

# Lab 02

*Alex Yang, John Kenney, Ram Balasubramanian*

*Oct 12, 2017*

## SECTION - 1 INTRODUCTION & KEY RESULTS

### Problem Introduction:

We have been hired by a Private University to identify who among their Alumni are most likely to contribute towards the University's foundation in future years. The university has provided us with data on past contributions from graduates - data includes some demographic information (like gender, marital status), university specific information (graduation year, major of studies), and some information on how "connected" an Alumnus is to the school (Alumni event attendance, historical contributions).

### 1.1 HIGH LEVEL DESCRIPTION OF MODELING APPROACH:

We have taken two approaches to the problem (named Beta-Hat and Y-Hat):

Approach "Beta-Hat": We will treat the problem as a "explanation" problem ( $\hat{\beta}$ ). The goal here is to figure out if and how much certain aspects of a person and their association with the university determines how much they will contribute to the university's foundation. We will develop a regression model that considers the 2016 contributions as a variable that depends on one or more of the other data elements that have been provided. The regression coefficients can then be interpreted as a measure of how much each aspect of a person influences their contributions.

Approach "Y-Hat":

We will treat the problem as a "prediction" problem ( $\hat{y}$  problem). Given all the data we have about a person and their past contributions, can we predict how much they will contribute in the future. We will develop a model that aims to predict the 2016 contribution amounts for each person. To evaluate the efficacy of our models, we will split the data into a "training" set and a "test" set. We will use the training data to estimate parameters for our prediction model and evaluate our model's prediction accuracy using the test set.

### 1.2 KEY RESULTS AND TECHNIQUES USED:

We will complete this section once we are done with the modeling work.

## SECTION 2 - DATA EXAMINATION AND EDA:

```
library(knitr)
opts_chunk$set(tidy.opts=list(width.cutoff=60),tidy=TRUE)

#Libraries required
library(car)
library(dplyr)
```

```
library(Hmisc)#Used by author for 3D plotting
library(ggplot2)
library(gridExtra)
library(effsize) #Used to calculate Cohen's D for T-Test
library(aod)      #Used for effect size of the logit model
library(mcpfile) #Used for confidence intervals
library(package = MASS) # Location of parcoord() function
library(vcd)
library(data.table)
library(stargazer)
library(caret) #Required for Confusion Matrix
library(ordinal)
library(GGally) #scatterplot matrix with jitter
library(reshape2) # For facet Grids
library(tidyr)
library(purrr)
```

```
dt <- fread("lab2data.csv")
# describe(dt) #This takes 2.5 pages on its own
```

## 2.1 Brief Description of Data Available:

We have data for 1000 past graduates of the University. There are 12 variables provided for each Alumnus. They are: 1. V1: Identifier for each record (Alumnus)

2. Gender: M/F, roughly 50/50 in the sample provided. 3. Class.Year: Appears like the “decade” of the graduating year. Goes from 1972 - 2012. We will assume 1972 represents students graduating from 1963-1972; 1982 represents students graduating from 1973 to 1982 etc.

4. Marital.Status: Has 4 categories - coded D,M,S,W. We will assume it means Divorced, Married, Single, Widowed with over 90% in the “married” or “single” categories.

5. Major: There are 45 majors represented in the sample. History, English, Biology & Economics are the top 4 representing about 37% of the sample.

6. Next.Degree: We assume this means what the alumnus went on to do after graduating from the university. 38% shows “None” implying they did not pursue another degree. The remainder (62%) seems rather high for this metric.

7. AttendanceEvent: Indicates whether the alumnus attended an alumni event between 2012 and 2015. If we choose to use this variable to model “Giving” we should probably not use it to model 2012-2014 Giving

8. FYGiving: There are 5 of these variables named FY12 - FY16 representing full year 2012 through full year 2016 contribution from the alumnus. There are some “outliers” (e.g. \$161,500 in 2013) in the data that we may need to be on the lookout for.

We do not have any missing values in the data; and there do not seem to be an obvious “data cleaning” that needs to be conducted. We will conduct an Exploratory Data Analysis next.

What are the important variables we want to include in our discussion here? What would we suppose would be meaningful? What can we omit?

```
# View the contents of Major and Next Degree - to identify if  
# there are any  
majortable = as.data.frame(round(prop.table(table(dt$Major)),  
  2))  
degtable = as.data.frame(round(prop.table(table(dt$Next.Degree)),  
  2))
```

