

# Jessica Ghai's Counselling Project

*Team C2*

*10/15/2019*

## Data Cleaning and EDA part

There are 3 files, Numeric.csv, ChoiceText.csv & Need to be excluded IDs.xlsx. The 1st two files are the results from the Survey, one in numeric form and the other lays out the actual choices the respondent made.

```
NumericDataRow <- read_csv("Numeric.csv")
TextDataRow <- read_csv("ChoiceText.csv")
```

Import "invalid" responses IDs, for the purpose of EDA, we will keep these response in the data and analyze them to determine the conditions or restrictions which would be post on our analysis if we take them out.

```
ExcludedIDs = readxl::read_xlsx(path = "Need to be excluded IDs.xlsx",
  sheet = 3, col_names = TRUE, col_types = "text") %>% select(1) %>% pull()
print(paste0("Number of IDs Client would like to exclude: ",
  as.character(length(ExcludedIDs))))
```

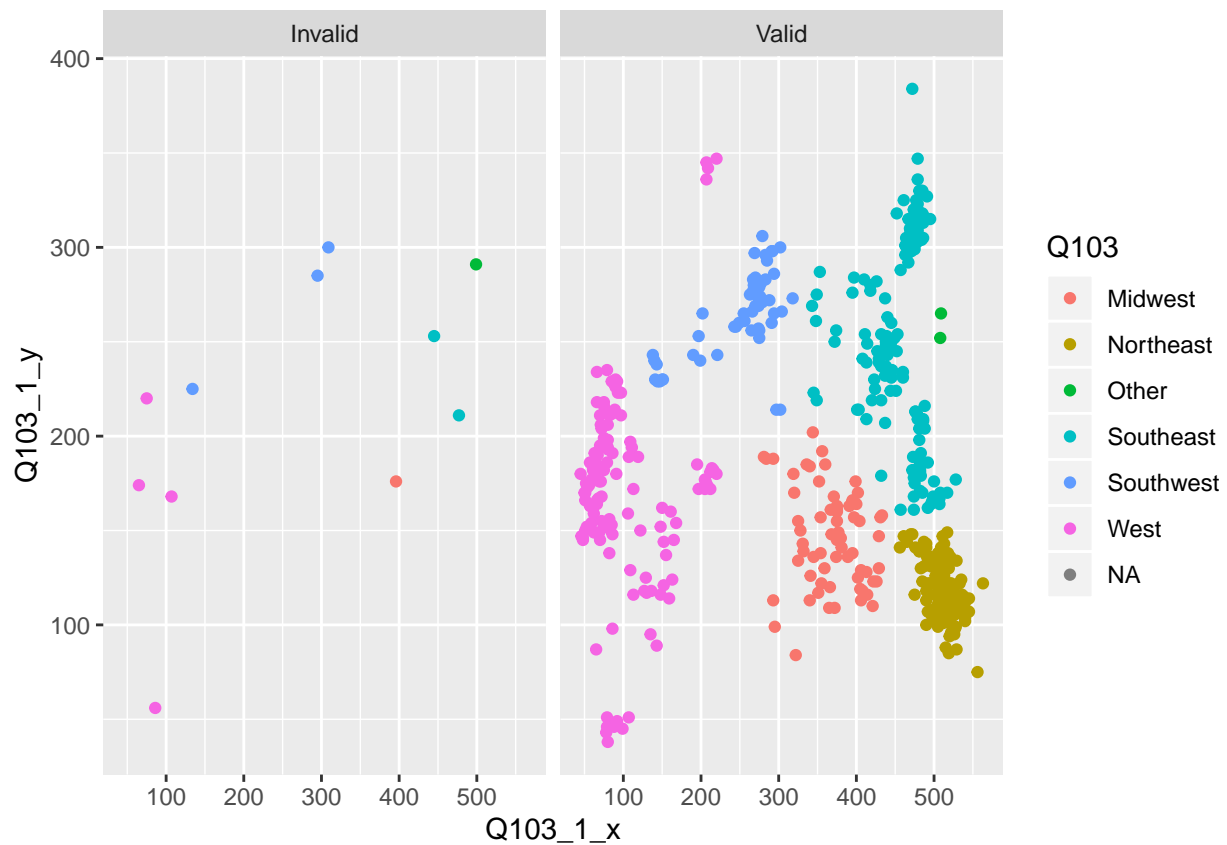
```
## [1] "Number of IDs Client would like to exclude: 87"
```

Import ends here.

Initial visualization of the demographic data

```
geo=select(TextDataRow, Q103_1_x:Q103, ValidResponse) %>%
  slice(-(1:2)) %>%
  mutate(Q103 = as.factor(Q103), Q103_1_x = as.numeric(Q103_1_x), Q103_1_y = as.numeric(Q103_1_y))

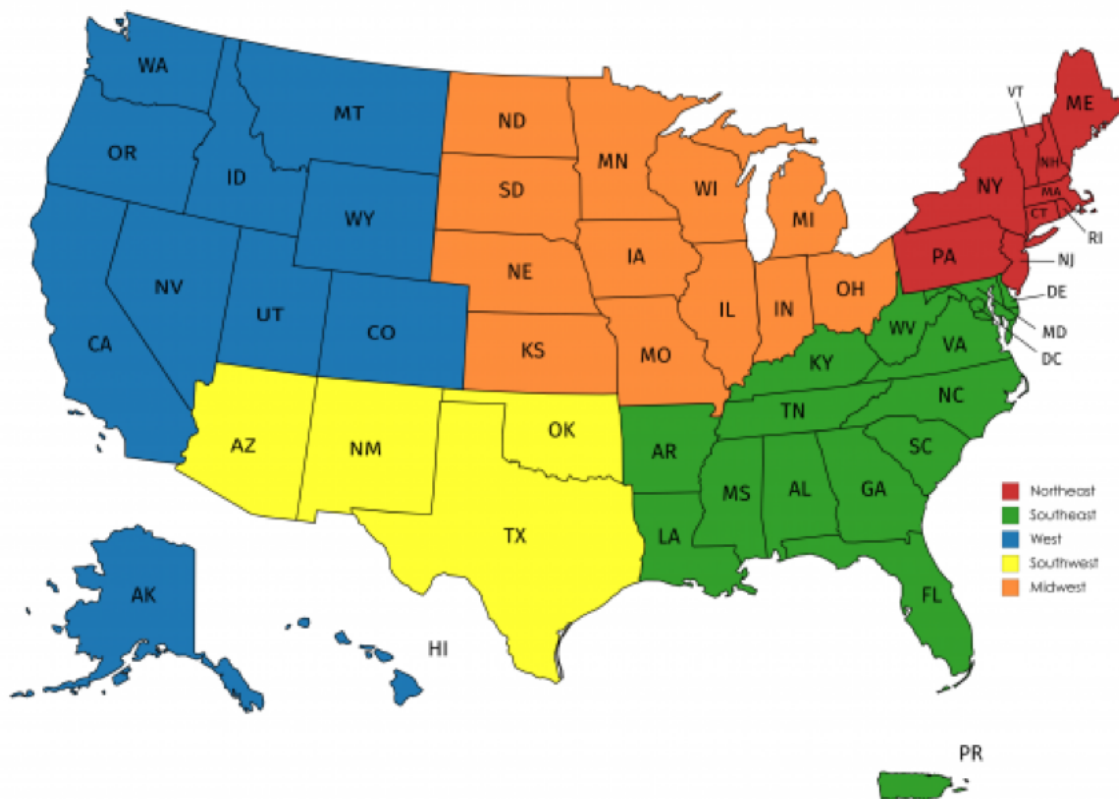
ggplot(geo) + aes(x = Q103_1_x, y = Q103_1_y, color = Q103) +
  geom_point() + facet_wrap(geo$ValidResponse)
```



