

Predicting Patient Survival Based on Physiological Variables using Logistic Regression

Joe Sanchez

January 20, 2020

Executive Summary

One of the many problems that hospital staff face is deciding which patients to take care of first. While ethicists and medical professionals don't always agree on who should get treatment first, the common procedure is to treat patients who arrive at a hospital in critical condition first. The goal of this study is to create a logistic regression model which will allow medical staff to quantify the severity of a patient's condition. Patients that are at a higher risk of dying would be given treatment first. The dataset used contains the physiological parameters of 112 critically ill patients which were collected in Southern California. Stepwise model selection was used in order to obtain the final logistic regression model. The model was then trained and tested on a 70/30 split of the data in order to gauge performance. The sensitivity achieved by the model was 0.8 and the specificity obtained was 0.68. The results of the analysis suggest that shock type, mean arterial pressure, mean central venous pressure, and body surface index are significant factors in determining a patient's chance of survival.

1 Introduction

It is a widely known that hospitals in the US face issues due to understaffing of registered nurses. When nurses cannot give patients the care that they require there are serious consequences. Many physicians have reported patient deaths in which they point to severe understaffing as the cause. It is critical for nurses to understand what level of care each patient needs in order to avoid patient deaths.

The purpose of this study is to create a logistic regression model which can accurately predict which patients are at a higher risk of dying. The dataset contains 21 fields and 224 observations, which includes 2 observations for each of the 112 critically ill patients (an initial measurement upon admission and a measurement just before death or discharge). Table 1 provides a brief summary of the variables. An accurate logistic regression model would allow medical staff to quickly identify higher risk patients which need more care than other patients.

Table 1: Descriptions of variables in the dataset.

Variable	Description	Units	Values
ID	ID	None	
AGE	Patient Age	yr	16-90
HT	Patient Height	cm	140-185
Sex		None	1=Male, 2=Female
SURVIVE	Survival	None	1=Survived, 3=Died
SHOCK_TYPE	Shock Type	None	2=Non-shock 3=Hypovolemic shock 4=Cardiogenic shock 5=Bacterial shock 6=Neurogenic shock 7=Other
SBP	Systolic Pressure	mmHg	26-171
MAP	Mean Arterial Pressure	mmHg	15-124
HR	Heart Rate	beats/min	25-217
DBP	Diastolic Pressure	mmHg	10-108
MCVP	Mean Central Venous Pressure	cm H ₂ O	2-302
BSI	Body Surface Index	m ²	109-224
CI	Cardiac Index	liters/min*m ²	17-763
AT	Appearance Time	sec	20-261
MCT	Mean Circulation Time	sec	81-590
UO	Urinary Output	ml/hr	0-510
PVI	Plasma Volume Index	ml/kg	207-1066
RCI	Red Cell Index	ml/kg	107-858
HG	Hemoglobin	gm/100 ml	66-180
HCT	Hematocrit	percent	20.0-54.0
RECORD	Card Sequence	None	1=Initial,2=Final

2 Methods

The data set contained two observations for each of the 112 patients. There were no missing values in the dataset, there was however, a discrepancy between the initial and final height measurements for one of the patients. It was assumed that this was a mistake and the initial height of 70 cm was changed to 170 cm to match the final height. Out of the 112 patients 43 (38.4%) patients died.

The goal was to create an accurate model which would allow hospital staff to quickly identify which incoming patients were in more severe condition (predicted to die by model). Therefore, only the initial records were used since, the final measurements may have either been taken shortly before a death or just before discharging a patient in which case they were irrelevant. The ID variable was useless for predicting survival therefore it was dropped along with the record variable once all the final measurements were removed. The Survival variable values were changed to 0 = survival and 1 = death. Since the response variable was binary, a logistic regression model was used in order to predict whether a patient in the dataset lived or died.

Exploratory data analysis was performed on the physiological variables and associations between survival and categorical tests were determined using χ^2 tests. Welch's T-test was used to determine differences in the means of continuous variables between the living and dead patients. The variables used in the model were selected by using stepwise model selection, χ^2 tests and T-tests. If the stepAIC function and tests for association both deemed a variable significant the term was kept if not the model was reduced and the process repeated until the final model was obtained. A highly reduced model was sought in order to avoid overfitting. The dataset was then split into training and testing sets using a ratio of 70/30 in order to evaluate the model's sensitivity and specificity. R was used for the analysis in this study.

3 Results

Exploratory Data Analysis

The data set contains both categorical and continuous predictor variables. Table 2a features the counts for survival and deaths among patients. Contingency tables were created for each of the predictor variables and the response and χ^2 -tests were used to test association. Table 2b is a contingency table of survival and Sex. A χ^2 -test determined that Sex is not associated with survival, however the P-value was 0.0638, meaning that association between Sex and survival was barely rejected by the test at the 0.05 alpha level. It is possible that the small sample size could have impacted the results. A χ^2 -test was also used to determine that there was a significant association shock type and survival.

Table 2: Contingency tables for survive and relation between survive and sex.

(a) Contingency table for survive.

Survived	Died
69	43

(b) Contingency table for survive and sex . A χ^2 -test resulted in no association with survive at the 0.05 significance level.

	Female	Male
Survived	28	41
Died	26	17

Figure 1 depicts boxplots of mean arterial pressure (MAP) and mean central venous pressure (MCVP) by survive. Welch’s two-sample t-tests were used to determine if there were differences in the means of continuous variables between the patients that survived and those that died. The t-tests determined that there were significant differences in MAP and MCVP between the two groups of patients.

