

Carvana: Predicting “Kicked” Purchase of Used Cars for An Online Dealership

DNSC 6279 Data Mining

May 5, 2018

Sahar Sohofi

Xuan Yang

Yinlu Wu

Problem Description

Carvana is a technology start-up based in Phoenix Arizona. It is an online-only used car dealership and it allows people to buy and sell their used cars through the website. Assuming a customer would like to purchase a used car, on the website they are able to search for a car according to body style, price range, financing, amount of monthly payments, year, color, interior features, engine transmission, cylinders, drive, fuel type.

As a way to improve their service and help their customers make more satisfactory purchases, Carvana would like to use their data on previous purchases to find out if each purchase of used car is a good or a bad purchase through supervised learning, and train the model built to predict good/bad purchases for future deals. A bad purchase, or a “kick”, refers to used vehicles with serious issues (e.g. tampered odometers) that prevent them from being resold to customers or requires heavy repair work. Kicked cars are costly to both dealerships and end buyers considering transportation cost, repairmen cost, and loss in market share.

Thus, we need to conduct predictive analysis using supervised learning techniques such as Logistics Regression, Decision Trees and Bootstrap Forest Model, etc. Then select the best model built to help Carvana better care for their customers. Besides, according to the output of the best model, we could realize other interesting observations, for instances, overall likelihood of a kick, any correlation between predictors, and which variables contribute most to our target variable.

Data Selection

The dataset was obtained from Kaggle as a past competition. Our target variable is “IsBadBuy”, a binary variable showing whether the purchased car is a “kick”. On Kaggle website, there are two datasets, one training set with 72,983 records and a test set with 48,707 records . Because only training dataset has the target variable “IsBadBuy”, we only focus on this dataset for building predictive models. Test set is for Kaggle participants to run models on and compete for scores. Besides that, there are 32 predictors, including 16 numerical variables and 16 categorical variables, describing vehicle age, make and model type, reference prices, auction information and other key information. A data dictionary describing all the variables is provided in the appendix.

Data Preprocessing

1. Correct data types as necessary: we found there were lots of predictors false classified. For example, our target variable should be categorical but in the original dataset it was recognized as continuous, so we needed make sure all variables are correctly identified.
2. Missing values: In the original dataset, only the variable “trim” has missing values, but we later discovered that in other variables there are many observations coded with “NULL”. We identified them as missing values. Since SAS JMP supports “informative missing” in decision trees and neural networks, we are not worried about the existence of missing values except for logistic regression models.
3. Corrected typing errors: we only corrected errors when we are certain that they are merely due to typing, such as uppercase vs lowercase.
4. Splitted text for data: in the original dataset, variables “Model” and “Submodel” are long text containing information on engine, volume, number of doors, wheel drive, etc. We used *Excel* to extract such information and created new variables. Unfortunately, since “Model” and “Submodel” are capped at 20 characters, many records are missing such information, making the new variables not reliable. Therefore, we still used the original “Model” and “Submodel” for predictive analysis. The training set with extracted variables are provided as an attachment.
5. Variable selection: we removed some variables which were obviously not related to target variable, like Buyer ID, and also removed variables with repeated information, like WheelTypeID vs WheelType, vehicle_year vs vehicle_age. Besides, we removed variables which were too complicated to model while not adding much predictive power, like Zip code and Purchase date.