What's the map actually used in the survey??

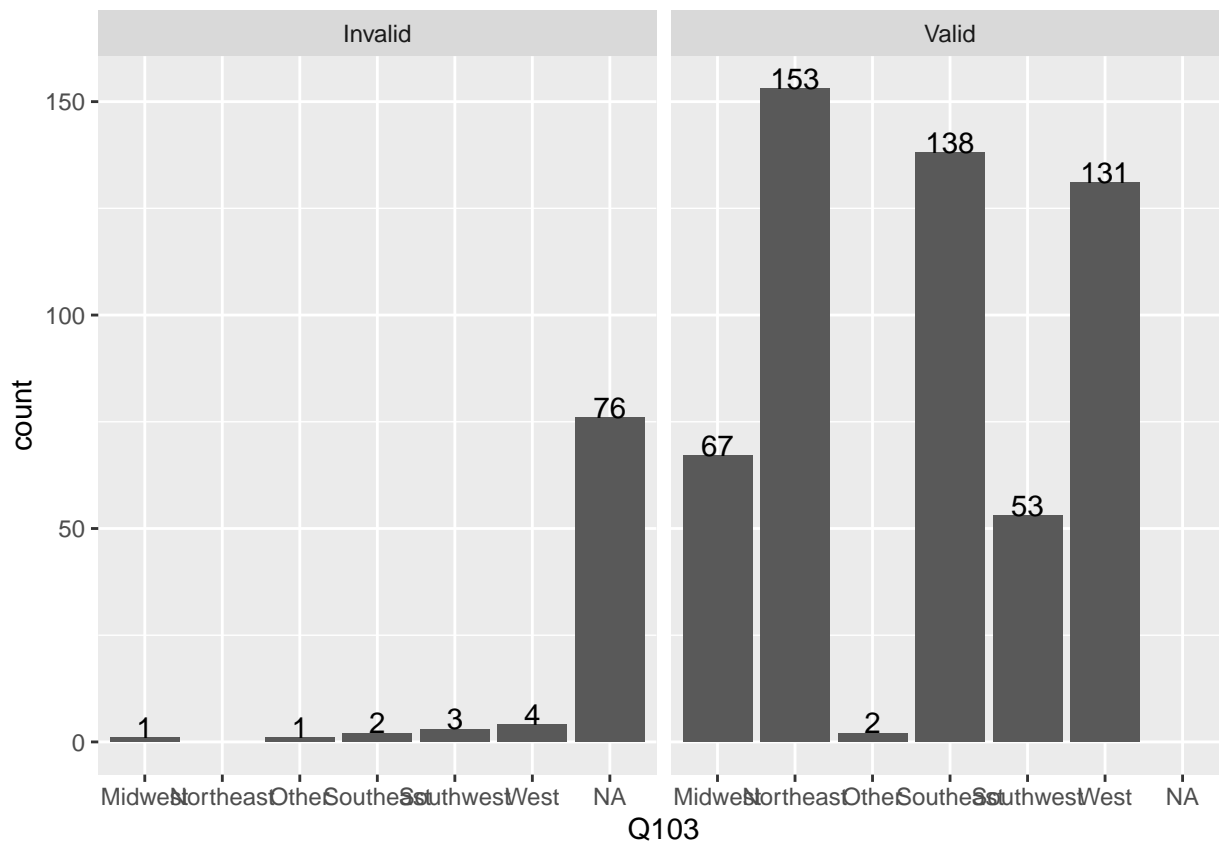
This should be the map from final survey version.

```
img <- readPNG("Picture1.png")
grid.raster(img)
```



Check counts in each region

```
ggplot(geo) + aes(x = Q103) + geom_bar()+geom_text(stat='count',
aes(label=..count..), vjust=0)+facet_wrap(geo$ValidResponse)
```



We have 76 NA values in geo questions. the population in this category is larger than Midwest and Southwest. Can we simply excluding them?