### 1.3 Create new variables:

Let's group the yearly contributions by the categories that the university is interested in; Classify the "next degree" variable into 0 (representing "none") and 1 (representing there was some next-degree). Create indicator variables for each year for giver(1) or not a giver(0). For each alumnus let's also count the number of years they have given between 2012 and 2015.

```
dt$FY16GivingCat <- cut(dt$FY16Giving, c(0, 1, 100, 250, 500,  
  2e+05), right = FALSE)  
dt$FY15GivingCat <- cut(dt$FY15Giving, c(0, 1, 100, 250, 500,  
  2e+05), right = FALSE)  
dt$FY14GivingCat <- cut(dt$FY14Giving, c(0, 1, 100, 250, 500,  
  2e+05), right = FALSE)  
dt$FY13GivingCat <- cut(dt$FY13Giving, c(0, 1, 100, 250, 500,  
  2e+05), right = FALSE)  
dt$FY12GivingCat <- cut(dt$FY12Giving, c(0, 1, 100, 250, 500,  
  2e+05), right = FALSE)  
  
# create an indicator for 'giver' and 'non giver' for each  
# year.  
dt$Giver16 = as.integer(dt$FY16Giving > 0)  
dt$Giver15 = as.integer(dt$FY15Giving > 0)  
dt$Giver14 = as.integer(dt$FY14Giving > 0)  
dt$Giver13 = as.integer(dt$FY13Giving > 0)  
dt$Giver12 = as.integer(dt$FY12Giving > 0)  
dt$YearsGiven = dt$Giver12 + dt$Giver13 + dt$Giver14 + dt$Giver15  
  
# Create identifier for next degree (1) or none (0)  
dt$NextDegCat = 1 - as.integer((dt$Next.Degree == "NONE"))  
  
# Group majors by broad categories  
  
dt$MajorCat = ifelse(dt$Major %in% c("American Studies", "Art",  
  "Chinese", "Classics", "Comparative Literature", "English",  
  "English-Journalism", "French", "German", "History", "Independent",  
  "Music", "Philosophy", "Philosophy-Religion", "Physical Education",  
  "Religious Studies", "Russian", "Spanish", "Speech (Drama, etc.)",  
  "Theatre"), "HUM_ART", ifelse(dt$Major %in% c("Biology",
```

```
"Chemistry", "Computer Science", "Engineering", "General Science",
"General Science-Biology", "General Science-Chemistry", "General Science-Math",
"General Science-Physics", "Mathematics", "Mathematics-Physics",
"Physics", "Zoology"), "STEM", ifelse(dt$Major %in% c("Economics",
"Economics-Regional Stds.", "Sociology", "Psychology", "Pol. Sci.-Regional Stds.",
"Sociology-Anthropology", "Political Science", "Anthropology",
"Economics-Business"), "SOCIAL_SCIENCE", "OTHER"))))
```

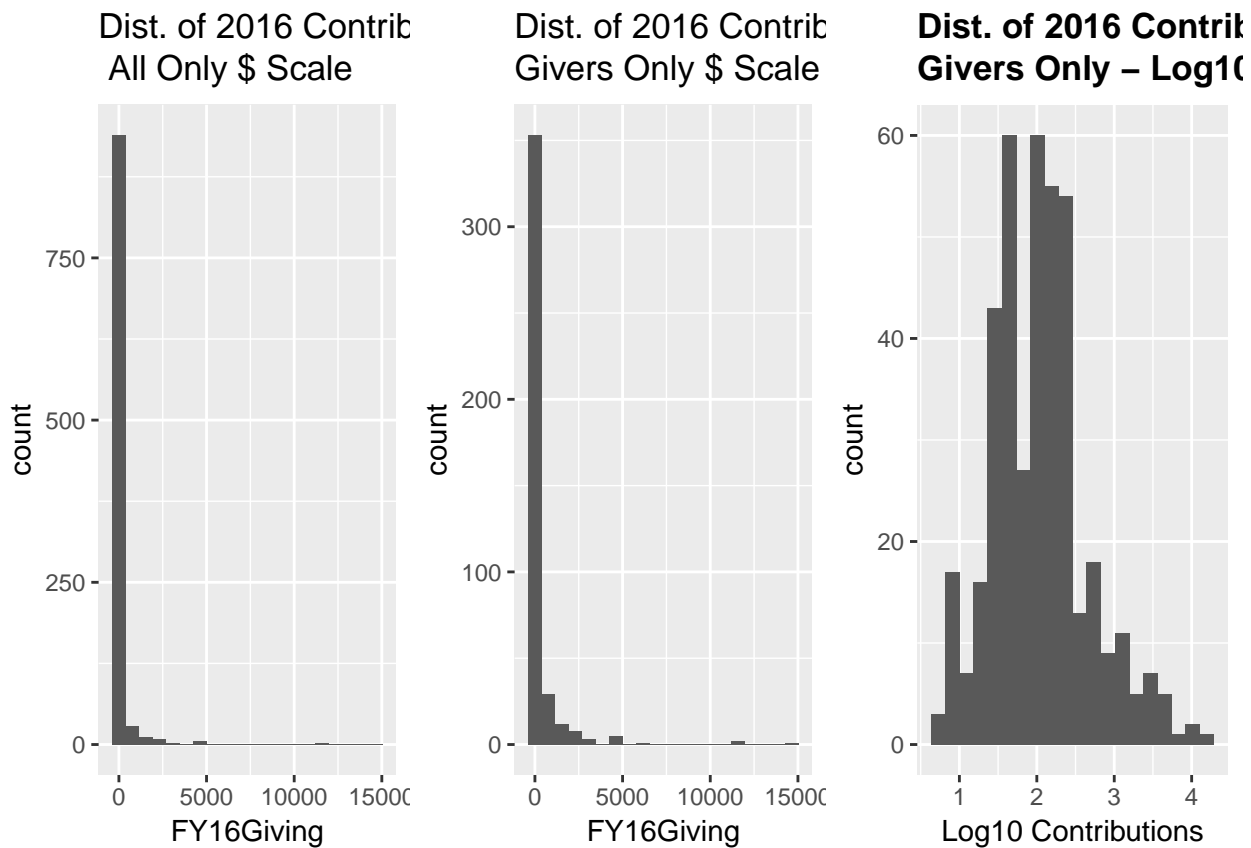
```
dt$MaritalStatusCat = factor(dt$Marital.Status)
dt$ClassYearCat = factor(dt$Class.Year)
```

## 2.2 Exploratory Data Analysis:

### 2.2.1 Univariate Analysis:

Let's examine each variable first starting with the "variable of interest" - 2016 contributions.