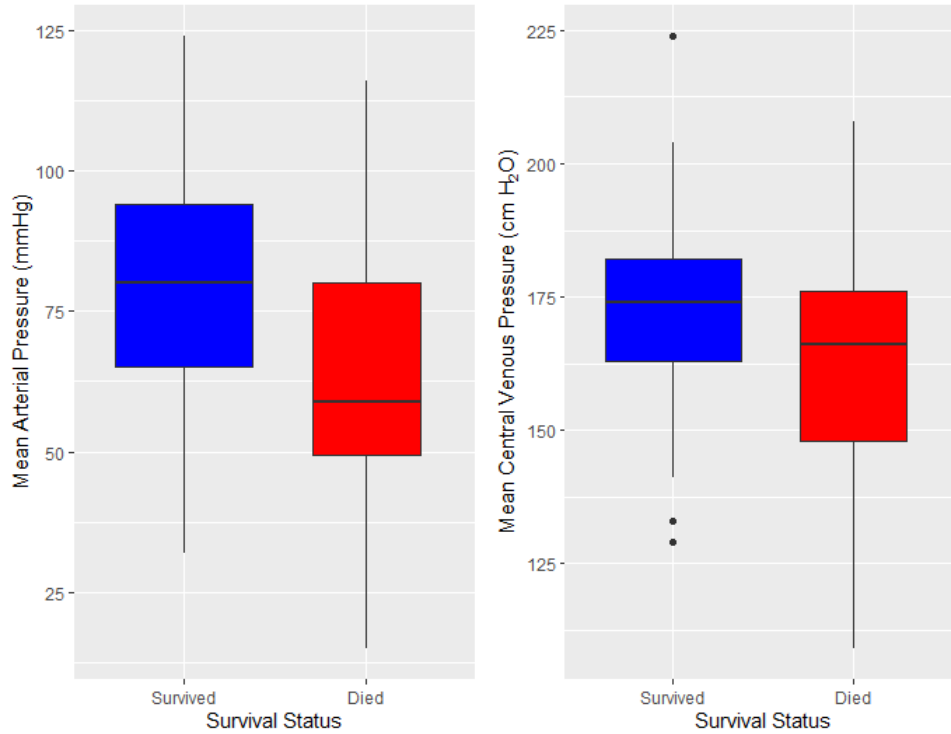


Figure 1: Boxplots for mean arterial pressure (left) and mean central venous pressure (right)

Cohen’s d was also calculated for each of the continuous variables. MAP had a value cohen’s d estimate of 4.7 indicating roughly a medium level of effect on survival. While MCVP had a cohen’s d estimate of 2.2 indicating a small level of effect on survival. Figure 2 shows a correlation plot of the variables. MAP, SBP, and DBP are highly correlated with each other. AT and MCT are also highly correlated with each other and HG and HCT form another duo of highly correlated variables.

Model Selection

In order to obtain the final model, stepwise AIC was used in conjunction with the results from t-tests and χ^2 -tests of association. Each variable was used in the initial model, however variables were removed if stepAIC and other tests did not agree that a term was significant. This was done iteratively until the final model was obtained. The results were that shock type, MAP, MCVP and BSI were deemed significant. The same procedure was repeated on model involving an interaction term between sex and shock type (a contingency table showed a possible association with survival) however, the interaction term was deemed insignificant.

Although the first iteration of stepwise AIC initially recommended a model which also included UO and RCI, these were dropped because the results of t-tests in the model summary deemed them

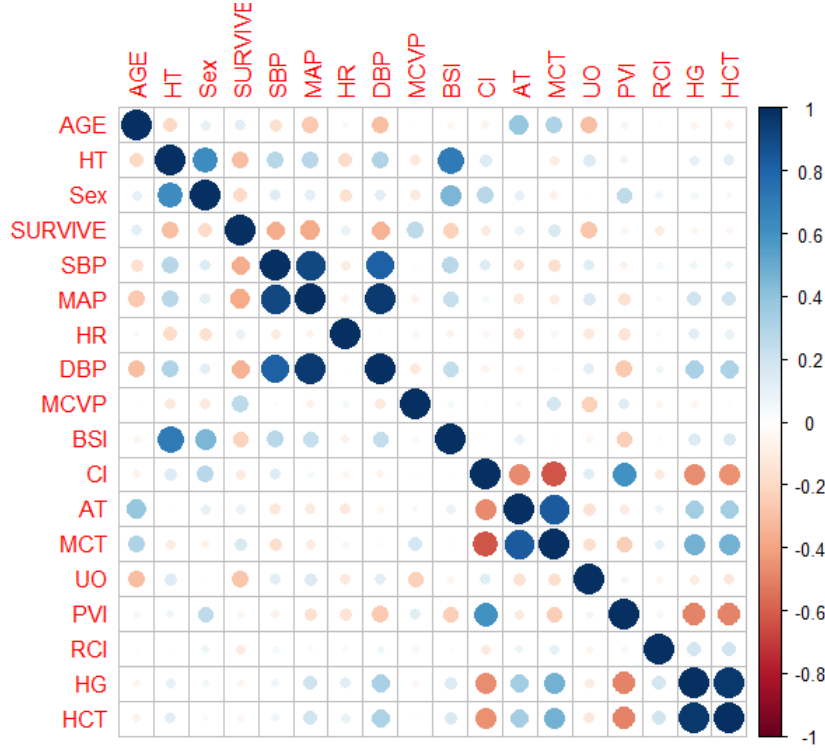


Figure 2: Correlation plot of the variables in the dataset

insignificant at the 0.05 significant level. The AIC score of the model including UO and RCI was 119.1 and the AIC score of the final model was 120.8. The increase in AIC was marginal while the final model benefited from further simplification in order to reduce overfitting.

Diagnostics

Transformation of the covariates into explanatory variable patterns (EVPs) was necessary in order to perform model diagnostics. Due to the small size of the dataset, the continuous variables used in the final model were split into 2 intervals based on percentiles. The data was aggregated into sets of similar covariates and the result of the transformation was reduction of the 112 observations into 39 EVPs.

