

Assignment 4: Ensemble Model Development 1 Using SAS Enterprise Miner

Daanish Ahmed

DATA 640 9041

Fall 2017

dsahmed2334@yahoo.com

Professor Sounak Chakraborty

UMUC

November 19, 2017

Introduction

Ensemble models are some of the most accurate predictive models to use for classification problems. They involve creating several distinct models and averaging their outputs to arrive at a single result. Ensemble models are typically more accurate than individual models, but they can only achieve this when each model is distinct from one another (“Ensemble models,” n.d.). These models work best on binary targets, and some of the most effective ensemble methods include bagging, boosting, gradient boosting, and random forests. In this assignment, I will create several ensemble models using SAS Enterprise Miner on a dataset about automobile purchases. I will focus on categorizing vehicles based on whether they are bad purchases, which will help my organization provide higher quality vehicles for consumers. I will build one model for each of the four ensemble types mentioned. I also plan on creating additional diagrams that use sampling, transformations, and variable selection—allowing me to study their impact on accuracy. I will then examine the effects of preprocessing by excluding data preparation from one of the models. Finally, I will compare the results to see which method is the most accurate, as well as which method can identify the highest number of bad vehicle purchases.

My dataset contains 72,983 observations and 34 variables. The target variable, "is bad buy," indicates whether the vehicle was a bad purchase. It is a binary variable with the values 0 (not bad) and 1 (bad). Some of the inputs include make, model, vehicle year, purchase date, warranty cost, and transmission (see Figure 4 in Appendix). The dataset contains 15 numeric variables and 19 character variables. Of the numeric variables, 2 are binary and 13 are interval. The binary variables are “is bad buy” and “is online sale.” Interval inputs include vehicle year, vehicle cost, and warranty cost. All character inputs are nominal, and they include wheel type, transmission, and model. I set the variables “ref ID,” “wheel type ID,” and “BYRNO” to rejected

because they are IDs and are not useful for the analysis. I also found that vehicle year is strongly correlated with vehicle age. Vehicle age is a better indicator of a vehicle's actual age, while vehicle year is just a "model year" label. Thus, I kept vehicle age and set vehicle year to rejected.

According to the descriptive statistics (see Figure 4 in Appendix), there are six variables with missing values. The input with the highest percentage of missing values is trim (4.37%). This figure also shows that warranty cost has the highest skewness of 2.07 (excluding BYRNO since it was rejected). I found that the target variable is skewed as well, since there are only 8976 vehicles that are labeled as bad purchases. Furthermore, some variables appear to have outliers. For instance, warranty cost, odometer reading, and vehicle cost have values more than three standard deviations from the mean. Additionally, some variables such as current retail clean price are listed as nominal character variables when they are supposed to be numeric interval inputs. It is impossible to set these variables to the interval type. Inspecting the dataset reveals that some of their values are listed as "NULL," which is written in text format (see Figure 6 in Appendix). Thus, I will have to create new variables that exclude these "NULL" cases.

Data Preparation

I began data preparation by handling issues such as outliers, missing values, and errors. I used the Replacement node to eliminate outliers by setting the interval variable default limits method to "standard deviations from the mean." When using the default cutoff value of 3.0, this node will replace all values more than 3 standard deviations from the mean. I also used this node to handle errors, namely the "NULL" issue described earlier. I used the replacement editor to find all "NULL" cases and set their replacement value to unknown, and then I set the class variable

unknown levels replacement to “missing value.” Next, I used the Transform Variables node to create new variables from the inputs which previously included “NULL” values. I used the formula editor to create copies of variables such as “current retail clean price” (see Figure 5 in Appendix). Since these new variables no longer contain any character entries, they are now labeled correctly as numeric interval inputs. Finally, I used the Impute node to handle missing values. I set the class variable default input method to “tree surrogate,” and I set the interval variable default input method to “mean” to compute a replacement for all missing values.

Next, I used the StatExplore node to create a bar graph showing the importance of each input variable (see Figure 7 in Appendix). According to this image, the variables with the highest worth are wheel type, model, vehicle age, sub model, and current auction clean price. Likewise, the least important variables are online sale, transmission, nationality, color, and auction. When looking at the variable distributions, we see that several inputs are skewed and may benefit from transformations (see Figure 8 in Appendix). For instance, warranty cost has a clear positive skew, suggesting that most vehicles have cheaper warranties. Using the “max normal” transform will reduce this skewness. Furthermore, size and VNST have spikes in their distributions—which can be resolved by creating dummy variables. However, this would create hundreds of dummy variables since some variables contain very high numbers of levels. Thus, I will use variable selection to limit the number of inputs for all models that use transformations.

In this analysis, I will create three similar diagrams to explore the effects of sampling, transformations, and variable selection. My first diagram will involve creating the four ensemble models after completing the steps described so far (see Figure 1 in Appendix). My second diagram will iterate on my first diagram by using sampling prior to each model (see Figure 2 in Appendix). My third diagram will expand on my second diagram by using transformations and variable

selection nodes after sampling (see Figure 3 in Appendix). Sampling will be done using the Sample node, and it serves two purposes. Since my dataset contains over 70,000 cases, sampling will firstly reduce the training time. Secondly, creating a weighted sample will balance the cases of my target variable and reduce its skewness. These weighted samples were created by selecting the “level based” sampling criterion and setting level selection to the rarest level. I set the rarest level to take up 40% of all cases, resulting in samples with 22,440 cases and 8976 cases of bad buys. Each Sample node uses a different random seed to ensure that every model uses different data. This will avoid any issues from using the same data in all models. In the end, I will compare the models in each diagram to see if sampling has any effect on the overall accuracy.

