

Tumor Penetration Rates for Prostate Cancer

November 12, 2020

Executive Summary: With the common prostate cancer diagnostic among men, the goal of this analysis is to find the variable that could predict the chances of a tumor has penetrated the prostatic capsule. A dataset was provided by the Ohio State University Comprehensive Cancer Center to analyze how demographic and screening results could cause the likelihood of penetration. A tool was built using multiple logistic regression to predict the chances of a patient getting diagnosed with prostate cancer using three variables: results of digital rectal exam, Prostatic Specific Antigen value, and total Gleason score. This model was proved to be the best model with an AUC of 0.82, an accuracy of 0.77, Sensitivity of 0.75, and Specificity of 0.79. Some interesting findings include a 269% increase in chances of having prostate cancer when a patient has a unilobar nodule from the left on their results of digital rectal exam compared to patient with no nodule. Furthermore, patients who have bilobar nodule results after their digital rectal exam will have a higher chance of having prostate cancer compared to patient with no nodule, 260% more to be exact. In addition, patients could gain 133% more likelihood of being diagnosed with prostate cancer when they have a higher Gleason score. Knowing these three variables (results of digital rectal exam, Prostatic Specific Antigen value, and total Gleason score) are significant factors of the likelihood of penetration rate in prostate cancer, health care professionals would be more cautious in treating each patient. In addition, patients could detect their prostate cancer stages earlier and seek help accordingly.

Introduction: Prostate cancer is one of the most common cancers among men. The survival rate is very dependent on the stages of the cancer. Prostate cancer has a high probability of being treated if detected at early stages. There are many important characteristics that can predict whether a tumor has penetrated the prostatic capsule. With the results of penetration rate, doctors and nurses would prioritize the patients with the high risk level. On the other hand, if the patients have low rate in penetration, health care workers could find the best treatments for them to slow down prostate growth rate or even treat the patients. With this being said, hospitals can plan their budget and supplies to focus more on high risk patients and minimize supplies to low risk patients. Knowing penetration rates does not only benefit doctors and nurses, it would also help all men out there to detect if their tumors have penetrated the prostatic capsule. This way, patients would seek professional help at an earlier stage to increase their survival rate. Many questions have been proposed in favor of these issues including 1) What variables cause the penetration rate to increase or decrease? 2) What are the likelihoods of penetration rate? 3) How accurate is the model proposed in this analysis? 4) How can we implement the model for everyone to use? and 5) How can we improve the current model? This report will answer all of the questions mentioned above by crucial data analysis processes such as data exploratory analysis process to explore the relationship of each variable, model development using multiple logistic regression to build model, and diagnostics process to check the quality of the model. This model will serve the purpose of predicting penetration rate for prostate cancer in men to minimize human errors and provide an effective tool for health care workers.

Methods: A subset of data was provided from the Ohio State University Comprehensive Cancer Center to find the variables to predict if a tumor has penetrated the prostatic capsule. This data set contain a total of 380 subjects, out of which, 153 subjects had a cancer that penetrated the capsule. There are three observations contain missing values. After eliminating missing values to build more quality models, the dataset now have a total of 377 observations. These data contains one identification column with numbers of each subject, two demographic columns including age and race of the subjects, and four health measurements such as the results of digital rectal exam, detection of capsular involvement, prostatic specific antigen value, and total gleason score. Health measurement variables will be explained in details in the exploratory data analysis section. The goal of this analysis is to find the average tumor penetration of prostatic capsule rate. To accomplish this goal, we will first analyze the relationship between each predictor and target variable, build models using bootstrapping method, consider interaction terms to be in the final model, check for the quality of the models, and conclude with an interpretaion of the model output. All analysis will be done in R Studio with R version 3.6.2.

Exploratory Data Analysis: Before picking variables for our model development process, it is important to access the relationship of each response variable with the target variable. Table A2 shows an descriptive statistics of each independent variable. As seen in this table, the maximum value of Prostatic Specific Antigen value are relatively high compare to the median and mean values suggesting there might be outliers. In addition, there are disproportion weights of each factor in race of the subjects and detection of capsular involvement. This indicates that there might be biased in predicting the outcome if we use these two variables in the models. The rest of the variables seems proportionally balanced and no skewness were detected among them. In terms of the tumor penetration of prostatic capsule, the response variable has a balance proportion of 40% of tumor has penetrated the prostatic capsule out of 377 subjects. Figure 1.1 show the histogram of chances of the tumor penetration of prostatic capsule (1 = penetration, 0 = no penetration).