Initial visual ends here

## Seperate Demographics

Seperate the Demographic part, setting them as adequate type

```
DemoText = select(TextDataRow, Q10:Q28, ValidResponse, ResponseId) %>% slice(-(1:2)) %>%
  mutate_at(.vars = c(1:4, 8:9, 12:13, 16:21), .funs = as.factor) %>%
  mutate_at(.vars = 5:6, .funs = as.numeric) %>% mutate(Q13 = as.list(strsplit(Q13, ",")),
    Q18 = (as.list(strsplit(Q18, ","))))
```

Combine “Day care” with “Daycare”

```
DemoText = DemoText %>% mutate(Q13_6_TEXT = if_else(DemoText$Q13_6_TEXT == "Daycare", "Day care",
  DemoText$Q13_6_TEXT))
```

Right now, two lists are in the result, need to separate them into columns

## Correcting column names and factor codings for DemoText

```
## check levels:
print(levels(DemoText$Q10))
```

```
## [1] "Board Certified Assistant Behavior Analyst (BCaBA)"
## [2] "Board Certified Behavior Analyst (BCBA)"
```

```
## [3] "Board Certified Behavior Analyst- Doctoral (BCBA-D)"
## [4] "Not certified"
## [5] "Registered Behavior Technician (RBT)"

## Relevel
DemoText$Q10 = relevel(DemoText$Q10, "Registered Behavior Technician (RBT)")
## Do the same to all demo questions
DemoText$Q104 <- factor(DemoText$Q104, levels = rev(levels(DemoText$Q104)))

DemoText$Q58_1 <- factor(DemoText$Q58_1, levels = c("< 2 years", "2 to 5 years", "6 to 10 years", "11 to 15 years"))
DemoText$Q58_2 <- relevel(DemoText$Q58_2, "Yes")

DemoText$Q19 <- relevel(DemoText$Q19, "Yes")

DemoText$Q17 <- factor(DemoText$Q17, levels = c("1", "2", "3-5", "+5", "None"))

DemoText$Q105 <- factor(DemoText$Q105, levels = c("Never", "Seldom", "Sometimes", "Often"))
DemoText$Q22 <- factor(DemoText$Q22, levels = c("Never", "Seldom", "Sometimes", "Often"))
DemoText$Q29 <- relevel(DemoText$Q29, "Yes")

DemoText$Q30 <- factor(DemoText$Q30, levels = c("Never", "Seldom", "Sometimes", "Often"))
DemoText$Q24 <- relevel(DemoText$Q24, "Yes")
DemoText$Q28 <- relevel(DemoText$Q28, "Yes")

## Rename the columns so the number match the number on the survey
colnames(DemoText) <- c("Q1", "Q2", "Q3", "Q3_text", "Q4_x", "Q4_y", "Q4_Region", "Q5a", "Q5b", "Q6", "Q6_text", "Q7", "Q8", "Q9", "Q10", "Q11", "Q12", "Q13", "Q14", "Q15", "Q16", "Q17", "Q18", "Q19", "Q20", "Q21", "Q22", "Q23", "Q24", "Q25", "Q26", "Q27", "Q28", "Q29", "Q30", "Q31", "Q32", "Q33", "Q34", "Q35", "Q36", "Q37", "Q38", "Q39", "Q40", "Q41", "Q42", "Q43", "Q44", "Q45", "Q46", "Q47", "Q48", "Q49", "Q50", "Q51", "Q52", "Q53", "Q54", "Q55", "Q56", "Q57", "Q58_1", "Q58_2", "Q59", "Q60", "Q61", "Q62", "Q63", "Q64", "Q65", "Q66", "Q67", "Q68", "Q69", "Q70", "Q71", "Q72", "Q73", "Q74", "Q75", "Q76", "Q77", "Q78", "Q79", "Q80", "Q81", "Q82", "Q83", "Q84", "Q85", "Q86", "Q87", "Q88", "Q89", "Q90", "Q91", "Q92", "Q93", "Q94", "Q95", "Q96", "Q97", "Q98", "Q99", "Q100", "Q101", "Q102", "Q103", "Q104", "Q105", "Q106", "Q107", "Q108", "Q109", "Q110", "Q111", "Q112", "Q113", "Q114", "Q115", "Q116", "Q117", "Q118", "Q119", "Q120", "Q121", "Q122", "Q123", "Q124", "Q125", "Q126", "Q127", "Q128", "Q129", "Q130", "Q131", "Q132", "Q133", "Q134", "Q135", "Q136", "Q137", "Q138", "Q139", "Q140", "Q141", "Q142", "Q143", "Q144", "Q145", "Q146", "Q147", "Q148", "Q149", "Q150", "Q151", "Q152", "Q153", "Q154", "Q155", "Q156", "Q157", "Q158", "Q159", "Q160", "Q161", "Q162", "Q163", "Q164", "Q165", "Q166", "Q167", "Q168", "Q169", "Q170", "Q171", "Q172", "Q173", "Q174", "Q175", "Q176", "Q177", "Q178", "Q179", "Q180", "Q181", "Q182", "Q183", "Q184", "Q185", "Q186