In my third diagram, I will perform transformations and variable selection using the Transform Variables and Variable Selection nodes. For the bagging and boosting models, I will use the “max normal” method to minimize the skewness of interval inputs. For the gradient boosting and random forest models, I will use the “best” transform for interval inputs—which has often produced more accurate models in past assignments. All models use dummy indicators for class inputs to allow handling of nominal variables. I will then attach Variable Selection nodes using default properties to reduce the number of dummy variables. These models will determine the impact of transformations and variable selection on accuracy. One disadvantage of variable selection is that it may limit the number of available inputs for models such as random forest.

Model Development

With data preparation complete, I can now create my ensemble models. The first four models use one of the following methods: bagging, boosting, gradient boosting, and random forest.

Bagging, or “Bootstrap Aggregating,” involves building many different decision trees where each tree uses a random sample of the data and replaces those cases from the original dataset (Knode, 2016a). This process results in some cases being used in multiple trees while others are not used at all. Boosting is an iterative process that adjusts the weights of subsequent trees such that incorrect classifications are given more weight while correct predictions have less weight (Srivastava, 2015). Gradient boosting, a variation of boosting, involves computing the error from the previous tree and attempting to correct that error in later trees (Ihler, 2012). Random forests consist of many decision trees made from random data samples, where each iteration of the tree uses a random subset of the available input variables (Knode, 2016a).

The bagging and boosting models were created using three nodes: Start Groups, Decision Tree, and End Groups. The Start Groups node had its mode set to “bagging” and “boosting” for each of the respective models. For the index count (i.e. number of iterations), I kept the default value of 10 due to extremely long training times lasting several minutes. I then attached a Decision Tree node to each model followed by an End Groups node, using the default properties for each. For the gradient boosting and random forest models, I used the Gradient Boosting and HP Forest nodes respectively—keeping the default properties for each node. These four models are the same across all three diagrams, with the only difference being the inclusion of sampling, transformation, and variable selection nodes for the appropriate diagrams. My fifth model will be a random forest model that does not use any preprocessing (see Figure 3 in Appendix). Its purpose is to evaluate whether data preparation will improve the results. I will not use data partitioning in any diagram because ensemble models are usually evaluated on the training data (Knode, 2016a).

I will now describe several important measures for evaluating the results. The ROC curve shows the sensitivity of each model compared to its specificity (see Figure 9 in Appendix). Since

we want to maximize the number of true positives, we prefer our model to have a higher sensitivity. Next, the classification chart shows a bar graph with the number of correct and incorrect classifications for the majority and minority classes (see Figure 10 in Appendix). This is useful for visualizing each model's ability to classify 1's and 0's. The fit statistics include several measures, of which the most important is misclassification rate (see Figure 11 in Appendix). A more accurate model will have a low misclassification rate. Finally, the classification table shows the number of true positives, false positives, true negatives, and false negatives in each model (see Figure 12 in Appendix). We want to minimize the error from false positives and negatives while maximizing the number of correct predictions—especially from true positives. Unfortunately, the results provide little information for finding the most effective inputs. Ensemble models are black box technologies which are difficult to interpret (“Ensemble models,” n.d.), which is a major disadvantage when compared to individual models. However, it is likely that “wheel type” will have the greatest impact on the target due to having the highest worth (see Figure 7 in Appendix).

Results

I will now assess the results of my analysis to determine which method is the most effective at classifying bad vehicle purchases. I will use Model Comparison nodes to analyze all models in each diagram, beginning with my first diagram (see Figure 1 in Appendix). According to the ROC curve, boosting has the highest sensitivity per specificity, followed by random forest, bagging, and gradient boosting (see Figure 9 in Appendix). This suggests that boosting is the most effective model for identifying the minority class, while gradient boosting is the weakest. The classification chart reveals that boosting identified the highest number of true positives, while gradient boosting identified the fewest (see Figure 10 in Appendix). However, boosting also misclassified a massive

number of 0's. And even though this model identified the highest number of true positives, it still failed to classify a significant amount of 1's. Gradient boosting accurately identified more 0's than any other model, yet it hardly identified any 1's. The bagging and random forest models have more balanced results, but they are still very weak at identifying true positives.

The fit statistics indicate that the random forest is the most accurate model with a misclassification of 9.88%, while boosting is by far the least accurate with a misclassification of 36.3% (see Figure 11 in Appendix). This number is significantly worse than the second least accurate model, gradient boosting, which has a misclassification of only 12.3%. The classification table (see Figure 12 in Appendix) corresponds with my findings from the classification chart. This figure shows that boosting identified the highest number of true positives while having the lowest number of false negatives. But it performed poorly at identifying 0's, as evidenced by having 22,492 false positives. Gradient boosting classified every value as a 0—meaning that it correctly identified all 0's but failed to classify any 1's. I would not recommend using the gradient boosting model due to its inability to identify bad vehicle purchases. And though boosting can identify the highest number of 1's, I cannot recommend it either due to its terrible overall accuracy. Of the remaining models, I would suggest using the random forest model since it has the highest accuracy. But when classifying true positives, there is significant room for improvement.

Next, I will examine the models from my second diagram that uses weighted samples (see Figure 2 in Appendix). Its ROC curve is very different from the previous diagram's ROC curve, since it lists random forest as the strongest model and boosting as the weakest (see Figure 13 in Appendix). The fit statistics show that boosting has the highest misclassification rate of 50.0%, while random forest has the lowest rate of 27.5% (see Figure 14 in Appendix). These results resemble those from the previous diagram, but the accuracy rates are now significantly worse.

