

W271 Section 3 Lab 1

Daghan Atlas, Zhaoning Yu, Hoang Phan

9/23/2017

Problem statement

In this lab, we are going to model the relationship between age and voters' preference for Bernie Sanders over Hillary Clinton.

Dataset

The dataset comes from the 2016 American National Election Survey.

```
library(dplyr)
library(ggplot2)
library(Hmisc)
library(GGally)
library(data.table)
library(stargazer)

if (dir.exists("/Users/daghanaltas/Hacking/Berkeley/W271/Labs/w271_lab1/")) {
  setwd("/Users/daghan/Hacking/Berkeley/W271/Labs/w271_lab1/")
} else if (dir.exists("/Users/daghan/Hacking/Berkeley/W271/Labs/w271_lab1/")) {
  setwd("/Users/daghan/Hacking/Berkeley/W271/Labs/w271_lab1/")
} else {
  print("add yor local directory path here")
}

df <- read.csv("./public_opinion.csv")
dt <- data.table(df)
head(dt)
```

```
##      sanders_preference party race_white gender birthyr
## 1:                1      1           1      1      1960
## 2:                0      2           1      2      1957
## 3:                1      3           1      1      1963
## 4:                1      1           1      1      1980
## 5:                1      2           1      1      1974
## 6:                1      2           1      1      1958
```

```
describe(dt)
```

```
## dt
##
## 5 Variables      1200 Observations
## -----
## sanders_preference
##      n missing distinct      Info      Sum      Mean      Gmd
## 1191      9         2    0.733    686    0.576    0.4889
```

```
##
## -----
## party
##      n missing distinct      Info      Mean      Gmd
##    1200      0        3    0.875    1.851    0.8309
##
## Value      1      2      3
## Frequency  459   461   280
## Proportion 0.382 0.384 0.233
## -----
## race_white
##      n missing distinct      Info      Sum      Mean      Gmd
##    1200      0        2    0.592      875    0.7292    0.3953
##
## -----
## gender
##      n missing distinct      Info      Mean      Gmd
##    1200      0        2    0.748    1.525    0.4992
##
## Value      1      2
## Frequency  570   630
## Proportion 0.475 0.525
## -----
## birthyr
##      n missing distinct      Info      Mean      Gmd      .05      .10
##    1200      0        73      1    1968    19.53    1940    1946
##      .25      .50      .75      .90      .95
##    1955    1968    1982    1991    1994
##
## lowest : 1921 1924 1925 1926 1927, highest: 1993 1994 1995 1996 1997
## -----
```

Description of the data

The dataset contains 5 variables with 1200 samples:

- **sanders__preference**: A categorical variable with 2 levels, denoting whether the voter prefers Bernie Sanders (=1) or Hillary Clinton (=0).
- **party**: A categorical variable with 3 levels, denoting whether the voter prefers is affiliated with the Democratic Party (=1), Independent (=2), or Republican Party (=3).
- **race__white**: A categorical variable with 2 levels, denoting wheter the voter is White (=1), or not (=0).
- **gender**: A categorical variable with 2 levels, denoting whether the voter is male (=1), or female (=2).
- **birthyr**: A numerical variable, denoting the birthyear of the voter.

Observations:

- There are 9 missing values (NAs) for the **sanders__preference** variable.
- There is no direct **age** variable. We'll derive it from the **birthyr** variable.
- None of the other variables have missing values.

Clean-up

- We are going to create a new variable, **age** based on the **birthyr** variable through the following formula:

$$age = 2015 - birthyr$$

We will use 2015, as the data was collected in 2016, since we are interested in the age of the voters *when* the data was collected. (All people born in 2015 will be counted as age=1, but this is only true if the survey was taken on 12/31/2016; suppose the survey was taken on 01/01/2016, then none of the people born in 2015 actually reached age 1, they are all actually age 0 when being surveyed, in this case “age = 2015 - birthyr” is the correct formula. We don’t know when the survey was taken, but by changing the calculation to “age = 2015 - birthyr”, the minimum age is 18 now)

- We’ll convert all categorical variables to R factor variables.

```
dim(dt)

## [1] 1200    5

dt$sanders_preference <- as.factor(dt$sanders_preference)
dt$party <- factor(x = dt$party, levels = c("1", "2", "3"), labels = c("D",
  "I", "R"))
dt$race_white <- as.factor(dt$race_white)
dt$gender <- as.factor(dt$gender)
dt$age <- 2015 - dt$birthyr
dim(dt)

## [1] 1200    6

head(dt)

##      sanders_preference party race_white gender birthyr age
## 1:                   1      D          1      1    1960  55
## 2:                   0      I          1      2    1957  58
## 3:                   1      R          1      1    1963  52
## 4:                   1      D          1      1    1980  35
## 5:                   1      I          1      1    1974  41
## 6:                   1      I          1      1    1958  57
```

- We are going to remove the 9 observations without the **sanders_preference** value. A possible way to impute these **NA** values could be to use a logistic regression but for this work, we’ll simply opt to remove these observations (9 out of 1200 observations is less than 1% of the data), especially given that the other values for each observation are close to the mean of their respective columns. A table of the NA observations are detailed below:

```
## Hoang to add table
dt_missing = dt[is.na(sanders_preference)]
dt_notmissing = dt[!is.na(sanders_preference)]

table = data.table(Data = c("NA's", "Remaining Data"), PartyDem = c(round(nrow(dt_missing[party ==
  "D"])/nrow(dt_missing), 2), round(nrow(dt_notmissing[party ==
  "D"])/nrow(dt_notmissing), 2)), PartyInd = c(round(nrow(dt_missing[party ==
  "I"])/nrow(dt_missing), 2), round(nrow(dt_notmissing[party ==
  "I"])/nrow(dt_notmissing), 2)), PartyRep = c(round(nrow(dt_missing[party ==
  "R"])/nrow(dt_missing), 2), round(nrow(dt_notmissing[party ==
  "R"])/nrow(dt_notmissing), 2)), Race_white = c(round(nrow(dt_missing[race_white ==
  1])/nrow(dt_missing), 2), round(nrow(dt_notmissing[race_white ==
  1])/nrow(dt_notmissing), 2)), GenderFemale = c(round(nrow(dt_missing[gender ==
  2])/nrow(dt_missing), 2), round(nrow(dt_notmissing[gender ==
```

```
2))/nrow(dt_notmissing), 2)), MeanAge = c(round(dt_missing[,
mean(age)], 2), round(dt_notmissing[, mean(age)], 2)))
table
```

```
##           Data PartyDem PartyInd PartyRep Race_white GenderFemale
## 1:           NA's      0.44      0.33      0.22      0.67      0.89
## 2: Remaining Data      0.38      0.38      0.23      0.73      0.52
##      MeanAge
## 1:    49.89
## 2:    47.04
```

```
dt <- dt[!is.na(dt$sanders_preference), ]
```

We observe that 8 out of 9 NAs are female. This requires a careful review. Even if they all preferred Bernie Sanders, it doesn't make a big difference in the overall ratios. Furthermore, our analysis (below) show that gender is not a meaningful factor for our final model. All other variables deviate very little between the NAs subset and the rest of the data. We justify our decision to remove NA data based on these observations.

Explotary Data Analysis

Univariate analysis

Sanders vs. Hillary Preference

```
# xtabs( ~ sanders_preference, data=dt)
c.table <- array(data = c(sum(dt$sanders_preference == 1), sum(dt$sanders_preference ==
1)/length(dt$sanders_preference), sum(dt$sanders_preference ==
0), sum(dt$sanders_preference == 0)/length(dt$sanders_preference),
sum((dt$sanders_preference == 0) | (dt$sanders_preference ==
1)), sum((dt$sanders_preference == 0) | (dt$sanders_preference ==
1))/length(dt$sanders_preference)), dim = c(2, 3), dimnames = list(Count = c("Voter Count",
"pi.hat"), Preference = c("Bernie", "Hillary", "Total")))

round(c.table, 2)
```

```
##           Preference
## Count      Bernie Hillary Total
## Voter Count 686.00  505.00 1191
## pi.hat      0.58    0.42    1
```

In this survey, 58% of voters prefer Hillary over Bernie. While there is a larger than expected Bernie preference (since Hillary Clinton won the Democratic nomination, one would expect to see a higher ratio for Clinton than Bernie), there isn't a substantial tilt in one direction or the other. One possible explanation is that Bernie is more popular among the Independents and Republicans. So we are going to assume that this sample does not exhibit a meaningful selection bias (sample set is random).

Party affiliations

```
c.table <- array(data = c(sum(dt$party == "D"), sum(dt$party ==
"D)/length(dt$party), sum(dt$party == "I"), sum(dt$party ==
"I)/length(dt$party), sum(dt$party == "R"), sum(dt$party ==
"R)/length(dt$party), sum((dt$party == "D" | (dt$party ==
```

```

"I") | (dt$party == "R")), sum((dt$party == "D") | (dt$party ==
"I") | (dt$party == "R"))/length(dt$party)), dim = c(2, 4),
dimnames = list(Count = c("Voter Count", "pi.hat"), Party_affiliation = c("Democrat",
"Independent", "Republican", "Total")))

round(c.table, 2)

```

```

##           Party_affiliation
## Count      Democrat Independent Republican Total
## Voter Count  455.00      458.00      278.00  1191
## pi.hat       0.38       0.38       0.23    1

```

We observe that 38% of the voters in the dataset are affiliated with the Democratic Party, whereas only 23% are affiliated with the Republican Party. The 0.6 Democratic to Republican ratio is noteworthy. The data appears to be skewed towards Democratic voters (perhaps a specific region of the country). Our model may not be applicable to the entire country. A further analysis of how the dataset was sampled from the entire population would be very useful.