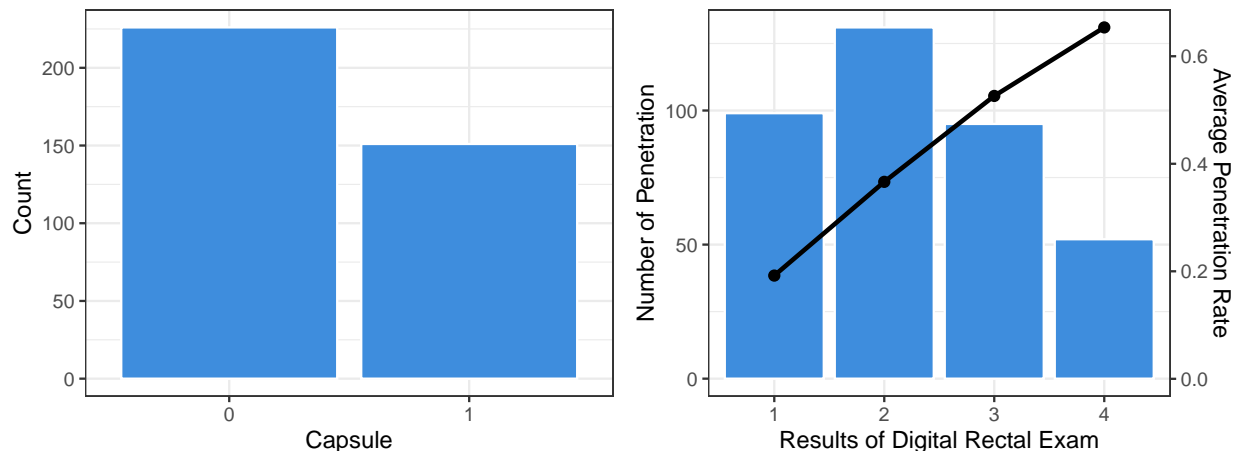


Figure 1: Histogram of 1 (Left): Tumor Penetration of prostatic capsule (1 = penetration) and 2 (Right): Results of Digital Rectal Exam (1 = no nodule, 2 = unilobar nodule left, 3 = unilobar nodule right, 4 = bilobar nodule) with average penetration rate in black.

The results of digital rectal exam is a simple exam to check the prostate to determine the size of the prostate and any abnormalities. Clearly, different abnormalities would have different chances

of penetration of prostatic capsule. Figure 1.2 shows four categories of the results of digital rectal exam including no nodule, unilobar nodule (left), unilobar nodule (right), and bilobar nodule not using number from 1 to 4, respectively. Here, this variable is identified as categorical variables with four levels instead of a continuous variable. As seen in plot, these four categories obviously have different average penetration rates. For example, a subject with no nodule would be less likely to penetrate compare to a person with bilobar nodule. It makes sense because if there is no nodule detected in the prostate area, then there is less chances of the tumor will penetrate the prostatic capsule. In addition, Figure 1.2 also depicts the number of subject in each category. Unilobar nodule (left) is the most common and bilobar nodule is the least common among prostate cancer. However, the proportion of each factor for digital rectal exam are not too different. Thus, this variables would provide significant impact for the models.

Detection of capsular involvement can describe if the tumor involving or extending beying the prostate capsule. The tumor is involved when there is capsular involment present. Clearly, the patients who have capsular involment present can be more risky than patient who do not have capsular involment. Figure 2.1 shows the count and proportion of penetration for the detection of capsular involment. This variable has two levels (1 = no, 2 = yes) to suggest if a subject has capsular involvement or not. As seen in the figure, patients with capsule involment would be much more likely to have tumor penetrated of prostatic capsule compare to patients that do not have capsule involvement. However, the count for each level is not proportional. The majority of the subject have level 1, no capsular involment. Therefore, this variable could be a good variable for our model since it shows a strong signal, however, cautious about the disproportional need to be consider before finalizing the model.

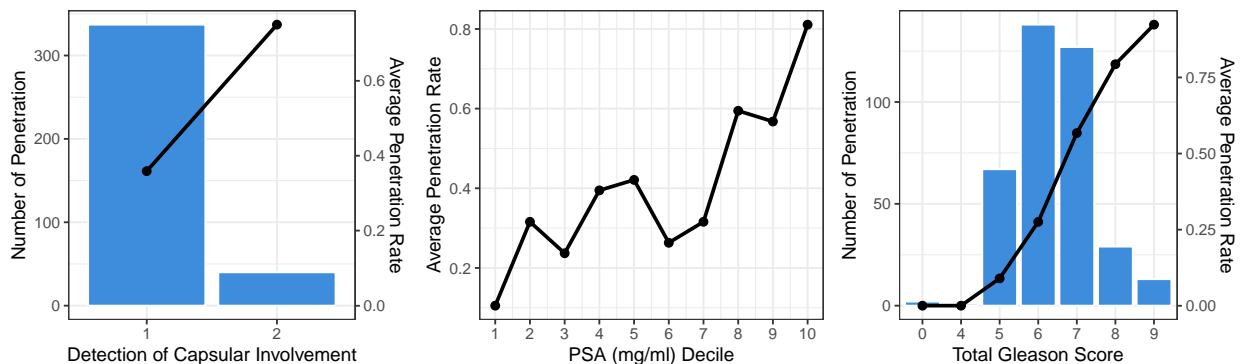


Figure 2: Plots of each variable with average penetration rate in black of 1 (Left): Dectation of capsular involvement along with the frequency of each factor, 2 (Center): Prostatic Specific Antigen value with mg/ml unit by decile, and 3 (Right): Total Gleason score along with its histogram.

