

Midterm Project



Warm Up: Predict Blood Donations

Name: Young, Brent

DrivenData Name: bdy3176

DrivenData Public Score: 0.4371

MSDS 454 Section #: 55

Quarter: Summer 2018

Introduction

Problem

The purpose of the midterm project is to analyze blood donation data (*includes information about each donor's history*) from a mobile blood donation vehicle in Taiwan using machine learning methods to predict the probability that a donor made a donation in March 2007. The Blood Transfusion Service Center drives to various universities in the area and collects blood as part of a blood drive. As a result, we ultimately want to predict whether or not a donor will give blood the next time the vehicle comes to campus. However, predicting the probability a person will donate blood can be difficult given that the data provided only contains 4 predictor variables (e.g., Months since Last Donation, Number of Donations, Total Volume Donated, and Months since First Donation) and has only 576 records in the training dataset and 200 in the test dataset.

Significance

The problem is significant/interesting because good data-driven systems for tracking and predicting blood donations and supply needs can help improve the supply chain, ensuring that more patients receive the blood transfusions they need. This is particularly important given that blood remains a critical resource during emergencies, for sickle cell patients, cancer patients, and car accident victims. In fact, according to the American Red Cross, every two seconds someone in the U.S. needs blood, more than 41,000 blood donations are needed every day, and 30 million blood components are transfused each year in the United States.

Applicability to Data Scientists

This real-world data problem is applicable to data scientists because it serves as a great resource for practicing our data science skills and testing out various machine learning algorithms. For instance, since this is a classification problem, it'll provide data scientists the opportunity to practice classification modeling techniques such as Logistic Regression, LDA, QDA, Decision Trees, Bagging, Random Forest, Gradient Boosting Machines, and Neural Network.

Data Exploration

Structure and Description of Training, Validation, and Test Datasets

To help streamline the analysis, both the train and test datasets have been consolidated into one dataset called total. The combined dataset prior to feature engineering consists of 6 variables (includes identification and response variable) and 776 observations. The data was then split into a 48/26/26: train/validation/test so that the validation and test sets have the same amount of observations. For instance, there are 376 training observations (48%), 200 validation observations (26%), and 200 test observations (26%). The training set will be used to fit the models, the validation set will be used to estimate prediction error for model selection, and the test set will be used for assessment of the prediction error of the final chosen model.

Data Mining/ Cleaning

Cleaning

In order to make the data more R friendly, we then added ID to help identify distinct records.

ID	MSLD	NUM	VOLUME	MSFD	TARGET_FLAG	DPH	TENRAT
0	0	0	0	0	200	0	0
DF	MSLD_bin	REPEAT	SQRT_MSLD	SQRT_NUM	SQRT_VOLUME	SQRT_MSFD	SQRT_TENRAT
0	0	0	0	0	0	0	0
SQRT_DF	LOG_NUM	LOG_VOLUME	LOG_MSFD				
0	0	0	0				

We also renamed Months since Last Donation, Number of Donations, Total Volume Donated, Months since First Donation, and Made Donation in March 2007 (response variable) to MSLD, NUM, VOLUME, MSFD, and TARGET_FLAG, respectively. After conducting summary statistics on the dataset. The results showed that there were no missing values, except the 200 TARGET_FLAG NA's in the test set (see top right). As a result, missing value imputation was not conducted.

Summary of Variables, Feature Creation, Measurement Levels, and Standardization

ID, which is used for identification purposes will be ignored. TARGET_FLAG represents our classification response variable of whether a person made a donation in March 2007 (1 = Yes, 0 = No). There is a total of 90 people who donated in March 2007 in the training dataset and 48 who donated in March 2007 in the validation dataset. The other 4 variables are numeric and consist of Months since Last Donation, Number of Donations, Total Volume Donated, and Months since First Donation. Through feature creation, we also created three new numeric variables: Donations per Month (DPM = MSFD/NUM), Ratio of Months since Last Donation to First Donation (TENRAT = MSLD/MSFD), and Donation Frequency (DF = MSFD-MSLD/NUM)). We also created two qualitative variables: MSLD_bin and REPEAT. MSLD_bin creates groups for Months since Last Donation of 0 to 4, 5 to 8, 9 to 12, and 13+, while REPEAT identifies repeat donors (e.g., if a donor donated more than once = 1, else 0). Furthermore, SQRT transformations were conducted on MSLD, NUM, VOLUME, MSFD, TENRAT, and DF, while LOG transformations were conducted on NUM, VOLUME, and MSFD due to the skewness of the variables. Lastly, all variables, except, ID, TARGET_FLAG, MSLD_bin, AND REPEAT have been standardized to have a mean of 0 and standard deviation of 1 in the training, validation, and test datasets. Note: Standardization was conducted after EDA.