Race

```

c.table <- array(data = c(sum(dt$race_white == 1), sum(dt$race_white ==
1)/length(dt$race_white), sum(dt$race_white == 0), sum(dt$race_white ==
0)/length(dt$race_white), sum((dt$race_white == 1) | (dt$race_white ==
0)), sum((dt$race_white == 1) | (dt$race_white == 0))/length(dt$race_white)),
dim = c(2, 3), dimnames = list(Count = c("Voter Count", "pi.hat"),
Race = c("White", "Non White", "Total")))

round(c.table, 2)

```

```

##           Race
## Count      White Non White Total
## Voter Count 869.00    322.00  1191
## pi.hat      0.73     0.27    1

```

The 73% white / non-white ratio is inline with the overall US population (according to the 2016 US Census, whites made up 72.4% of the population). The dataset does not appear to have a selection bias with respect to the voter race

Gender

```

c.table <- array(data = c(sum(dt$gender == 1), sum(dt$gender ==
1)/length(dt$gender), sum(dt$gender == 2), sum(dt$gender ==
2)/length(dt$gender), sum((dt$gender == 1) | (dt$gender ==
2)), sum((dt$gender == 1) | (dt$gender == 2))/length(dt$gender)),
dim = c(2, 3), dimnames = list(Count = c("Voter Count", "pi.hat"),
Gender = c("Male", "Female", "Total")))

round(c.table, 2)

```

```

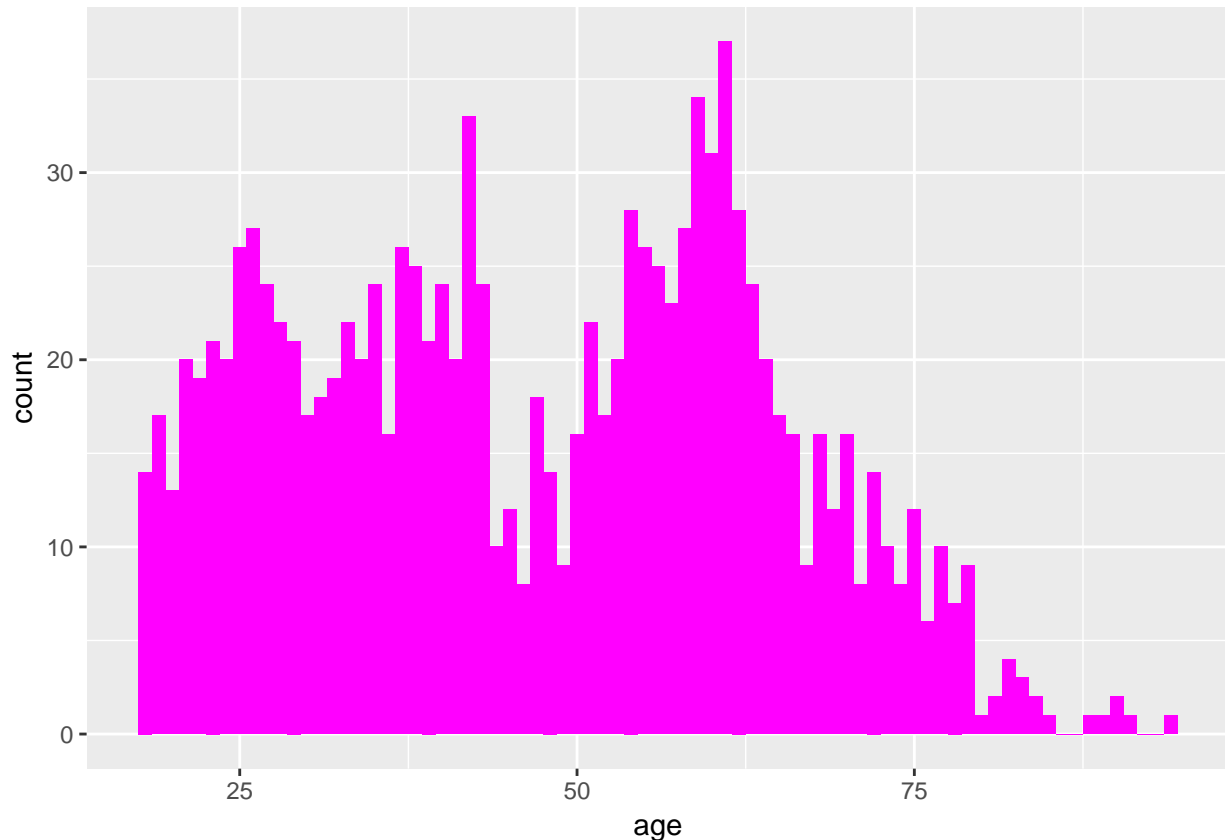
##           Gender
## Count      Male Female Total
## Voter Count 569.00 622.00  1191
## pi.hat     0.48  0.52    1

```

The female / male ration is 1.10 on the dataset and the population ratio (according to Wikipedia) is 1.05. So the sample data doesn't appear to have a meaningful skew.

Age

```
ggplot(dt, aes(age)) + geom_histogram(binwidth = 1, fill = "magenta")
```



```
c.table <- array(data = c(range(dt$age)), dim = c(1, 2), dimnames = list(c("Voter Age"),  
  c("Youngest", "Oldest")))  
round(c.table, 2)
```

```
##           Youngest Oldest  
## Voter Age       18     94
```

We observe that age data for the voters in the survey is within the norms (i.e, there is no one below 18) as expected, as the age variable moves beyond 70, there is a rapid decline. However the data appears to be tri-model. There is not an obvious explanation for that either. This may point to a selection bias (i.e, the sample isn't trully random). We'll note the shortcoming in our final analysis as a caution.

Multivariate analysis

At this section, we are going to focus on the relationship between the outcome (**sanders__preference**) and possible explanatory variables including:

- preference vs. party

- preference vs. race
- preference vs. gender
- age vs. party
- age vs. race
- age vs. gender

preference vs. party

```
c.tabs <- xtabs(~sanders_preference + party, data = dt)
c.tabs <- rbind(c.tabs, colSums(c.tabs))
c.tabs <- cbind(c.tabs, rowSums(c.tabs))

c.table <- array(data = as.array(c.tabs), dim = c(3, 4), dimnames = list(Preference = c("Hillary",
  "Bernie", "Total"), Party = c("Democratic", "Independent",
  "Republican", "Total")))

round(c.table, 2)
```

```
##           Party
## Preference Democratic Independent Republican Total
##   Hillary          249          156          100   505
##   Bernie           206          302          178   686
##   Total            455          458          278  1191
```

```
round(c.table/dim(dt)[1], 2)
```

```
##           Party
## Preference Democratic Independent Republican Total
##   Hillary          0.21          0.13          0.08  0.42
##   Bernie           0.17          0.25          0.15  0.58
##   Total            0.38          0.38          0.23  1.00
```

We have further proof that Bernie is popular with the wrong group (i.e, Independents and Republicans), which is a point we touched on the univariate analysis section for the preference variable. While he enjoys roughly 2x the popularity of Hillary Clinton among the Independents and Republicans, he is less popular among the Democrats. **Our intuition is to include Independents to the target audience** as their lack of enthusiasm for the Democratic party may be offset by their support for Bernie.

preference vs. race__white

```
c.tabs <- xtabs(~sanders_preference + race_white, data = dt)
c.tabs <- rbind(c.tabs, colSums(c.tabs))
c.tabs <- cbind(c.tabs, rowSums(c.tabs))

c.table <- array(data = as.array(c.tabs), dim = c(3, 3), dimnames = list(Preference = c("Hillary",
  "Bernie", "Total"), Race = c("Whites", "Non Whites", "Total")))

round(c.table, 2)
```

```
##           Race
## Preference Whites Non Whites Total
##   Hillary      190        315   505
##   Bernie       132        554   686
```

```
##      Total      322      869 1191
```

```
round(c.table/dim(dt)[1], 2)
```

```
##           Race
## Preference Whites Non Whites Total
##   Hillary   0.16     0.26 0.42
##   Bernie    0.11     0.47 0.58
##   Total     0.27     0.73 1.00
```

Bernie enjoys 4 times more support from Non White voters than White voters. This is definitely a strong signal for our model, so we will explore adding race as a dependent variable to our model.

preference vs. gender

```
t <- table(dt[, .(Gender = ifelse(gender == 1, "Female", "Male"),
  Preference = ifelse(sanders_preference == 1, "Sanders", "Hillary"))])
t
```

```
##           Preference
## Gender   Hillary Sanders
##   Female     242     327
##   Male      263     359
```

```
round(t/sum(t), 2)
```

```
##           Preference
## Gender   Hillary Sanders
##   Female     0.20     0.27
##   Male      0.22     0.30
```