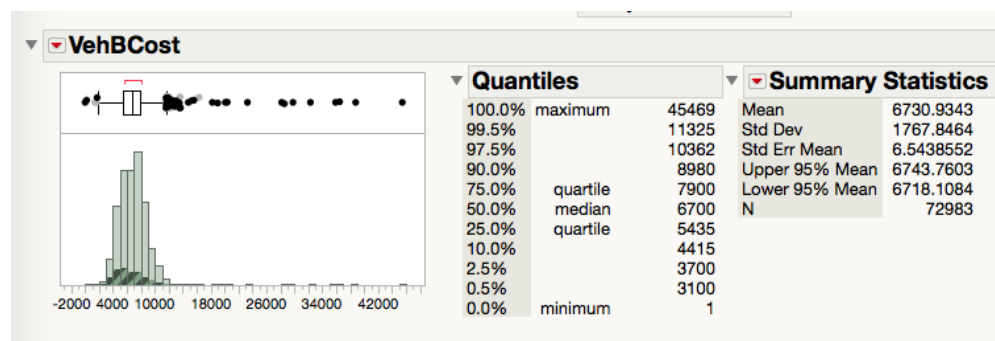
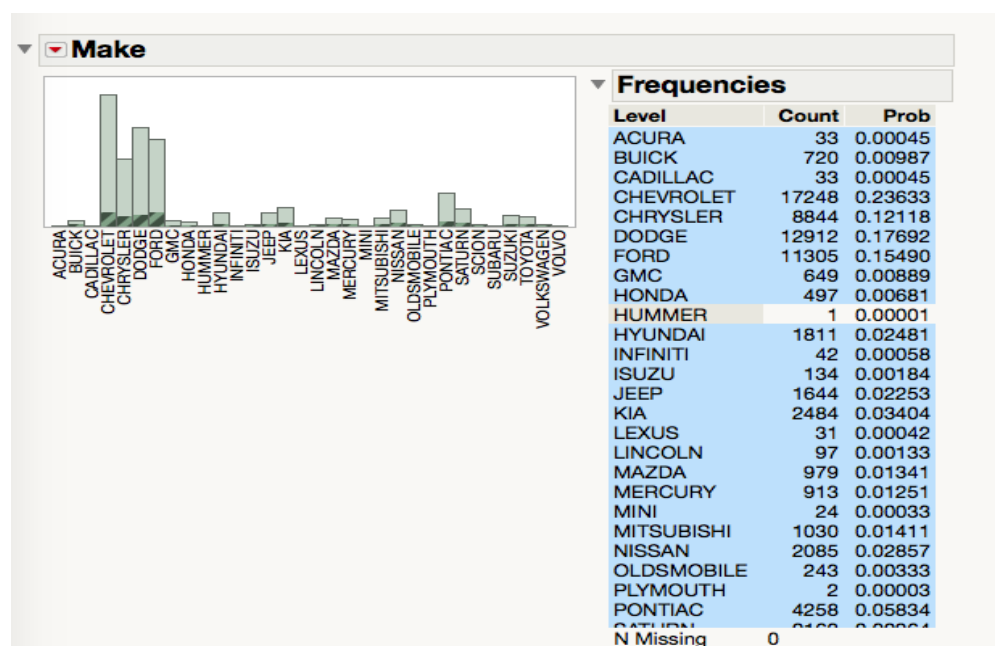
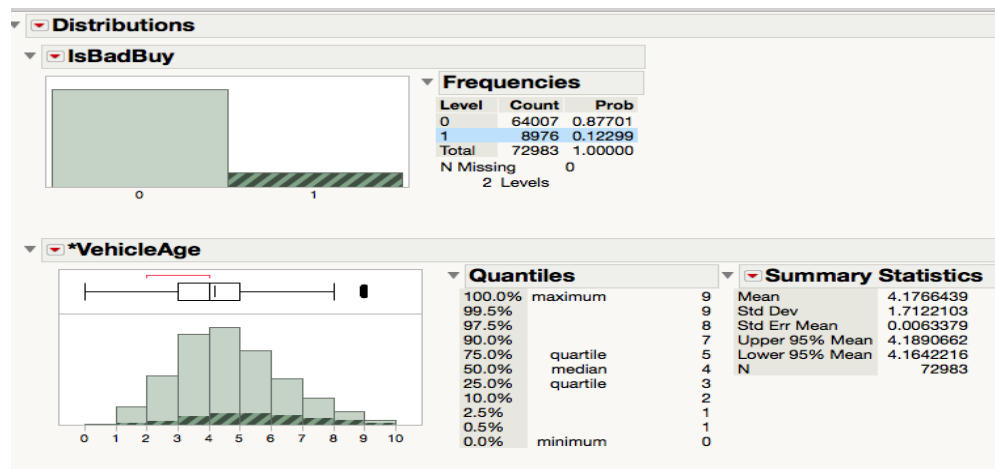
Below is the data summary after pre-processing.

Columns	N	N Missing	N Categories	Min	Max	Mean	Std Dev
IsBadBuy	72983	0	2
**PurchDate	72983	0	517
Auction	72983	0	3
*VehYear	72983	0	.	2001	2010	2005.3430524917	1.7312516057818
VehicleAge	72983	0	.	0	9	4.1766438759711	1.7122102878797
Make	72983	0	32
**Model	72983	0	1064
**Trim	70623	2360	134
**SubModel	72975	8	863
Color	72881	102	15
Transmission	72974	9	2
*WheelTypeID	69809	3174	3
WheelType	69809	3174	3
VehOdo	72983	0	.	4825	115717	71499.995916857	14578.913128203
**Nationality	72978	5	4
Size	72978	5	12
**TopThreeAmericanName	72978	5	4
**PRIMEUNIT	3419	69564	2
**AUCGUART	3419	69564	2
*BYRNO	72983	0	.	835	99761	26345.842155022	25717.351218728
*VNZIP1	72983	0	.	2764	99224	58043.059945467	26151.640414915
VehBCost	72983	0	.	1	45469	6730.934326213	1767.8464352309
IsOnlineSale	72983	0	2
WarrantyCost	72983	0	.	462	7498	1276.5809846129	598.84678820937
Prin1	72668	315	.	-7.767877407126	29.558634404353	0.0011367955648	2.7156883170312
Prin2	72668	315	.	-4.372561241603	4.0153932335083	0.0014192650965	0.5424591535496
Prin3	72668	315	.	-10.07727880101	6.8346784931143	-0.000676700531	0.5002597639041
WheelDrive	22516	50467	5
Cylinder	48242	24741	8
Volume	40994	31989	.	1.5	8.1	3.3096770259064	0.9580790808608
Doors	63322	9661	4
Type	68772	4211	13