	1979	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100	2101	2102	2103	2104	2105	2106	2107	2108	2109	2110	2111	2112	2113	2114	2115	2116	2117	2118	2119	2120	2121	2122	2123	2124	2125	2126	2127	2128	2129	2130	2131	2132	2133	2134	2135	2136	2137	2138	2139	2140	2141	2142	2143	2144	2145	2146	2147	2148	2149	2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	2160	2161	2162	2163	2164	2165	2166	2167	2168	2169	2170	2171	2172	2173	2174	2175	2176	2177	2178	2179	2180	2181	2182	2183	2184	2185	2186	2187	2188	2189	2190	2191	2192	2193	2194	2195	2196	2197	2198	2199	2200	2201	2202	2203	2204	2205	2206	2207	2208	2209	2210	2211	2212	2213	2214	2215	2216	2217	2218	2219	2220	2221	2222	2223	2224	2225	2226	2227	2228	2229	2230	2231	2232	2233	2234	2235	2236	2237	2238	2239	2240	2241	2242	2243	2244	2245	2246	2247	2248	2249	2250	2251	2252	2253	2254	2255	2256	2257	2258	2259	2260	2261	2262	2263	2264	2265	2266	2267	2268	2269	2270	2271	2272	2273	2274	2275	2276	2277	2278	2279	2280	2281	2282	2283	2284	2285	2286	2287	2288	2289	2290	2291	2292	2293	2294	2295	2296	2297	2298	2299	2300	2301	2302	2303	2304	2305	2306	2307	2308	2309	2310	2311	2312	2313	2314	2315	2316	2317	2318	2319	2320	2321	2322	2323	2324	2325	2326	2327	2328	2329	2330	2331	2332	2333	2334	2335	2336	2337	2338	2339	2340	2341	2342	2343	2344	2345	2346	2347	2348	2349	2350	2351	2352	2353	2354	2355	2356	2357	2358	2359	2360	2361	2362	2363	2364	2365	2366	2367	2368	2369	2370	2371	2372	2373	2374	2375	2376	2377	2378	2379	2380	2381	2382	2383	2384	2385	2386	2387	2388	2389	2390	2391	2392	2393	2394	2395	2396	2397	2398	2399	2400	2401	2402	2403	2404	2405	2406	2407	2408	2409	2410	2411	2412	2413	2414	2415	2416	2417	2418	2419	2420	2421	2422	2423	2424	2425	2426	2427	2428	2429	2430	2431	2
--	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	---

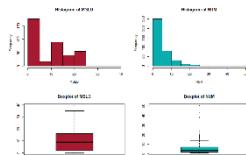


Figure 2

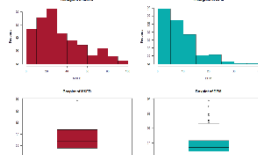


Figure 3

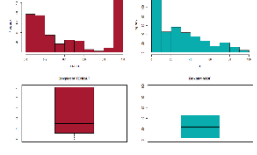
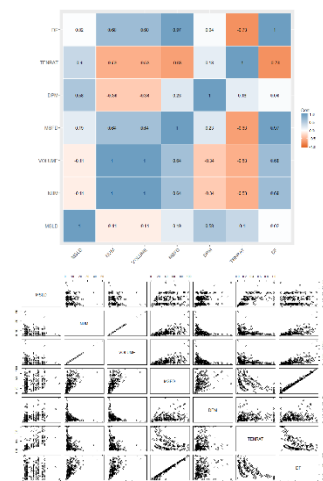


Figure 4

Multivariate Plots for Numeric Variables

Observations: Figure 5 shows a correlation and scatterplot matrix of the numeric variables DF, TENRAT, DPM, MSFD, VOLUME, NUM, and MSLD. This allows us to see which variables may be correlated with each other so that we can glean interesting insights. The plots show that VOLUME and NUM are perfectly positive correlated, which validates what we saw in our summary statistics. As a result, given the duplicative nature of these two variables, VOLUME was removed from the dataset. Furthermore, the plots revealed strong positive correlations between MSFD vs. DF, which makes sense given that the longer someone has been a donor for, the more frequently he/she will donate.