According to the classification table and classification chart (see Figures 15 and 16 in Appendix), boosting once again identifies the highest number of true positives (6212) and the lowest number of true negatives by far (5004). However, gradient boosting no longer classifies every value as a 0 and has successfully identified 4196 true positives. Thus, it is possible that the gradient boosting model has a better fit on the current data sample than on the full dataset. But overall, all four models perform better at classifying 1's but are significantly weaker at classifying 0's. This causes every model to have a much lower accuracy, and thus I do not recommend using this sampling method on my dataset. Sampling may have weakened the results by reducing the number of cases available in each model. Based on these results, I would recommend the random forest model if focusing on overall accuracy. But if focusing on bad vehicle purchases, I would suggest using gradient boosting since it has the highest number of true positives.

Now I will evaluate the models from my third diagram that uses sampling, transformations, and variable selection (see Figure 3 in Appendix). The ROC curve once again shows that the random forest has the highest sensitivity per specificity while boosting has the lowest (see Figure 17 in Appendix). The fit statistics show that boosting has the highest misclassification of 46.9% while random forest has the lowest rate of 27.4% (see Figure 18 in Appendix). These models are slightly more accurate than those in the second diagram, except for gradient boosting which has a higher misclassification of 35.8%. Gradient boosting used the "best" transform, suggesting that this transformation may not be a good fit for this model. It is also possible that variable selection lowered its accuracy by limiting the number of input variables. Based on the classification table and classification chart (see Figures 19 and 20 in Appendix), boosting still has the highest number of true positives (5167) and the lowest number of true negatives by far (6743). Overall, these models are slightly better at identifying 0's compared to the second diagram, but they are also

slightly worse at identifying 1's. Still, the accuracy improved when using transformations and variable selection in most cases. But these models still succumb to the shortcomings from sampling, making them weaker than my original set of models. The random forest has the highest overall accuracy, and it identified the highest number of true positives aside from the boosting model (the least accurate model). Thus, I still recommend using the random forest model.

Finally, I will evaluate my last model—a random forest without any data preparation (see Figure 3 in Appendix). This model will determine whether preprocessing improves the classification accuracy. The fit statistics reveal that this model has a misclassification rate of 9.96% (see Figure 21 in Appendix), which is slightly worse than my original random forest's misclassification of 9.88% (see Figure 11 in Appendix). The classification table shows that this model has a slightly higher number of true positives, but its increase of false positives outweighs its increase of true positives (see Figure 22 in Appendix). These findings show that preprocessing has slightly improved the accuracy, but it has also slightly reduced the sensitivity.

Conclusion

In my analysis, I created five ensemble models across three diagrams to classify bad vehicle purchases. My goal was to obtain not only the most accurate model, but also the model that can identify the highest number of bad buys. I found that the random forest consistently had the highest accuracy regardless of sampling, transformations, or variable selection. When focusing on true positives, boosting always had the highest sensitivity by far. However, it also had the highest misclassification rate due to its ineffectiveness at classifying 0's. Due to its terrible accuracy, I do not recommend the boosting model for classifying bad vehicle purchases. Excluding this model,

the random forest often had one of the highest sensitivity rates—though gradient boosting had higher sensitivity when using weighted data samples. But since gradient boosting fails to classify 1's when using the full dataset, I ultimately recommend the random forest model.

Through this analysis, I found that my sampling method significantly worsened the results. Although it improved each model's sensitivity, it also drastically lowered the specificity and accuracy. As such, I would suggest using the entire dataset instead. Transformations and variable selection, on the other hand, slightly improved each model's accuracy. However, these nodes slightly reduced the number of true positives identified. I may recommend using transformations and variable selection to maximize accuracy, but this should be done without sampling. Finally, I found that data preparation helps to improve accuracy but with a minimal loss of sensitivity. My fifth model—a random forest with no preprocessing nodes—had a slightly lower accuracy than my original random forest but also had a slightly higher number of true positives.

The main limitation of my analysis is the poor sensitivity of my models. Even my most accurate models failed to achieve sensitivity rates of at least 80%—regardless of data preparation, sampling, transformations, or variable selection. Although ensemble models are usually more accurate than individual models, it is possible that my dataset is not the best fit for using bagging, boosting, gradient boosting, or random forests. It is likely that other models, such as SVM or neural networks, may produce stronger results. Thus, one suggestion for further analysis would be to use combine other types of models using the Ensemble node. These heterogeneous models might provide a closer balance between accuracy and sensitivity. Another idea for future analysis would involve experimenting with a validation or test set to evaluate the results. Though ensemble models are typically evaluated on the training set, it is possible that using this approach may yield helpful insights to help my organization effectively classify bad vehicle purchases.