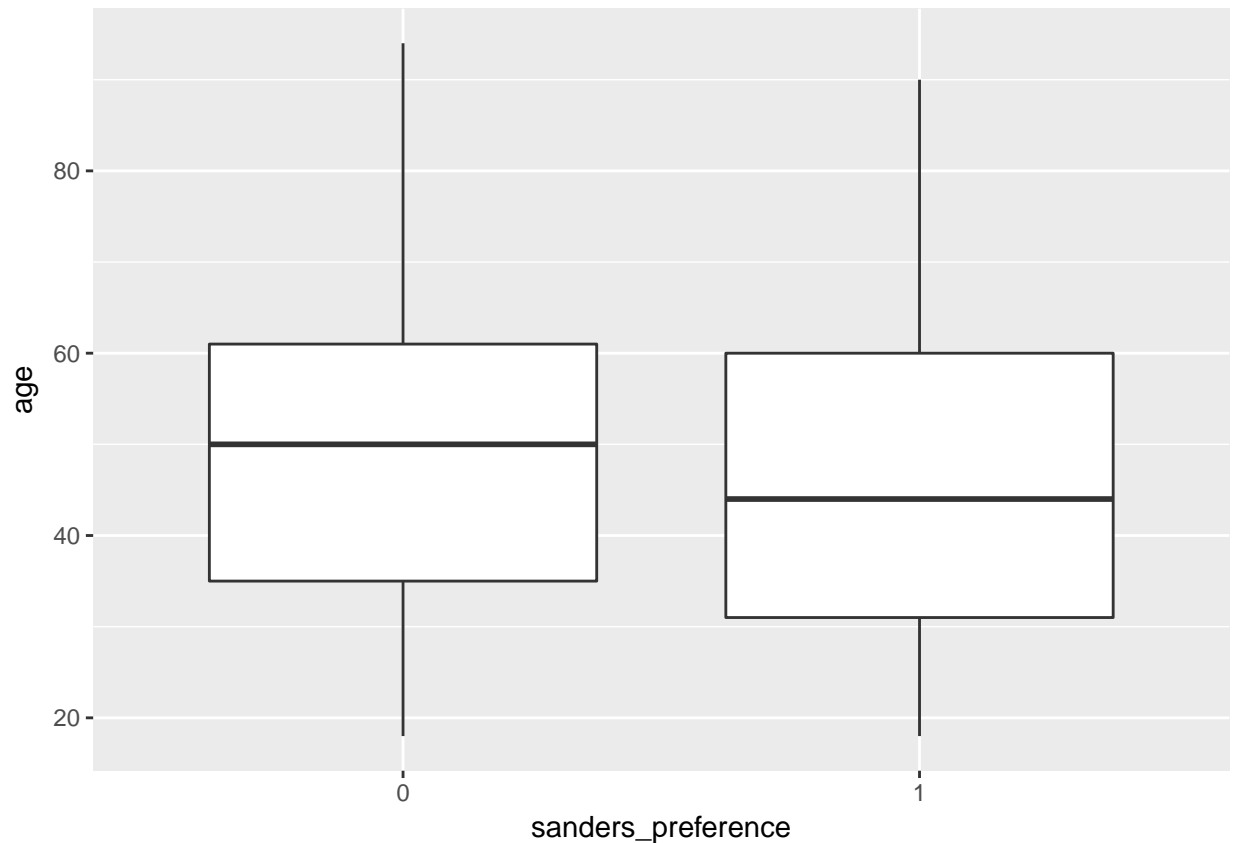
Both males and females are more likely to prefer Sanders. $P(\text{prefers_sanders}|\text{female}) = \frac{0.27}{0.47} = 0.58$ and $P(\text{prefers_sanders}|\text{male}) = \frac{0.30}{0.52} = 0.58$. So there doesn't seem to be a strong direct relationship between preference for Bernie Sanders and gender. We'll revisit the gender variable as part of our interaction analysis.

Preference vs. age

```
dt[, .(age, Preference = ifelse(sanders_preference == 1, "Sanders",
  "Hillary"))], .(`Mean Age` = mean(age)), by = Preference] # Average age of those who prefer Sanders
```

```
##      Preference Mean Age
## 1:      Sanders 46.11953
## 2:      Hillary 48.29307
```

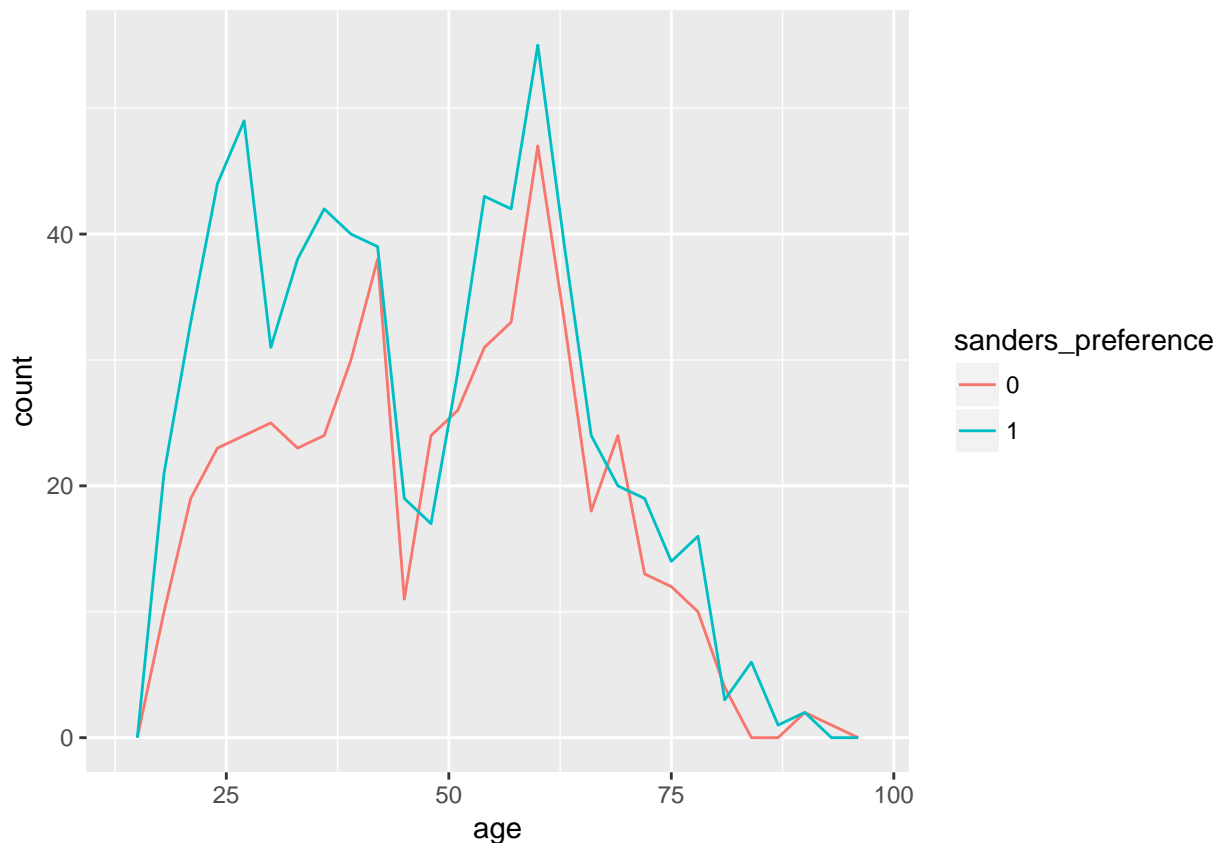
```
ggplot(dt, aes(x = sanders_preference, y = age, group = sanders_preference)) +
  geom_boxplot()
```

Looking at the boxplot, it is clear that mean age for people who prefer Bernie Sanders is younger, but not that much. Also the 1st quartile is lower, but also not that much. So the boxplot doesn't provide a strong visual evidence for either adding or excluding the age as an explanatory variable.

We are going to plot a frequency polygon for age conditioned on sander_preference to have a closer look.

```
ggplot(dt, aes(age, color = sander_preference)) + geom_freqpoly(binwidth = 3)
```



The frequency polygon indicates that there is a clear difference in favor of Bernie Sanders for voters below the age of 40. We'll use 40 years as a cutoff point for old vs young age groups in our model.

```
dt$younger = as.factor(dt$age < 40)
t <- table(dt[, .(Age = ifelse(younger == 1, "Younger", "Older"),
  Preference = ifelse(sanders_preference == 1, "Sanders", "Hillary"))])
t
```

```
##      Preference
## Age   Hillary Sanders
## Older    505    686
```

```
round(t/sum(t), 2)
```

```
##      Preference
## Age   Hillary Sanders
## Older    0.42    0.58
```

We confirm that $P(\text{preference_bernie}|\text{Younger}) = \frac{0.24}{0.38} \approx 0.63$, whereas $P(\text{preference_bernie}|\text{Older}) = \frac{0.34}{0.62} \approx 0.55$. So Bernie is more popular among the voters under the age of 40 than those 40 years or older.

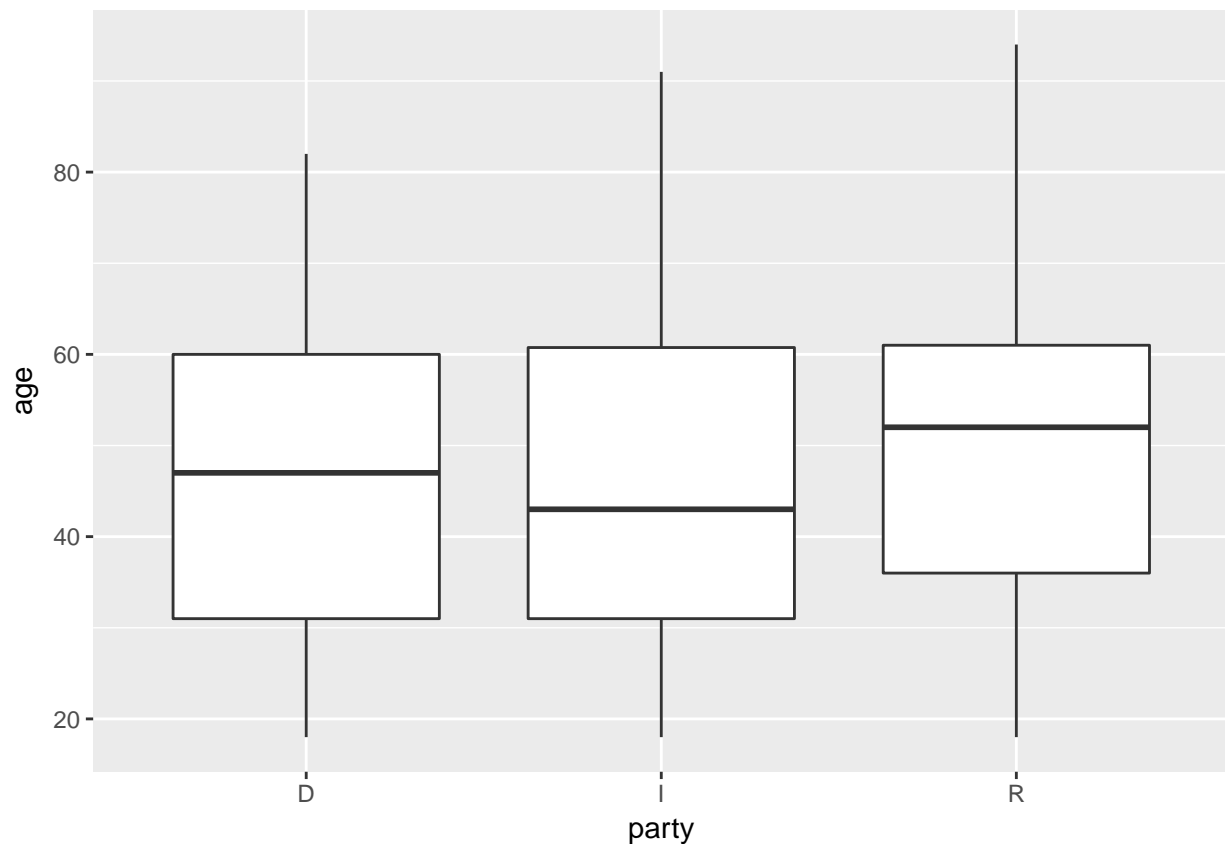
age vs. party

```
dt[, .(`Mean Age` = mean(age)), by = party] # 48:47:52 Age of Democrat:Independent:Republican
```

```
##      party Mean Age
## 1:      D 46.34945
```

```
## 2:      I 45.90393
## 3:      R 50.04676
```

```
ggplot(dt, aes(x = party, y = age, group = party)) + geom_boxplot()
```