Figures 6a to 6c: Boxplot of MSLD, NUM, DF vs. TARGET_FLAG

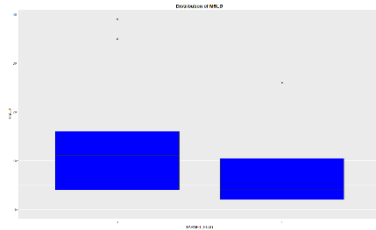


Figure 6a

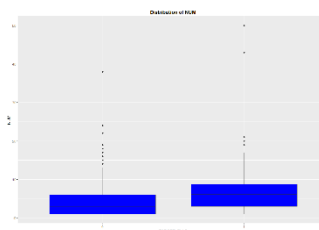


Figure 6b

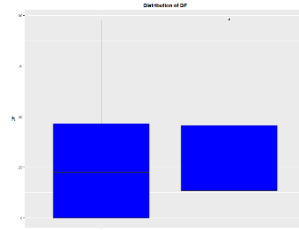


Figure 6c

Observations: Figure 6a shows a boxplot of MSLD vs. TARGET_FLAG, so that we can compare the median differences and variability between the numeric variable and TARGET_FLAG. The results show that the lower the number of months since the donor's most recent donation, the more likely he/she made a donation in March 2007 (vice versa). Figure 6b shows a boxplot of NUM vs. TARGET_FLAG. The results show that as the number of donations that the donor has made increases, the more likely he/she made a donation in March 2007 (vice versa). Figure 6c shows a boxplot of DF vs. TARGET_FLAG. The results show that as donation frequency increases, the more likely he/she made a donation in March 2007 (vice versa). As a result, this provides evidence that MSLD, NUM, and DF are strong predictors to include in our models since the median difference between whether a person made a donation in March 2007 (1 = Yes, 0 = No) is wide. I also conducted additional boxplots of TARGET_FLAG (x-axis) vs. the other numeric variables (y-axis) and found that DPM and TENRAT have median differences between TARGET_FLAG, while MSFD did not. However, these variables do not seem as strong as the predictors mentioned above since they either had a lot of variability or there were subtle differences between whether a person made a donation in March 2007 (1 = Yes, 0 = No). This was also validated using the variable importance feature in the Caret package.

Multivariate Plots for Qualitative Variables

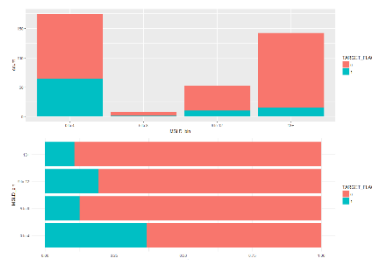


Figure 7a

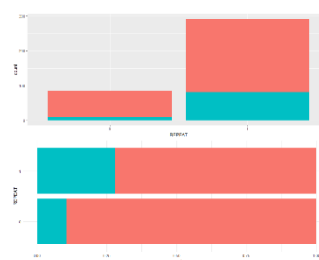


Figure 7b

Observations: Figures 7a and 7b shows bar plots of MSLD_bin and REPEAT vs. donors who made a donation in March 2007 (blue) and those who did not (red). The data shows that donors who fell between 0 to 4 months in regards to the number of months since their most recent donation, donated in March 2007 more than the other categories ranges (36.8%) (e.g., 5 to 8, 9 to 12, and 13+ months). For instance, those who were in the 13+ month range donated only 10.6% of the time. Furthermore, the data shows that those who were repeat customers (27.8%) donated in March 2007 more than those who were not (10.6%). For example, 81 out of the 90 people who donated in March 2007 were repeat customers.

Reviews of Literature & Formulation of Models

Reviews of Literature (see last page for references)

There were many peer reviewed journals in the NU library database that used logistic regression, LDA, QDA, decision trees etc. to predict the probability that a donor made a donation within a specified time period or predicted the probability on a closely related subject area. For instance, in the *Transfusion Medicine* (2000), Flegel, Besenfelder, and Wagner use logistic regression to calculate the probability that someone will donate blood within a preselected time frame (e.g., 6-9 months after an index donation). Interestingly, first-time donors had a donation probability of 33% and were more likely to return than repeat donors. Second, in the *Environmental Pollution* (2018), Wang, Li, Ma, Li, Wang, Huang, Xu, Chenzi; and An, Yi use QDA, logistic regression, and decision trees to predict the probability of cadmium pollution (Cd) in rice in China. The results showed that the accuracy rate of 74% with QDA was significantly higher than the decision tree (67%) and logistic regression (68%) models. Lastly, Bhoi, Sherpa, and Khandelwal use LDA and decision trees to predict the probability of cardiovascular diseases such as myocardial ischemia and cardiac arrhythmias.