References

- Ensemble models and portioning algorithms in SAS Enterprise Miner. (n.d.). Retrieved November 6, 2017, from <http://www.sas.com/apps/webnet/video-sharing.html?bcid=4363855671001>
- Ihler, A. (Performer) (2012, November 7). *Ensembles (3): Gradient boosting* [Web]. Retrieved November 7, 2017, from <https://www.youtube.com/watch?v=sRktKszFmSk>
- Knodel, S. (2016a, November 1). *Ensemble Models (bagging, boosting, random forest)*. Lecture presented at UMUC. Retrieved November 6, 2017.
- Knodel, S. (2016b, November 3). *SAS Enterprise Miner Ensemble Models (bagging, boosting, random forest)*. Lecture presented at UMUC. Retrieved November 7, 2017.
- Srivastava, T. (2015, August 2). Basics of Ensemble Learning Explained in Simple English. Retrieved November 6, 2017, from <https://www.analyticsvidhya.com/blog/2015/08/introduction-ensemble-learning/>

Appendix

Relevant SAS Enterprise Miner Output Images

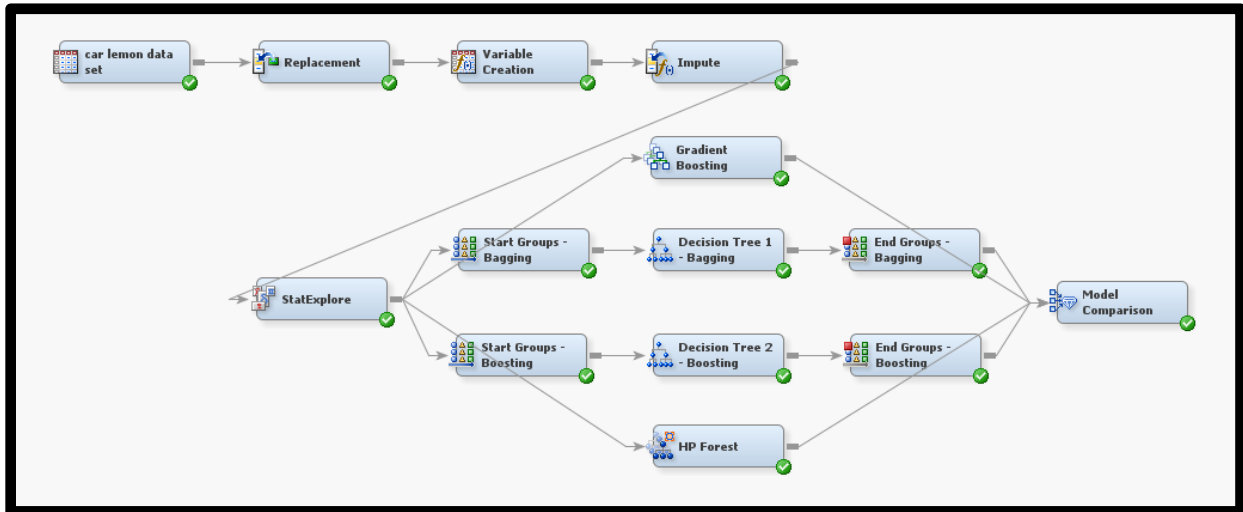


Figure 1. Initial Process Flow Diagram for Car Lemon Ensemble Models.

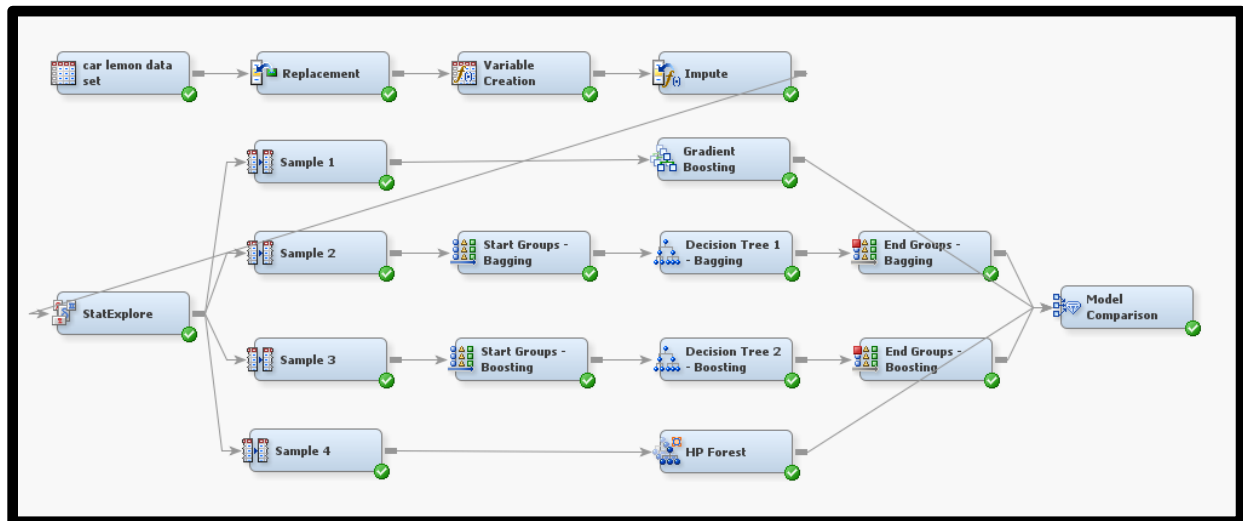


Figure 2. Process Flow Diagram for Ensemble Models with Sampling.