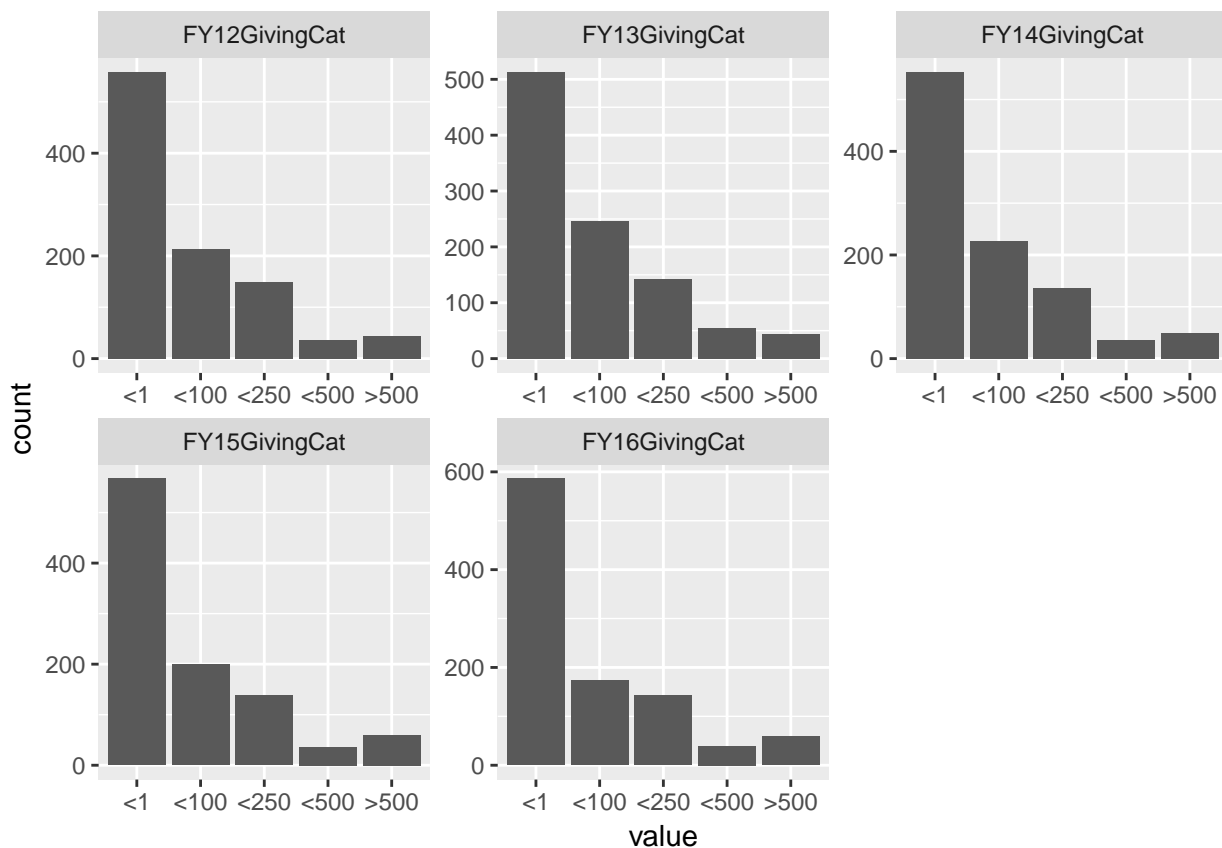
```
h1 = ggplot(data = dt, aes(x = FY16Giving)) + geom_histogram(bins = 20) +
  ggtitle("Dist. of 2016 Contributions \n All Only $ Scale") +
  theme(plot.title = element_text(lineheight = 1))
h2 = ggplot(data = dt[FY16Giving > 0], aes(x = FY16Giving)) +
  geom_histogram(bins = 20) + ggtitle("Dist. of 2016 Contributions \nGivers Only $ Scale") +
  theme(plot.title = element_text(lineheight = 1))
h3 = ggplot(data = dt[FY16Giving > 0], aes(x = log10(FY16Giving))) +
  geom_histogram(bins = 20) + ggtitle("Dist. of 2016 Contributions\nGivers Only - Log10 $ Scale") +
  xlab("Log10 Contributions") + theme(plot.title = element_text(lineheight = 1,
  face = "bold"))
grid.arrange(h1, h2, h3, ncol = 3)
```



Most of the alumni contributed \$0; so we see a big spike at 0 and then it tapers off quickly. When we filter for alumni who donated more than 0, we again see a very similar pattern-most alumni who did contribute contributed very little. The log transform shows a something approaching a normal distribution, but skewed to the left.