The studentized residual plot of the EVPS shows two points which are potential outliers since they are just below 3 standard deviations away from 0. The cook's distance vs leverage plot of the EVPs shows four potential outliers. The same two points have a somewhat higher cook's distance along with another point while an additional point has a high leverage compared to all of the other points. These values are displayed in Table 3 and Diagnostic plots are shown in Appendix C. None of the possible outliers were dropped since none of the points had a cook's distance above $4/n$ (where n is the number of EVPs) while also simultaneously having a high leverage and being over three standard deviations from 0. The Generalized Variance Inflation Factor (VIF) was also calculated for each of the variables included in the model with all results being low, indicating that there was little multicollinearity in the model.

Training and Testing

Table 3: The EVPs which are potential outliers are shown

	SHOCK_TYP	MAP_interval	MCVP_interval	SURVIVE	BSI_interval	std.res	cookd	h
15	5	(15,72.5]	(80,302]	4.00	(109,169]	-0.15	0.00	0.63
16	6	(15,72.5]	(80,302]	1.00	(109,169]	-2.90	0.70	0.43
32	2	(15,72.5]	(80,302]	2.00	(169,224]	2.44	0.28	0.30
38	6	(72.5,124]	(80,302]	2.00	(169,224]	1.59	0.14	0.34

Table 4: Generalized VIF scores and Degrees of freedom

	GVIF	Df	GVIF ^{1/(2*Df)}
SHOCK_TYP	1.27	5.00	1.02
MAP	1.12	1.00	1.06
MCVP	1.24	1.00	1.11
BSI	1.08	1.00	1.04

The dataset was split into training and testing sets using a 70/30 ratio. The final model was then fit to the training data and predictions were made using the testing data. If the probability predicted was greater than 0.40 a patient was classified as dead if not the patient was classified as a survivor. The cutoff of 0.40 was used since there was an imbalance in the number of patients which lived and died in the dataset with approximately 40% of patients dying.

The results of the training and testing were that a sensitivity of 80% was achieved and a specificity of 68% was also achieved. The results are not quite ideal but they can still be somewhat useful to medical staff. Since the goal is to create a model which can predict which patients are at a higher risk of death, it is better to have higher sensitivity than specificity if one must be lower than the other. As the saying goes "it's better to be safe than sorry" and in this case it is better to give a lower risk patient more care than to give a high risk patient less care in order to avoid patient deaths.

Inferences

In order to make sense of the final model, odds ratios must be used in order to describe relationships between the covariates and survival. Changes in odds ratios specify the effects of an increase or decrease in a variable while all other variables are held constant. According to the model, there is a 1450% increase in the odds of death if a patient has hypovolemic shock. If a patient has cardiogenic shock there is a 827% increase in odds of death. Bacterial shock increases odds of death by 279% and neurogenic shock results in a 598% increase in odds of death. Patients exhibiting "other" shock type have a 2940% increase in odds of death. A one unit increase (mmHg) in MAP reduces odds of death by 2.9%. A one unit increase (mmHg) in MCVP increases odds of death by 2% a one unit increase (m²) in BSI reduces odds of death by 2.8%. A summary of the model estimates is provided in table 5.

4 Conclusion

The goal of developing a model which could identify patients likely to die was achieved with shock type, MAP, MCVP, and BSI being identified significant predictor variables. The sensitivity

Table 5: Summary of model estimates including odds ratios, standard errors, p-values, and confidence intervals.

	Coefficient	Odds ratio	SE	p-value	95% CI on OR
Intercept	3.10	23.000	2.500	0.21	(0.195 , 3950.000)
Hypovolemic Shock	2.70	15.500	0.935	0.00	(2.750 , 115.000)
Cardiogenic Shock	2.20	9.270	0.835	0.01	(1.980 , 56.100)
Bacterial Shock	1.30	3.790	0.897	0.14	(0.683 , 24.900)
Neurogenic Shock	1.90	6.980	0.871	0.03	(1.360 , 44.400)
Other	3.40	30.400	1.070	0.00	(4.270 , 308.000)
MAP	-0.03	0.971	0.013	0.02	(0.945 , 0.995)
MCVP	0.01	1.020	0.005	0.00	(1.010 , 1.030)
BSI	-0.03	0.972	0.014	0.05	(0.944 , 0.999)

achieved by the final model was 80% and the specificity achieved was 68.4%. While these values could use some improvement, the model could still be used to predict which patients are at a high risk of dying and would allow medical staff to allocate resources and care in order to minimize deaths.

One limitation of this study is that a small sample size was used. This could have affected which terms were deemed significant and could also have affected the coefficients in the model. Additionally, the results from training and testing may be somewhat unreliable because the training set included only a small amount of observations. Another limitation of this study is that nothing is known about a patient's history. Such information could prove useful in assessing whether or not a patient is at a higher risk of dying. Another limitation is that the data was collected in Southern California. It is possible that the model would not be accurate for people in different regions.

In order to produce a more accurate model future studies could include a larger sample size, include a patients history, and use data that was obtained from all over the nation or possibly world. Perhaps researchers could also include variables which describe conditions at the hospital and describe medical staff in order to get a better picture of whether or not a patient will survive.