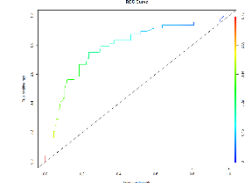
Modeling Strategy

Given that the goal of the DrivenData competition Warm Up: Predict Blood Donations is to predict the probability that that a donor made a donation in March 2007, I began my analysis using linear classification techniques such as logistic regression (with stepAIC) and LDA to serve as initial baselines prior to conducting more sophisticated modeling techniques. Additionally, given that this data set is the smallest and least complex dataset on DrivenData (e.g., only 4 variables such as Months since Last Donation, Number of Donations, Total Volume Donated, and Months since First Donation are included), I felt like linear classification modeling techniques would perform admirably. After building logistic regression and LDA models, I then moved to QDA, Decision Trees, Bagging, Random Forest, Gradient Boosting Machines, and Neural Network. The next section provides a summary of my results for each modeling technique.

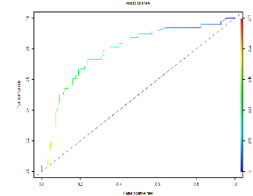
Application of Tools

Logistic Regression: Logistic regression models the probability that the response variable belongs to a specific category and assumes a linear decision boundary (James, Witten, Hastie, & Tibshirani, 2013). For instance, it models the probabilities of the K classes using linear functions in x, while also ensuring that they sum to 1 and remain in-between 0 and 1 (Hastie, Tibshirani, & Friedman, 2009). This is accomplished using the logistic function and maximum likelihood, which is used to fit the model (James, et al., 2013). As a result, using the glm function, we produced a logistic regression model of $\text{TARGET_FLAG} \sim \text{MSLD} + \text{LOG_NUM} + \text{DF}$. These variables were chosen using stepwise regression (stepAIC). The Analysis of Deviance table showed that all the variables that were included in this model were statistically significant, which illustrates that these variables improved the model. Additionally, the Analysis of Deviance table and varimp showed that LOG_NUM, followed by MSLD, and DF impacted the model the most. In regards to the coefficients of the model, variables such as LOG_NUM, which

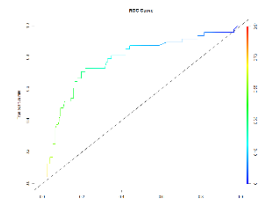
has a positive coefficient and MSLD, which has a negative coefficient make intuitive sense and were statistically significant. For instance, as the number of donations that the donor has made increases, the more likely he/she will donate (vice versa). Additionally, the lower the number of months since the donor's most recent donation, the more likely he/she will donate (vice versa). The model produced the following performance metrics on the training dataset: AIC: 361.6423, BIC: 377.3606 and the following accuracy metrics and evaluation criteria on the validation dataset: accuracy: 0.77, AUC: 0.785636 (see ROC curve on the right), and LogLoss: 0.4608874. Note: We also tried logistic regression GAM, but the performance was exactly the same as a standard logistic regression model (e.g., `model.gam1 <- glm(TARGET_FLAG ~ MSLD + s(NUM,1) + s(MSFD,50)+s(TENRAT,1))`).



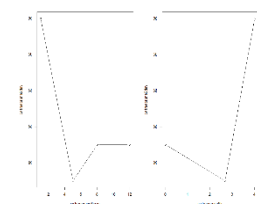
Linear Discriminant Analysis: LDA is very similar in form to logistic regression (distributions are assumed to be normal), except it models the distribution of the predictors separately in each of the response classes and then applies Bayes theorem (James, et al., 2013). This model also uses Gaussian densities, which arises when we assume that the classes have a common covariance matrix, and assumes a linear decision boundary (Hastie, et al., 2009). Using the `lda` function, we then produced a linear discriminant analysis model of MSLD + LOG_NUM + DF. These variables were chosen using the same variables from the logistic regression model (with `stepAIC`). The model also produced the following accuracy metrics and evaluation criteria: accuracy: 0.77, AUC: 0.785636 (see ROC curve on the right) and LogLoss: 0.4601341. LDA performed similarly to logistic regression.



