

# MA478 TEE Report

CDT Joshua L. Blackmon

## 1 Research Question

A charity organization has essentially provided us with two requests to try to complete:

1. *Provide a model that will maximize the expected net profit.*
2. *Provide a model that will predict the expected gift amounts from those that donate.*

This is in an effort to fix the overall response rate based on the and maintain the most cost-effective way to mail everyone.

## 2 Data

### 2.1 Data Exploration/Preparation

The data we were given includes all that are in Figure 1 and 2.

ID: unique donor number [Do NOT use this as a predictor variable in any models]  
REG1, REG2, REG3, REG4: Region (There are five geographic regions; only four are needed for analysis since if a potential donor falls into none of the four he or she must be in the other region. Inclusion of all five indicator variables would be redundant and cause some modeling techniques to fail. A "1" indicates the potential donor belongs to this region.)  
HOME: (1 = homeowner, 0 = not a homeowner)  
CHLD: Number of children  
HINC: Household income (7 categories)  
GENF: Gender (0 = Male, 1 = Female)  
WRAT: Wealth Rating (Wealth rating uses median family income and population statistics from each area to index relative wealth within each state. The segments are denoted 0-9, with 9 being the highest wealth group and 0 being the lowest.)  
AVHV: Average Home Value in potential donor's neighborhood in \$ thousands  
INCM: Median Family Income in potential donor's neighborhood in \$ thousands  
INCA: Average Family Income in potential donor's neighborhood in \$ thousands  
PLOW: Percent categorized as "low income" in potential donor's neighborhood  
NPRO: Lifetime number of promotions received to date  
TGIF: Dollar amount of lifetime gifts to date  
LGIF: Dollar amount of largest gift to date  
RGIF: Dollar amount of most recent gift  
TDON: Number of months since last donation  
TLAG: Number of months between first and second gift  
AGIF: Average dollar amount of gifts to date  
DONR: Classification Response Variable (1 = Donor, 0 = Non-donor)  
DAMT: Prediction Response Variable (Donation Amount in \$).

Figure 1: Variables given in Charity Dataset.

Immediately, we notice a few things:

- REG1, REG2, REG3, and REG4 are all used to identify location of the potential donor.
- WRAT is an organized category.
- AVHV, INCM, and INCA are all in the thousands.

Regarding the Region used, we chose to merge them all into one variable that has 5 factors being each region. It may make the model faster and help ascertain a better prediction outcome.

For WRAT, since it is variable that is in order, we decided to ascribe an order to it correlating positively to the value. As 9 is the highest wealth, 9 is also the highest order in the variable.

Finally, for the three economic variables, we sought to explore the state at which the variables are impacted by each other. In Figure 2, you can see the correlation heatmap for all present quantitative variables to see how much they overlap.

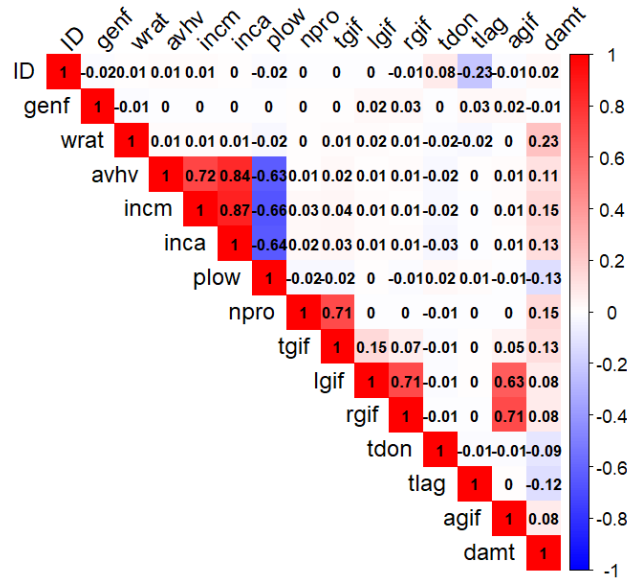


Figure 2: Correlation Heatmap with each Quantitative Variable

Figure 2 shows that the three are all correlated positively to some higher degree but this also show how the percentage of low income is also related to the three negatively. Anecdotaly this makes sense as those with similar economic statuses tend to live near each other.

In Figure 3, we can see the direct degree that AVHV, INCM, and INCA are all correlated. This also shows how many outliers are present in the dataset regarding the data. For this reason, we elected to raise them to their log to minimize the impact of all outliers. This is in order to make the model more accurate.

Through further exploration, there was not missing data so we had no reason to impute or add anything from that.