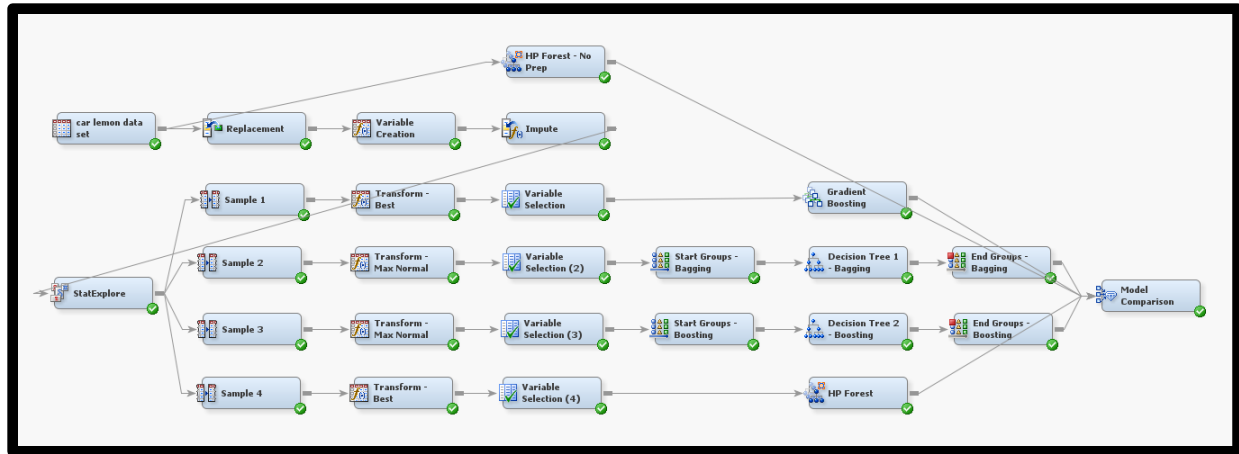


Figure 3. Process Flow Diagram with Sampling, Transformations, and Variable Selection.

Name	Role	Level	Type	Number of Levels	Percent Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis
AUCGUART	Input	Nominal	Character	3	0
Auction	Input	Nominal	Character	3	0
BYRNO	Rejected	Interval	Numeric	.	0	835	99761	26345.84	25717.35	2.129225	3.474184
Color	Input	Nominal	Character	17	0
IsBadBuy	Target	Binary	Numeric	2	0
IsOnlineSale	Input	Binary	Numeric	2	0
MMRAcquisitionAuctionAveragePric	Input	Interval	Numeric	.	0.024663	0	35722	6128.909	2461.993	0.463641	1.593728
MMRAcquisitionAuctionCleanPrice	Input	Interval	Numeric	.	0.024663	0	36859	7373.636	2722.492	0.466501	1.651147
MMRAcquisitionRetailAveragePrice	Input	Interval	Numeric	.	0.024663	0	39080	8497.034	3156.285	0.209214	0.680599
MMRAcquisitionRetailCleanPrice	Input	Interval	Numeric	.	0.024663	0	41482	9850.928	3385.79	0.1763	0.924827
MMRCurrentAuctionAveragePrice	Input	Nominal	Character	21	0
MMRCurrentAuctionCleanPrice	Input	Nominal	Character	21	0
MMRCurrentRetailAveragePrice	Input	Nominal	Character	21	0
MMRCurrentRetailCleanPrice	Input	Nominal	Character	21	0
Make	Input	Nominal	Character	21	0
Model	Input	Nominal	Character	21	0
Nationality	Input	Nominal	Character	5	0
PRIMEUNIT	Input	Nominal	Character	3	0
PurchDate	Time ID	Interval	Numeric	.	0
Refid	Rejected	Interval	Numeric	.	0	1	73014	36511.43	21077.24	-0.000203	-1.19996
Size	Input	Nominal	Character	13	0
SubModel	Input	Nominal	Character	21	0
TopThreeAmericanName	Input	Nominal	Character	5	0
Transmission	Input	Nominal	Character	3	0.00137
Trim	Input	Nominal	Character	20	4.367863
VNST	Input	Nominal	Character	21	0
VNZIP1	Input	Interval	Numeric	.	0	2764	99224	58043.06	26151.64	-0.10353	-1.68869
VehBCost	Input	Interval	Numeric	.	0	1	45469	6730.934	1767.846	0.715931	8.144378
VehOdo	Input	Interval	Numeric	.	0	4825	115717	71500	14578.91	-0.45315	-0.19874
VehYear	Rejected	Interval	Numeric	.	0	2001	2010	2005.343	1.731252	-0.33736	-1.28193
VehideAge	Input	Interval	Numeric	.	0	0	9	4.176644	1.71221	0.393616	-0.20927
WarrantyCost	Input	Interval	Numeric	.	0	462	7498	1276.581	598.8468	2.070831	9.964808
WheelType	Input	Nominal	Character	4	0
WheelTypeID	Rejected	Nominal	Character	5	0

Figure 4. Descriptive Statistics of Car Lemon Dataset.

Name ▲	Type	Length	Format	Level	Formula	Label	Role	Report
NEW_CurrentAuctionAveragePric	Numeric	8		Interval	input(REP_MMRCurrentAuctionAveragePric, best32.)		Input	No
NEW_CurrentAuctionCleanPrice	Numeric	8		Interval	input(REP_MMRCurrentAuctionCleanPrice, best32.)		Input	No
NEW_CurrentRetailAveragePrice	Numeric	8		Interval	input(REP_MMRCurrentRetailAveragePrice, best32.)		Input	No
NEW_CurrentRetailCleanPrice	Numeric	8		Interval	input(REP_MMRCurrentRetailCleanPrice, best32.)		Input	No

Figure 5. Interval Variable Creation Using Transform Variables Formula Editor.



Figure 8. Distributions of Variables in Car Lemon Dataset.