Prostatic Specific Antigen (PSA) can be found in the prostate gland cells. This screening test measure the amount of prostate-specific antigen in patients' blood using miligram per milliliter (mg/ml). Higher PSA level could lead to a high posibility of diagnosing with prostate cancer. Figure 2.2 show a plot of the relationship between PSA and the average penetration of prostatic capsule rate. Here, the data point is being sorted from smallest to highest, divided into 10 equal buckets, and assinged a decide for visualization purposes. As on can see, as the decile number increases by one unit, the chances of penetration rate increases, on average. Since PSA in higher buckets are generally the largest, they have higher risk of getting prostate cancer than PSA in lower

buckets. Therefore, this proves that PSA could provide significant relationship for the likelihood of tumor penetrated of prostatic capsule.

Last but not least, total gleason score can measure the abnormality of the cells to see how the cells are arranged using a scale of 1 to 10. Higher scores lead to higher possibility of have penetration in prostatic capsule. On the other hand, cancer cells that looks similar to normal cells can be consider as low risk and have lower scores. Figure ??3 show a histogram of total gleason score overlayed by the average risk of having cancer. Total gleason score seems to have a normal distribution with less subjects having low or high scores and most subjects having median score of 6. In addition, the average penetration rate is increasing as the total gleason score is higher. This mean that if a patient's cancer cells looks similar to normal cells, they have less chances of being diagnosed with prostate cancer. On the other hand, if the gleason scores are high, they would have high risk of being diagnosed with prostate cancer. Therefore, this variable can also provide excellent signal for the response variable. Other variables did not mention above have little to no relationship with tumor penetration including age and race of the subjects.

After accessing the relationship of predictors and response variable, it is important to check the correlation between predictors. Fortunately, all correlation between predictors are less than 40%. Total gleason score and PSA values have the highest correlation score of 38.60% which is considered as low. This indicate that all variables can be in the same models without raising the multicollinearity issues. However, to make sure, we will check VIF scores after the modeling process. Therefore, we will go ahead and build our model using four variables including results of digital rectal exam ("dpros"), detection of capsule involvement ("dcaps"), PSA value ("psa"), and total Gleason score ("gleason").

Model Fitting/Inferences: After four variables were detected to be considered as significant, logistic regression will be used in model development process. Bootstrapping method is a common method for variable selection using random partitions. Here, the models are generated using a combination of each variable using different random partition of training and testing data. An average AUC and AIC of each model was computed after to compare. Table A3 shows the top 10 models with the highest AUC scores. Models with three variables results of digital rectal exam, PSA values, and total Gleason scores seems to be the best model with overall AUC above 80% and low AIC scores compares to other models. Therefore, the model development process will continue using the three variables mentioned previously with a train set of 70% and hold-out test set of 30%.

Before finalizing the model, all combination of interaction terms were added to the model with three variables including results of digital rectal exam, PSA values, and total Gleason score. The results of this model are shown in Table A4. The coefficient of all three variables seems to increase in the model with interaction terms. This proves that the model with interaction terms coefficient has steeper slope than the model without interaction terms (see Table 1). However, p-values for all interaction terms are quite high suggesting that interaction terms are not needed in the model. To further prove this argument, a stepwise model selection was run using the model with interaction terms using AIC as a validation metric. The best model seems to be the model with the original three variables: results of digital rectal exam, PSA values, and total Gleason score. As seen in Table 1, the final model's coefficients are positive consistent with the exploratory analysis. In

addition, all variables in this model have low p-values of less than 0.05 and odd ratio confident interval does not contain 1 suggesting all variables are significant. Finally, this model has a low AIC of 248.56, high AUC of 0.82, and high accuracy of 0.77 consistent with high probability of producing accurate predictions.

term	estimate	std.error	statistic	p.value	OR	2.5 %	97.5 %
(Intercept)	-7.254	1.296	-5.595	0.000	0.001	0.000	0.008
dpros2	0.822	0.461	1.785	0.074	2.276	0.944	5.835
dpros3	1.306	0.483	2.700	0.007	3.689	1.466	9.892
dpros4	1.280	0.577	2.220	0.026	3.598	1.176	11.452
psa	0.041	0.015	2.760	0.006	1.042	1.015	1.075
gleason	0.846	0.200	4.231	0.000	2.331	1.599	3.514

Table 1: Summary regression of final model with three variables: results of digital rectal exam (dpros), PSA scores (psa), and total Gleason scores (gleason). The columns represent the features names, coefficient estimated values, standard errors, test statistics, significant p-value, odd ratios, 95 percent lower confident interval of odd ratios, and 95 percent upper confident interval of odd ratios, respectively from left to right.