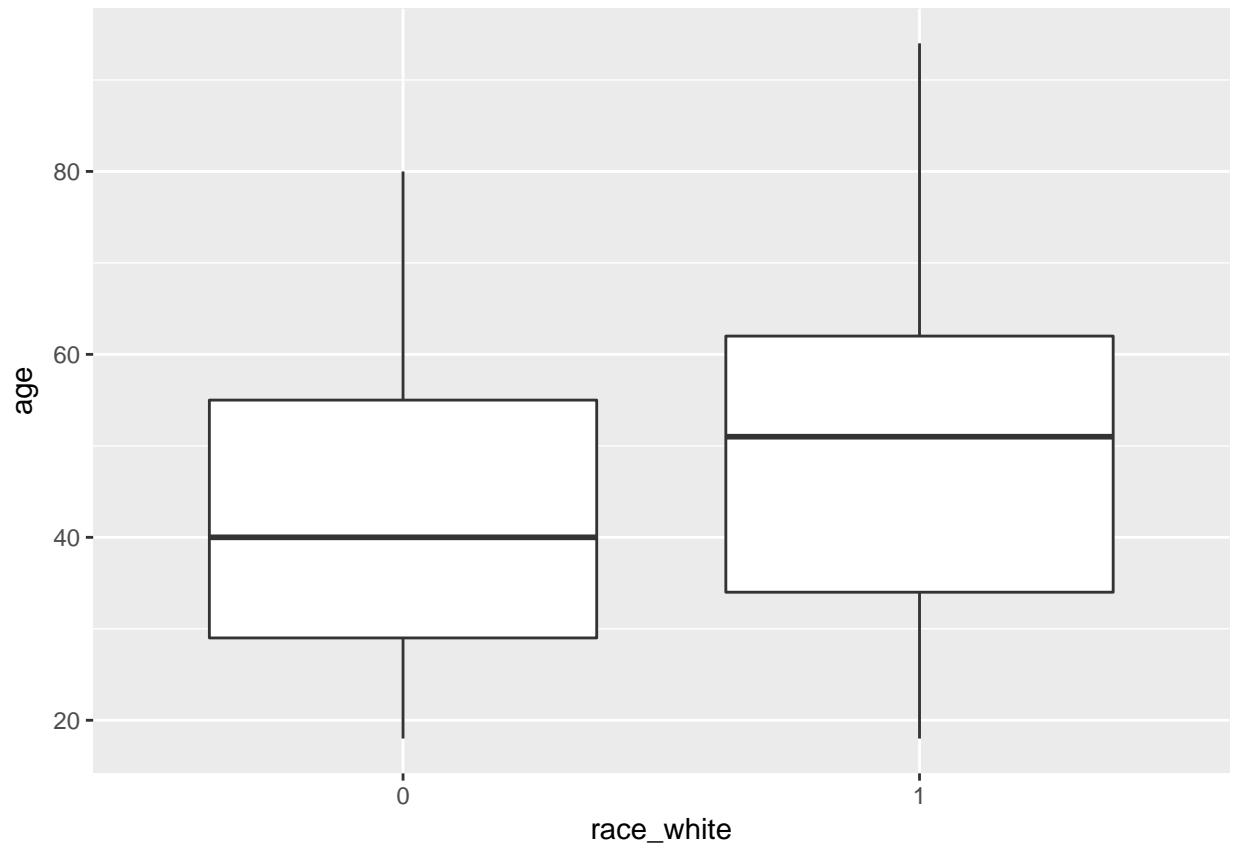
We observe that Independents and Democrats are relatively younger than Republicans, which is another indicator that adding age in the model is a good idea.

age vs. race_white

```
dt[, .(age, race_white = ifelse(race_white == 1, "White", "Non-white"))[,  
  .(`Mean Age` = mean(age)), by = race_white] # 50 vs 44 for race = white
```

```
##      race_white Mean Age  
## 1:      White 48.78941  
## 2:   Non-white 42.32298
```

```
ggplot(dt, aes(x = race_white, y = age, group = race_white)) +  
  geom_boxplot()
```



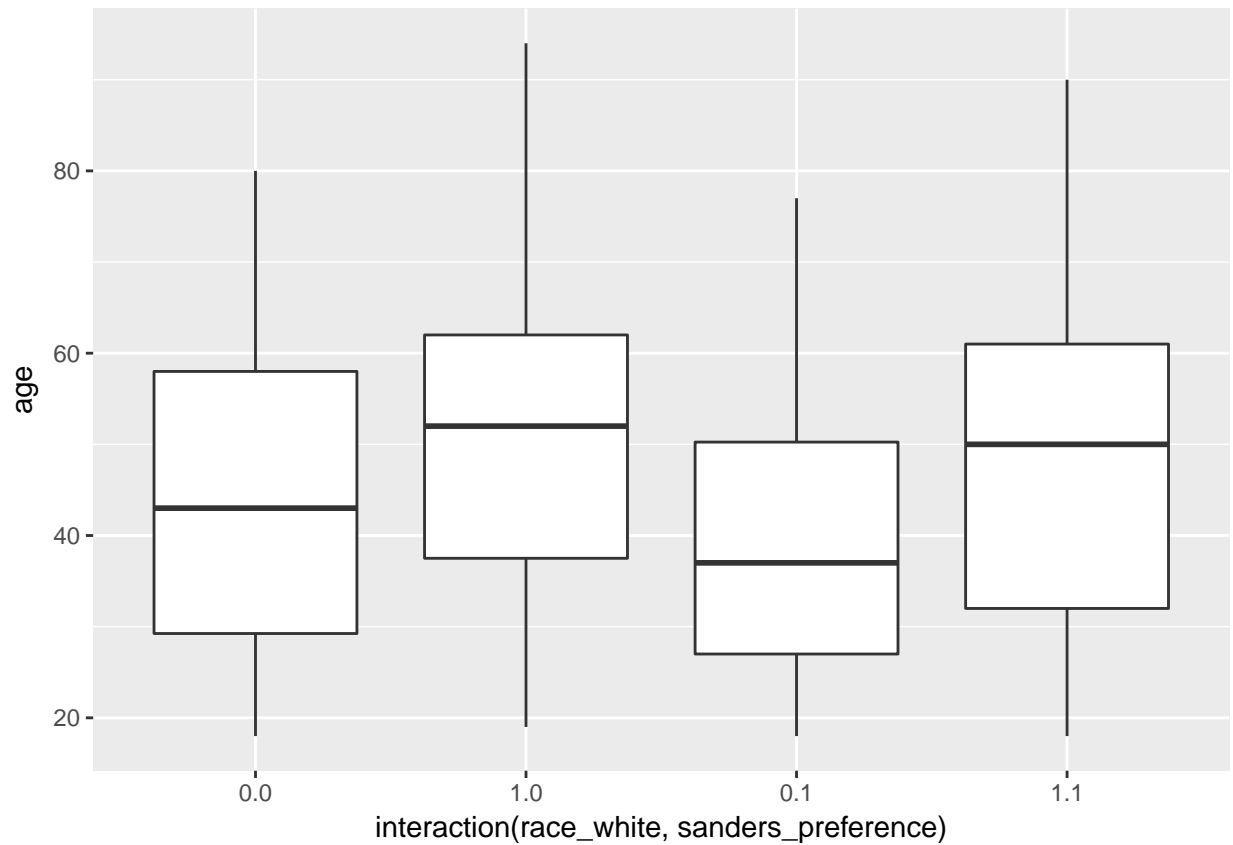
Non-white people prefer Bernie to Hillary Clinton ($\frac{\hat{\pi}_{Bernie|non_white}}{\hat{\pi}_{Hillary|non_white}} \approx 4$). The boxplot provides evidence that non-white voters are younger, so that is further evidence that age will be a strong candidate as an explanatory variable in our model.