Quadratic Discriminant Analysis: QDA is similar to LDA, in which the QDA classifier results from assuming the observations in each class are drawn from a Gaussian distribution and then Bayes theorem is applied to perform prediction (James, et al., 2013). However, unlike LDA, QDA assumes that each class has its own covariance matrix and assumes a quadratic decision boundary (James, et al., 2013). Using the `qda` function, we then produced a quadratic discriminant analysis model of MSLD + LOG_NUM + DF. These variables were chosen using the same variables from the logistic regression model (with `stepAIC`). The model also produced the following accuracy metrics and evaluation criteria: accuracy: 0.795, AUC: 0.7783717 (see ROC curve on the right), and LogLoss: N/A. QDA had higher accuracy than logistic and LDA, but slightly lower AUC.



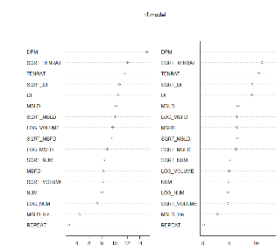
Decision Tree: A decision tree is a tree-based method that involves stratifying or segmenting the predictor space into a number of simple regions (James, et al., 2013). For instance, predictions are made by assigning an observation in a given region to the most common occurring class of training observations in that region (James, et al., 2013). Using the `tree` function, we produced a decision tree model with 5 predictor variables, after cross-validation helped eliminate 12 predictor variables (see plot on the right). The model produced the following accuracy metrics and evaluation criteria: accuracy: 0.76,



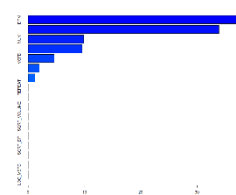
AUC: 0.6761239, and LogLoss: 0.5129595 on the validation dataset. The performance was worse than logistic regression, LDA, and QDA.

Bagging: Bagging is a technique for reducing the variance of an estimation prediction function. For classification, a committee of trees each cast a vote for the predicted class (aka: majority vote) (Hastie, et al., 2009). As a result, using the randomForest function (mtry=17, ntree=100), a bagged decision tree model was produced using all 17 predictor variables. The model produced the following accuracy metrics and evaluation criteria: accuracy: 0.755, AUC: 0.6057429, and LogLoss: N/A on the validation dataset. Interestingly, the performance was worse than logistic regression, LDA, QDA, and a basic decision tree.

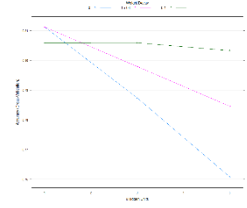
Random Forest: Random forest provides an improvement over bagged trees by incorporating a small tweak that decorrelates the trees and then averages them (e.g., forces each split to only consider a subset of predictors and will not consider strong predictors so that other predictors will have more of a chance (James, et al., 2013)). As a result, using the randomForest function (mtry=4, ntree=400), a random forest model was produced using all 17 predictor variables. Using the variable importance function (which measures prediction strength), the plot on the right shows that DPM, SQRT_TENRAT, and TENRAT are the most important variables. This is different than what we saw earlier using stepAIC. The model produced the following accuracy metrics and evaluation criteria: accuracy: 0.78, AUC: 0.683114, and LogLoss: N/A on the validation dataset. This was an improvement over a basic decision tree and bagging. The accuracy score was also the second highest compared to all the models we've fitted thus far. However, AUC was dramatically lower than logistic regression, LDA, and QDA.



Gradient Boosting Machines: Boosting provides another approach for improving the predictions resulting from a decision tree by fitting each tree on an altered version of the original dataset (James, et al., 2013). In other words, trees are grown sequentially (e.g., each tree is grown using information from previously grown trees). As a result, using the gbm function, a boosted decision tree model was produced using all 17 predictor variables. I also incorporated n.trees =50, shrinkage=0.1, and depth=1, which was determined using a grid search. Relative importance showed that DPM and TENRAT are the most important variables, similar to what was seen in randomForest (see plot on the right). The model produced the following accuracy metrics and evaluation criteria: accuracy: 0.765, AUC: 0.6387061, and LogLoss: 0.4992688 on the validation dataset. However, according to these metrics and evaluation criteria, boosting performed inadequately compared to the other models.