", "Q187", "Q188", "Q189", "Q190", "Q191", "Q192", "Q193", "Q194", "Q195", "Q196", "Q197", "Q198", "Q199", "Q200", "Q201", "Q202", "Q203", "Q204", "Q205", "Q206", "Q207", "Q208", "Q209", "Q210", "Q211", "Q212", "Q213", "Q214", "Q215", "Q216", "Q217", "Q218", "Q219", "Q220", "Q221", "Q222", "Q223", "Q224", "Q225", "Q226", "Q227", "Q228", "Q229", "Q230", "Q231", "Q232", "Q233", "Q234", "Q235", "Q236", "Q237", "Q238", "Q239", "Q240", "Q241", "Q242", "Q243", "Q244", "Q245", "Q246", "Q247", "Q248", "Q249", "Q250", "Q251", "Q252", "Q253", "Q254", "Q255", "Q256", "Q257", "Q258", "Q259", "Q260", "Q261", "Q262", "Q263", "Q264", "Q265", "Q266", "Q267", "Q268", "Q269", "Q270", "Q271", "Q272", "Q273", "Q274", "Q275", "Q276", "Q277", "Q278", "Q279", "Q280", "Q281", "Q282", "Q283", "Q284", "Q285", "Q286", "Q287", "Q288", "Q289", "Q290", "Q291", "Q292", "Q293", "Q294", "Q295", "Q296", "Q297", "Q298", "Q299", "Q300", "Q301", "Q302", "Q303", "Q304", "Q305", "Q306", "Q307", "Q308", "Q309", "Q310", "Q311", "Q312", "Q313", "Q314", "Q315", "Q316", "Q317", "Q318", "Q319", "Q320", "Q321", "Q322", "Q323", "Q324", "Q325", "Q326", "Q327", "Q328", "Q329", "Q330", "Q331", "Q332", "Q333", "Q334", "Q335", "Q336", "Q337", "Q338", "Q339", "Q340", "Q341", "Q342", "Q343", "Q344", "Q345", "Q346", "Q347", "Q348", "Q349", "Q350", "Q351", "Q352", "Q353", "Q354", "Q355", "Q356", "Q357", "Q358", "Q359", "Q360", "Q361", "Q362", "Q363", "Q364", "Q365", "Q366", "Q367", "Q368", "Q369", "Q370", "Q371", "Q372", "Q373", "Q374", "Q375", "Q376", "Q377", "Q378", "Q379", "Q380", "Q381", "Q382", "Q383", "Q384", "Q385", "Q386", "Q387", "Q388", "Q389", "Q390", "Q391", "Q392", "Q393", "Q394", "Q395", "Q396", "Q397", "Q398", "Q399", "Q400", "Q401", "Q402", "Q403", "Q404", "Q405", "Q406", "Q407", "Q408", "Q409", "Q410", "Q411", "Q412", "Q413", "Q414", "Q415", "Q416", "Q417", "Q418", "Q419", "Q420", "Q421", "Q422", "Q423", "Q424", "Q425", "Q426", "Q427", "Q428", "Q429", "Q430", "Q431", "Q432", "Q433", "Q434", "Q435", "Q436", "Q437", "Q438", "Q439", "Q440", "Q441", "Q442", "Q443", "Q444", "Q445", "Q446", "Q447", "Q448", "Q449", "Q450", "Q451", "Q452", "Q453", "Q454", "Q455", "Q456", "Q457", "Q458", "Q459", "Q460", "Q461", "Q462", "Q463", "Q464", "Q465", "Q466", "Q467", "Q468", "Q469", "Q470", "Q471", "Q472", "Q473", "Q474", "Q475", "Q476", "Q477", "Q478", "Q479", "Q480", "Q481", "Q482", "Q483", "Q484", "Q485", "Q486", "Q487", "Q488", "Q489", "Q490", "Q491", "Q492", "Q493", "Q494", "Q495", "Q496", "Q497", "Q498", "Q499", "Q500", "Q501", "Q502", "Q503", "Q504", "Q505", "Q506", "Q507", "Q508", "Q509", "Q510", "Q511", "Q512", "Q513", "Q514", "Q515", "Q516", "Q517", "Q518", "Q519", "Q520", "Q521", "Q522", "Q523", "Q524", "Q525", "Q526", "Q527", "Q528", "Q529", "Q530", "Q531", "Q532", "Q533", "Q534", "Q535", "Q536", "Q537", "Q538", "Q539", "Q540", "Q541", "Q542", "Q543", "Q544", "Q545", "Q546", "Q547", "Q548", "Q549", "Q550", "Q551", "Q552", "Q553", "Q554", "Q555", "Q556", "Q557", "Q558", "Q559", "Q560", "Q561", "Q562", "Q563", "Q564", "Q565", "Q566", "Q567", "Q568", "Q569", "Q570", "Q571", "Q572", "Q573", "Q574", "Q575", "Q576", "Q577", "Q578", "Q579", "Q580", "Q581", "Q582", "Q583", "Q584", "Q585", "Q586", "Q587", "Q588", "Q589", "Q590", "Q591", "Q592", "Q593", "Q594", "Q595", "Q596", "Q597", "Q598", "Q599", "Q600", "Q601", "Q602", "Q603", "Q604", "Q605", "Q606", "Q607", "Q608", "Q609", "Q610", "Q611", "Q612", "Q613", "Q614", "Q615", "Q616", "Q617", "Q61
```