Data Exploration

Summary and Visualization

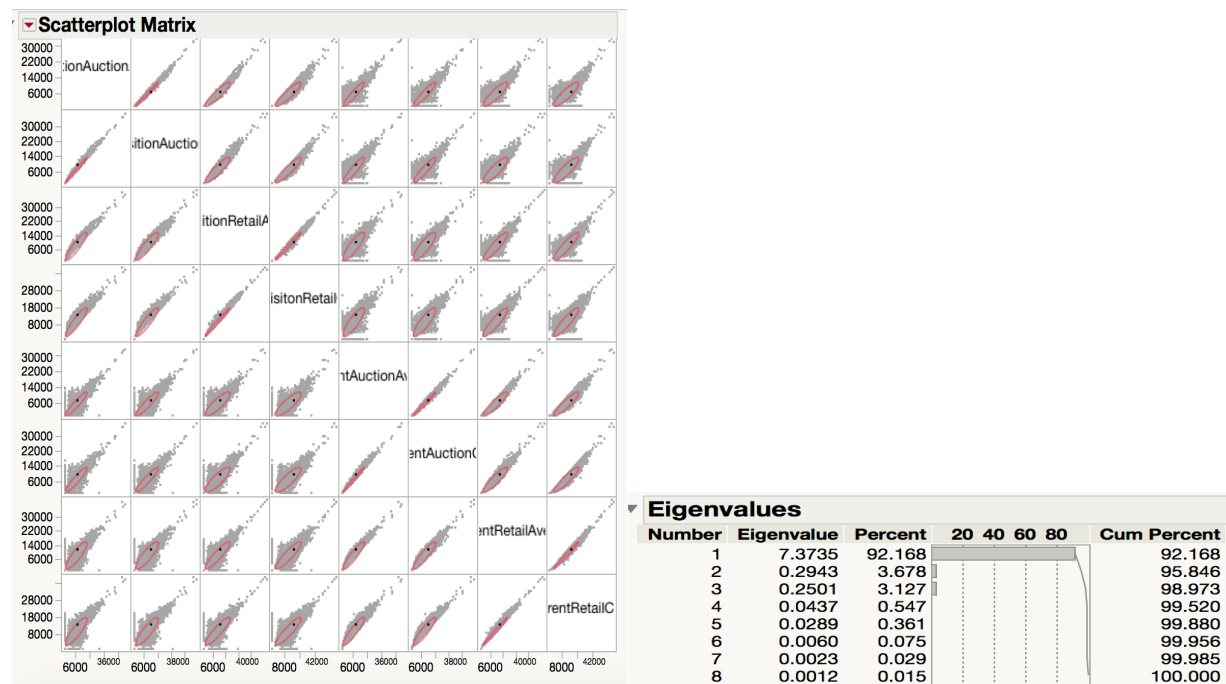
As we explored the data we noticed that there is a small proportion of cars labeled as a bad purchase. In the below graphs these cars are highlighted. You can see that these cars even though scattered all over they are mostly aged 4-6 years old. Looking at the 'Make' distribution we can see which brands of cars are mainly to be labeled as a bad buy. Some of these cars are Dodge, Chrysler, Ford and Chevrolet. Looking at the Vehicle Cost column distribution we notice how much are most of these bad purchases valued at. These values mostly vary between \$4,000-\$10,000.