To access how good the model fit, predictions were generated using the hold-out test set. Figure 3.1 show a lift chart of empirical and indicated likelihood of penetration of prostate capsule. Here, the predicted proportion were sorted from smallest to largest, divided into 10 equal buckets, and assigned a decile number from 1 to 10. The lift chart shows that the average predicted proportion of each bucket are similar to the average empirical values. This indicates that the model is predicting well.

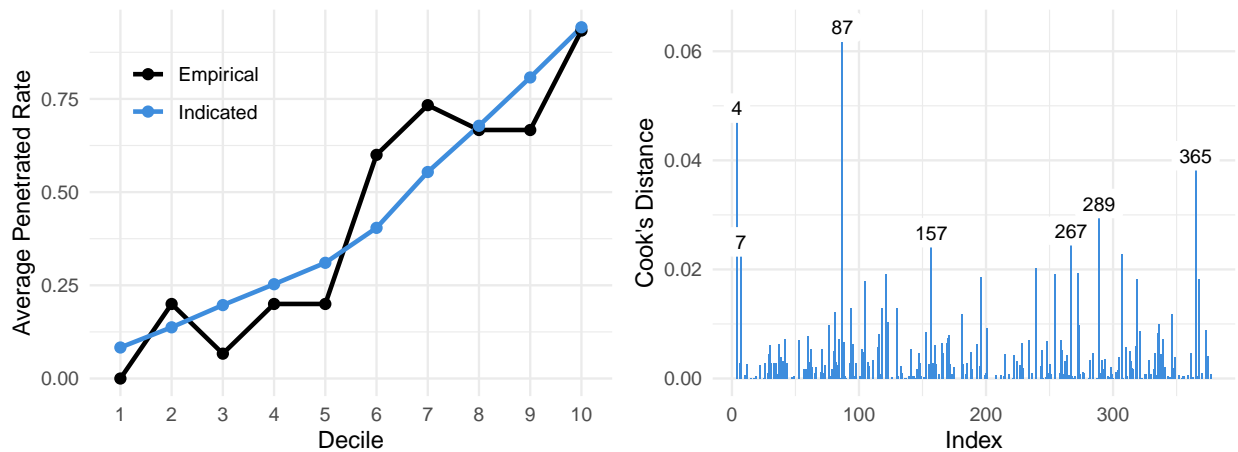


Figure 3: 1 (Left): Lift chart of empirical (black) and indicated (blue) penetration rate sorted from smallest to largest and assigned by decile and 2 (Right): Cook's Distance plot with labeled outlier observations.

Checking the quality of the model is one of the most important step in a data analysis. First, multicollinearity would not be an issue in this analysis since VIF scores are smaller than 2 for each variable. Secondly, Figure 3.2 shows the Cook's distance of each observation point. A few observations did not meet the cutoff point of 0.02 including rows 4, 7, 87, 121, 157, 196, 239, 254, 267,

272, 289, 307, and 365. These points influence all fitted values. In other words, these subjects have huge impact on the chances of getting penetrated of prostate capsule.

The student residuals are also very important to validate while accessing the quality of the model. Figure 4.1 portrays a residual plot of the model. The plots is color coded by two section, one for underestimate points and one for over estimate subjects. The number of overestimte values and overestimate values are approximately the same which constructed a mean value for resudials is approximately 0 or -0.07 to be exact. Furthermore, there are no pattern to be detected in the residuals consistent with constant variance. Most residuals are between -2 and 2 which is ideal for all residuals. However, there are a few points that clearly have large residuals including rows 87, 239, and 289. This indicates that these three subjects are outliers.

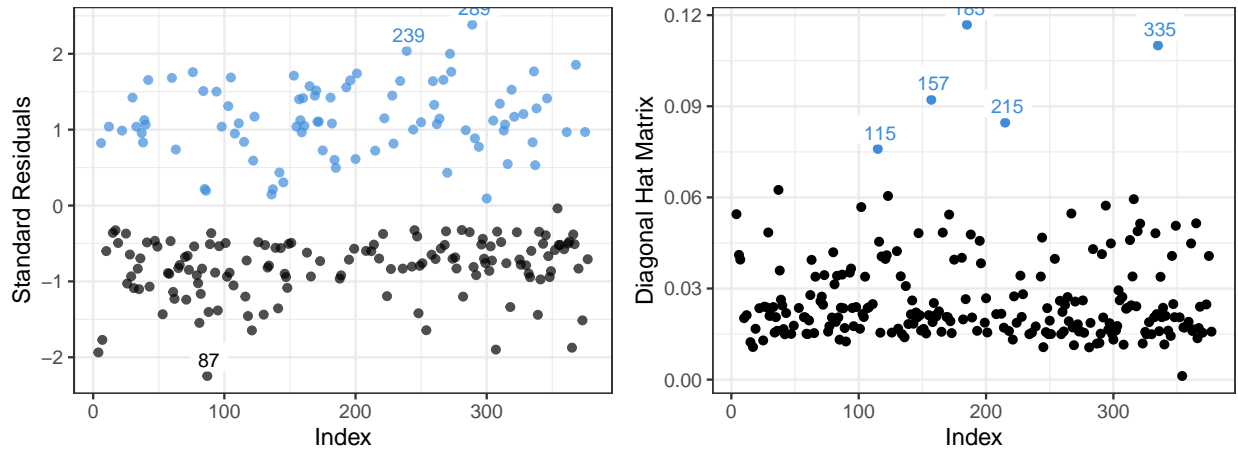


Figure 4: 1 (Left): Studendized residuals plot with blue color are positive residuals (overestimate values) and black are negatie residuals (underestimate values), and 2 (Right): Diagonal hat matrix plot with influential points in colored in blue.