Interactions

preference ~ age & gender

preference ~ age & race_white

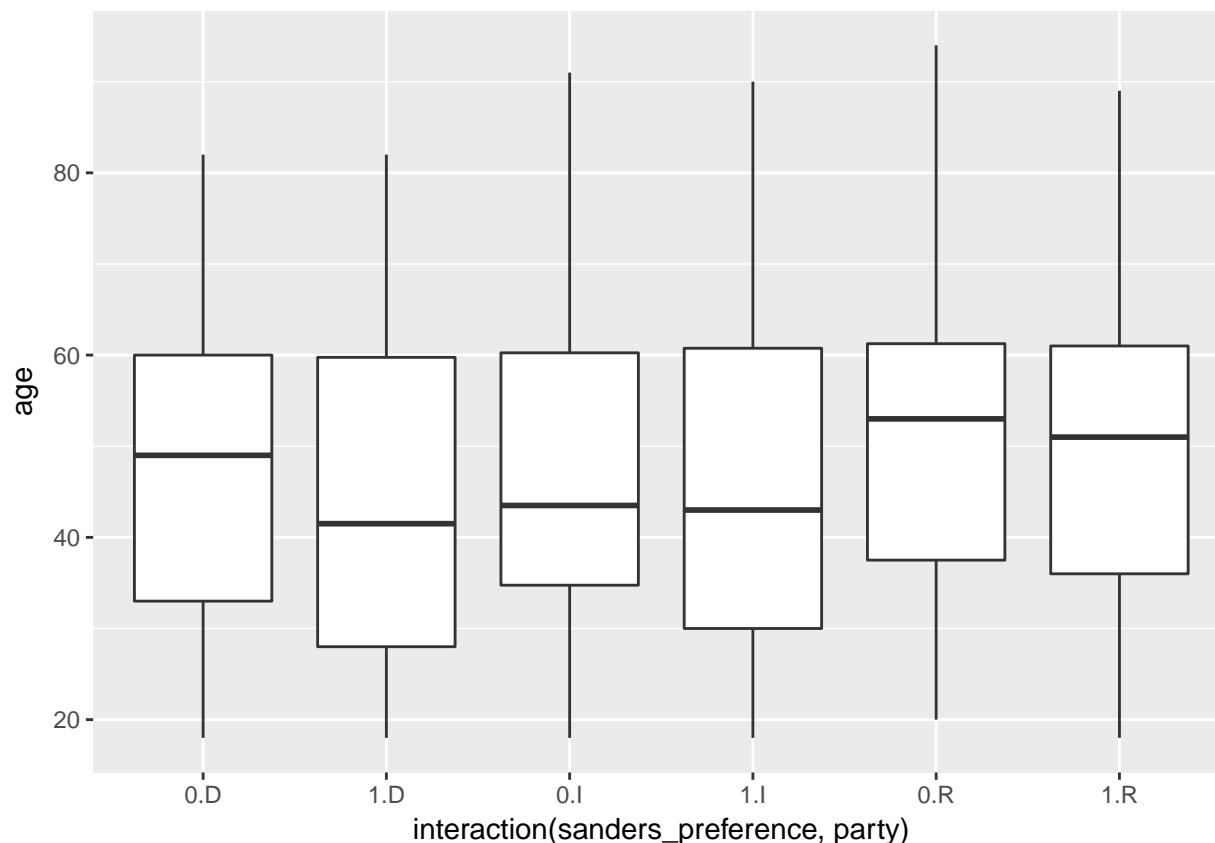
```
ggplot(dt, aes(x = interaction(race_white, sanders_preference),  
  y = age)) + geom_boxplot()
```



Non-white voters who prefer Bernie are the youngest group (based on their 1st, 3rd quartile and mean age). We'll explore the interaction between age and race during our model exploration.

preference ~ age & party

```
ggplot(dt, aes(x = interaction(sanders_preference, party), y = age)) +  
  geom_boxplot()
```



For Republicans and Independents, the mean and 3rd quartiles for age don't change that much based on sanders_preference. But 1st quartile for Independents as well as the mean age for Democrats changes noticeably based on sanders_preference. **We are going to explore the interaction between the age and party variables.**

Models

Based on the exploratory data analysis, our model exploration strategy is as follows:

- Investigate whether gender is a significant explanatory variable or not.
- Investigate that **younger** (age < 40) is a better explanatory variable than **age**.
- Investigate that both **party** and **race_white** are meaningful explanatory variables.
- Investigate additional interactions:
 - Interaction between age and race
 - Interaction between age and party

Is gender important?

Null Hypothesis

In our null hypothesis we are going to assume that $\beta_{gender} = 0$

```
model1.H0 = glm(sanders_preference ~ age, dt, family = binomial(link = "logit"))
summary(model1.H0)
```

```
##
## Call:
## glm(formula = sanders_preference ~ age, family = binomial(link = "logit"),
##      data = dt)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.4074  -1.2842   0.9841   1.0620   1.1840
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.661384   0.173958   3.802 0.000144 ***
## age         -0.007522   0.003457  -2.176 0.029566 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1623.5  on 1190  degrees of freedom
## Residual deviance: 1618.7  on 1189  degrees of freedom
## AIC: 1622.7
##
## Number of Fisher Scoring iterations: 4
```

Alternative hypothesis

In the alternative hypothesis, we are going to assume that $\beta_{gender} \neq 0$

```
model1.Ha = glm(sanders_preference ~ age + gender, dt, family = binomial(link = "logit"))
summary(model1.Ha)
```

```
##
## Call:
## glm(formula = sanders_preference ~ age + gender, family = binomial(link = "logit"),
##      data = dt)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.4111  -1.2869   0.9838   1.0625   1.1816
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.653207   0.183064   3.568 0.000359 ***
## age         -0.007535   0.003458  -2.179 0.029346 *
## gender2      0.016857   0.117657   0.143 0.886079
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1623.5  on 1190  degrees of freedom
## Residual deviance: 1618.7  on 1188  degrees of freedom
## AIC: 1624.7
##
```

```
## Number of Fisher Scoring iterations: 4
```

Comparing H0 to Ha

```
anova(model1.H0, model1.Ha, test = "Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: sanders_preference ~ age
## Model 2: sanders_preference ~ age + gender
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1      1189      1618.7
## 2      1188      1618.7  1  0.020525   0.8861
```

There is no evidence supporting the alternative hypothesis that **gender** is an explanatory variable given that **age** is in the model. We fail to reject the null hypothesis that $\beta_{gender} \neq 0$.

Age (numeric) or younger (binary, <40) variable?

We are going to compare the 2 models and decide which variable provides a better model for our *specific task*

```
model2.age = glm(sanders_preference ~ age, dt, family = binomial(link = "logit"))
model2.younger = glm(sanders_preference ~ younger, dt, family = binomial(link = "logit"))
stargazer(model2.age, model2.younger, type = "text")
```

```
##
## =====
##                               Dependent variable:
##                               -----
##                               sanders_preference
##                               (1)           (2)
## -----
## age                        -0.008**
##                               (0.003)
##
## younger                                0.349***
##                               (0.122)
##
## Constant                    0.661***      0.176**
##                               (0.174)      (0.074)
##
## -----
## Observations                1,191        1,191
## Log Likelihood              -809.356     -807.623
## Akaike Inf. Crit.          1,622.712     1,619.246
## =====
## Note:                *p<0.1; **p<0.05; ***p<0.01
```

While there is not a major difference in the AIC criterion between the the models, the 2nd. model providers a more practical interpretation (i.e., providing a clear indication to focus on voters younger than 40).

Adding race and party affiliation

We are now going to include race and party. To recap:

- Null hypothesis H_0 :
 - $\text{logit}(\pi_{\text{preference_sanders}}) = \beta_0 + \beta_1 \text{younger}$
- Alternative hypothesis H_a :
 - $\text{logit}(\pi_{\text{preference_sanders}}) = \beta_0 + \beta_1 \text{younger} + \beta_2 \text{partyI} + \beta_3 \text{partyR} + \beta_4 \text{race_white}$

```
model3.H0 = glm(sanders_preference ~ younger, dt, family = binomial(link = "logit"))
model3.Ha = glm(sanders_preference ~ younger + party + race_white,
  dt, family = binomial(link = "logit"))
summary(model3.Ha)
```

```
##
## Call:
## glm(formula = sanders_preference ~ younger + party + race_white,
##      family = binomial(link = "logit"), data = dt)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6899  -1.1605   0.7406   0.9550   1.5725
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.8932     0.1446  -6.179 6.44e-10 ***
## youngerTRUE    0.4746     0.1287   3.687 0.000227 ***
## partyI         0.7190     0.1405   5.117 3.10e-07 ***
## partyR         0.5886     0.1629   3.612 0.000303 ***
## race_white1    0.8533     0.1409   6.054 1.41e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1623.5  on 1190  degrees of freedom
## Residual deviance: 1531.1  on 1186  degrees of freedom
## AIC: 1541.1
##
## Number of Fisher Scoring iterations: 4
```

```
# Comparing H0 and Ha
```

```
anova(model3.H0, model3.Ha, test = "Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: sanders_preference ~ younger
## Model 2: sanders_preference ~ younger + party + race_white
##   Resid. Df Resid. Dev Df Deviance  Pr(>Chi)
## 1      1189      1615.2
## 2      1186      1531.1  3    84.175 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The p-values for the **party(I|R)** and **race_white1** explanatory variables are very low. In addition, the

p-value for H_a (model3.Ha) is also very low. We conclude that there is strong empirical evidence for including race and party affiliation in our model.

Exploring models with interactions

We are going to construct 3 additional models to explore the interactions:

- Between race and age
- Between party and age
- Between party and race

```
model4.H0 = glm(sanders_preference ~ younger + party + race_white,
  dt, family = binomial(link = "logit")) # also = H3.HA
model4.Ha1 = glm(sanders_preference ~ younger + party + race_white +
  younger:race_white, dt, family = binomial(link = "logit"))
model4.Ha2 = glm(sanders_preference ~ younger + party + race_white +
  younger:party, dt, family = binomial(link = "logit"))
model4.Ha3 = glm(sanders_preference ~ younger + party + race_white +
  party:race_white, dt, family = binomial(link = "logit"))

stargazer(model4.H0, model4.Ha1, model4.Ha2, model4.Ha3, type = "text",
  report = ("vc*p"))
```