Multicollinearity Issue and PCA

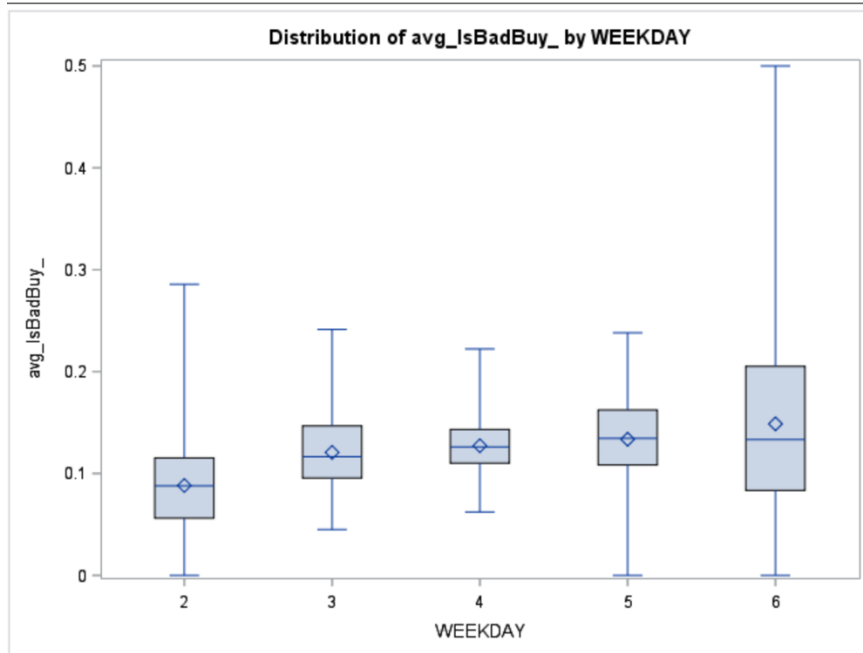
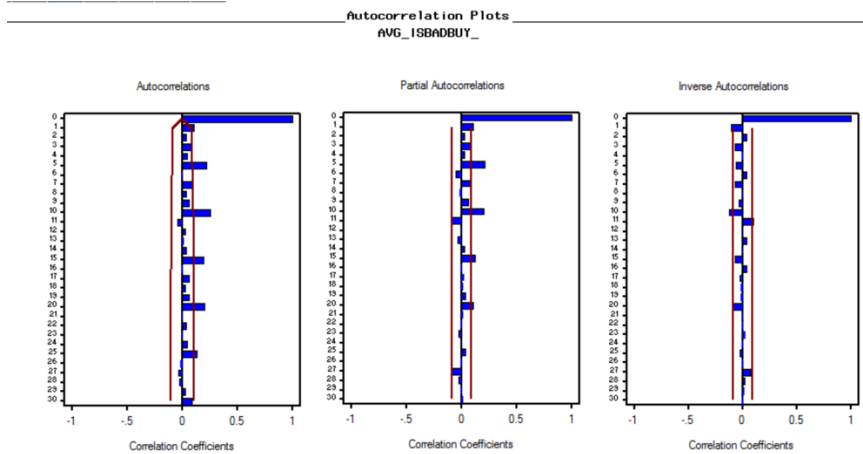
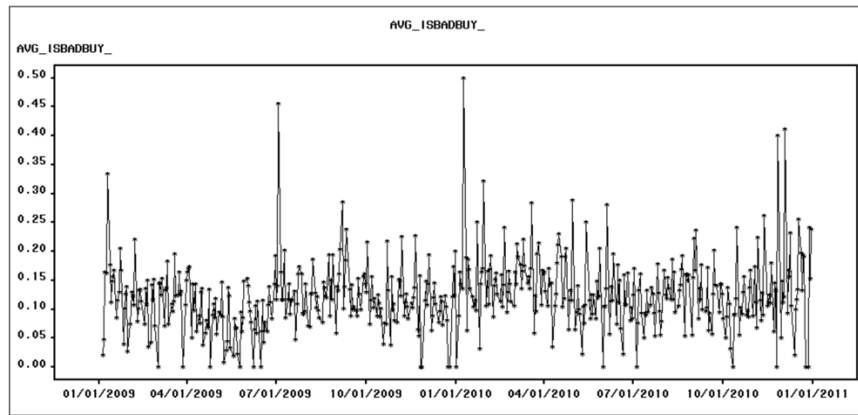
The 8 variables related to auction price and retail price in the dataset are highly correlated as you

can see below. A principal component analysis shows that the first PC alone explains 92% of the variation, and the first 3 components explain 99% of total variation. As a result we used the first three principle components instead of the 8 original variables for model building.



Time Series Pattern

We are interest in whether purchase date has any impact on the likelihood of a kick, and whether time series modeling is necessary. We used R to aggregate the data into daily proportion of bad buys for each day in the dataset, and performed time series analysis using SAS forecasting system. The plot of the series does not show any obvious trend of daily kick rate. There is a strong autocorrelation at lags of multiples of 5. Since the dataset mainly consists of workdays, strong autocorrelation at lag 5 means a weekly pattern. A boxplot by weekday also suggests our suspect of weekly seasonality. Specifically, bad purchases are least likely to be found on Mondays. Since this is beyond the scope of this course and dramatically complicates modeling, we will consider including time series analysis in the next stage. R code and SAS code are included in Appendix B and the aggregated daily kicked rates are included as attachments.



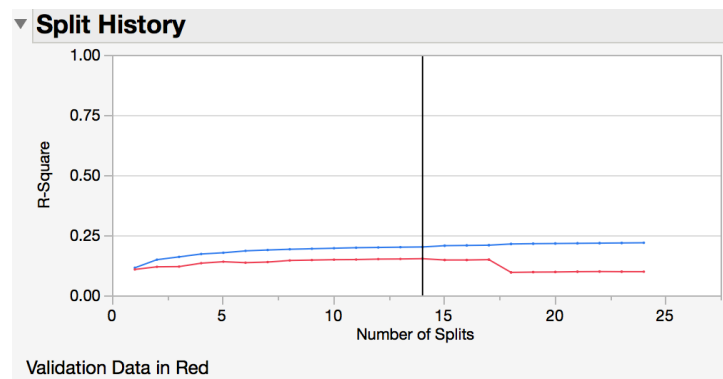
Data Mining Techniques selected

We tried 7 different types of supervised data mining techniques which we share here the significant results we got through some of these techniques. The model comparison section of this document goes into details as to which of the models are better and how. We partitioned the training set of 72,983 records into 60% training and 40% validation. Unless otherwise specified, all modeling is done using SAS JMP. The seven methods are:

- Logistic Regression
- Decision Tree — Partition
- Neural Networks
- Bootstrap Forest
- Boosted Tree
- K Nearest Neighbors
- Naïve Bayes

Decision Tree

The decision tree we got ended in 14 splits. The fit details show a misclassification rate of less than 10% and an RMSE of 0.29. We are focus on predicting IsBadBuy=1 and this model predicts IsBadBuy=0 very well.