Lastly, large hat values show potential outlying observation with respect to each predictors. As seen in Figure 4, rows 115, 215, 185, and 335. These subjects potentially have really high/low values consisted in one or more of their predictors. To further access the outliers and influential points, Table A5 was constructed to show the details of each observation that were detected as outliers or influential points. Most subjects noticeably have really high PSA scores compare to the mean PSA scores of 15.26. However, all subjects were decided to be kept in the model since there are no valid reasons to remove them from the model.

After finishing model development and validation, the model is ready for interpretation. All interpretaion values used in this section are from Table 1. The results of digital exam has great impact on penetration. For instant, a patient with a unilobar nodule (left) would be 2.28 times or 128% riskier than a patient without a nodule. A person with a unilobar nodule (right) would be 3.69 times or 269% more likely to have their tumor penetrated with prostate capsule. Lastly, a person with bilobar nodule would be more likely to be diagnosed with prostate cancer 3.60 times compared to patients with no nodule. In addition, a patient as PSA scores increases by 1 mg/ml, on average, the patient will have 4% increase in the risk of having a prostate cancer. Last but not least, having a low Gleason score would decrease the chances of having cancer since an increase of 1-unit in Gleason score would cause the patient to be 2.33 times riskier.

Conclusion: Generally, there are protocols for doctors and nurses to follow when it comes to important diagnosis like prostate cancer. Manually, it is quite difficult to avoid human error when providing a diagnostic to a patient. With this model, health professionals can be able to minimize predicting errors and generate the likelihood of a men getting prostate cancer efficiently. The three main factors that change the chance of a tumor have penetrated of prostate capsule are the results of digital rectal exam, PSA scores, and total Gleason scores. Penetration rate are likely to increase as the results of digital rectal exam are in a critical nodule, PSA are higher, and total Gleason scores are larger. On the other hand, having no nodule in the results of digital rectal exam, low PSA scores, and low Gleason scores would cause the likelihood of penetration to be very low. This leads to lower risk of having prostate cancer. With an AUC of 0.82 and accuracy of 0.77, we can consider this model to be the best model that can provide the likelihood of penetration rate in men. There are many limitations for this models. First, the sample size in this dataset is relatively small. The subset of this small sample size could cause the results to be not as accurate as the full dataset. If we have access to the full dataset from the Ohio State University Comprehensive Cancer Center, then our results might be different. In addition, there were only 6 predictors provided to use in the model. If there is more variables, the model could provide better results. For example, we can consider genetic could also be a factor that cause the risk of having prostate cancer in the family of the patient to improve the current model.

Appendix A: Supplemental Tables

Table 2: Summary Statistics for all independent features

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
age	377	66.066	6.425	47	62	71	79
psa	377	15.256	19.869	0.300	5.000	16.800	139.700
gleason	377	6.379	1.092	0	6	7	9

	random_partition	nrow_train	nrow_test	features	AIC	auc
1	0.900	339	38	dpros,psa,gleason,dcaps	298.077	0.822
2	0.800	301	76	dpros,psa,gleason,dcaps	278.940	0.820
3	0.600	226	151	dpros,psa,gleason	241.799	0.819
4	0.700	263	114	dpros,psa,gleason	260.832	0.818
5	0.600	226	151	dpros,psa,gleason,dcaps	240.567	0.818
6	0.700	263	114	dpros,psa,gleason,dcaps	258.675	0.815
7	0.900	339	38	dpros,psa,gleason	297.050	0.812
8	0.800	301	76	dpros,psa,gleason	277.654	0.809
9	0.600	226	151	psa,gleason,dcaps	251.167	0.805
10	0.900	339	38	dpros,gleason	301.355	0.805

Table 3: Iteration log of model bootstrapping process. Showing the top 10 models that have the highest AUC scores. The columns represent random partition for train set, number of rows in the training data, number of rows in the testing data, predictors, AIC scores, and AUC scores, respectively from left to right.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-9.7465	3.8714	-2.52	0.0118
dpros2	1.7503	4.1518	0.42	0.6733
dpros3	2.9024	4.2654	0.68	0.4962
dpros4	4.8459	4.7046	1.03	0.3030
psa	0.2090	0.1227	1.70	0.0886
gleason	1.0748	0.6161	1.74	0.0811
dpros2:psa	-0.0986	0.0638	-1.55	0.1221
dpros3:psa	-0.1225	0.0801	-1.53	0.1260
dpros4:psa	-0.0601	0.0699	-0.86	0.3902
dpros2:gleason	0.0468	0.6584	0.07	0.9434
dpros3:gleason	-0.0354	0.6803	-0.05	0.9585
dpros4:gleason	-0.4452	0.7431	-0.60	0.5491
psa:gleason	-0.0130	0.0168	-0.77	0.4396