Appendices

Appendix A: R Code

```
1 library(dplyr)
2 library(vioplot)
3 library(ggplot2)
4 library(corrplot)
5 library(MASS)
6 library(gridExtra)
7 library(xtable)
8 library(car)
9 library(effsize)
10
11 column.names<-c('ID','AGE','HT','Sex','SURVIVE','SHOCK_TYP',
12 'SBP','MAP','HR','DBP','MCVP','BSI',
13 'CI','AT','MCT','UO','PVI','RCI','HG',
14 'HCT','RECORD')
15
16 physio.data<-read.csv('DATA-FILEsp2020.csv',header=FALSE,
17 col.names=column.names)
18 physio.data$SURVIVE<-ifelse(physio.data$SURVIVE==1,0,1) #changing
    survive to 0 and death to 1
19 physio.data$Sex<-ifelse(physio.data$Sex == 2,0,1) #changed female
    to 0 kept male as 1
20 physio.data$HCT<-physio.data$HCT/10
21 #3rd digit in HCT was assumed to specify tenths since it is
22 #not possible to have an hct above 100
23 physio.data[physio.data$ID==539 & physio.data$RECORD == 1,'HT
    ']<-170 #error in height measurement fixed
24 physio.data[physio.data$ID==539,]
25
26
27 #removing ID column
28 physio.data$ID<-NULL
29 #getting rid of second observations for each patient
30 physio.data<-filter(physio.data,RECORD==1)
31 physio.data2<-physio.data
32 #removing record column
33 physio.data$RECORD<-NULL
34 physio.data2$RECORD<-NULL
35
36 #changing categorical columns to factors
37 physio.data$SHOCK_TYP<-as.factor(physio.data$SHOCK_TYP)
38 physio.data2$SURVIVE<-as.factor(physio.data2$SURVIVE)
```



```

39 physio.data2$Sex<-as.factor(physio.data2$Sex)
40 physio.data2$SHOCK_TYP<- as.factor(physio.data2$SHOCK_TYP)
41
42 head(physio.data2)
43
44 #112 observations and 19 variables after removing final
    measurements
45 dim(physio.data2)
46 sum(is.na(physio.data2)) #no missing values
47 summary(physio.data2)
48 attach(physio.data2)
49
50
51
52 #####variable distributions
53
54 table(SURVIVE)
55 signif(prop.table(table(SURVIVE)),digits=3) #38.4% of patients did
    not survive
56
57
58 #AGE distribution
59 plot(1, 1, xlim = c(0, 2), ylim = range(AGE), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
60 vioplot(AGE, at=1, add=T)
61 axis(2, at = 50, pos=-0.2, tck=0, labels="Age")
62
63 #height distribution
64 plot(1, 1, xlim = c(0, 2), ylim = range(HT), type = 'n', xlab = '',
    ylab = '', xaxt = 'n')
65 vioplot(HT, at=1, add=T)
66 axis(2, at = 165, pos=-0.2, tck=0, labels="Height (cm)")
67
68 table(Sex) # 58 Males and 54 Females
69
70 table(SHOCK_TYP)
71 # 2 Non-shock patients
72 # 3 Hypovolemic
73 # 20 Cardiogenic
74 # 5 Bacterial
75 # 16 Neurogenic
76 # 7 Other
77
78 #SBP distribution
79 plot(1, 1, xlim = c(0, 2), ylim = range(SBP), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
80 vioplot(SBP, at=1, add=T)

```

```

81 axis(2, at = 100, pos=-0.2, tck=0, labels="Systolic Blood Pressure
    (mmHg)")
82
83 #MAP distribution
84 plot(1, 1, xlim = c(0, 2), ylim = range(MAP), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
85 vioplot(MAP, at=1, add=T)
86 axis(2, at = 70, pos=-0.2, tck=0, labels="Mean Arterial Pressure (
    mmHg)")
87
88 #HR distribution
89 plot(1, 1, xlim = c(0, 2), ylim = range(HR), type = 'n', xlab = '',
    ylab = '', xaxt = 'n')
90 vioplot(HR, at=1, add=T)
91 axis(2, at = 120, pos=-0.2, tck=0, labels="Heart Rate (bpm)")
92
93 #DBP distribution
94 plot(1, 1, xlim = c(0, 2), ylim = range(DBP), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
95 vioplot(DBP, at=1, add=T)
96 axis(2, at = 60, pos=-0.2, tck=0, labels="Diastolic Blood Pressure
    (mmHg)")
97
98 #MCVP distribution
99 plot(1, 1, xlim = c(0, 2), ylim = range(MCVP), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
100 vioplot(MCVP, at=1, add=T)
101 axis(2, at = 150, pos=-0.2, tck=0, labels="Mean Central Venous (cm
    H2O pressure)")
102
103 #BSI distribution
104 plot(1, 1, xlim = c(0, 2), ylim = range(BSI), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
105 vioplot(BSI, at=1, add=T)
106 axis(2, at = 170, pos=-0.2, tck=0, labels="Body Surface Index (m^2)
    ")
107
108 #CI distribution
109 plot(1, 1, xlim = c(0, 2), ylim = range(CI), type = 'n', xlab = '',
    ylab = '', xaxt = 'n')
110 vioplot(CI, at=1, add=T)
111 axis(2, at = 380, pos=-0.2, tck=0, labels="Cardiac index (liters/
    min m^2)")
112
113 #AT distribution
114 plot(1, 1, xlim = c(0, 2), ylim = range(AT), type = 'n', xlab = '',
    ylab = '', xaxt = 'n')

```

```

115 vioplot(CI, at=1, add=T)
116 axis(2, at = 380, pos=-0.2, tck=0, labels="Cardiac index (liters/
    min m^2)")
117
118 #MCT distribution
119 plot(1, 1, xlim = c(0, 2), ylim = range(MCT), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
120 vioplot(MCT, at=1, add=T)
121 axis(2, at = 350, pos=-0.2, tck=0, labels="Mean circulation time (
    sec)")
122
123 #UO distribution
124 plot(1, 1, xlim = c(0, 2), ylim = range(UO), type = 'n', xlab = '',
    ylab = '', xaxt = 'n')
125 vioplot(UO, at=1, add=T)
126 axis(2, at = 250, pos=-0.2, tck=0, labels="Urinary Output (ml/hr)")
127
128 #PVI distribution
129 plot(1, 1, xlim = c(0, 2), ylim = range(PVI), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
130 vioplot(PVI, at=1, add=T)
131 axis(2, at = 600, pos=-0.2, tck=0, labels="Plasma volume index (ml/
    kg)")
132
133 #RCI distribution
134 plot(1, 1, xlim = c(0, 2), ylim = range(RCI), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
135 vioplot(RCI, at=1, add=T)
136 axis(2, at = 500, pos=-0.2, tck=0, labels="Red Cell Index ml/kg")
137
138 #HG distribution
139 plot(1, 1, xlim = c(0, 2), ylim = range(RCI), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
140 vioplot(RCI, at=1, add=T)
141 axis(2, at = 500, pos=-0.2, tck=0, labels="Hemoglobin (gm/100ml)")
142
143 #HCT distribution
144 plot(1, 1, xlim = c(0, 2), ylim = range(HCT), type = 'n', xlab =
    '', ylab = '', xaxt = 'n')
145 vioplot(HCT, at=1, add=T)
146 axis(2, at = 375, pos=-0.2, tck=0, labels="Hematocrit (percent)")
147
148
149
150
151 #####comparing covariates to survival status
152 #continuous covariates are compared using boxplots