### Separation For Q6

```
Q6Text = select(DemoText, Q6, ValidResponse, ResponseId) %>%
  unnest(Q6) %>% spread(key = Q6, value = Q6) %>%
  mutate_at(3:9, ~replace(., !is.na(.), 1)) %>%
  mutate_at(3:9, ~replace(., is.na(.), 0))
```

### Similarly Separation for Q9

```
Q9Text = select(DemoText, Q9, ValidResponse, ResponseId) %>%
  unnest(Q9) %>% spread(key = Q9, value = Q9) %>% select(-3) %>%
  mutate_at(3:25, ~replace(., !is.na(.), 1)) %>%
  mutate_at(3:25, ~replace(., is.na(.), 0))
# Bird's colname is buggy replace it
colnames(Q9Text)[colnames(Q9Text)=="Bird (i.e.)"] <- "Bird"
```

At this point, we will have DemoText as the demographic part cleaned. Q6Text and Q9Text as the suppliment table to those two questions. What's missing here is the EDA on the valid and invalid groups to see if they belongs to the same population or not. In other words, what and how much bias will be introduced into the sample if we exclude them?

## Last Question Data Cleaning

Preparing the data

```
LastQText = select(TextDataRaw, Q32_1:Q32_4, ValidResponse, ResponseId) %>% slice(-(1:2))

Q32_1Text = select>LastQText, Q32_1, ValidResponse, ResponseId) %>% unnest(Q32_1) %>% spread(key = Q32_1, value = Q32_1Text[, c(1, 2, 6, 3, 5, 4, 7, 8)])

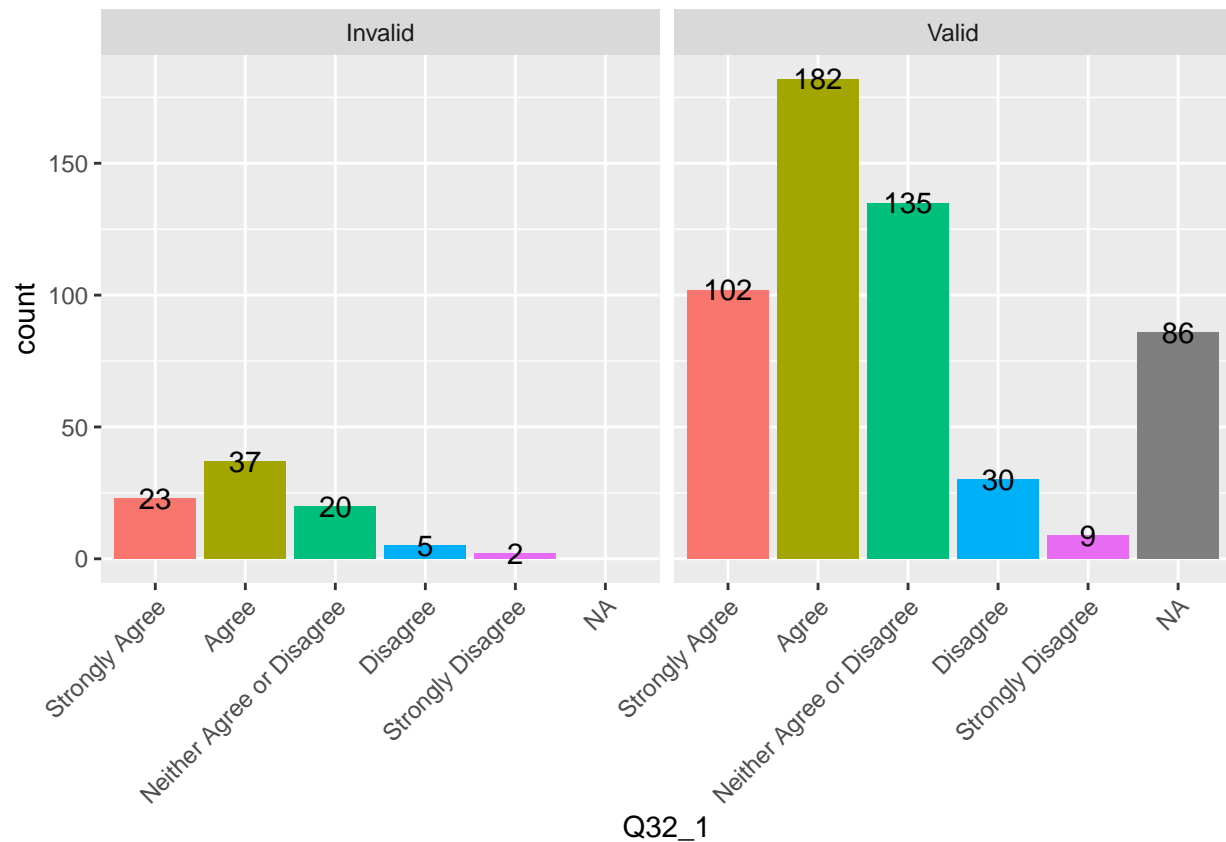
Q32_2Text = select>LastQText, Q32_2, ValidResponse, ResponseId) %>% unnest(Q32_2) %>% spread(key = Q32_2, value = Q32_2Text[, c(1, 2, 6, 3, 5, 4, 7, 8)])

Q32_3Text = select>LastQText, Q32_3, ValidResponse, ResponseId) %>% unnest(Q32_3) %>% spread(key = Q32_3, value = Q32_3Text[, c(1, 2, 6, 3, 5, 4, 7, 8)])

Q32_4Text = select>LastQText, Q32_4, ValidResponse, ResponseId) %>% unnest(Q32_4) %>% spread(key = Q32_4, value = Q32_4Text[, c(1, 2, 6, 3, 5, 4, 7, 8)])
```

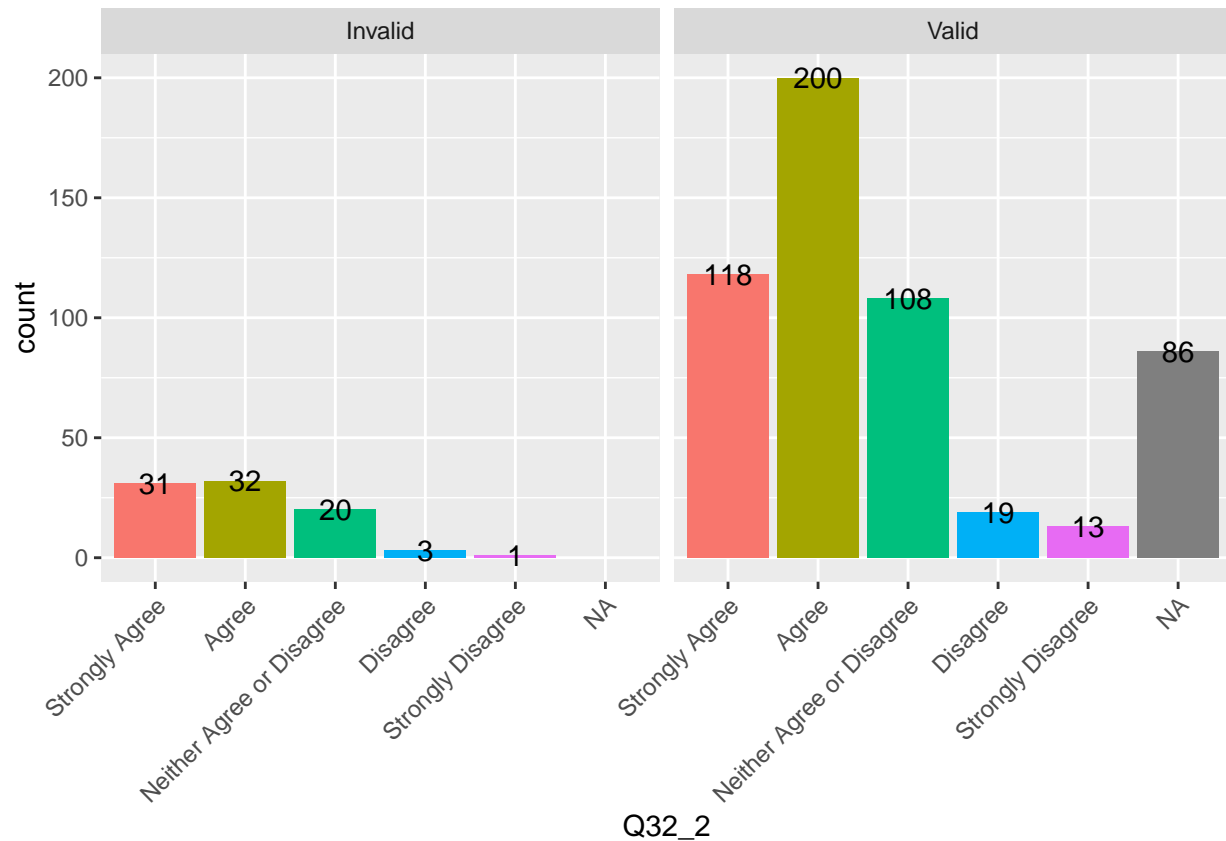
**Q55a.1.** I believe a majority of the target population outlined above would respond positively to therapeutic services (e.g., ABA services) that incorporate animals.

```
LastQText$Q32_1 <- factor>LastQText$Q32_1, levels = c("Strongly Agree", "Agree", "Neither Agree or Disagree", "Disagree", "Strongly Disagree", "NA")
ggplot>LastQText) + aes(Q32_1, fill = Q32_1) + geom_bar(show.legend = FALSE) + geom_text(stat = 'count', aes(label = count))
```