Table 4: Summary regression of the final model with all combination of interaction terms. The columns represent the variables, coefficient estimated values, standard error, test statistics, and significant p-values, respectively from left to right.

	capsule	dpros	psa	gleason	rownames
1	0	2	51.20	7	4
2	0	4	31.90	7	7
3	0	4	17.10	9	87
4	1	4	45.30	6	115
5	0	4	25.20	7	121
6	1	1	44.40	6	157
7	1	1	85.40	7	185
8	1	3	17.70	5	196
9	1	4	53.90	6	215
10	1	1	6.70	6	239
11	0	4	25.10	7	254
12	1	4	18.70	5	267
13	1	1	8.90	6	272
14	1	1	6.80	5	289
15	0	3	18.10	8	307
16	1	2	58.00	6	335
17	0	2	48.00	7	365

Table 5: Table of all potential outlier and influential values. The columns represent the tumor penetration of prostatic capsule, results of digital rectal exam, PSA scores, total Gleason scores, and row number of each subject, respectively from left to right.

Appendix B: R Code

```

1 ##### Load packages #####
2 library(MASS)
3 library(magrittr)
4 library(tidyverse)
5 library(ggcorrplot)
6 library(car)
7 library(broom)
8
9
10 ##### Functions #####
11
12 #' Explore variables function, returning plots of variables count and average
   proportion
13 #'
14 #' @param df dataframe
15 #' @param xvar independent variable
16 #' @param count return histogram of count, default is TRUE
17 #' @param x_axis x axis name
18 #'
19 #' @return
20 #' @export
21 #'
22 #' @examples
23 ExploreVariable <- function(df, xvar, count = TRUE, x_axis){
24
25   group_df <- df %>%
26     group_by({{xvar}}) %>%
27     summarise(avg_prop_capsule = mean(capsule),
28               count = n(),
29               .groups = "drop")
30
31   ratio <- max(group_df$count) / max(group_df$avg_prop_capsule)
32
33   if(count == TRUE){
34     p <- ggplot(group_df, aes(x = factor({{xvar}}), group = 1)) +
35       geom_bar(aes(y = count), stat = "identity", fill = "#3e8ddd", col = "white") +
36       geom_point(aes(y = avg_prop_capsule * ratio), size = 2, color = "black") +
37       geom_line(aes(y = avg_prop_capsule * ratio), size = 1, color = "black") +
38       scale_y_continuous(sec.axis = sec_axis(~./ratio, name = "Average Penetration
   Rate"))
39
40     y_lab <- "Number of Penetration"
41
42   } else {
43     p <- ggplot(group_df, aes(x = {{xvar}}, group = 1)) +
44       geom_point(aes(y = avg_prop_capsule), size = 2, color = "black") +
45       geom_line(aes(y = avg_prop_capsule), size = 1, color = "black")
46
47     y_lab <- "Average Penetration Rate"
48   }
49   p +
50     labs(y = y_lab, x = x_axis) +
51     theme_bw()
52 }
53

```

[illegible]

```

110         features = paste(xvars, collapse = ","),
111         AIC = mean(all_aic),
112         auc = mean(all_auc))
113     index <- index + 1
114 }
115 all_iteration %<>% bind_rows()
116
117 return(all_iteration)
118 }
119
120 SafelyBinaryFit <- safely(BinaryFit)
121
122
123 #' Function to loop all the models through each variable
124 #'
125 #' @param rp
126 #' @param nrounds
127 #' @param list_of_xvars
128 #' @param number_of_xvars
129 #'
130 #' @return
131 #' @export
132 #'
133 #' @examples
134 LoopAllVars <- function(rp, nrounds, list_of_xvars, number_of_xvars){
135   # group variables
136   vars <- combn(list_of_xvars, number_of_xvars)
137
138   # how many group do we have?
139   no_vars <- dim(vars)[[2]]
140
141   # initiate print list and index
142   print_list <- list()
143   k <- 1
144
145   for (i in 1:no_vars) {
146     iteration_log <- SafelyBinaryFit(rp, xvars = vars[,i], nrounds)
147     print_list[[k]] <- iteration_log
148     k <- k + 1
149   }
150   # map all separate log to bind them together into a dataframe
151   iteration_log <- map(print_list, "result") %>% bind_rows()
152
153   return(iteration_log)
154 }
155
156
157 ##### Load data #####
158
159 ## Load in prostate data set
160 prostate <- read.table("data/prostate.txt", header = T)
161
162
163 ##### Clean data & Data Pre-processing #####
164
165 ## check for NA