```

```

153 #contingency tables are used for discrete covariates
154
155 ggplot(data=physio.data2,aes(x=SURVIVE,y=AGE)) +
156 geom_boxplot(notch=FALSE,varwidth=T) +
157 scale_x_discrete('Survival Status',breaks=c(0,1),
158 labels=c("Survived","Died"))
159
160
161 ggplot(data=physio.data2,aes(x=SURVIVE,y=HT)) +
162 geom_boxplot(notch=FALSE,varwidth=T) +
163 scale_x_discrete('Survival Status',breaks=c(0,1),
164 labels=c("Survived","Died"))
165
166 table(SURVIVE,Sex)
167
168
169 table(SURVIVE,SHOCK_TYP)
170
171 p3<-ggplot(data=physio.data2,aes(x=SURVIVE,y=SBP)) +
172 geom_boxplot(notch=FALSE,varwidth=T,fill=c("blue","red")) +
173 scale_x_discrete('Survival Status',breaks=c(0,1),
174 labels=c("Survived","Died")) +
175 labs(y="Systolic Blood Pressure (mmHg)")
176
177 p1<-ggplot(data=physio.data2,aes(x=SURVIVE,y=MAP)) +
178 geom_boxplot(notch=FALSE,varwidth=T,fill=c('blue','red')) +
179 scale_x_discrete('Survival Status',breaks=c(0,1),
180 labels=c("Survived","Died")) +
181 labs(y="Mean Arterial Pressure (mmHg)")
182
183 ggplot(data=physio.data2,aes(x=SURVIVE,y=HR)) +
184 geom_boxplot(notch=FALSE,varwidth=T) +
185 scale_x_discrete('Survival Status',breaks=c(0,1),
186 labels=c("Survived","Died"))
187
188
189 p4<-ggplot(data=physio.data2,aes(x=SURVIVE,y=DBP)) +
190 geom_boxplot(notch=FALSE,varwidth=T,fill=c("blue","red")) +
191 scale_x_discrete('Survival Status',breaks=c(0,1),
192 labels=c("Survived","Died")) +
193 labs(y="Diastolic Blood Pressure (mmHg)")
194
195 ggplot(data=physio.data2,aes(x=SURVIVE,y=MCVP)) +
196 geom_boxplot(notch=FALSE,varwidth=T,fill=c('blue','red')) +
197 scale_x_discrete('Survival Status',breaks=c(0,1),
198 labels=c("Survived","Died")) +
199 labs(y="Mean Arterial Pressure (mmHg)")

```

```

200
201 p2 <- ggplot(data=physio.data2,aes(x=SURVIVE,y=BSI)) +
202 geom_boxplot(notch=FALSE,varwidth=T,fill=c('blue','red')) +
203 scale_x_discrete('Survival Status',breaks=c(0,1),
204 labels=c("Survived","Died")) +
205 labs(y=expression(paste("Mean Central Venous Pressure (cm ",H[2],"
    0)"))))
206
207 ggplot(data=physio.data2,aes(x=SURVIVE,y=CI)) +
208 geom_boxplot(notch=FALSE,varwidth=T) +
209 scale_x_discrete('Survival Status',breaks=c(0,1),
210 labels=c("Survived","Died"))
211
212 ggplot(data=physio.data2,aes(x=SURVIVE,y=AT)) +
213 geom_boxplot(notch=FALSE,varwidth=T) +
214 scale_x_discrete('Survival Status',breaks=c(0,1),
215 labels=c("Survived","Died"))
216
217 ggplot(data=physio.data2,aes(x=SURVIVE,y=MCT)) +
218 geom_boxplot(notch=FALSE,varwidth=T) +
219 scale_x_discrete('Survival Status',breaks=c(0,1),
220 labels=c("Survived","Died"))
221
222 p5<-ggplot(data=physio.data2,aes(x=SURVIVE,y=UO)) +
223 geom_boxplot(notch=FALSE,varwidth=T,fill=c("blue","red")) +
224 scale_x_discrete('Survival Status',breaks=c(0,1),
225 labels=c("Survived","Died")) +
226 labs(y="Urinary Output (ml/hr)")
227
228 ggplot(data=physio.data2,aes(x=SURVIVE,y=PVI)) +
229 geom_boxplot(notch=FALSE,varwidth=T) +
230 scale_x_discrete('Survival Status',breaks=c(0,1),
231 labels=c("Survived","Died"))
232
233 p6<-ggplot(data=physio.data2,aes(x=SURVIVE,y=RCI)) +
234 geom_boxplot(notch=FALSE,varwidth=T,fill=c("blue","red")) +
235 scale_x_discrete('Survival Status',breaks=c(0,1),
236 labels=c("Survived","Died")) +
237 labs(y="Red Cell Index (ml/kg)")
238
239 ggplot(data=physio.data2,aes(x=SURVIVE,y=HG)) +
240 geom_boxplot(notch=FALSE,varwidth=T) +
241 scale_x_discrete('Survival Status',breaks=c(0,1),
242 labels=c("Survived","Died"))
243
244
245 ggplot(data=physio.data2,aes(x=SURVIVE,y=HCT)) +