```
##
## =====
##                               Dependent variable:
##                               -----
##                               sanders_preference
##                               (1)         (2)         (3)         (4)
## -----
## younger                    0.475***    0.592**    0.630***    0.474***
##                             p = 0.0003    p = 0.011    p = 0.002    p = 0.0003
##
## partyI                     0.719***    0.717***    0.810***    0.753***
##                             p = 0.00000    p = 0.00000    p = 0.00001    p = 0.004
##
## partyR                     0.589***    0.585***    0.700***    0.565
##                             p = 0.0004    p = 0.0004    p = 0.0005    p = 0.128
##
## race_white1                0.853***    0.929***    0.853***    0.867***
##                             p = 0.000    p = 0.00000    p = 0.000    p = 0.00002
##
## youngerTRUE:race_white1          -0.170
##                                p = 0.543
##
## youngerTRUE:partyI                    -0.232
##                                p = 0.417
##
## youngerTRUE:partyR                    -0.323
##                                p = 0.341
##
## partyI:race_white1                                -0.048
##                                                    p = 0.876
##
```

```
## partyR:race_white1                                0.024
##                                                     p = 0.955
##
## Constant          -0.893***    -0.950***    -0.955***    -0.902***
##                   p = 0.000    p = 0.00000    p = 0.000    p = 0.00000
##
## -----
## Observations          1,191          1,191          1,191          1,191
## Log Likelihood        -765.535        -765.350        -764.974        -765.517
## Akaike Inf. Crit.      1,541.071      1,542.700      1,543.948      1,545.034
## =====
## Note:                                     *p<0.1; **p<0.05; ***p<0.01
```

We observe that none of the alternative models have any new explanatory variable with a low p-value. We'll confirm our findings with a pairwise anova table $H_0 vs. H_{a1}$, $H_0 vs. H_{a2}$, and $H_0 vs. H_{a3}$

```
model4.1.anova <- anova(model4.H0, model4.Ha1, test = "Chisq")
model4.2.anova <- anova(model4.H0, model4.Ha2, test = "Chisq")
model4.3.anova <- anova(model4.H0, model4.Ha3, test = "Chisq")

model4.anova.results <- data.frame(c(model4.1.anova$`Pr(>Chi)`[2],
  model4.2.anova$`Pr(>Chi)`[2], model4.3.anova$`Pr(>Chi)`[2]))
colnames(model4.anova.results) <- "Pr(>Chi)"
rownames(model4.anova.results) <- c("model4.Ha1", "model4.Ha2",
  "model4.Ha3")
model4.anova.results
```

```
##           Pr(>Chi)
## model4.Ha1 0.5424986
## model4.Ha2 0.5703348
## model4.Ha3 0.9815497
```

None of the 3 interaction models yields a test statistic satisfactory to reject the null hypothesis. **There is no empirical evidence of any interaction between younger and race, younger and party.**

Model selection conclusion

Based on our model study, we conclude that the most appropriate model for the specific question we want to answer (*whether it is a good idea to focus on younger voters*) given the dataset is:

$\text{logit}(\pi_{\text{preference_sanders}}) = \beta_0 + \beta_1 \text{younger} + \beta_2 \text{partyI} + \beta_3 \text{partyR} + \beta_4 \text{race_white}$

Probability Plots

the base model (FOR TESTING PLOTS)

```
logit.mod.base <- glm(formula = sander_preference ~ age, family = binomial(link = logit),
  data = dt)
```

the full model (FOR TESTING PLOTS)

```
logit.mod.full <- glm(formula = sander_preference ~ age + race_white +
  party + gender + age:party, family = binomial(link = logit),
  data = dt)
```

function for calculating C.I.

```
# Wald confidence interval
wald.ci.pi <- function(newdata, mod.fit.obj, alpha) {
  linear.pred <- predict(object = mod.fit.obj, newdata = newdata,
    type = "link", se = TRUE)
  CI.lin.pred.lower <- linear.pred$fit - qnorm(p = 1 - alpha/2) *
    linear.pred$se # Wald interval
  CI.lin.pred.upper <- linear.pred$fit + qnorm(p = 1 - alpha/2) *
    linear.pred$se # Wald interval
  CI.pi.lower <- exp(CI.lin.pred.lower)/(1 + exp(CI.lin.pred.lower))
  CI.pi.upper <- exp(CI.lin.pred.upper)/(1 + exp(CI.lin.pred.upper))
  list(lower = CI.pi.lower, upper = CI.pi.upper)
}

# test case for logit.mod.1 at age = 50 wald.ci.pi(newdata =
# data.frame(age = 50), mod.fit.obj = logit.mod.1, alpha =
# 0.05)
```

Bubble plot of the base model for all observations

aggregate the data by age for symbols

In this case, the data include all observations.

```
w <- aggregate(formula = as.numeric(dt$sanders_preference) ~
  1 ~ dt$age, FUN = sum) # sanders supporters at each age
n <- aggregate(formula = as.numeric(dt$sanders_preference) ~
  dt$age, FUN = length) # total voters at each age

names(w)[1] <- "age"
names(w)[2] <- "preference"
names(n)[1] <- "age"
names(n)[2] <- "preference"

w.n <- data.frame(age = w$age, n = n$preference, w = w$preference,
  ratio = round(w$preference/n$preference, 4))
# head(w.n) # ratio = (sanders_supporters/total number of
# voters)
```

The estimated logistic model with the 95% Wald confidence intervals for the base model $sanders_{preference} = \beta_0 + \beta_1 age$, the bubble plot also shows the observed ratio of voters who prefer Sanders at each age, with the plotting size being proportional to the number of observations at that age.

```
par(mfrow = c(1, 1))
# Plot data points
# #####
symbols(x = w$age, y = w$preference/n$preference, circles = sqrt(n$preference),
  inches = 0.12, xlab = "age", ylab = "Estimated probability",
  xlim = c(15, 100), panel.first = grid(col = "gray", lty = "dotted"))

# Plot model fit
curve(expr = predict(object = logit.mod.base, newdata = data.frame(age = x),
  type = "response"), col = "red", add = TRUE, xlim = c(15,
  95))
```

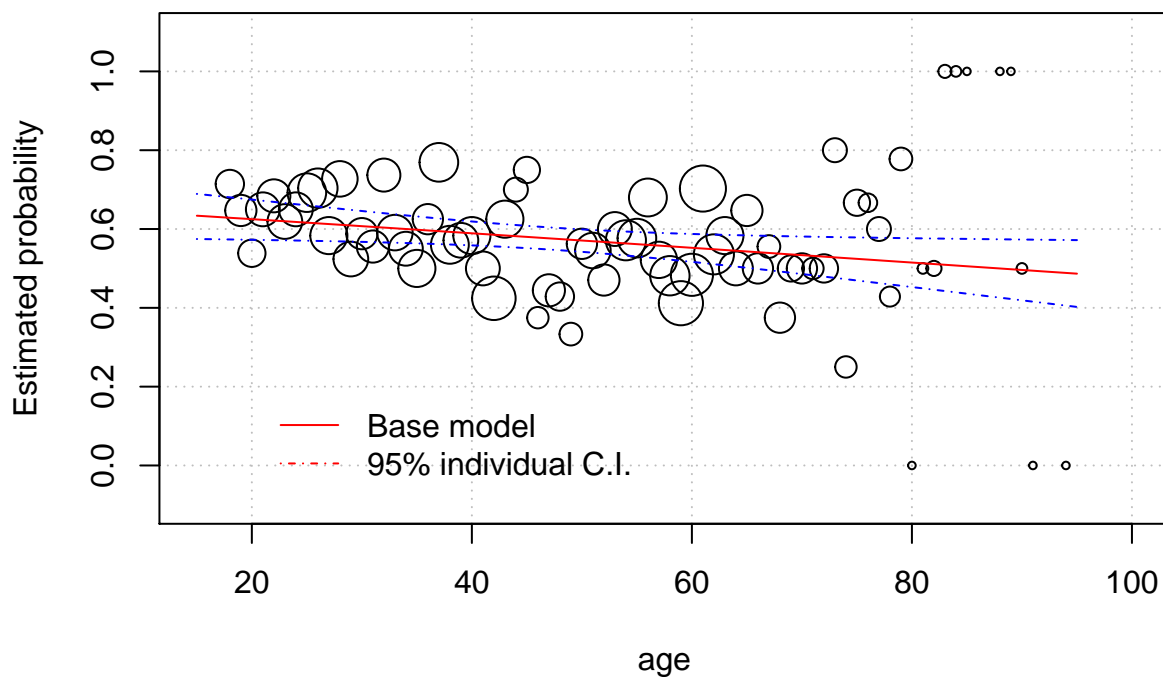
```

# Plot C.I. bands
curve(expr = wald.ci.pi(newdata = data.frame(age = x), mod.fit.obj = logit.mod.base,
  alpha = 0.05)$lower, col = "blue", lty = "dotdash", add = TRUE,
  xlim = c(15, 95))

curve(expr = wald.ci.pi(newdata = data.frame(age = x), mod.fit.obj = logit.mod.base,
  alpha = 0.05)$upper, col = "blue", lty = "dotdash", add = TRUE,
  xlim = c(15, 95))

# Legend
legend(x = 20, y = 0.2, legend = c("Base model", "95% individual C.I."),
  lty = c("solid", "dotdash"), col = c("red", "red"), bty = "n")

```



Comparing different voter groups using the full model

For comparing different groups of voters using the full model, we will plot the observed data points for the subgroup of voters each model represents with the estimated probability and 95% confidence intervals.

Define function `plotsubgroup.pi.ci()` for subgroup probability plot