In 2016, roughly half the alumni population did not give anything. 20-25%% gave less than \$100; Around 14% gave \$100-250; About 4% gave \$250-500 and less than 1% give more than \$500. We see from the figures that the alumni contributions in all years followed similar distributions. In 2013, we had the highest percentage of the Alumni contributing (about 49%) and 2015 marked the lowest % giving (43%)

```
lcat <- c("FY16GivingCat", "FY15GivingCat", "FY14GivingCat",
         "FY13GivingCat", "FY12GivingCat")
dt[, lcat, with = FALSE] %>% na.omit() %>% gather() %>% ggplot(aes(value)) +
  facet_wrap(~key, scales = "free") + scale_x_discrete(labels = c("<1",
    "<100", "<250", "<500", ">500")) + geom_bar()
```



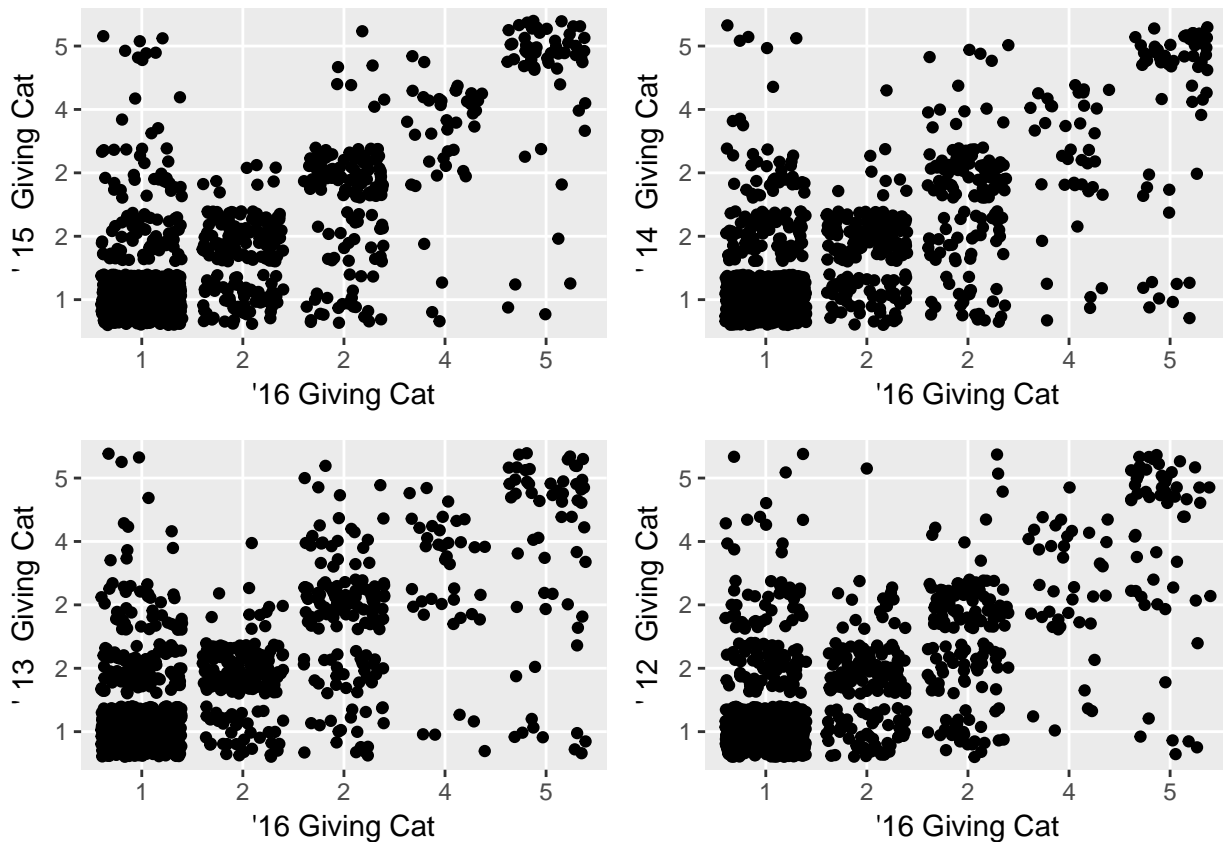
### 2.2.2 Bivariate Analysis:

#### Giving in 2016 vs. 2015:

Let's look at how 2016 giving relates to 2015. The plot shows that there is a reasonable correlation between the two (note the density of dots along the diagonal). We also conducted the test for each of the other years. Each of the tests show that there is a dependence between 2016 giving and past years' giving category. Essentially, what we are observing is that someone's giving category in 2016 is most likely to be the same as their past; it is also interesting to note that in most instances the second highest category is the "[0,1)" category - so basically either they give like they have given in the past or not give at all!

```
ggScatterJitter <- function(yearcat) {
  yeargive <- paste("FY", toString(yearcat), "GivingCat", sep = "")
  # print(yeargive)
  p <- ggplot(dt, aes_string(x = "FY16GivingCat", y = yeargive))
  p <- p + geom_jitter()
  p <- p + scale_x_discrete(name = "'16 Giving Cat", labels = c("1",
    "2", "3", "4", "5")) + scale_y_discrete(name = paste("",
    toString(yearcat), " Giving Cat"), labels = c("1", "2",
    "3", "4", "5"))
  return(p)
}
sj1 <- ggScatterJitter(15)
sj2 <- ggScatterJitter(14)
sj3 <- ggScatterJitter(13)
```

```
sj4 <- ggScatterJitter(12)
grid.arrange(sj1, sj2, sj3, sj4, ncol = 2)
```



```
# #Let's comment on assocstats
GenXtab = function(dframe, x1, x2, nlist) {
  x1vsx2 = xtabs(formula = ~x1 + x2, data = dframe)
  names(dimnames(x1vsx2)) = nlist
  print(x1vsx2)
  print("Percentage of Column Totals Shown Below")
  print(round(prop.table(x1vsx2, 2), 2))
  a.s = assocstats(x1vsx2)
  a.s
  if (is.null(a.s$phi) | is.na(a.s$phi)) {
    print("Phi: Not Applicable")
  } else if (abs(a.s$phi) > 0.5) {
    print("Phi: Large Effect")
  } else if (abs(a.s$phi) > 0.3) {
    print("Phi: Medium Effect")
  } else if (abs(a.s$phi) > 0.1) {
    print("Phi: Small Effect")
  } else {
    print("Phi: Negligible Effect")
  }

  if (is.null(a.s$contingency)) {
    print("Contingency Coef: Not Applicable")
  }
}
```

```

} else if (abs(a.s$contingency) > 0.5) {
  print("Contingency Coef: Large Effect")
} else if (abs(a.s$contingency) > 0.3) {
  print("Contingency Coef: Medium Effect")
} else if (abs(a.s$contingency) > 0.1) {
  print("Contingency Coef: Small Effect")
} else {
  print("Contingency Coef: Negligible Effect")
}

if (is.null(a.s$cramer)) {
  print("Cramer's V: Not Applicable")
} else if (abs(a.s$cramer) > 0.5) {
  print("Cramer's V: Large Effect")
} else if (abs(a.s$cramer) > 0.3) {
  print("Cramer's V: Medium Effect")
} else if (abs(a.s$cramer) > 0.1) {
  print("Cramer's V: Small Effect")
} else {
  print("Cramer's V: Negligible Effect")
}
}