```

```

246 geom_boxplot(notch=FALSE,varwidth=T) +
247 scale_x_discrete('Survival Status',breaks=c(0,1),
248 labels=c("Survived","Died"))
249
250 table(Sex,SHOCK_TYP,SURVIVE)
251 #May be worth including interaction term Sex:SHOCK_TYP
252
253 grid.arrange(p1,p2,ncol=2)
254 grid.arrange(p3,p4,ncol=2)
255 grid.arrange(p5,p6,ncol=2)
256
257 #####correlation plot
258
259
260 corrplot(cor(physio.data[apply(physio.data, is.numeric)]))
261
262
263 #####chisq test for associations
264 chisq.test(SHOCK_TYP,SURVIVE)
265 #P-value =0.0007616
266
267 chisq.test(Sex,SURVIVE)
268 #P-value is 0.06376
269
270 #####t-tests for differences in means
271
272 t.test(AGE[SURVIVE==0],AGE[SURVIVE==1])
273 #P-value is 0.1882
274
275 t.test(HT[SURVIVE==0],HT[SURVIVE==1])
276 #P-values is 0.002579
277
278 t.test(SBP[SURVIVE==0],SBP[SURVIVE==1])
279 #P-value is 0.000116
280
281 t.test(MAP[SURVIVE==0],MAP[SURVIVE==1])
282 #p-value is 9.991e-005
283
284 t.test(HR[SURVIVE==0],HR[SURVIVE==1])
285 #p-value is 0.3006
286
287 t.test(DBP[SURVIVE==0],DBP[SURVIVE==1])
288 #P-value is 0.0003469
289
290 t.test(MCVP[SURVIVE==0],MCVP[SURVIVE==1])
291 #P-value is 0.004267
292

```

```

293 t.test(BSI[SURVIVE==0],BSI[SURVIVE==1])
294 #P-value is 0.02435
295
296 t.test(CI[SURVIVE==0],CI[SURVIVE==1])
297 #P-value is 0.2308
298
299 t.test(AT[SURVIVE==0],AT[SURVIVE==1])
300 #P-value is 0.3418
301
302 t.test(MCT[SURVIVE==0],MCT[SURVIVE==1])
303 #P-value is 0.08888
304
305 t.test(UO[SURVIVE==0],UO[SURVIVE==1])
306 #P-value is 0.0004551
307
308 t.test(PVI[SURVIVE==0],PVI[SURVIVE==1])
309 #P-value is 0.5683
310
311 t.test(RCI[SURVIVE==0],RCI[SURVIVE==1])
312 #P-value is 0.2342
313
314 t.test(HG[SURVIVE==0],HG[SURVIVE==1])
315 #P-value is 0.5729
316
317 t.test(HCT[SURVIVE==0],HCT[SURVIVE==1])
318 #P-value is 0.6216
319
320 t.test(HCT[SURVIVE==0],HCT[SURVIVE==1])
321 #P-value is 0.6216
322
323
324 #####cohen's d
325
326 cohen.d(AGE,SURVIVE)
327 #estimate is -0.2593205
328
329 cohen.d(HT,SURVIVE)
330 #estimate is 0.6589045
331
332 cohen.d(SBP,SURVIVE)
333 #estimate is 0.7908941
334
335 cohen.d(MAP,SURVIVE)
336 #estimate is 4.684463
337
338 cohen.d(HR,SURVIVE)
339 #estimate is -0.1999015

```

```

340
341 cohen.d(DBP,SURVIVE)
342 #estimate is 0.7417424
343
344 cohen.d(MCVP,SURVIVE)
345 #estimate 2.187079
346
347 cohen.d(BSI,SURVIVE)
348 #estimate is 0.4664643
349
350 cohen.d(CI,SURVIVE)
351 #estimate is 0.2275672
352
353 cohen.d(AT,SURVIVE)
354 #estimate is -0.182227
355
356 cohen.d(MCT,SURVIVE)
357 #estimate is -0.349621
358
359 cohen.d(UO,SURVIVE)
360 #estimate is0.5933606
361
362 cohen.d(PVI,SURVIVE)
363 #estimate is 0.1100068
364
365 cohen.d(AT,SURVIVE)
366 #estimate is -0.182227
367
368 cohen.d(MCT,SURVIVE)
369 #estimate is -0.349621
370
371 cohen.d(RCI,SURVIVE)
372 #estimate is .2094894
373
374 cohen.d(HG,SURVIVE)
375 #estimate is 0.1076231
376
377 cohen.d(HCT,SURVIVE)
378 #estimate is 0.09268039
379
380 #####model selection
381
382 fit1<-glm(SURVIVE~.,data=physio.data2,family='binomial')
383 summary(fit1)
384
385 stepAIC(fit1,direction='both')
386 #stepAIC did not choose SBP and DBP even though they both