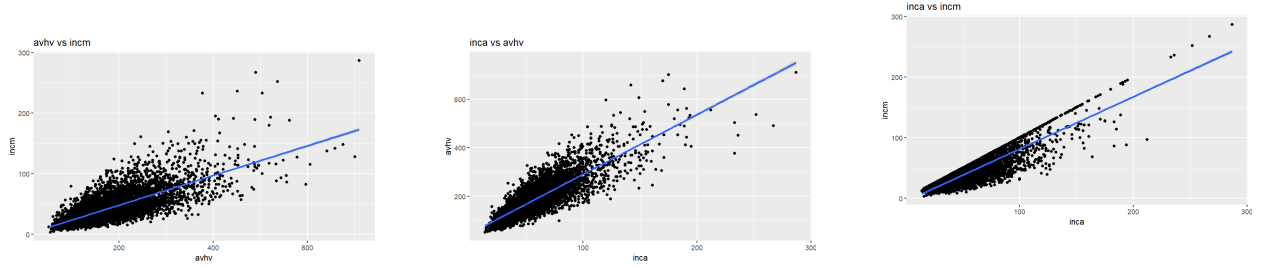


Figure 3: Overall caption for the group of images

### 3 Model Performance and Results

For the 3 models used, we settled on one full binomial model, one binomial model that uses step regression based on AIC to determine the right variables and one that used the same variables but through quasi-binomial. Model 2 and 3 had the same results regardless so my best model but Model 2 is easier to interpret. For this reason, we chose Model 2 as our model of choice.

Based on this model we will be able to mail 1433 of our highest posterior probabilities.

For the models used for the Amount Donated, we settled on our linear model using the stepwise function that gave us the lowest MSE: 89.79938 with the standard deviation of 3.451564.

### 4 Conclusions

The actual error rate of the model is quite high due to a lack of standardization which is one limit of the model. We also couldn't account for the effects that couldn't be covered (Location, Space, etc.) Outside of that, this is the best model possible.

## 5 Appendix A: Code Used

```
1 ---
2 title: "TEE MA478 Blackmon"
3 author: "CDT Joshua Blackmon"
4 date: "r Sys.Date()"
5 output: word_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 ```
11
12 ```{r}
13 library(CARBayesST)
14 library(lmtest)
15 library(rstan)
16 library(Matrix)
17 library(dplyr)
18 library(tidyverse)
19 library(lubridate)
20 library(knitr)
21 library("spdep")
22 library("lme4")
23 library("spatstat")
24 library("sp")
25 # library("maptools")
26 library("lattice")
27 library(caTools)
28 library(pROC)
29 library(ggplot2)
30 library(MASS)
31 library(pscl)
32 library(nnet)
33 library(geepack)
34 library(INLA)
35 library(GGally)
36 library(MASS)
37 library(dplyr)
38 library(ggplot2)
39 library(lme4)
40 library(corrplot)
41
42 ```
43
44 ```{r}
45
46 charity_test <- read.csv(file.choose())
47 charity_overall <- read.csv(file.choose())
48
49 ```
50
51
52 <!-- Split the data into training and validation based on if the 'part' column
53     says train or valid. -->
54 <!-- ```{r} -->
```

```

54
55 <!-- charity_train <- charity_overall[charity_overall$part == "train",] -->
56 <!-- charity_valid <- charity_overall[charity_overall$part == "valid",] -->
57
58 <!-- ''' -->
59 The purpose of the train and validate splits are to train the model on the
    training data and then validate the model on the validation data. This is done
    to ensure that the model is not overfitting the training data and is
    generalizable to new data.
60
61 '''{r}
62
63 summary(charity_overall)
64
65 #Make wrat column ordinal categorical since it is a scale of 0-9, 9 being the
    highest.
66
67 charity_overall$wrat <- ordered(charity_overall$wrat, levels = c("0", "1", "2", "3",
    "4", "5", "6", "7", "8", "9"))
68 charity_test$wrat <- ordered(charity_test$wrat, levels = c("0", "1", "2", "3", "4",
    "5", "6", "7", "8", "9"))
69
70 # Make home, chld, and hinc into Categorical Variables
71
72
73 charity_overall$home <- as.factor(charity_overall$home)
74 charity_overall$chld <- as.factor(charity_overall$chld)
75 charity_overall$hinc <- as.factor(charity_overall$hinc)
76
77 charity_test$home <- as.factor(charity_test$home)
78 charity_test$chld <- as.factor(charity_test$chld)
79 charity_test$hinc <- as.factor(charity_test$hinc)
80
81 charity_overall$donr <- as.factor(charity_overall$donr)
82 charity_overall$damt <- as.numeric(charity_overall$damt)
83
84 '''
85 reg1, reg2, reg3, and reg4 are 4 different columns that are all categorical
    variables that represent 5 different regions. If there is a 1 in reg1, then a
    new column named region will say 1, if there is a 1 in reg2, then region will
    say 2, and so on. If there is a 0 in all of the reg columns, then the region
    will be 5.
86
87 '''{r}
88
89 charity_overall <- charity_overall %>%
90   pivot_longer(cols = starts_with("reg"), names_to = "region", values_to = "region
    _value") %>%
91   group_by(ID) %>%
92   mutate(region = ifelse(region_value == 1, region, "reg5")) %>%
93   select(-region_value) %>%
94   mutate(region = gsub("reg", "", region))
95
96 charity_overall <- charity_overall %>% group_by(ID) %>% filter(region == min(
    region))
97