```

### Giving in 2016 vs. Alumni Event Attendance:

Attendance at an alumni event is the factor (other than previous years' contributions), that our intuitions say will have the strongest relationship with 2016 give.

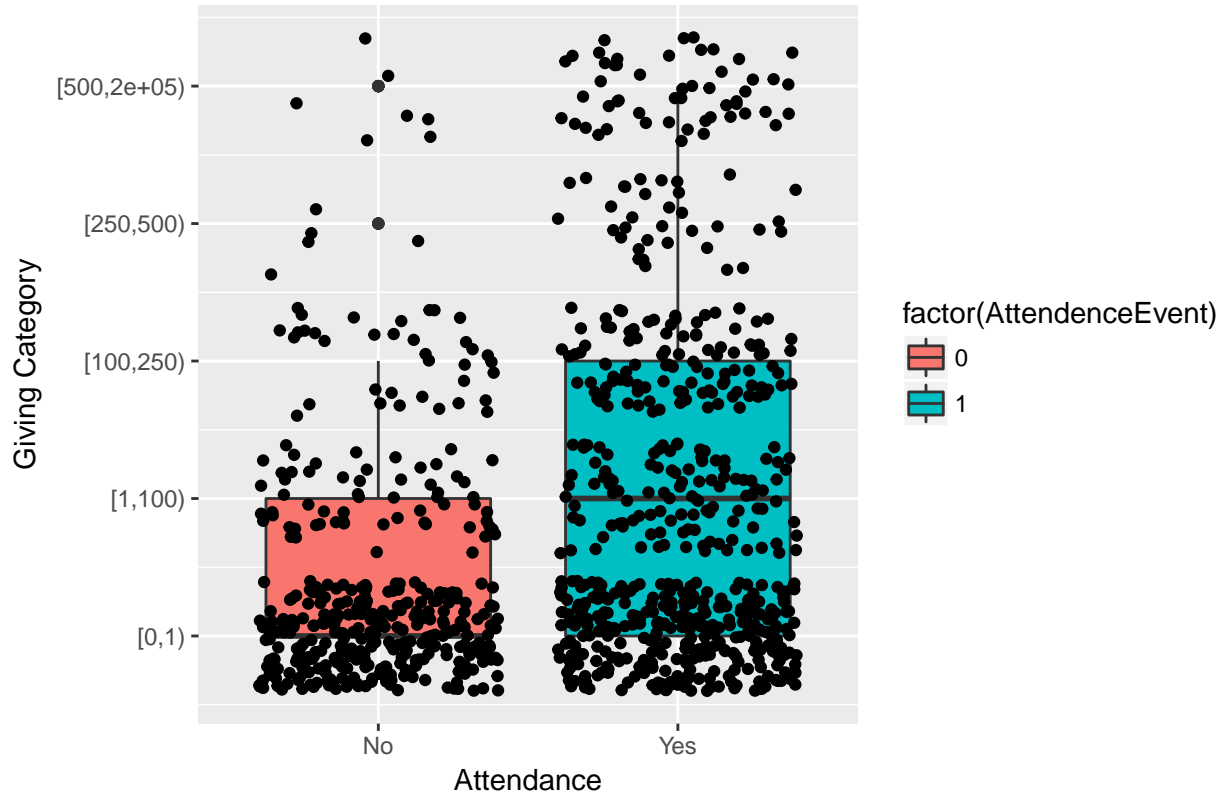
```

ggplot(dt, aes(factor(AttendanceEvent), as.numeric(FY16GivingCat))) +
  geom_boxplot(aes(fill = factor(AttendanceEvent))) + ggtitle("Giving Category by Attendance at
  geom_jitter() + scale_x_discrete(name = "Attendance", labels = c("No",
  "Yes")) + scale_y_continuous(name = "Giving Category", breaks = 1:5,
  labels = c("[0,1)", "[1,100)", "[100,250)", "[250,500)",
  "[500,2e+05)")) + theme(plot.title = element_text(lineheight = 1,
  face = "bold"))