Fit Details

Measure	Training	Validation	Definition
Entropy RSquare	0.2043	0.1555	$1 - \text{Loglike}(\text{model}) / \text{Loglike}(0)$
Generalized RSquare	0.2690	0.2080	$(1 - (L(0)/L(\text{model}))^{2/n}) / (1 - L(0)^{2/n})$
Mean -Log p	0.2978	0.3131	$\sum -\text{Log}(\rho_{ij})/n$
RMSE	0.2901	0.2960	$\sqrt{\sum (y_{ij} - \rho_{ij})^2/n}$
Mean Abs Dev	0.1683	0.1723	$\sum y_{ij} - \rho_{ij} /n$
Misclassification Rate	0.0995	0.0996	$\sum (\rho_{ij} \neq \rho_{\text{Max}})/n$
N	43873	29110	n

Confusion Matrix

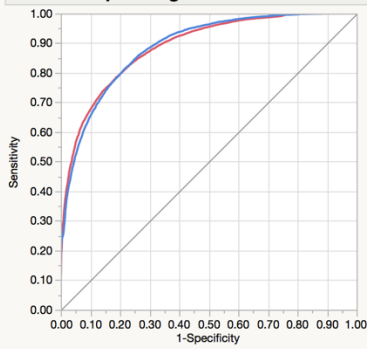
Training			Validation		
Actual	Predicted Count		Actual	Predicted Count	
IsBadBuy	0	1	IsBadBuy	0	1
0	38244	203	0	25405	155
1	4163	1263	1	2744	806

Column Contributions				
Term	Number of Splits	G^2		Portion
WheelType	1	3861.52314	<div></div>	0.5758
**Model	3	1748.00388	<div></div>	0.2606
Auction	1	394.82038	<div></div>	0.0589
VehicleAge	4	340.286521	<div></div>	0.0507
VNST	2	187.829799	<div></div>	0.0280
**PRIMEUNIT	2	101.892624	<div></div>	0.0152
Prin3	1	72.0613985	<div></div>	0.0107

Neural Networks

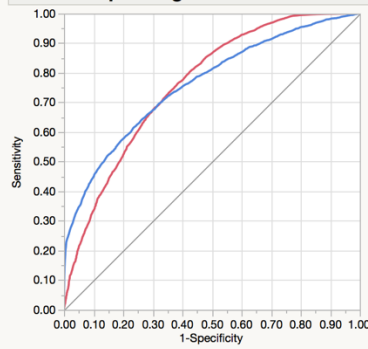
After trying different layers and nodes we reached a more significant result by applying 2 layers and 3 nodes in each layer. Below shows the confusion matrix we got in the results. This model predicts true negatives really well and if we were interested in our model predicting IsBadBuy=0 this would have been a good model but we are focusing on a true positive prediction which is IsBadBuy=1.

▼ Receiver Operating Characteristic



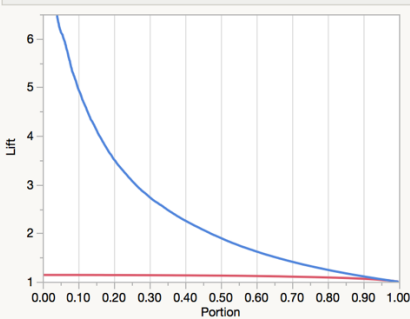
IsBadBuy	Area
0	0.8887
1	0.8887

▼ Receiver Operating Characteristic on Validation Data



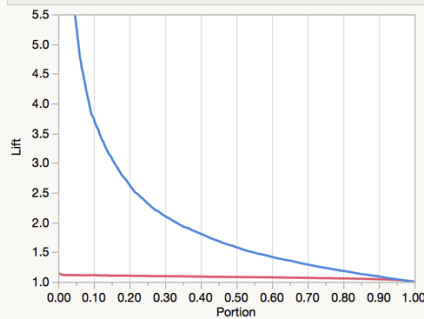
IsBadBuy	Area
0	0.7613
1	0.7613

▼ Lift Curve



IsBadBuy
0
1

▼ Lift Curve on Validation Data



IsBadBuy
0
1

▼ Specifications

Target Column:	IsBadBuy	Training Rows:	43813
Validation Column:	Validation	Validation Rows:	29170
		Test Rows:	0
Number of Trees in the Forest:	23	Number of Terms:	17
Number of Terms Sampled per Split:	4	Bootstrap Samples:	43813
		Minimum Splits per Tree:	10
		Minimum Size Split:	72

▼ Overall Statistics

Measure	Training	Validation	Definition
Entropy RSquare	0.2811	0.1743	$1 - \text{Loglike}(\text{model}) / \text{Loglike}(0)$
Generalized RSquare	0.3599	0.2317	$(1 - (L(0)/L(\text{model}))^{2/n}) / (1 - L(0)^{2/n})$
Mean -Log p	0.2683	0.3074	$\sum -\text{Log}(p[i]) / n$
RMSE	0.2787	0.2949	$\sqrt{\sum (y[i] - p[i])^2 / n}$
Mean Abs Dev	0.1679	0.1779	$\sum y[i] - p[i] / n$
Misclassification Rate	0.0985	0.1012	$\sum (p[i] \neq p_{\text{Max}}) / n$
N	43813	29170	n

▼ Confusion Matrix

Training			Validation		
Actual	Predicted Count		Actual	Predicted Count	
IsBadBuy	0	1	IsBadBuy	0	1
0	38203	214	0	25365	225
1	4103	1293	1	2728	852

Below we can see the details on our column contributors. Wheel type plays a significant role in our model and the second most significant column is Model. Overall 17 predictors contribute to our Bootstrap Forest model.