```

```

166 sum(is.na(prostate)) # missing data: 3 in race
167
168 ## omit NA
169 prostate %>%
170   na.omit()
171
172 ## change categorical variables to class factor
173 prostate %>%
174   modify_at(c("dcaps", "dpros", "race"), as.factor)
175
176 ## add bins for psa
177 prostate %>%
178   mutate(psa_bins = ntile(psa, 10))
179
180
181 ##### Data Exploratory Analysis #####
182
183 ## Target Variable — capsule — histogram
184 ggplot(prostate, aes(factor(capsule))) +
185   geom_histogram(stat = "count", fill = "#3e8ddd", col = "white") +
186   theme_bw() +
187   labs(x = "Capsule",
188        y = "Count")
189
190 ## Predictors
191 ExploreVariable(prostate, age, count = FALSE, x_axis = "Age")
192 ExploreVariable(prostate, race, x_axis = "Race")
193 ExploreVariable(prostate, dpros, x_axis = "Results of Digital Rectal Exam")
194 ExploreVariable(prostate, dcaps, x_axis = "Detection of Capsular Involvement")
195 ExploreVariable(prostate, psa_bins, count = FALSE,
196                 x_axis = "Prostatic Specific Antigen Value (mg/ml)") +
197   scale_x_continuous(breaks = seq(1, 10, 1))
198 ExploreVariable(prostate, gleason, x_axis = "Total Gleason Score")
199
200 # after the variable exploratory, we can see that
201 # — age don't have a relationship with capsule
202 # — there is little to no relationship with capsule, and unbalance factors
203 # — dpros has a strong effect of getting penetration by different levels
204 # — dcaps also has strong effect of getting penetration
205 # — psa has a strong positive relationship with capsule
206 # — gleason also has a strong positive relationship with capsule
207 # we can continue the analysis using 4 variables: dpros, dcaps, psa, and gleason
208
209
210 ##### Check Correlation for Continuous variables #####
211 corr_table <- prostate %>%
212   select(age, psa, gleason) %>%
213   cor()
214
215 ggcorrplot(corr_table, hc.order = TRUE,
216            outline.col = "white",
217            colors = c("blue", "white", "red"))
218
219 # no high correlation between continuous predictors
220 # all variables can be in the same model
221

```

```

222
223 ##### Model Development #####
224
225 ## parameter for bootstrapping
226
227 xvars <- c("dpros", "psa", "gleason", "dcaps")
228
229 rp_par <- c(0.6, 0.7, 0.8, 0.9)      # random partition parameters
230
231 nrounds <- 10                        # number of rounds per model
232
233 ## Bootstrapping
234
235 # initiate print list
236 print_list <- list()
237
238 # loop
239 for (n_features in 1:4){
240   results <- LoopAllVars(rp = rp_par,
241                           nrounds = nrounds,
242                           list_of_xvars = xvars,
243                           number_of_xvars = n_features)
244   print_list[[n_features]] <- results
245 }
246
247 # bind iteration log into a dataframe
248 iteration_log <- print_list %>% bind_rows()
249
250 # after running the iteration logs, models with variables dpros, psa, and gleason
251 # has the high AUC and low AIC, using train set of 60% prop and test set of 40%
252
253
254 ##### Select Model after bootstrap #####
255
256 ## Split data to train/test set using random partition with proportion .7/.3
257 set.seed(1234)
258
259 n <- floor(.6 * nrow(prostate))
260
261 train_ind <- sample(seq_len(nrow(prostate)), size = n)
262
263 train <- prostate[train_ind, ]
264
265 test <- prostate[-train_ind, ]
266
267
268 ## start with initial fit using 4 variables
269 fit1 <- glm(capsule ~ dpros + psa + gleason,
270             family = binomial(link = logit),
271             data = train)
272
273
274 ## Checking for Interaction term
275
276 fit2 <- glm(capsule ~ .*,
277             family = "binomial",