```



```

387 #appeared significant it may be because MAP is very highly
388 #correlated with SBP and DBP and it is slightly more
389 #correlated with survival status therefore stepAIC
390 #only picked MAP
391 #model chosen by stepAIC below
392 #Call: glm(formula = SURVIVE ~ SHOCK_TYP + MAP + MCVP + BSI + UO +
      RCI,
393
394 fit2<-glm(SURVIVE~.+Sex:SHOCK_TYP,data=physio.data2,family='
      binomial')
395 stepAIC(fit2,direction='both')
396 summary(fit2)
397 #interaction term deemed insignificant
398
399 fit3<-glm(SURVIVE~SHOCK_TYP+MAP+MCVP+BSI+UO+RCI,data=physio.data2,
      family='binomial')
400 stepAIC(fit3,direction='both')
401 summary(fit3)
402 #decided to get rid of UO and RCI based on results from
403 #summary(fit3). Kept MAP because there was a very large
404 #difference in means of MAP between the living and dead patients
405
406 ###final model
407 fit4<-glm(SURVIVE~SHOCK_TYP+MAP+MCVP+BSI,data=physio.data2,family='
      binomial')
408 stepAIC(fit4,direction="both")
409 summary(fit4)
410 extractAIC(fit4)
411
412 #####diagnostics
413 one.fourth.root=function(x){
414 x^0.25
415 }
416
417 source("examine.logistic.reg.R")
418
419
420 physio.data2$SURVIVE<-as.numeric(as.character(physio.data2$SURVIVE)
      )
421 SURVIVE<-as.numeric(as.character(SURVIVE))
422
423 g<-2
424 MAP_interval<-cut(MAP,quantile(MAP,0:g/g),include.lower=TRUE)
425 MCVP_interval<-cut(MCVP,quantile(MCVP,0:g/g),include.lower=TRUE)
426 BSI_interval<-cut(BSI,quantile(BSI,0:g/g),include.lower=TRUE)
427
428 w <- aggregate(formula = SURVIVE ~ SHOCK_TYP+MAP_interval+

```

```

    MCVP_interval+BSI_interval,
429 data = physio.data2, FUN = sum)
430 n<- aggregate(formula= SURVIVE ~SHOCK_TYP+MAP_interval+
    MCVP_interval+BSI_interval,
431 data=physio.data2,FUN=length)
432 w.n<-data.frame(w, trials=n$SURVIVE, prop=round(w$SURVIVE/n$SURVIVE
    ,2))
433 w.n
434
435 mod.prelim1<- glm(formula=SURVIVE/trials~SHOCK_TYP+MAP_interval+
    MCVP_interval+BSI_interval,
436 family=binomial(link=logit), data=w.n,
437 weights=trials)
438
439
440 save1<-examine.logistic.reg(mod.prelim, identify.points = T,
441 scale.n=one.fourth.root, scale.cookd=sqrt)
442
443 w.n.diag1<- data.frame(w.n, pi.hat=round(save1$pi.hat, 2),
444 std.res=round(save1$stand.resid, 2),
445 cookd=round(save1$cookd, 2),
446 h=round(save1$h, 2))
447
448 p<-length(mod.prelim1$coef)
449 ck.out<-abs(w.n.diag1$std.res)>2 |
450 w.n.diag1$cookd>4/nrow(w.n)|w.n.diag1$h>0.5
451 extract.EVPs<-w.n.diag1[ck.out,]
452 extract.EVPs
453 #a higher cutoff was used for leverage since there were so few data
    points
454 #and most data points were above the 3p/n cutoff
455
456 xtable(extract.EVPs[,c('SHOCK_TYP', 'MAP_interval', 'MCVP_interval', '
    SURVIVE', 'BSI_interval', 'std.res', 'cookd', 'h')])
457
458 vif(fit4)
459 xtable(vif(fit4))
460
461 #####training and testing
462 set.seed(1)
463 n<- nrow(physio.data2)
464 p<-0.7
465 c<-0.4
466 train<- sample(n, p*n)
467 train_set<-physio.data2[train,]
468 test_set<-physio.data2[-train,]
469

```

```

470 fit_train <- glm(SURVIVE~SHOCK_TYP+MAP+MCVP+BSI ,
471 data=train_set ,
472 family=binomial(link="logit"))
473 nrow(physio.data2[-train,])
474 test_probs <- predict.glm(fit_train,newdata=test_set ,
475 type='response')
476
477 predicted_class<- as.numeric(test_probs > c)
478 length(predicted_class)
479 length(test_set$SURVIVE)
480 test_confusion_matrix<-table(predicted_class ,test_set$SURVIVE ,
481 dnn=c("Predicted","Actual"))
482 test_confusion_matrix
483
484 sensitivity<-test_confusion_matrix[2,2]/(test_confusion_matrix
485 [2,2]+test_confusion_matrix[1,2])
486 sensitivity
487 specificity<-test_confusion_matrix[1,1]/(test_confusion_matrix
488 [1,1]+test_confusion_matrix[2,1])
489 specificity
490
491 #####model inferences
492
493 betahat<-formatC(signif(fit5$coeff,digits=2),digits=2,format='f',
494 flag='#')
495 OR<-formatC(signif(exp(fit5$coeff),digits=3),digits=3,format='f',
496 flag='#')
497 SE<-formatC(signif(summary(fit5)$coeff[,2],digits=3),digits=3,
498 format='f',flag='#')
499 cibounds<-formatC(signif(exp(confint(fit5)),digits=3),digits=3,
500 format='f',flag='#')
501 pval<-formatC(signif(summary(fit5)$coeff[,4],digits=2),digits=2,
502 format='f',flag='#')
503 x<-cbind(betahat,OR,SE,pval,matrix(paste("(",cibounds[,1],
504 " ",cibounds[,2],")")))
505 colnames(x)<-cbind("Coefficient","Odds ratio","SE","p-value",
506 "95% CI on OR")
507 x
508 xtable(x)

```

Appendix B: Additional Boxplots of covariates by survival

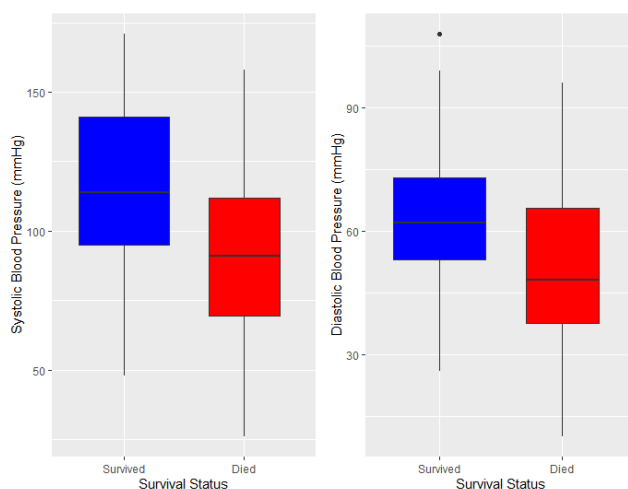


Figure 3: Boxplot of SBP by survival status (left) and Boxplot of DBP by survival status (right) these variables were highly correlated with MAP and were ultimately not used in the final model.