```

```

98 charity_overall <- charity_overall %>% distinct(ID, .keep_all = TRUE)
99
100
101 charity_test <- charity_test %>%
102   pivot_longer(cols = starts_with("reg"), names_to = "region", values_to = "region
    _value") %>%
103   group_by(ID) %>%
104   mutate(region = ifelse(region_value == 1, region, "reg5")) %>%
105   select(-region_value) %>%
106   mutate(region = gsub("reg", "", region))
107
108 charity_test <- charity_test %>% group_by(ID) %>% filter(region == min(region))
109
110 charity_test <- charity_test %>% distinct(ID, .keep_all = TRUE)
111
112 charity_overall$region <- as.factor(charity_overall$region)
113 charity_test$region <- as.factor(charity_test$region)
114
115 summary(charity_overall)
116
117 '''
118
119
120 Find if there are any missing variables in the data set.
121
122 '''{r}
123
124 sum(is.na(charity_overall))
125
126
127 sum(is.na(charity_test))
128
129 '''
130
131 No missing data so no need to impute. Time to check correlation between variables.
132
133
134
135 '''{r}
136
137 #create a df from the training set that only includes the columns that are
    quantitative
138
139 charity_overall_quant <- charity_overall[, sapply(charity_overall, is.numeric)]
140 corr_matrix <- cor(charity_overall_quant)
141 corrplot(corr_matrix, method = "color", type = "upper",
142   tl.col = "black", tl.srt = 45,
143   addCoef.col = "black", number.cex = 0.7,
144   col = colorRampPalette(c("blue", "white", "red"))(200))
145
146
147 '''
148 Here are the plots for each significant correlated variable.
149
150 '''{r}
151

```

```

152 #Plot for lgif and agif
153
154 ggplot(charity_overall, aes(x = lgif, y = agif)) +
155   geom_point() +
156   geom_smooth(method = "lm") +
157   labs(title = "lgif vs agif", x = "lgif", y = "agif")
158
159 #Plot for rgif and agif
160
161 ggplot(charity_overall, aes(x = rgif, y = agif)) +
162   geom_point() +
163   geom_smooth(method = "lm") +
164   labs(title = "rgif vs agif", x = "rgif", y = "agif")
165
166 #plot for avhv and incm, inca and avhv, and inca and incm all on one plot
167
168 ggplot(charity_overall, aes(x = avhv, y = incm)) +
169   geom_point() +
170   geom_smooth(method = "lm") +
171   labs(title = "avhv vs incm", x = "avhv", y = "incm")
172
173 ggplot(charity_overall, aes(x = inca, y = avhv)) +
174   geom_point() +
175   geom_smooth(method = "lm") +
176   labs(title = "inca vs avhv", x = "inca", y = "avhv")
177
178 ggplot(charity_overall, aes(x = inca, y = incm)) +
179   geom_point() +
180   geom_smooth(method = "lm") +
181   labs(title = "inca vs incm", x = "inca", y = "incm")
182
183 '''
184 The plots show many outliers as well for the variables. We will attempt to
185   minimize the impact of outliers through standardizing the income and value
186   variables.
187
188 '''{r}
189
190 charity_overall$incm <- log(charity_overall$incm)
191 charity_overall$inca <- log(charity_overall$inca)
192 charity_overall$avhv <- log(charity_overall$avhv)
193
194 charity_test$incm <- log(charity_test$incm)
195 charity_test$inca <- log(charity_test$inca)
196 charity_test$avhv <- log(charity_test$avhv)
197
198 '''
199
200 We seek to determine the expected profit from each mailing. We k
201 Log Regression with everything in the model
202
203 '''{r}
204