Neural Network: Neural network (aka: single hidden layer back-propagation network) is a nonlinear statistical model that is basically a nonlinear generalization of a linear model (Hastie, et al., 2009). It contains inputs, a hidden layer, and outputs that are typically represented by a network diagram (Hastie, et al., 2009). Additionally, neural network has unknown parameters called weights that introduce nonlinearities where needed (Hastie, et al., 2009). Using the nnet function, we then produced a neural network model of MSLD + LOG_NUM + DF. These variables were chosen using the same variables from the logistic regression model (with stepAIC). Additionally, using the variables from the logistic regression model makes sense given that neural network is essentially a bunch of logistic regressions, fed into a multinomial logit model. I also incorporated 1 hidden layer into the model with a decay=1e-04 and maxit=1000, which was determined using a grid search (see plot on the right). The model produced the following accuracy metrics and evaluation criteria: accuracy: 0.77, AUC: 0.7871436, and LogLoss: 0.4579446 on the validation dataset. As a result, neural network performed better than all the tree-based methods and had similar performance to logistic regression, LDA, and QDA.



Performance/Accuracy of Classification Models on Validation Set & DrivenData

Model Name	Accuracy	AUC	LogLoss
Logistic Regression	0.77	0.785636	0.4609
Linear Discriminant Analysis	0.77	0.785636	0.4601
Quadratic Discriminant Analysis	0.795	0.778372	N/A
Decision Tree	0.76	0.676124	0.5130
Bagging	0.755	0.605743	N/A
Random Forest	0.78	0.683114	N/A
Gradient Boosting Machines	0.765	0.638706	0.4993
Neural Network	0.77	0.787144	0.4579

Figure 8

Submissions

BEST	CURRENT RANK	# COMPETITORS	SUBS. TODAY
0.4371	158	4592	0 / 3

EVALUATION METRIC

Log loss = $-\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$
The metric used for this competition is logarithmic loss. \hat{y}_i is the probability that $y = 1$. Logarithmic loss provides a steep penalty for predictions that are both confident and wrong. The goal is to minimize the log loss.

Figure 9

Observations: Figure 8 shows the performance/accuracy metrics and evaluation criteria for the following models in the validation dataset: logistic regression, LDA, QDA, decision trees, bagging, random forest, boosting, and neural networks. The results show that QDA had the highest accuracy, while bagging had the lowest accuracy out of all the models. Additionally, neural network, logistic regression, LDA, and QDA had similar AUC around 0.77 to 0.78. Given that QDA had the highest accuracy and a high AUC, I applied the model to the test dataset, and then submitted it to DrivenData. The LogLoss score was 0.4371, which currently places me in the top 3% out of 4,592 competitors (as of 7/18/18) (figure 9).

Conclusion

Future Work

In regards to future work, there are three areas that could help improve my models. First, it would be beneficial to obtain additional data and variables such as gender, age, ethnicity, education level, job status, job type, college/university, donation region/city, climate/temperature variables, income (individual and family/parents), average time spent waiting prior to donating at the university, and average dollars spent on advertising at the university. Second, it could be helpful exploring different mixing of models using an ensemble approach or model averaging since model diversity can help increase accuracy and performance. Lastly, it could be helpful to explore different transformations of the variables, interactions, and other R packages.

Learnings

In the end, I learned three primary things from building these models. First, I learned how to build different classification and prediction models using various methods. Second, I learned that it's really important to conduct a thorough EDA and that a lot can be learned from it. Lastly, I learned that trying different modeling approaches can result in better performance and to never settle on a model due to gut instinct. For instance, I found it interesting that vanilla approaches such as logistic regression, LDA, and QDA performed better than more sophisticated modeling techniques such as Bagging, Random Forest and GBM.

References

1. Flegel, W.A., Besenfelder, W., & Wagner, F.F. (2000). Predicting a donor's likelihood of donating within a preselected time interval. *Transfusion medicine*, 10(3), 181-92. [Peer Reviewed Journal].
2. Wang, X., Li, X., Ma, R., Li, Y., Wang, W., Huang, H., Xu, C., and An, Y. (2018). Quadratic discriminant analysis model for assessing the risk of cadmium pollution for paddy fields in a county in China. *Environmental Pollution*, 236, 366-372. [Peer Reviewed Journal].
3. Bhoi, A.K., Sherpa, K.S., and Khandelwal, B. (2015). Classification Probability Analysis for Arrhythmia and Ischemia Using Frequency Domain Features of QRS Complex. *International Journal Bioautomation*, 19(4), 531-542. [Peer Reviewed Journal].
4. James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning with Applications in R*. New York: Springer Science + Business Media.
5. Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*. New York: Springer Science + Business Media.