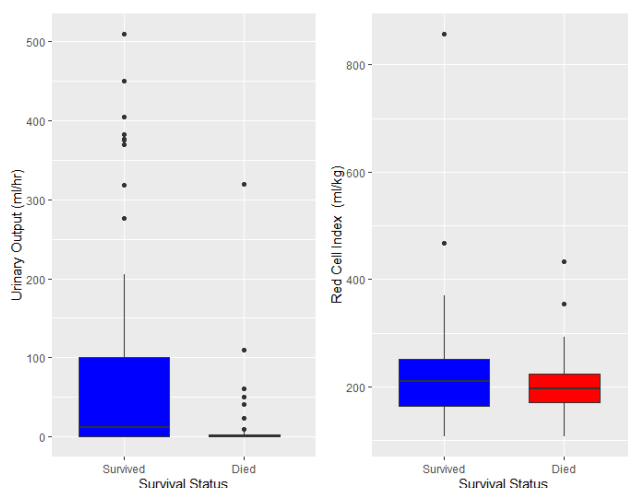
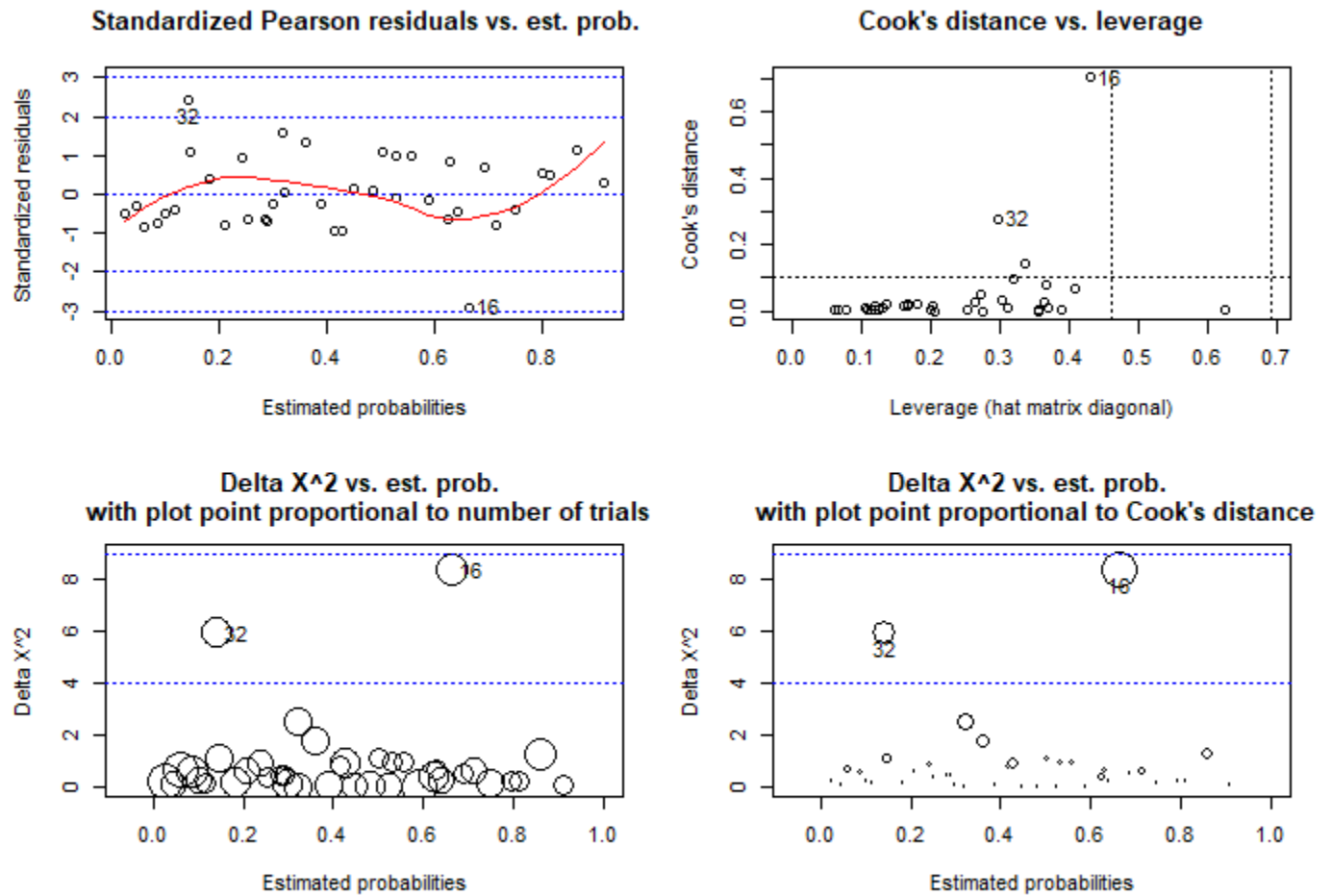


Figure 4: Boxplot of urinary output by survival status (left) and Boxplot of red cell index by survival status (right) removal of these variables increased AIC score by a marginal amount while simplifying the model.

Appendix C: Diagnostic Plots



Deviance/df = 0.94; GOF thresholds: 2 SD = 1.52, 3 SD = 1.77

Figure 5: Residual plot (top left) Cook's distance vs leverage plot (top right) and Delta X^2 plots (bottom)