```

```

205 #split overall data into train and validation sets based on if the part column
    says train or valid
206
207 charity_train <- charity_overall[charity_overall$part == "train",]
208 charity_valid <- charity_overall[charity_overall$part == "valid",]
209 '''
210
211 '''{r}
212
213 x_train <- charity_train[, !(names(charity_train) %in% c("damt", "part", 'ID', 'donr
    ')))]
214
215 c_train <- charity_train[, (names(charity_train) %in% c('donr'))]# donr
216 n_train_c <- 3984 # 3984
217
218 #damt for rows with donr=1
219
220 y_train <- charity_train$damt[charity_train$donr == 1]# damt for observations with
    donr=1
221 n_train_y <- 1995
222
223 '''
224
225 '''{r}
226
227
228 x_valid <- charity_valid[, !(names(charity_valid) %in% c("damt", "part", 'ID', 'donr
    ')))]
229
230 c_valid <- charity_valid[, (names(charity_valid) %in% c('donr'))]# donr
231 n_valid_c <- 2018
232
233 #damt for rows with donr=1
234
235 y_valid <- charity_valid$damt[charity_valid$donr == 1]# damt for observations with
    donr=1
236 n_valid_y <- 999
237
238 '''
239
240 '''{r}
241
242 n_test <- 2007
243 x_test <- charity_test[, !(names(charity_test) %in% c("damt", "part", 'ID', 'donr'))
    ]
244
245 '''
246
247 '''{r}
248
249 x_train_mean <- apply(x.train, 2, mean)
250 x_train_sd <- apply(x.train, 2, sd)
251 x_train_std <- t((t(x.train)-x.train.mean)/x.train.sd) # standardize to have zero
    mean and unit sd
252 apply(x_train_std, 2, mean) # check zero mean
253 apply(x_train_std, 2, sd) # check unit sd

```



```

254 data_train <- data.frame(x.train.std, donr=c.train) # to classify donr
255 data_train_y <- data.frame(x.train.std[c.train==1,], damt=y.train) # to predict
    damt when donr=1
256
257 x_valid <- t((t(x.valid)-x.train.mean)/x.train.sd) # standardize using training
    mean and sd
258 data_valid <- data.frame(x.valid.std, donr=c.valid) # to classify donr
259 data_valid_y <- data.frame(x.valid.std[c.valid==1,], damt=y.valid) # to predict
    damt when donr=1
260
261 x_test<- t((t(x.test)-x.train.mean)/x.train.sd) # standardize using training mean
    and sd
262 data_test <- data.frame(x.test.std)
263
264 '''
265
266 Cycling through the best binomial model based on Accuracy
267
268 '''{r}
269
270 best_model <- NULL
271 best_mse <- Inf
272
273 data.train <- data.frame(x_train, donr = c_train)
274 data.train$donr <- as.factor(data.train$donr)
275
276 # create full model
277
278 full_model <- glm(donr ~ ., data = data.train, family = binomial(link = "logit"))
279 summary(full_model)
280
281 #model2 based on AIC
282
283 #model2 <- step(full_model, direction = "both", trace = 0)
284 model2 <- glm(formula = donr ~ home + chld + hinc + wrat + incm + plow +
285     npro + tgif + tdon + tlag + region, family = binomial(link = "logit"),
286     data = data.train)
287 summary(model2)
288
289 #model3 using quasi-binomial
290
291 model3 <- glm(formula = donr ~ home + chld + hinc + wrat + incm + plow +
292     npro + tgif + tdon + tlag + region, family = quasibinomial(link = "logit"),
293     data = data.train)
294 summary(model3)
295
296 '''
297 '''{r}
298
299 # Calculate ordered profit function using average donation = $14.50 and mailing
    cost = $2
300
301 # Calculate profit for each model
302
303 pred1 <- predict(full_model, newdata = data_vlau, type = "response")
304 profit.log1 <- cumsum(14.5*c.valid[order(pred1, decreasing=T)]-2)