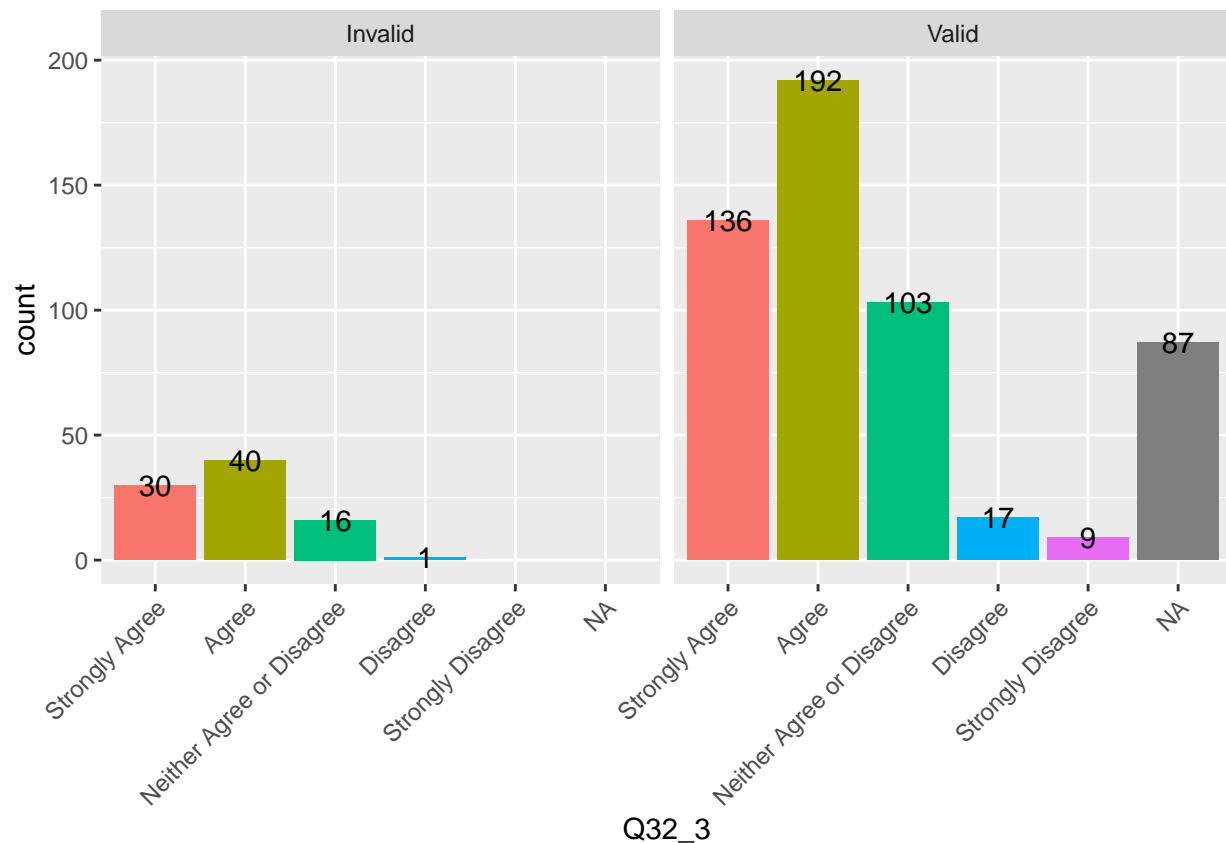
55b.2.I believe the presence of an animal can have a calming effect on many individuals within the target population outlined above

```
LastQText$Q32_2 <- factor>LastQText$Q32_2,levels = c("Strongly Agree", "Agree", "Neither Agree or Disagree", "Disagree", "Strongly Disagree", "NA")
ggplot>LastQText) + aes(Q32_2,fill = Q32_2)+geom_bar(show.legend = FALSE)+geom_text(stat='count', aes(l
```



55c.3.I believe the presence of an animal can increase the frequency of social interaction opportunities between individuals within the target population outlined above and other people.