```

```

278         data = train %>% dplyr::select(dpros, psa, gleason, capsule))
279
280 test$preds2 <- predict(fit2, newdata = test, type = "response")
281
282 test_roc <- pROC::roc(response = test$capsule, predictor = test$preds2)
283 test_auc <- as.data.frame(pROC::auc(test_roc))
284 aic <- fit2$aic
285
286 summary(fit2)
287
288 stepAIC(fit2, direction = c("both"))
289
290 # although auc for this model is high and AIC is low,
291 # each coefficient in the model is not statistically significant
292 # so we're not considering this model to be the final model
293 # also using model selection function stepAIC,
294 # interaction terms are not needed in the model
295
296 ##### Final Model #####
297
298 fit_final <- fit1
299
300 ## generate predictions using test set
301 test$preds <- predict(fit_final, newdata = test, type = "response")
302
303 ## create percentile to generate lift chart
304 test$percentile <- ntile(test$preds, 10)
305
306 lift_plot_df <- test %>%
307   group_by(percentile) %>%
308   summarise("Empirical" = mean(capsule),
309             "Indicated" = mean(preds)) %>%
310   gather("key", "value", -percentile)
311
312 ## Lift chart to visualize the quality of the model
313 ggplot(lift_plot_df, aes(x = factor(percentile), y = value, color = key, group = key
314   )) +
315   geom_line(size = 1) +
316   geom_point(size = 2) +
317   theme_minimal() +
318   labs(x = "Percentile",
319        y = "Average Proportion Penetrated") +
320   theme(legend.title = element_blank()) +
321   scale_color_manual(values = c("black", "#3e8ddd"))
322
323 ## compute AIC and AUC of final model
324 aic <- fit_final$aic
325 test_roc <- pROC::roc(response = test$capsule, predictor = test$preds)
326 test_auc <- as.data.frame(pROC::auc(test_roc))
327
328 tidy <- broom::tidy(fit_final) %>%
329   mutate(OR = exp(fit_final$coefficients)) %>%
330   cbind(exp(confint(fit_final))) %>%
331   modify_if(is.numeric, round, 4)
332
333 ## wald test

```

```

333 aod::wald.test(b = coef(fit_final), Sigma = vcov(fit_final), Terms = 4)
334
335
336 ## Generate confusion matrix and accuracy
337 test %>%
338   mutate(preds = case_when(preds < 0.5 ~ 0.5, TRUE ~ 1.1))
339
340 confusion_matrix <- ftable(test$capsule, test$preds)
341
342 accuracy <- sum(diag(confusion_matrix))/nrow(test)*100
343
344 accuracy
345
346
347 ##### Diagnostics #####
348
349 ##### Outliers test #####
350 ## Bonferonni p-value for most extreme obs
351 outlierTest(fit_final)
352
353
354 ## Influential Observations
355 influencePlot(fit_final,
356               id.method="identify",
357               main="Influence Plot",
358               sub="Circle size is proportional to Cook's Distance" )
359
360 diag <- augment(fit_final) %>%
361   mutate(Index = 1:nrow(.))
362
363 cutoff <- 4/((nrow(train) - length(fit_final$coefficients) - 2))
364
365 diag %>%
366   mutate(high_cooksd = case_when(
367     .cooks > cutoff ~ 1, TRUE ~ 0),
368     col_stdresid = case_when(
369       .std.resid > 0 ~ 1,
370       .std.resid < 0 ~ 0),
371     high_hat = case_when(
372       .hat > .07 ~ 1,
373       TRUE ~ 0))
374
375 ## cook's distant plot
376 ggplot(diag, aes(x = Index, y = .cooks)) +
377   geom_bar(stat = "identity", fill = "#3e8ddd", col = "white") +
378   labs(y = "Cook's Distance") +
379   theme_minimal() +
380   geom_label(data = diag %>% filter(.cooks > cutoff + .005),
381             aes(label = .rownames), label.size = NA, size = 3)
382
383
384 ## standard residual plot
385 ggplot(diag, aes(x = Index, y = .std.resid, color = factor(col_stdresid))) +
386   geom_point(alpha = 0.7) +
387   labs(y = "Standard Residuals") +
388   theme_bw() +

```



```

389   scale_color_manual(values = c("black", "#3e8ddd")) +
390   theme(legend.position = "none") +
391   geom_label(data = diag %>% filter(.std.resid > 2 | .std.resid < -2),
392             aes(label = .rownames), label.size = NA, size = 3, hjust=0.45, vjust
              =-.15)
393
394 ## diagonal hat matrix plot
395 ggplot(diag, aes(x = Index, y = .hat, color = factor(high_hat))) +
396   geom_point() +
397   labs(y = "Diagonal Hat Matrix") +
398   theme_bw() +
399   scale_color_manual(values = c("black", "#3e8ddd")) +
400   theme(legend.position = "none") +
401   geom_label(data = diag %>% filter(.hat > .07),
402             aes(label = .rownames), label.size = NA, size = 3, hjust=0.45, vjust
              =-.15)
403
404 ## student residual distribution
405 sresid <- studres(fit_final)
406 hist(sresid, freq=FALSE,
407       main="Distribution of Studentized Residuals")
408 xfit<-seq(min(sresid),max(sresid),length=40)
409 yfit<-dnorm(xfit)
410 lines(xfit, yfit, col = "red")
411 # looks normal
412
413 ##### Collinearity #####
414 vif(fit_final) # variance inflation factors
415 # all VIF are low = no multicollinearity
416
417
418 ##### Find influential point in data #####
419
420 inf_points <- (diag %>%
421               filter(.std.resid > 2 | .std.resid < -2 | .cooksd > cutoff | .hat >
422                     .07))$.rownames %>%
423   as.numeric()
424
425 prostate[inf_points, ] %>%
426   as.data.frame() %>%
427   dplyr::select(capsule, dpros, psa, gleason) %>%
428   mutate(rownames = inf_points) %>%
429   arrange(rownames)

```

Listing 1: Appendix of Code