```

plotsubgroup.pi.ci <- function(newdata, mod.fit.obj, Race, Party,
  Gender) {
  race_white.G = Race

```

```

party.G = Party
gender.G = Gender

# Creating subgroup data frames for the LEFT and RIGHT SIDE
# PLOT
dt.G <- newdata[newdata$race_white == race_white.G & newdata$party ==
  party.G & newdata$gender == gender.G, ]

# Creating aggregated data for PLOT
# #####
w.G <- aggregate(formula = as.numeric(dt.G$sanders_preference) ~
  1 ~ dt.G$age, FUN = sum) # sanders supporters at each age
n.G <- aggregate(formula = as.numeric(dt.G$sanders_preference) ~
  dt.G$age, FUN = length) # total voters at each age

names(w.G)[1] <- "age"
names(w.G)[2] <- "preference"
names(n.G)[1] <- "age"
names(n.G)[2] <- "preference"

w.n.G <- data.frame(age = w.G$age, n = n.G$preference, w = w.G$preference,
  ratio = round(w.G$preference/n.G$preference, 4))

# PLOT
# #####

# Plot data points
symbols(x = w.G$age, y = w.G$preference/n.G$preference, circles = sqrt(n.G$preference),
  inches = 0.1, xlab = "age", ylab = "Estimated probability",
  xlim = c(15, 100), panel.first = grid(col = "gray", lty = "dotted"),
  main = paste("race:", race_white.G, " party:", party.G,
    ", gender:", gender.G))

# Plot model fit
curve(expr = predict(object = logit.mod.full, newdata = data.frame(age = x,
  race_white = race_white.G, party = party.G, gender = gender.G),
  type = "response"), col = "red", add = TRUE, xlim = c(15,
  100), ylim = c(0, 1), xlab = "age", ylab = expression(hat(pi)))

# Plot C.I. bands
curve(expr = wald.ci.pi(newdata = data.frame(age = x, race_white = race_white.G,
  party = party.G, gender = gender.G), mod.fit.obj = logit.mod.full,
  alpha = 0.05)$lower, col = "blue", lty = "dotdash", add = TRUE,
  xlim = c(15, 100))

curve(expr = wald.ci.pi(newdata = data.frame(age = x, race_white = race_white.G,
  party = party.G, gender = gender.G), mod.fit.obj = logit.mod.full,
  alpha = 0.05)$upper, col = "blue", lty = "dotdash", add = TRUE,
  xlim = c(15, 100))

legend(x = 15, y = 0.3, legend = c("model", "95% C.I."),
  lty = c("solid", "dotdash"), col = c("red", "blue"),
  bty = "n")

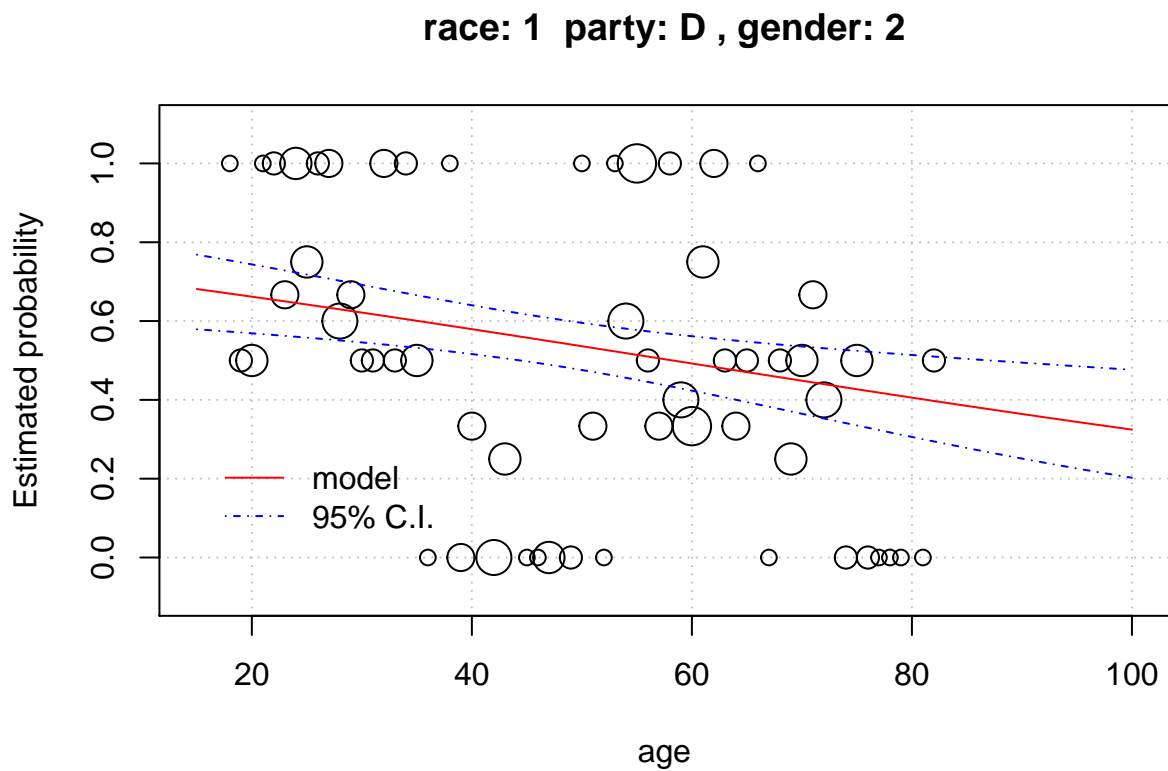
```

```
}
```

Plots comparing subgroups

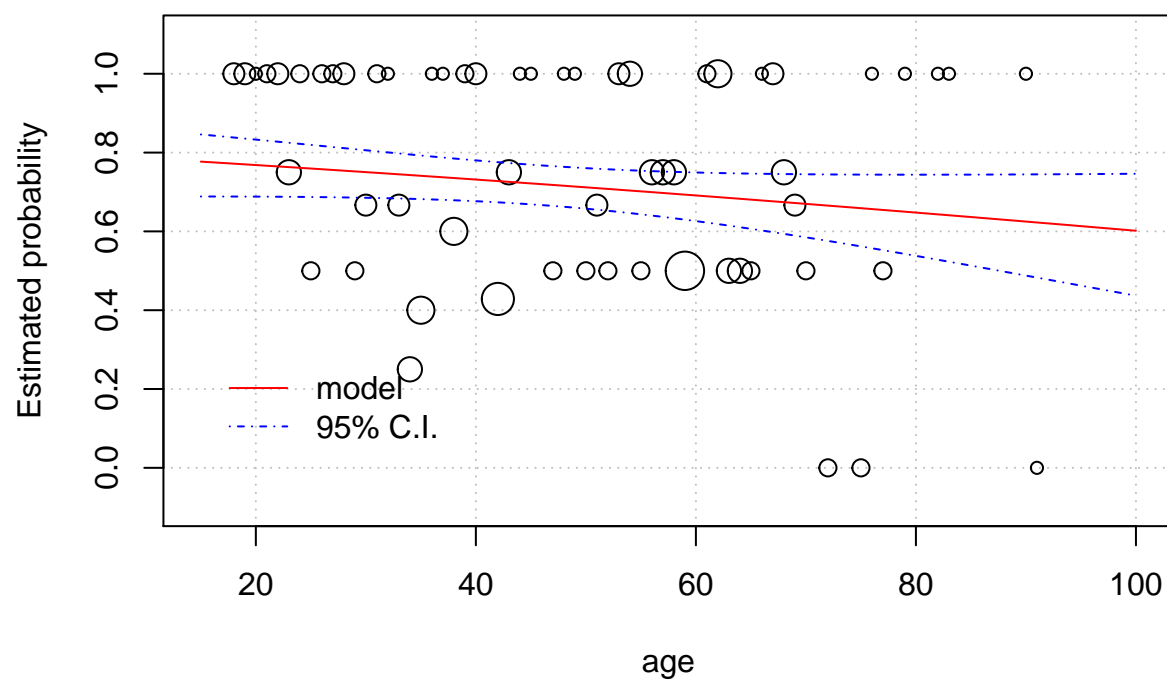
Example: white female, different parties

```
# White Male - Democratic  
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,  
  Race = "1", Party = "D", Gender = "2")
```



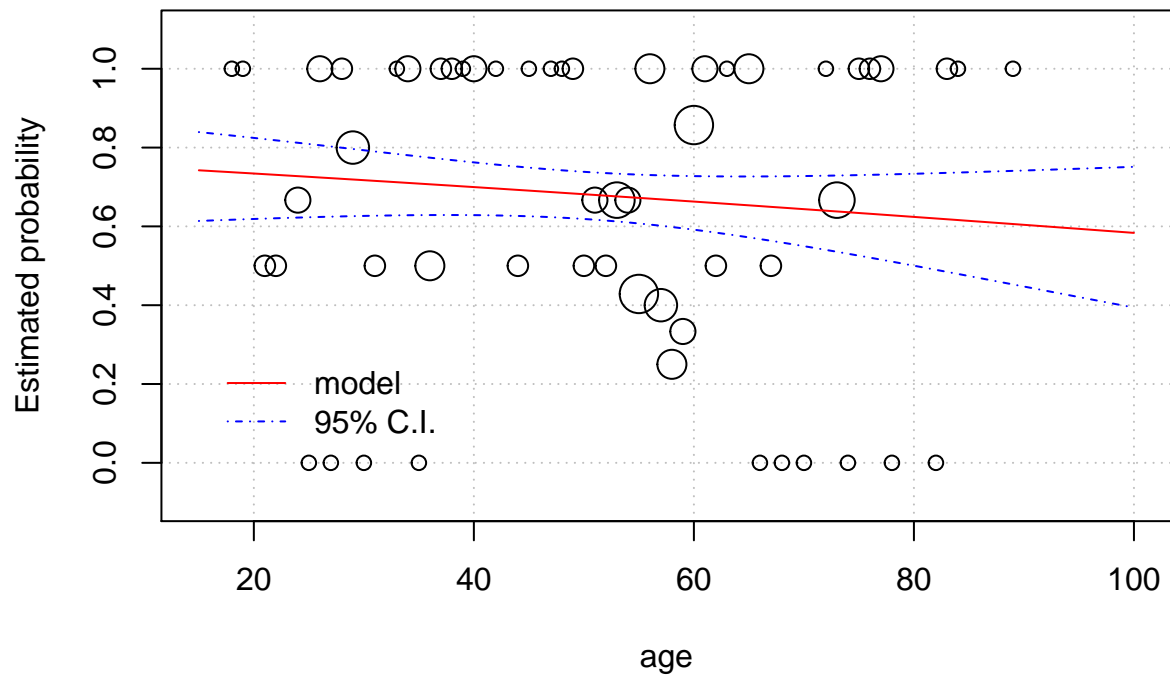
```
# White Male - Independent  
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,  
  Race = "1", Party = "I", Gender = "2")
```

race: 1 party: I , gender: 2