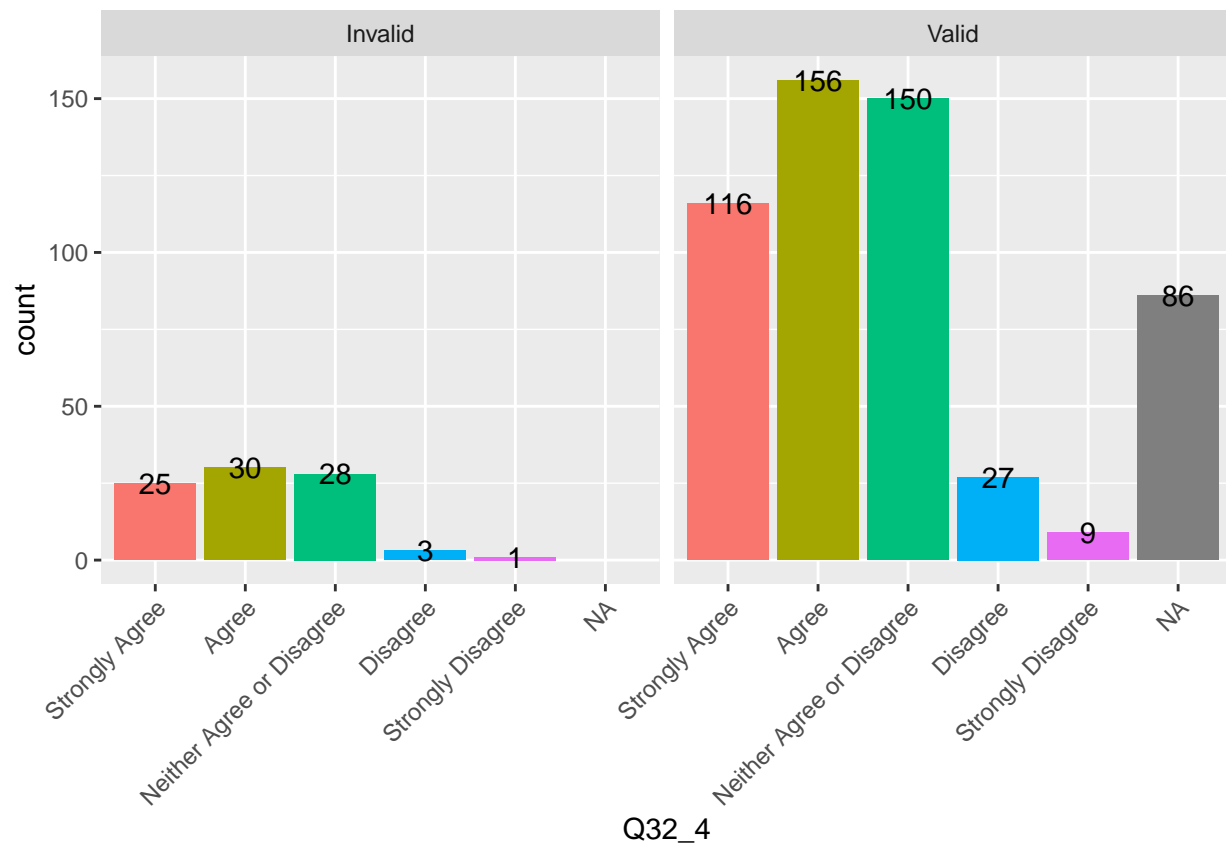
```
LastQText$Q32_3 <- factor>LastQText$Q32_3,levels = c("Strongly Agree", "Agree", "Neither Agree or Disagree", "Disagree", "Strongly Disagree", "NA")
ggplot>LastQText) + aes(Q32_3,fill = Q32_3)+geom_bar(show.legend = FALSE)+geom_text(stat='count', aes(l
```



55d.4.I believe animals can aid in effectively teaching social interaction and communication skills associated with ASD symptomatology for a majority of the target population outlined above.

```
LastQText$Q32_4 <- factor>LastQText$Q32_4,levels = c("Strongly Agree", "Agree", "Neither Agree or Disagree", "Disagree", "Strongly Disagree", "NA"))
ggplot>LastQText) + aes(Q32_4,fill = Q32_4)+geom_bar(show.legend = FALSE)+geom_text(stat='count', aes(l
```





## Testing the Invalid IDs' Responses

Check if the factor is adequately transferred into integers

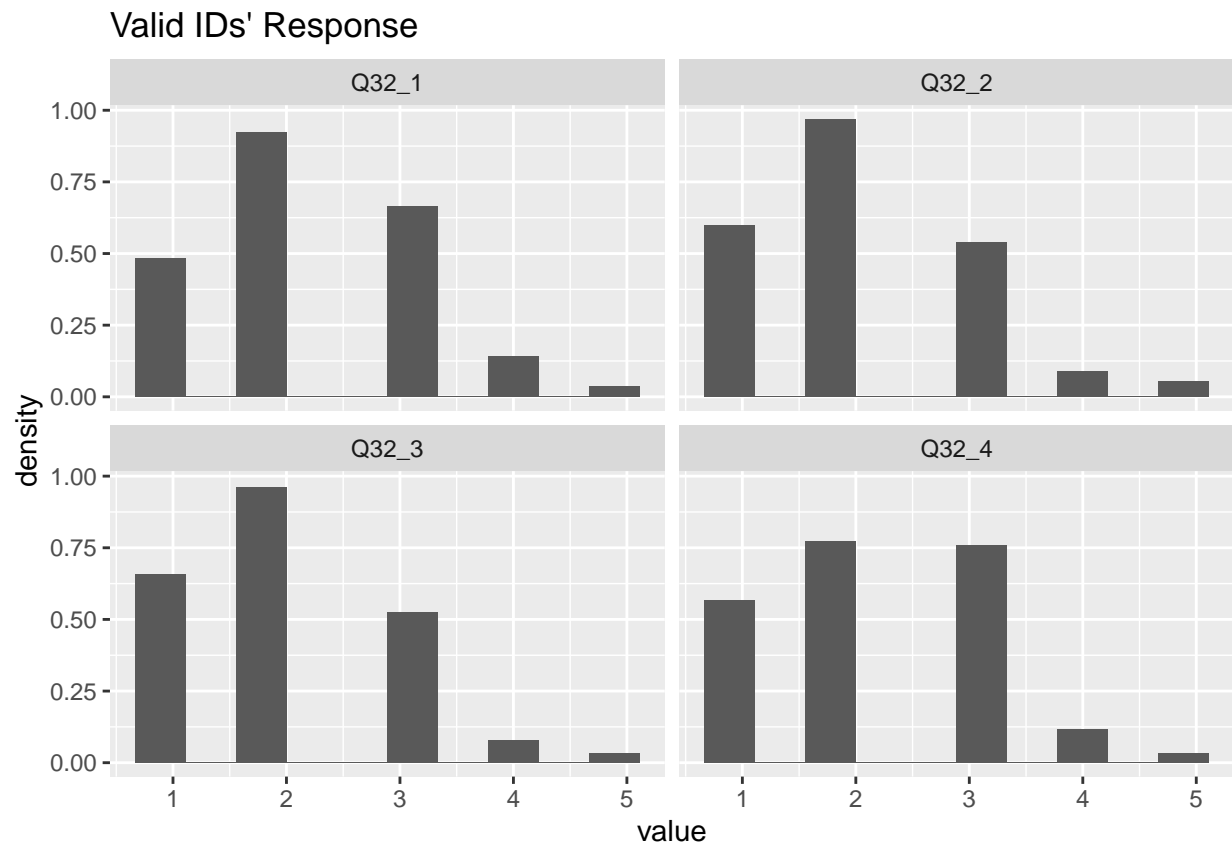
```
LastQTextTEST <- LastQText %>% mutate_at(.var = 1:4, .funs = as.factor) %>% mutate_at(.var = 1:4, .funs = as.integer)
```