```

```

305 plot(profit.log1)
306
307 pred2 <- predict(model2, newdata = x_valid, type = "response")
308 profit.log2 <- cumsum(14.5*c.valid[order(pred2, decreasing=T)]-2)
309 plot(profit.log2)
310
311 pred3 <- predict(model3, newdata = x_valid, type = "response")
312 profit.log3 <- cumsum(14.5*c.valid[order(pred3, decreasing=T)]-2)
313 plot(profit.log3)
314
315 '''
316 Report the number for maximum profit
317
318 '''{r}
319 which.max(profit.log1)
320 which.max(profit.log2)
321 which.max(profit.log3)
322 max(profit.log1)
323 max(profit.log2)
324 max(profit.log3)
325
326 '''
327
328 '''{r}
329 cutoff.log <- sort(pred2, decreasing=T)[which.max(profit.log2)+1] # set cutoff
    based on n.mail.valid
330 chat.valid.log <- ifelse(pred2>cutoff.log, 1, 0) # mail to everyone above the
    cutoff
331 table(chat.valid.log, c.valid)
332
333 '''
334
335
336 '''{r}
337
338
339 post.test <- predict(model2, x_test, type="response") # post probs for test data
340
341 # Oversampling adjustment for calculating number of mailings for test set
342 n.mail.valid <- which.max(profit.log1)
343 tr.rate <- .1 # typical response rate is .1
344 vr.rate <- .5 # whereas validation response rate is .5
345 adj.test.1 <- (n.mail.valid/n.valid.c)/(vr.rate/tr.rate) # adjustment for mail yes
346 adj.test.0 <- ((n.valid.c-n.mail.valid)/n.valid.c)/((1-vr.rate)/(1-tr.rate)) #
    adjustment for mail no
347 adj.test <- adj.test.1/(adj.test.1+adj.test.0) # scale into a proportion
348 n.mail.test <- round(n.test*adj.test, 0) # calculate number of mailings for test
    set
349
350 cutoff.test <- sort(post.test, decreasing=T)[n.mail.test+1] # set cutoff based on
    n.mail.test
351 chat.test <- ifelse(post.test>cutoff.test, 1, 0) # mail to everyone above the
    cutoff
352 table(chat.test)
353
354 '''

```

```

355
356 ''{r}
357
358 data_train_y <- charity_train[, !(names(charity_train) %in% c("part", 'ID', 'donr'))
359 ]
360
361 model1 <- lm(damt ~ ., data_train_y)
362
363 summary(model1)
364
365 #choose model 2 based on AIC
366
367 model2 <- step(model1, direction = "both", trace = 0)
368 summary(model2)
369
370 #choose model 3 using stepwise selection
371
372 model3 <- step(model1, direction = "both", trace = 0, k = log(n_train_y))
373 summary(model3)
374
375 ''
376 ''{r}
377 #compare the models using the validation set
378
379 data_valid_y <- charity_valid[, !(names(charity_train) %in% c("part", 'ID', 'donr'))
380 ]
381
382 pred.valid1 <- predict(model1, newdata = data_valid_y) # validation predictions
383 mean((y_valid - pred.valid1)^2) # mean squared error
384 # 1.867523
385 sd((y_valid - pred.valid1)^2)/sqrt(n_valid_y)
386
387 pred.valid2 <- predict(model2, newdata = data_valid_y) # validation predictions
388 mean((y_valid - pred.valid2)^2) # mean squared error
389 sd((y_valid - pred.valid2)^2)/sqrt(n_valid_y)
390
391 pred.valid3 <- predict(model3, newdata = data_valid_y) # validation predictions
392 mean((y_valid - pred.valid3)^2) # mean squared error
393 sd((y_valid - pred.valid3)^2)/sqrt(n_valid_y)
394
395 ''
396 ''{r}
397 yhat.test <- predict(model3, newdata = charity_test) # test predictions)
398
399 length(chat.test) # check length = 2007
400 length(yhat.test) # check length = 2007
401 chat.test[1:10] # check this consists of 0s and 1s
402 yhat.test[1:10]
403
404 ip <- data.frame(ID=data.test$ID, donr=chat.test, damt=yhat.test)
405 ip$damt[ip$donr == 0] <- 0
406
407 ''
408 ''{r}

```

```

409
410 idx.donors <- ip[(ip$donr > 0), ]
411 num_donors <- nrow(idx.donors)
412 est.profits <- sum(idx.donors$damt) - 2*num_donors
413 round(est.profits,2)
414
415 '''
416
417
418 '''{r}
419
420 submit <- data.frame(ID=ip$ID, damt=ip$damt) # data frame with two variables: ID
      and DAMT
421 write.csv(submit, file="Blackmon_submission.csv", row.names=FALSE)
422 '''
423
424
425
426 '''{r}
427
428 c.valid <- charity_valid[, (names(charity_valid) %in% c("donr"))]
429
430 profit.log1 <- cumsum(14.5*c.valid[order(log_donr_pred, decreasing=T)]-2)
431
432
433 '''

```