Column Contributions				
Term	Number of Splits	G ²		Portion
WheelType	461	2113.90643		0.3785
**Model	359	1163.28391		0.2083
VNST	707	372.342552		0.0667
VehBCost	636	280.531774		0.0502
*VehicleAge	710	271.727404		0.0487
Color	849	220.985423		0.0396
Make	617	193.472277		0.0346
Auction	517	192.382025		0.0344
WarrantyCost	866	182.694439		0.0327
Prin1	408	132.790766		0.0238
Size	562	128.411466		0.0230
Prin3	423	108.916939		0.0195
VehOdo	451	107.574215		0.0193
Prin2	378	83.9806732		0.0150
**Nationality	198	23.6918742		0.0042
Transmission	56	5.6594442		0.0010
IsOnlineSale	35	2.73955646		0.0005

Boosted Tree

We move on to Boosted Tree in hopes of a better result. As shown below the confusion matrix seems similar to Bootstrap Forest confusion matrix. Fit details show a misclassification rate of 10% and RMSE of 0.29. The lift curve also seems similar and it shows that the top 10% of our model does 3.4 times better than a random sample. ROC curve on the training dataset shows a much poorer performance than ROC curve in Bootstrap Forest.

▼ Boosted Tree for IsBadBuy

▼ Specifications

Target Column:	IsBadBuy	Number of training rows:	43813
Validation Column:	Validation	Number of validation rows:	29170
Number of Layers:	50		
Splits per Tree:	3		
Learning Rate:	0.1		
Overfit Penalty:	0.0001		

▼ Overall Statistics

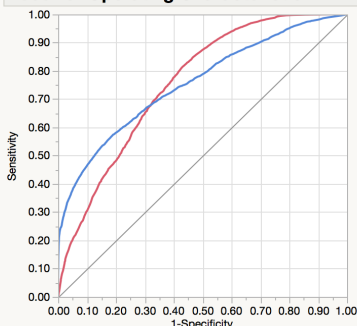
Measure	Training	Validation	Definition
Entropy RSquare	0.1851	0.1657	$1 - \text{Loglike}(\text{model}) / \text{Loglike}(0)$
Generalized RSquare	0.2454	0.2210	$(1 - (L(0)/L(\text{model}))^{(2/n)}) / (1 - L(0)^{(2/n)})$
Mean -Log p	0.3041	0.3106	$\sum -\text{Log}(p_{ij})/n$
RMSE	0.2922	0.2963	$\sqrt{\sum (y_{ij} - p_{ij})^2/n}$
Mean Abs Dev	0.1766	0.1797	$\sum y_{ij} - p_{ij} /n$
Misclassification Rate	0.0996	0.1026	$\sum (p_{ij} \neq p_{\text{Max}})/n$
N	43813	29170	n

▼ Confusion Matrix

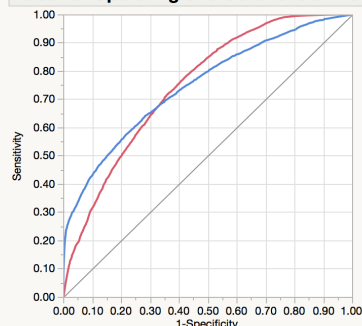
Training			Validation		
		Predicted Count			Predicted Count
Actual	IsBadBuy		Actual	IsBadBuy	
		0 1			0 1
0		38183 234	0		25369 221
1		4128 1268	1		2773 807

▼ Boosted Tree for IsBadBuy

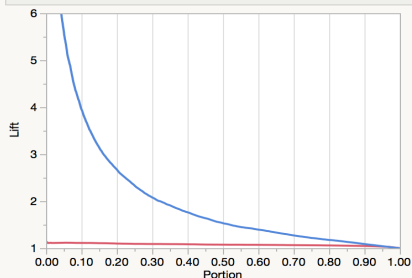
▼ Receiver Operating Characteristic



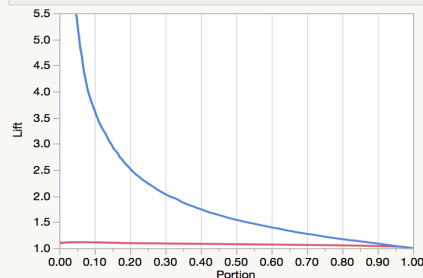
▼ Receiver Operating Characteristic on Validation Data



▼ Lift Curve



▼ Lift Curve on Validation Data



The column contributor seems somewhat different from our previous model. There are 11 contributors, the most important column is still wheel type but the second most important is no longer model, it is now vehicle age/year.

