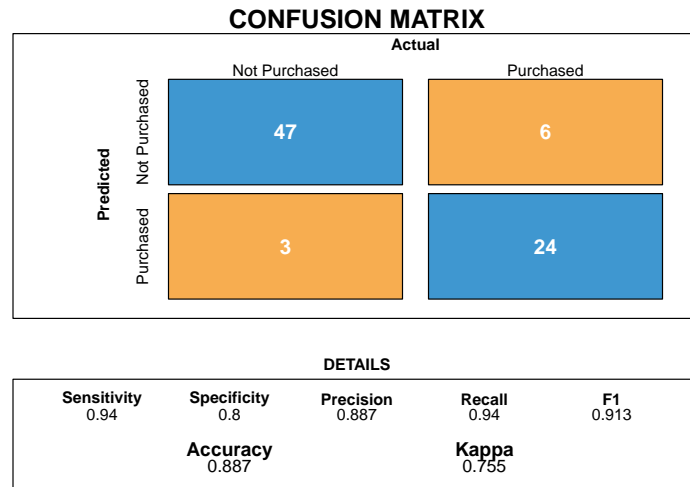
		Actual	
		Not Purchased	Purchased
Predicted	Not Purchased	47	3
	Purchased	3	27

DETAILS

Sensitivity 0.94	Specificity 0.9	Precision 0.94	Recall 0.94	F1 0.94
Accuracy 0.925		Kappa 0.84		

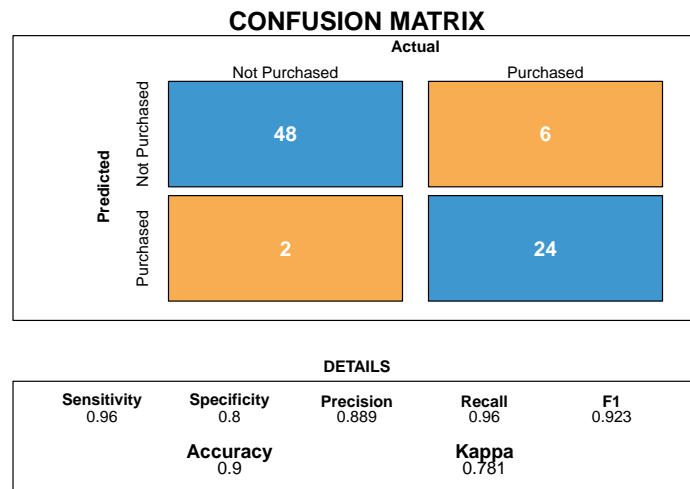
The new confusion matrix is given above and the test error is now 7.5%, which is significantly improved after pruning.

Here we try to improve the performance of the tree classifier by apply bagging, all three predictors are considered for each split of the tree and we choose 200 trees.



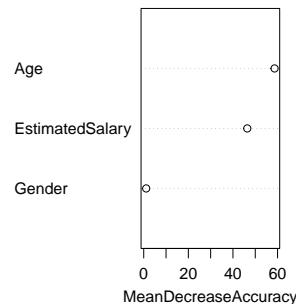
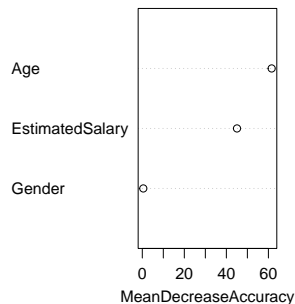
The out-of-bag observation test error is 10.94% and the test error is 11.25%.

For random forest model, we choose 200 trees and try 2 variables at each split. The out-of-bag observation test error is 10.31% and test error rate is 10%.