```
# White Male - Republican
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,
  Race = "1", Party = "R", Gender = "2")
```

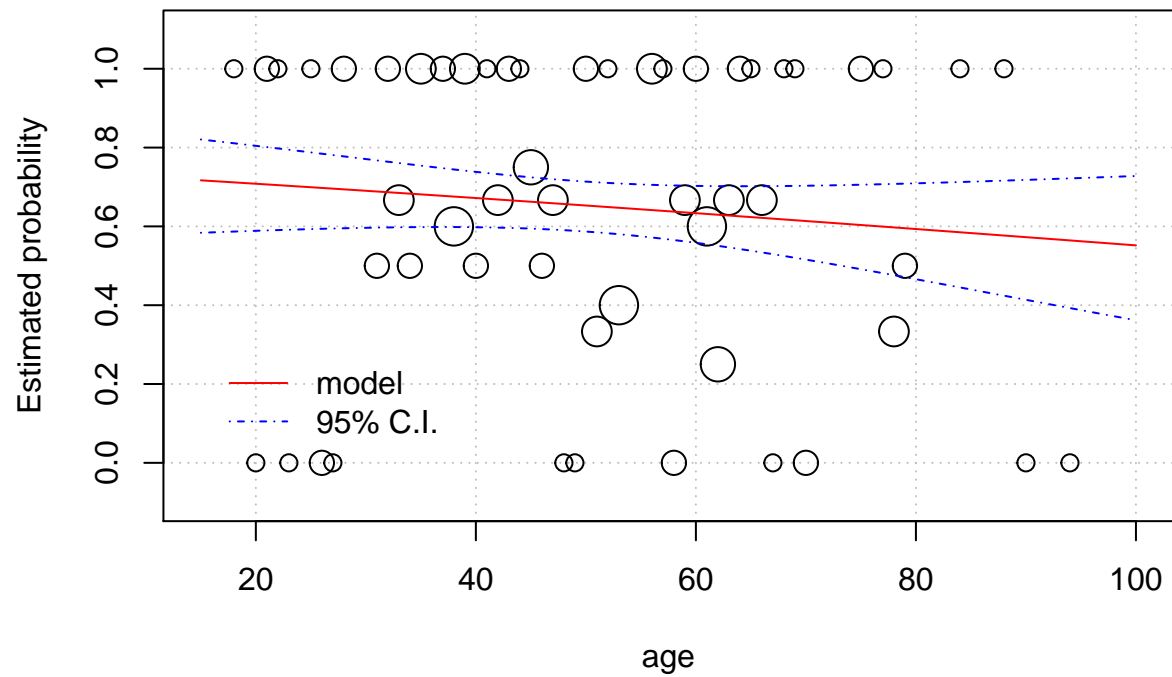

race: 1 party: R , gender: 2



Example: male Republican, white vs. non-white

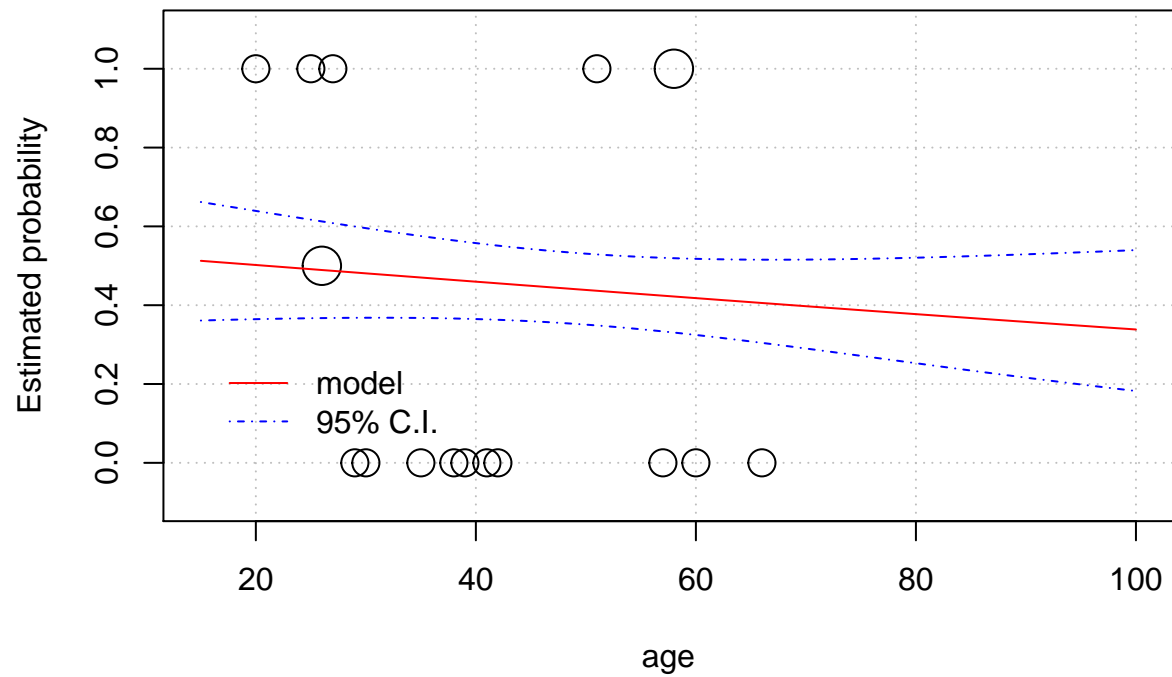
```
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,  
  Race = "1", Party = "R", Gender = "1")
```

race: 1 party: R , gender: 1



```
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,  
  Race = "0", Party = "R", Gender = "1")
```

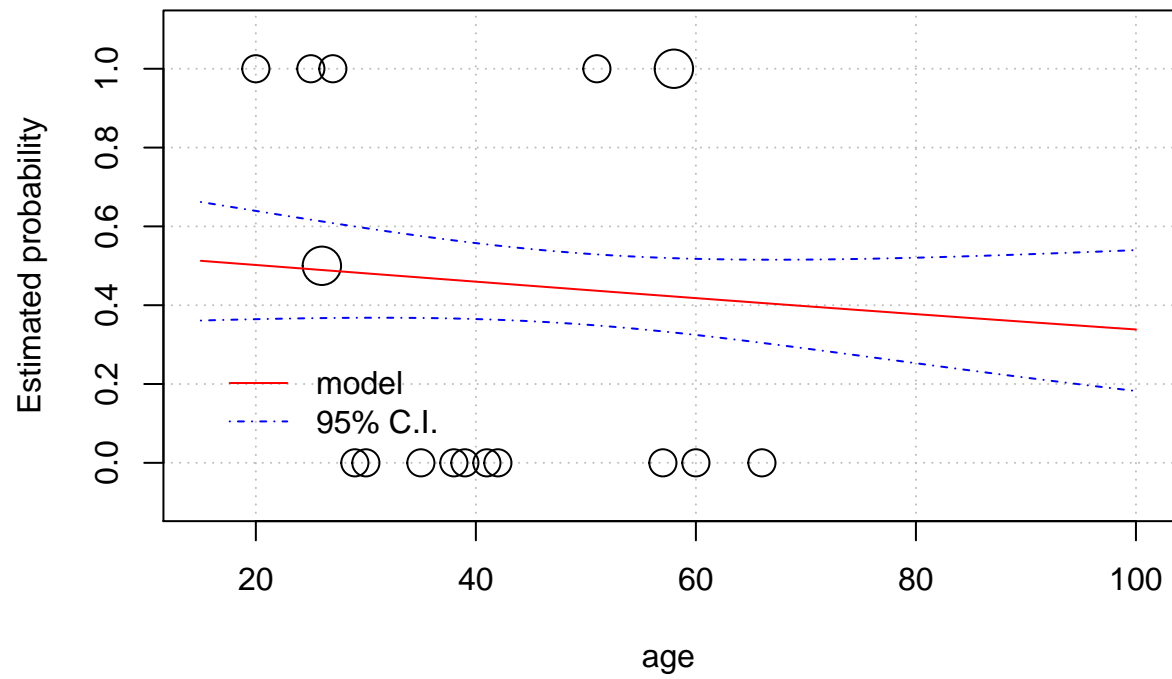
race: 0 party: R , gender: 1



Example: non-white Democratic, male vs. female

```
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,  
  Race = "0", Party = "R", Gender = "1")
```

race: 0 party: R , gender: 1



```
plotsugroup.pi.ci(newdata = dt, mod.fit.obj = logit.mod.full,  
  Race = "0", Party = "R", Gender = "2")
```

race: 0 party: R , gender: 2