Create separate tables for Valid and Invalid Responses

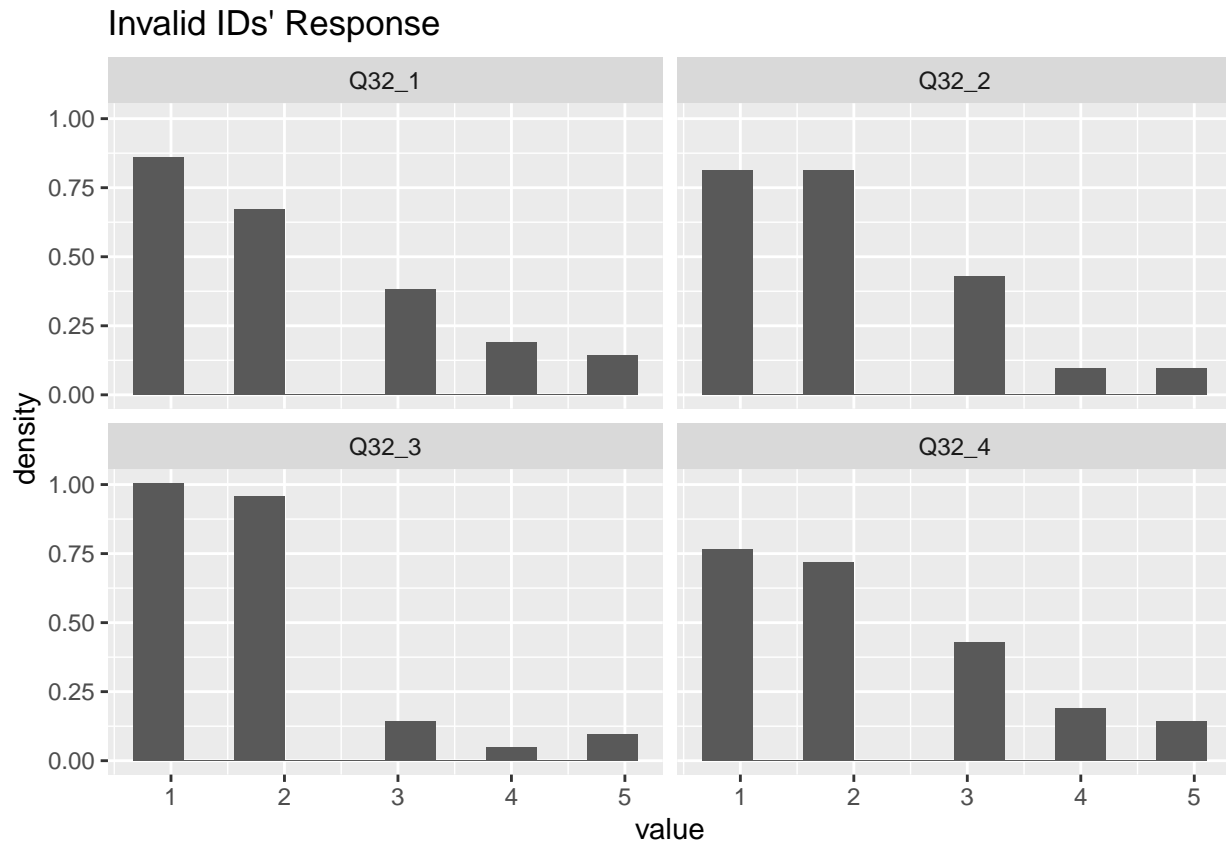
```
LQvalid <- LastQTextTEST %>% filter(ValidResponse == "Valid") %>% select(1:4) %>% drop_na()
LQInvalid <- LastQTextTEST %>% filter(ValidResponse == "Invalid") %>% select(1:4) %>% drop_na()
```

Plot the Histogram, but 1st needs to rearrange the data frame

```
LQvalid1 <- LQvalid %>% melt(measure.vars = 1:4)
ggplot(LQvalid1) + aes(x = value) + geom_histogram(bins = 10, aes(y = ..density..)) + facet_wrap(LQvalid1$variable)
```



```
LQInvalid1 <- LQInvalid %>% melt(measure.vars = 1:4)
ggplot(LQInvalid1)+aes(x = value)+geom_histogram(bins = 10,aes(y=..density..))+facet_wrap(LQInvalid1$va
```



The data does not seem to have a normal distribution, can we still use t test to test these two sample means?

In case we can, below is the t test results.

```
for(i in 1:4){
  print(t.test(select(LQInvalid,i)%>%pull(),select(LQvalid,i)%>%pull()))
}

##
##  Welch Two Sample t-test
##
## data:  select(LQInvalid, i) %>% pull() and select(LQvalid, i) %>% pull()
## t = -0.5746, df = 51.084, p-value = 0.5681
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##   -0.4699758  0.2608059
## sample estimates:
## mean of x mean of y
##  2.148936  2.253521
##
##
##  Welch Two Sample t-test
##
## data:  select(LQInvalid, i) %>% pull() and select(LQvalid, i) %>% pull()
## t = -0.52465, df = 52.917, p-value = 0.602
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
```

```

## -0.4061454 0.2377307
## sample estimates:
## mean of x mean of y
## 2.042553 2.126761
##
##
## Welch Two Sample t-test
##
## data: select(LQInvalid, i) %>% pull() and select(LQvalid, i) %>% pull()
## t = -1.7922, df = 53.435, p-value = 0.07877
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.5616919 0.03153223
## sample estimates:
## mean of x mean of y
## 1.787234 2.052314
##
##
## Welch Two Sample t-test
##
## data: select(LQInvalid, i) %>% pull() and select(LQvalid, i) %>% pull()
## t = -0.11495, df = 51.44, p-value = 0.9089
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.3809444 0.3396755
## sample estimates:
## mean of x mean of y
## 2.212766 2.233400

```

Judging from the t statistics, it seems that the id excluded were appeared to have high probability that they comes from the same population? Can we safely exclude them? Or does it matters if we exclude them or not?