```



## Giving Category by Attendance at Event



And

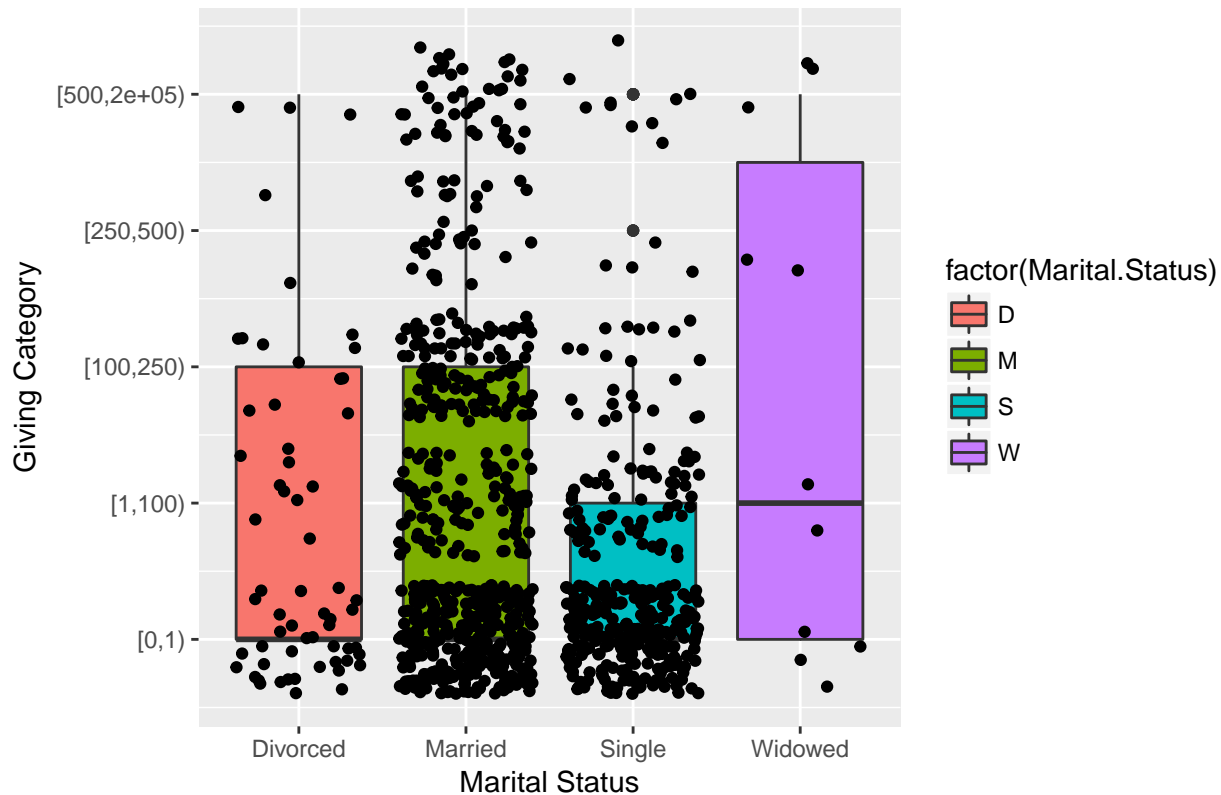
there is a relationship between attendance at alumni events in 2012-2015 and 2016 giving. 50% of those attending gave in 2016 while only 28% of those not attending gave in 2016.

## Giving in 2016 by Marital status,:

We do not expect a relationship between 2016 contributions and marital status

```
ggplot(dt, aes(factor(Marital.Status), as.numeric(FY16GivingCat))) +
  geom_boxplot(aes(fill = factor(Marital.Status))) + ggtitle("Giving Category by Marital Status") +
  geom_jitter() + scale_x_discrete(name = "Marital Status",
  labels = c("Divorced", "Married", "Single", "Widowed")) +
  scale_y_continuous(name = "Giving Category", breaks = 1:5,
  labels = c("[0,1)", "[1,100)", "[100,250)", "[250,500)", "[500,2e+05)")) + theme(plot.title = element_text(lineheight = 1,
  face = "bold"))
```