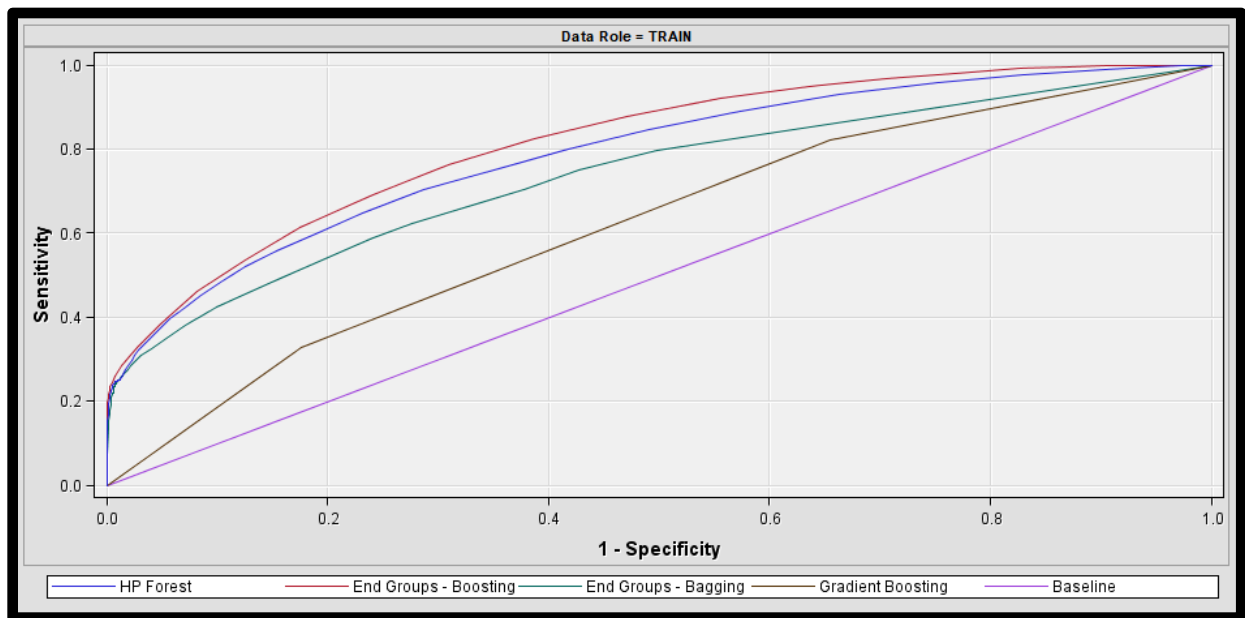


Figure 9. ROC Curve for Models in Initial Process Flow Diagram.

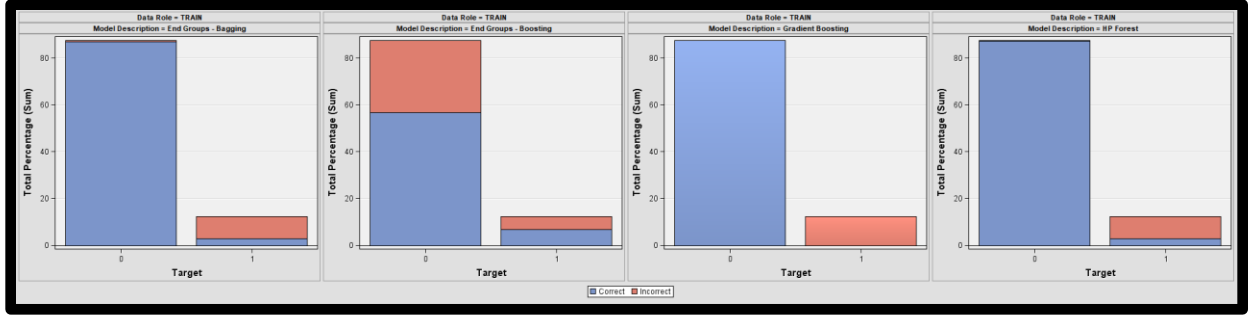


Figure 10. Classification Chart for Models in Initial Process Flow Diagram.

Selected Model	Predecessor Node	Model Node	Model Description ▲	Target Variable	Target Label	Selection Criterion: Train: Misclassification Rate	Train: Average Squared Error	Train: Divisor for ASE	Train: Maximum Absolute Error	Train: Sum of Frequencies
	EndGrp	EndGrp	End Groups - Bagging	REP_IsBad...	Replacement...	0.100571	0.087049	145966	0.929135	72983
	EndGrp2	EndGrp2	End Groups - Boosting	REP_IsBad...	Replacement...	0.363455	0.172631	145966	0.704396	72983
	Boost	Boost	Gradient Boosting	REP_IsBad...	Replacement...	0.122988	0.107366	145966	0.882822	72983
Y	HPDMForest	HPDMForest	HP Forest	REP_IsBad...	Replacement...	0.098804	0.085702	145966	0.974614	72983

Figure 11. Fit Statistics for Models in Initial Process Flow Diagram.

Event Classification Table				
Model Selection based on Train: Misclassification Rate (_MISC_)				
Model Node	Model Description	Data Role	Target	Target Label
EndGrp	End Groups - Bagging	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
EndGrp2	End Groups - Boosting	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
Boost	Gradient Boosting	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
HPDMForest	HP Forest	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
False Negative	True Negative	False Positive	True Positive	
6868	63535	472	2108	
4034	41515	22492	4942	
8976	64007	0	0	
6872	63668	339	2104	

Figure 12. Classification Table for Models in Initial Process Flow Diagram.