▼ Column Contributions				
Term	Number of Splits	G ²		Portion
WheelType	28	1230058.3	<div></div>	0.4093
VehYear	29	951455.132	<div></div>	0.3166
VehBCost	24	176419.068	<div></div>	0.0587
WarrantyCost	18	165425.099	<div></div>	0.0550
**Model	20	160411.237	<div></div>	0.0534
VNST	7	99379.8272	<div></div>	0.0331
Auction	11	98322.2086	<div></div>	0.0327
**TopThreeAmericanName	2	59508.7091	<div></div>	0.0198
Make	1	32765.2078	<div></div>	0.0109
Prin3	9	16174.181	<div></div>	0.0054
VehOdo	1	15635.317	<div></div>	0.0052
Color	0	0	<div></div>	0.0000
Transmission	0	0	<div></div>	0.0000
**Nationality	0	0	<div></div>	0.0000
Size	0	0	<div></div>	0.0000
IsOnlineSale	0	0	<div></div>	0.0000
Prin1	0	0	<div></div>	0.0000
Prin2	0	0	<div></div>	0.0000

K Nearest Neighbors

K Nearest neighbors is extremely slow to run and the classification result is not as satisfactory as bootstrap forest and boosted trees.

is also included because it is related to the measurement used in the Kaggle competition.

	Mis-classification	Sensitivity	AUC	Top Contributors
Logistic Regression	0.1073	0.2332	0.7383	<u>WheelType</u> , <u>Model</u> , <u>VehAge</u> , VNST, VehOdo, VehBCost
Decision Tree	0.0996	0.2270	0.7275	<u>WheelType</u> , <u>Model</u> , Auction, <u>VehAge</u>
Neural Networks	0.1250	0.283	0.7270	N/A
Bootstrap Forest	0.1012	0.24	0.7611	<u>WheelType</u> , <u>Model</u> , <u>VehAge</u> , VNST
Boosted Tree	0.1004	0.2251	0.7560	<u>VehAge</u> , <u>WheelType</u> , Primeunit, Auction, Top3American
K Nearest Neighbors	0.1244	0.0355	N/A	N/A
Naive Bayes	0.1596	0.3581	N/A	N/A

As shown above, there is no universal “best” model with best performance in all measurements. The best misclassification rate is provided by Decision tree, while best sensitivity is provided by Naïve Bayes and best AUC provided by Bootstrap forest. We need to further clarify the business goal to decide which measure to rely on and which data mining model to choose.

The Kaggle competition used Gini index as the scoring criterion, and the formula of Gini Index is:

$$\text{Gini Index} = (\text{AUC} - 0.5) * (1 - \text{probability}(\text{IsBadBuy}=1))^{[3]}.$$

Since our bootstrap forest gave the highest AUC, we ran bootstrap forest on the test set and submitted to Kaggle as a late submission. We received a score of 0.23732, ranking 129 among the total 570 submissions, or 22% in place, and the highest score is 0.26719. The test set pre-processed using the same methods, and the test set submission, are also attached to this report.



CARVANA

Don't Get Kicked!

Predict if a car purchased at auction is a lemon

\$10,000 · 570 teams · 6 years ago

[Overview](#)[Data](#)[Discussion](#)[Leaderboard](#)[Rules](#)[Team](#)[My Submissions](#)[Late Submission](#)

Your most recent submission

Name	Submitted	Wait time	Execution time	Score
Carvana-test-forest.csv	a few seconds ago	0 seconds	1 seconds	0.23732

Complete

[Jump to your position on the leaderboard](#) ▾

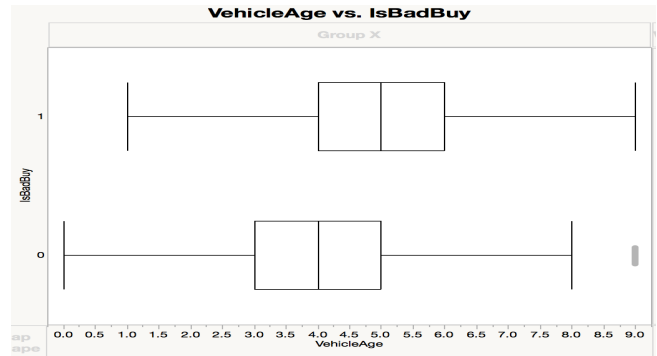
Conclusion

For the purpose of the Kaggle competition, our best model is bootstrap forest with AUC of 0.76. In reality, we expect businesses to care more about misclassification results and sensitivity rate. Since the cost of misclassifying a Kick and buys it outweighs the cost of misclassifying a good car and misses it, we recommend our Naïve Bayes model to Carvana, which has the highest sensitivity of capturing kicked purchases. However, even the highest sensitivity (36%) is not good enough to make improvement in business, and it is achieved by sacrificing misclassification rate.