## Giving Category by Marital Status



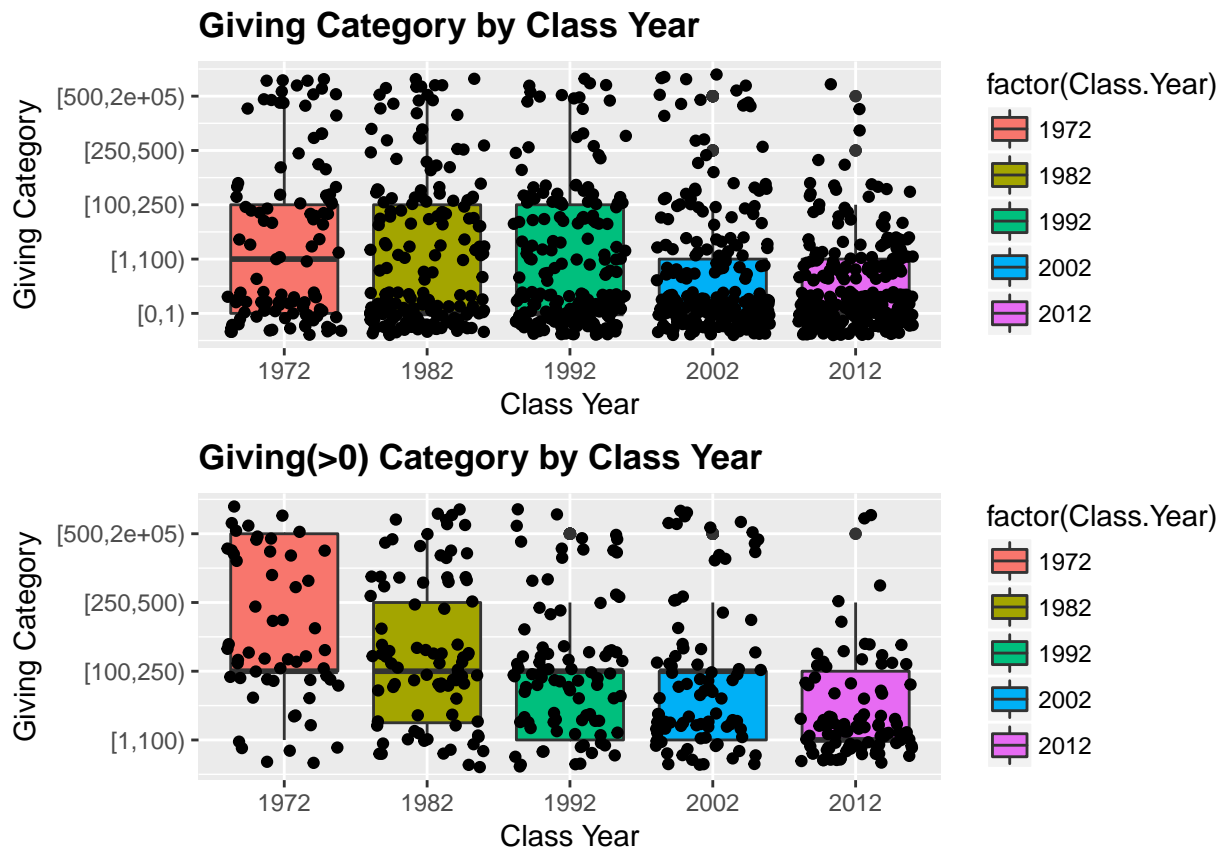
Most of the Alumni fall into either Married or Single category and Married Alumni are more likely to give than Single.

There are very few data points for Widowed and Divorced alumni - so we do not want to make broad conclusions, but it appears that there are both Divorced and Widowed Alumni that contribute high amounts (and there are those that contribute nothing too in these categories) One possible reason we are seeing Married giving more than Single might actually have to do with age. Older Alumni are more likely to be married and older alumni are also probably a bit more well established financially - so more likely to contribute to charitable causes. So Marital status vs. Giving might simply be capturing the relationship between Age and Giving. While age is not a variable that's available in the dataset, we have "Class Year" which is a good proxy for age.

## Giving 2016 vs. Class Year:

```
year1 <- ggplot(dt, aes(factor(Class.Year), as.numeric(FY16GivingCat))) +
  geom_boxplot(aes(fill = factor(Class.Year))) + ggtitle("Giving Category by Class Year") +
  geom_jitter() + scale_x_discrete(name = "Class Year") + scale_y_continuous(name = "Giving Category",
    breaks = 1:5, labels = c("[0,1)", "[1,100)", "[100,250)",
      "[250,500)", "[500,2e+05)")) + theme(plot.title = element_text(lineheight = 1,
    face = "bold"))
year2 <- ggplot(dt[dt$FY16Giving > 0, ], aes(factor(Class.Year),
  as.numeric(FY16GivingCat))) + geom_boxplot(aes(fill = factor(Class.Year))) +
  ggtitle(">0 Giving Category by Class Year") + geom_jitter() +
  scale_x_discrete(name = "Class Year") + scale_y_continuous(name = "Giving Category",
    breaks = 1:5, labels = c("[0,1)", "[1,100)", "[100,250)",
      "[250,500)", "[500,2e+05)")) + theme(plot.title = element_text(lineheight = 1,
```

```
face = "bold"))
grid.arrange(year1, year2, ncol = 1)
```