Variable Importance of Bagging

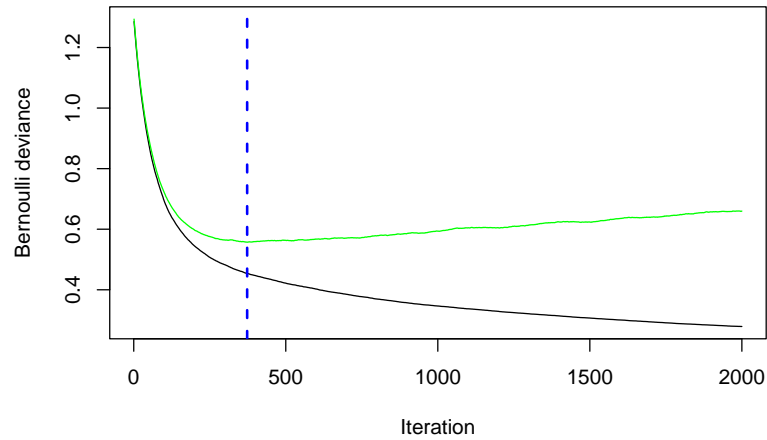
Variable Importance of Random Forest



The variable importance of bagging and random forest is drawn above. From the plots, we can see that the Age is the most influential factor and gender has the least importance which is consistent with our result in general linear model.

We try boosting to further optimize the result. We try 2000 trees, using shrinkage parameter 0.01 to control the learning rate of boosting. We set interaction depth to 2.

The result of relative influence is consistent. The optimal number of trees result from 5 fold cross validation is 373, and the test error rate is 10%.

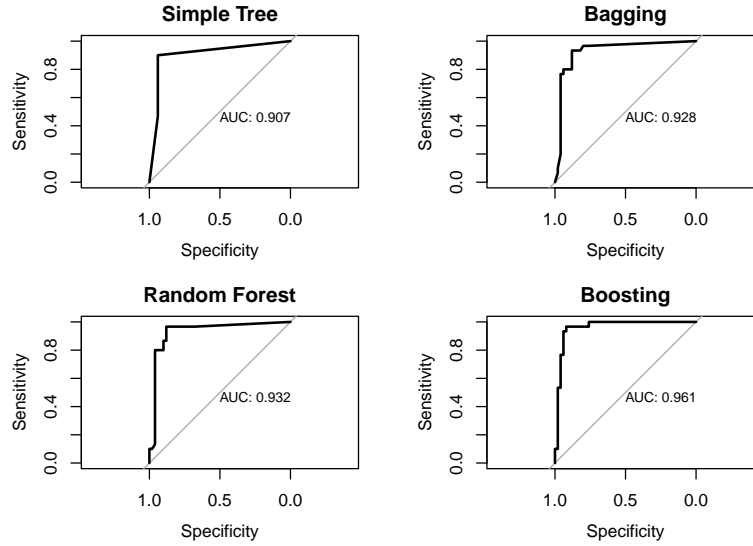


CONFUSION MATRIX

		Actual	
		Not Purchased	Purchased
Predicted	Not Purchased	47	5
	Purchased	3	25

DETAILS

Sensitivity 0.94	Specificity 0.833	Precision 0.904	Recall 0.94	F1 0.922
Accuracy 0.9		Kappa 0.784		



The ROC curve for single tree, bagging, random forest, and boosting are drawn above. The area under curves are 0.907, 0.928, 0.932, and 0.961 respectively. From the plots we can see that the boosting has the best AUC which shows it has the best overall performance.

Conclusion

To solve the problem of interest, we tried logistic regression, LDA, QDA, Naive Bayes Classifier, KNN, and Tree based classifiers (single tree with pruning, bagging, random forest, boosting). We summarize the test error of each model below.

Model_Names	Test_Error	AUC
Logistic Regression	16.25%	0.918
LDA	18.75%	0.909
QDA	12.5%	0.961
Naive Bayes	13.75%	0.951
KNN	20%	0.689
Single Tree with Pruning	7.5%	0.907
Bagging	11.25%	0.928
Random Forest	10%	0.932
Boosting	10%	0.961

Overall speaking, based on AUC and test errors, we decided that the Tree with boosting is the best classifier. It has AUC value of 0.961 and test error of 10%.

Contribution: Xiangru Tan, Yiran Zhu: Logistic regression, LDA, QDA, KNN, Naive Bayes; Ziyue Lin: Tree, random forest, bagging, boosting; Pranay Nagpal: Problem of interest, video, final check.

Reference

Arsh Anwara. (2021). Social Network Ads. <https://www.kaggle.com/datasets/d4rklucif3r/social-network-ads>