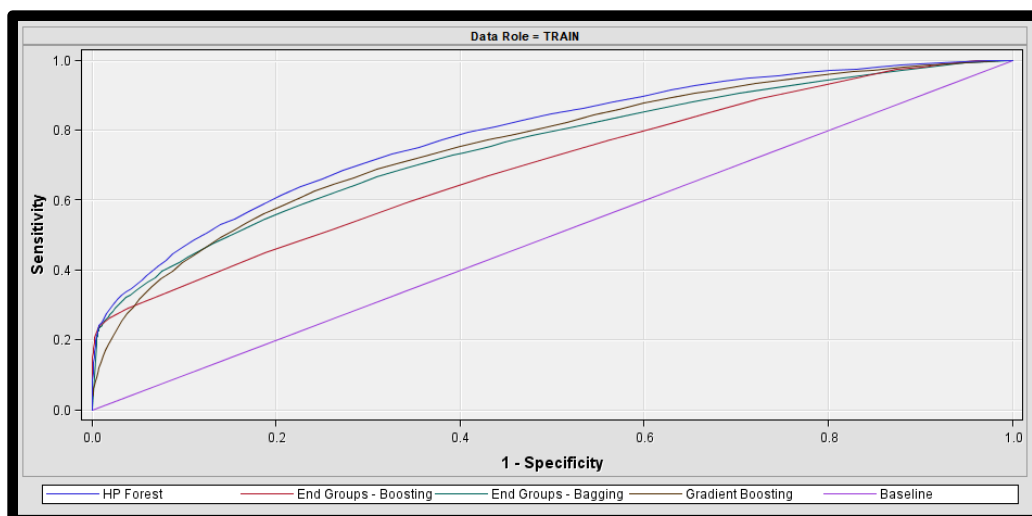


Figure 13. ROC Curve for Models with Sampling.

Selected Model	Predecessor Node	Model Node	Model Description ▲	Target Variable	Target Label	Selection Criterion: Train: Misclassification Rate	Train: Average Squared Error	Train: Divisor for ASE	Train: Maximum Absolute Error	Train: Sum of Frequencies
	EndGrp	EndGrp	End Groups - Bagging	REP_IsBad...	Replaceme...	0.30918	0.189837	44880	0.956063	22440
	EndGrp2	EndGrp2	End Groups - Boosting	REP_IsBad...	Replaceme...	0.500178	0.231253	44880	0.712571	22440
	Boost	Boost	Gradient Boosting	REP_IsBad...	Replaceme...	0.288547	0.195944	44880	0.876068	22440
Y	HPDMForest	HPDMForest	HP Forest	REP_IsBad...	Replaceme...	0.274911	0.182981	44880	0.893592	22440

Figure 14. Fit Statistics for Models with Sampling.

Event Classification Table				
Model Selection based on Train: Misclassification Rate (_MISC_)				
		Data		
Model Node	Model Description	Role	Target	Target Label
EndGrp	End Groups - Bagging	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
EndGrp2	End Groups - Boosting	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
Boost	Gradient Boosting	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
HPDMForest	HP Forest	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
	False	True	False	True
Negative	Negative	Positive	Positive	
4906	11432	2032	4070	
2764	5004	8460	6212	
4780	11769	1695	4196	
4981	12276	1188	3995	

Figure 15. Classification Table for Models with Sampling.

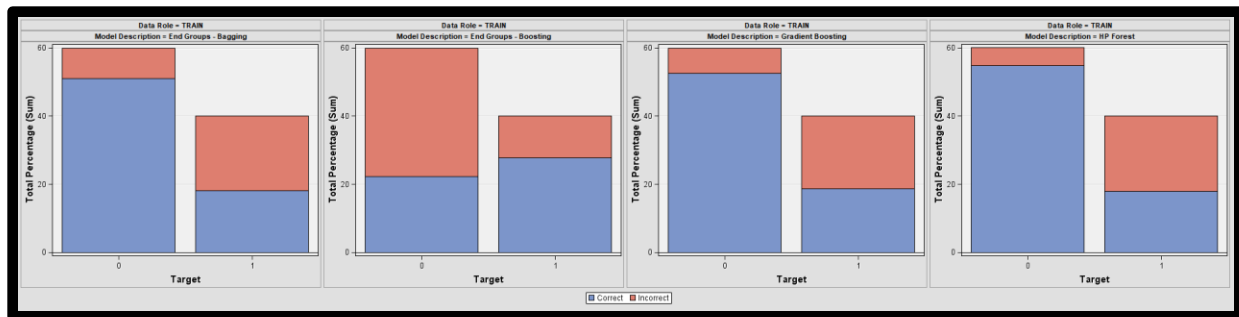


Figure 16. Classification Chart for Models with Sampling.