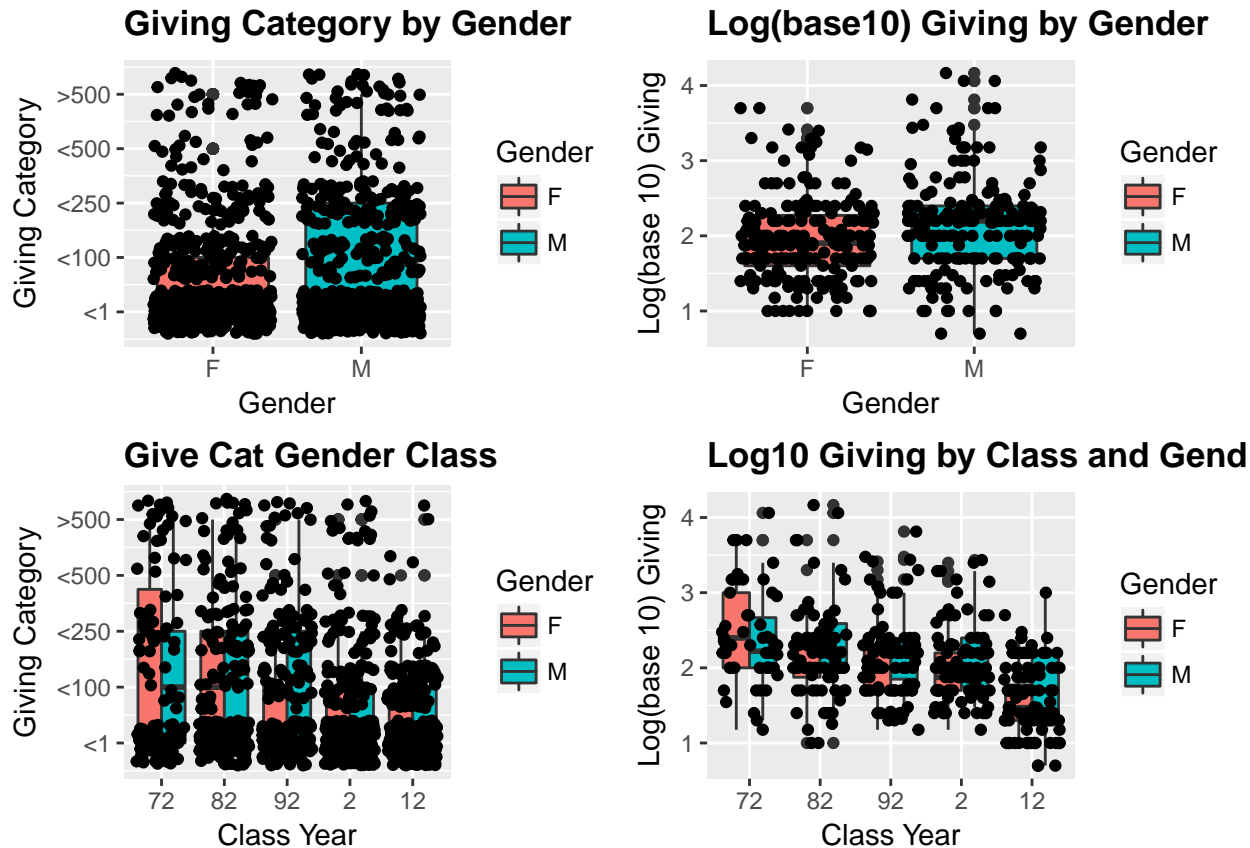


In the first chart, we see the more recent graduates tend to be more likely to have donated nothing. In the second chart, we filter out all the non-givers, and we see that older graduates tend to have donated more, given that they donated anything at all.

### Giving in 2016 by Gender:

```
gender1 <- ggplot(dt, aes(Gender, as.numeric(FY16GivingCat))) +
  geom_boxplot(aes(fill = Gender)) + ggtitle("Giving Category by Gender") +
  geom_jitter() + scale_x_discrete(name = "Gender") + scale_y_continuous(name = "Giving Category",
    breaks = 1:5, labels = c("<1", "<100", "<250", "<500", ">500")) +
  theme(plot.title = element_text(lineheight = 1, face = "bold"))
gender2 <- ggplot(dt[FY16Giving > 0], aes(Gender, log10(FY16Giving))) +
  geom_boxplot(aes(fill = Gender)) + ggtitle("Log(base10) Giving by Gender") +
  geom_jitter() + scale_x_discrete(name = "Gender") + scale_y_continuous(name = "Log(base 10) Giving",
    breaks = 1:5, labels = c("<1", "<100", "<250", "<500", ">500")) +
  theme(plot.title = element_text(lineheight = 1, face = "bold"))
yearGender1 <- ggplot(dt, aes(factor(Class.Year), as.numeric(FY16GivingCat))) +
  geom_boxplot(aes(fill = Gender)) + ggtitle("Give Cat Gender Class") +
  geom_jitter() + scale_x_discrete(name = "Class Year", labels = c(72, 82, 92, 2, 12)) + scale_y_continuous(name = "Giving Category",
    breaks = 1:5, labels = c("<1", "<100", "<250", "<500", ">500")) +
  theme(plot.title = element_text(lineheight = 1, face = "bold"))
yearGender2 <- ggplot(dt[FY16Giving > 0], aes(factor(Class.Year),
```

```
log10(FY16Giving))) + geom_boxplot(aes(fill = Gender)) +
ggtitle("Log10 Giving by Class and Gender") + geom_jitter() +
scale_x_discrete(name = "Class Year", labels = c(72, 82,
92, 2, 12)) + scale_y_continuous(name = "Log(base 10) Giving") +
theme(plot.title = element_text(lineheight = 1, face = "bold"))
grid.arrange(gender1, gender2, yearGender1, yearGender2, ncol = 2)
```



There is no statistical evidence that whether or not Alumni Give in 2016 varies by Gender (41% of Females Gave vs. 42% of Males). However when we look at the categories of contribution in 2016 by Gender, we see differences worth investigating later. There may be other factors at play here, for e.g. we know that older alumni tend to give more than younger alumni; it is possible that there are fewer “older” female alumni (fewer women attended college in 1972) than Male.

In the top two charts, we that of those that donate, men donate more. But we can’t be sure this is significant as the median man donated less than the 75th percentile woman.

In the bottom two, we examine the interaction between gender and class year.

Split by class year, gender does not seem to be a significant factor, and in fact older women donate more than older men. Limiting it to just alumni who donated in 2016, we do see that even split by age, men donated more than women, except for the oldest alumni