All our data mining models agree on the below list of variables that contribute most to identify a bad buy: Wheel Type, Vehicle Age, Model, State, Purchase Cost. Based on the below tabulates, we recommend Carvana when making used car purchases:

- Avoid vehicles with no Wheel Type information, and old vehicles!
- Be cautious in certain states
- Be cautious with cheap deals
- Be cautious with certain model types

	IsBadBuy			
	0		1	
WheelType	Row %	N	Row %	N
Alloy	88.95%	32065	11.05%	3985
Covers	91.96%	30349	8.04%	2655
Null	29.52%	937	70.48%	2237
Special	86.89%	656	13.11%	99



The biggest challenges in this data mining project are:

- Messy data, with large proportion of missing values, cropped model names, and missing important vehicle features like engine information.
- Some models are not stable because of randomness in nature: neural networks provide different results when the model is run repeatedly, even with identical settings. Bootstrap forest and boosted tree provides different column contribution information when run repeatedly.
- Trade-off between different measurements: since these measurements do not agree on the same model, choosing the most important measurement is crucial to business.
- Slow algorithms: logistic regression is extremely slow when one or more of the categorical variables contain too many levels, like our "model" variable with 1000 levels. K-nearest neighbors method is extremely computationally intensive, which took us 5 hours to run in SAS JMP.

In our next stage, we expect to further improve our models by:

- Trying different cut-off probabilities to balance between misclassification rate and sensitivity rate.
- For records with cropped model name, search for extra data and supplement with information like engine, volume, etc.
- Perform time series analysis using purchase date.

Appendix A: Data Dictionary

Field Name	Definition
RefID	Unique (sequential) number assigned to vehicles
IsBadBuy	Identifies if the kicked vehicle was an avoidable purchase
PurchDate	The Date the vehicle was Purchased at Auction
Auction	Auction provider at which the vehicle was purchased
VehYear	The manufacturer's year of the vehicle
VehicleAge	The Years elapsed since the manufacturer's year
Make	Vehicle Manufacturer
Model	Vehicle Model
Trim	Vehicle Trim Level
SubModel	Vehicle Submodel
Color	Vehicle Color
Transmission	Vehicles transmission type (Automatic, Manual)
WheelTypeID	The type id of the vehicle wheel
WheelType	The vehicle wheel type description (Alloy, Covers)

VehOdo	The vehicle's odometer reading
Nationality	The Manufacturer's country
Size	The size category of the vehicle (Compact, SUV, etc.)
TopThreeAmericanName	Identifies if the manufacturer is one of the top three American manufacturers
MMRAcquisitionAuctionAveragePrice	Acquisition price for this vehicle in average condition at time of purchase
MMRAcquisitionAuctionCleanPrice	Acquisition price for this vehicle in the above Average condition at time of purchase
MMRAcquisitionRetailAveragePrice	Acquisition price for this vehicle in the retail market in average condition at time of purchase
MMRAcquisitionRetailCleanPrice	Acquisition price for this vehicle in the retail market in above average condition at time of purchase
MMRCurrentAuctionAveragePrice	Acquisition price for this vehicle in average condition as of current day
MMRCurrentAuctionCleanPrice	Acquisition price for this vehicle in the above condition as of current day
MMRCurrentRetailAveragePrice	Acquisition price for this vehicle in the retail market in average condition as of current day
MMRCurrentRetailCleanPrice	Acquisition price for this vehicle in the retail market in above average condition as of current day
PRIMEUNIT	Identifies if the vehicle would have a higher demand than a standard purchase

AUCGUART	The level guarantee provided by auction for the vehicle (Green light - Guaranteed/arbitratable, Yellow Light - caution/issue, red light - sold as is)
BYRNO	Unique number assigned to the buyer that purchased the vehicle
VNZIP	Zipcode where the car was purchased
VNST	State where the the car was purchased
VehBCost	Acquisition cost paid for the vehicle at time of purchase
IsOnlineSale	Identifies if the vehicle was originally purchased online
WarrantyCost	Warranty price (term=36month and millage=36K)

Appendix B: Time Series Codes

R code:

```
car <- read.csv(file='/Users/florencewu/Desktop/DNSC 6279 Data Mining/Project/Carvana - clean and PCA 4-12.csv')
```

```
ts <- read.csv(file='/Users/florencewu/Desktop/DNSC 6279 Data Mining/Project/IsBadBuy-Time.csv')
```

```
library(sqldf)
```

```
ts$PurchDate = as.Date(ts$PurchDate, "%m/%d/%y")
```

```
ts1 <- sqldf("select avg(IsBadBuy), PurchDate from ts group by PurchDate")
```

```
ts3 = ts(ts1)
```

```
write.csv(ts3, file='/Users/florencewu/Desktop/DNSC 6279 Data Mining/Project/IsBadBuy-TSW.csv')
```

SAS Code:

```
DATA NEW;  
  
SET WORK.CAR;  
  
TIME= _N_;  
  
WEEKDAY=WEEKDAY(PURCHDATE);  
  
RUN;  
  
PROC SORT DATA=NEW;  
  
BY WEEKDAY;  
  
RUN;  
  
PROC BOXPLOT;  
  
PLOT AVG_ISBADBUY_ * WEEKDAY;  
  
RUN;
```

Reference

[1] <https://www.kaggle.com/c/DontGetKicked/data>

[2] <https://www.carvana.com/>

[3] <https://www.kaggle.com/c/DontGetKicked/discussion/925>