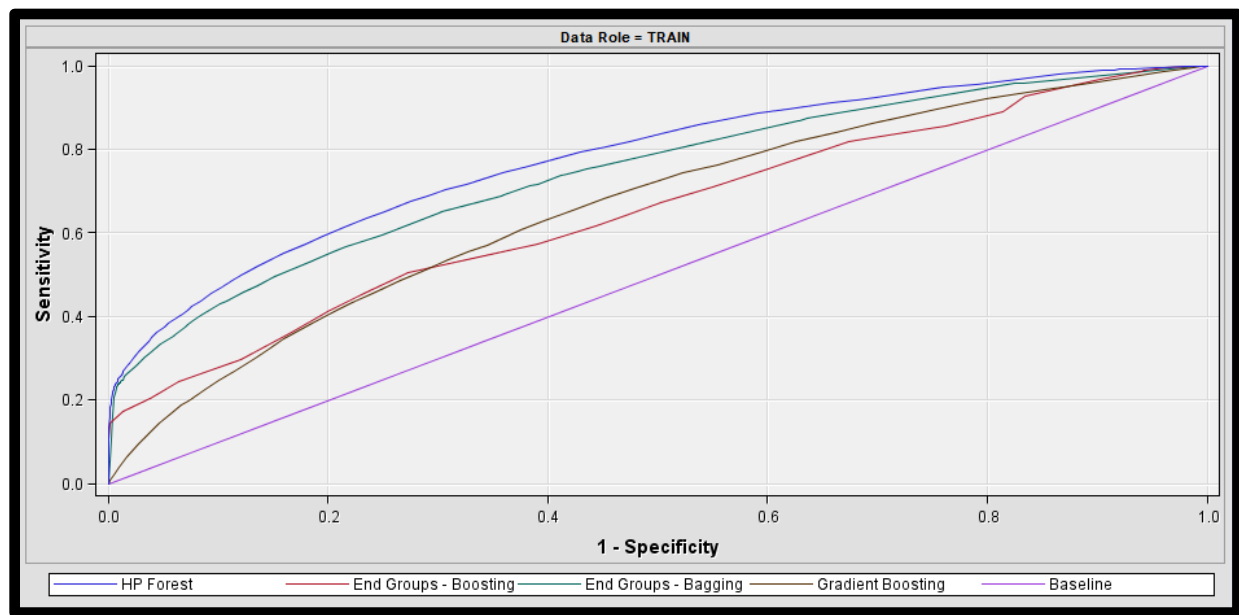


Figure 17. ROC Curve for Models with Sampling, Transforms, and Variable Selection.

Selected Model	Predecessor Node	Model Node	Model Description ▲	Target Variable	Target Label	Selection Criterion: Train: Misclassification Rate	Train: Average Squared Error	Train: Divisor for ASE	Train: Maximum Absolute Error	Train: Sum of Frequencies
	EndGrp	EndGrp	End Groups - Bagging	REP_IsBad...	Replaceme...	0.301248	0.191123	44880	0.946435	22440
	EndGrp2	EndGrp2	End Groups - Boosting	REP_IsBad...	Replaceme...	0.469251	0.235435	44880	0.661701	22440
	Boost	Boost	Gradient Boosting	REP_IsBad...	Replaceme...	0.358066	0.22198	44880	0.801505	22440
Y	HPDMForest	HPDMForest	HP Forest	REP_IsBad...	Replaceme...	0.274287	0.181057	44880	0.93545	22440

Figure 18. Fit Statistics for Models with Sampling, Transforms, and Variable Selection.

Event Classification Table				
Model Selection based on Train: Misclassification Rate (_MISC_)				
Model Node	Model Description	Data Role	Target	Target Label
EndGrp	End Groups - Bagging	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
EndGrp2	End Groups - Boosting	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
Boost	Gradient Boosting	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
HPDMForest	HP Forest	TRAIN	REP_IsBadBuy	Replacement: IsBadBuy
	False Negative	True Negative	False Positive	True Positive
	5941	12645	819	3035
	3809	6743	6721	5167
	5880	11313	2151	3096
	4889	12198	1266	4087

Figure 19. Classification Table for Models with Sampling, Transforms, and Variable Selection.

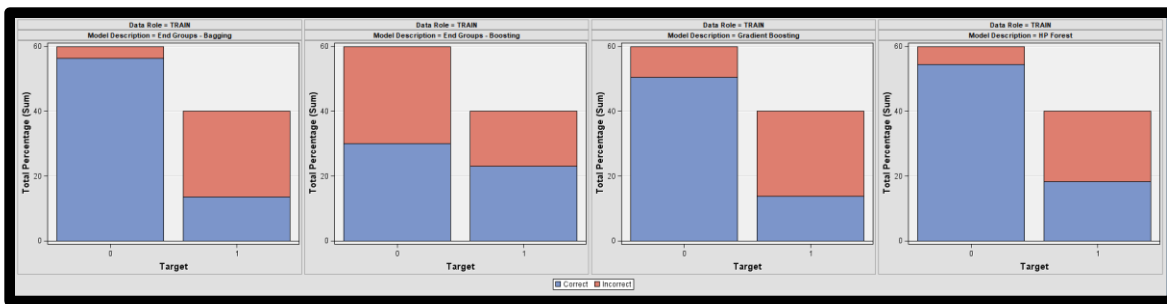


Figure 20. Classification Chart for Models with Sampling, Transforms, and Variable Selection.

Target	Target Label	Fit Statistics	Statistics Label ▲	Train	Validation	Test
IsBadBuy		_ASE_	Average Squared Error	0.085185		
IsBadBuy		_DIV_	Divisor for ASE	145966		
IsBadBuy		_DISF_	Frequency of Classified Cases	72983		
IsBadBuy		_MAX_	Maximum Absolute Error	0.982691		
IsBadBuy		_MISC_	Misclassification Rate	0.099626		
IsBadBuy		_WRONG_	Number of Wrong Classifications	7271		

Figure 21. Fit Statistics for Random Forest Model with no Data Preparation.

Event Classification Table				
Data Role=TRAIN Target=IsBadBuy Target Label=' '				
False Negative	True Negative	False Positive	True Positive	
6818	63554	453	2158	

Figure 22. Classification Table for Random Forest Model with no Data Preparation.