

# EDA

Anni Hong

9/15/2020

## Basic Info:

### Individual dataset

- The individual dataset contains 370 observations, and 10 variables :  
ID Participant ID number team.id ID number of the team this participant belonged to Age Age, in years  
Gender Gender (Male or Female) Ethnicity Ethnicity of the participant Cortisol Participant's cortisol levels, nMol/L Testosterone Participant's testosterone levels, pg/mL log.cortisol Natural logarithm of the participant's cortisol level log.testosterone Natural logarithm of the participant's testosterone level  
Country Country of citizenship of the participant
- 18 rows contain at least one missing value in one of the columns

### Team dataset

- The team dataset contains 74 teams and 14 variables:  
team.id Team ID number team.size Number of people on the team final.performance The team's final performance score time.of.day The time of day the team's hormone sample was collected (hh.mm)  
females Number of females in the group final.cash Total cash earned by the team final.contracts Total number of contracts won by the team final.reorders Total number of reorders won by the team  
final.rank Team's final rank at the end of the project, relative to other teams in their class section  
interim.performance Same as above, but measured at Day 5 of the study (missing for some teams)  
interim.cash  
interim.contracts  
interim.reorders  
interim.rank
- 22 teams have no interim.rank, interim.reorders, interim.contracts, interim.cash, and interim.performance

*#factorize certain variables and create diversity score for each group*

```
indi_dat <- indi_dat %>%
```

```
  mutate_if(sapply(indi_dat, is.character), as.factor) %>%
```

```
  group_by(team.id) %>%
```

```
  mutate(diversity_score = n_distinct(Gender, Ethnicity, Country))
```

*#intermediate step of calculating aggregated statistics by group*

```
agg_indi_dat <- indi_dat %>%
```

```
  group_by(team.id) %>%
```

```
  summarise(mean_testo = mean(Testosterone), mean_log_testo = mean(log.testosterone), sd_testo = sd(Tes
```

```
## `summarise()` ungrouping output (override with `.groups` argument)
```

```
#combing the info from the individual dataset to the team dataset, average/sd cortisol and testosterone.
combo_dat <- team_dat %>%
  inner_join(agg_indi_dat, by = "team.id")
```

Understanding the variables

Typical amount of diversity present:

```
diversity_vars <- indi_dat %>% select(Gender, Ethnicity, Country)
```

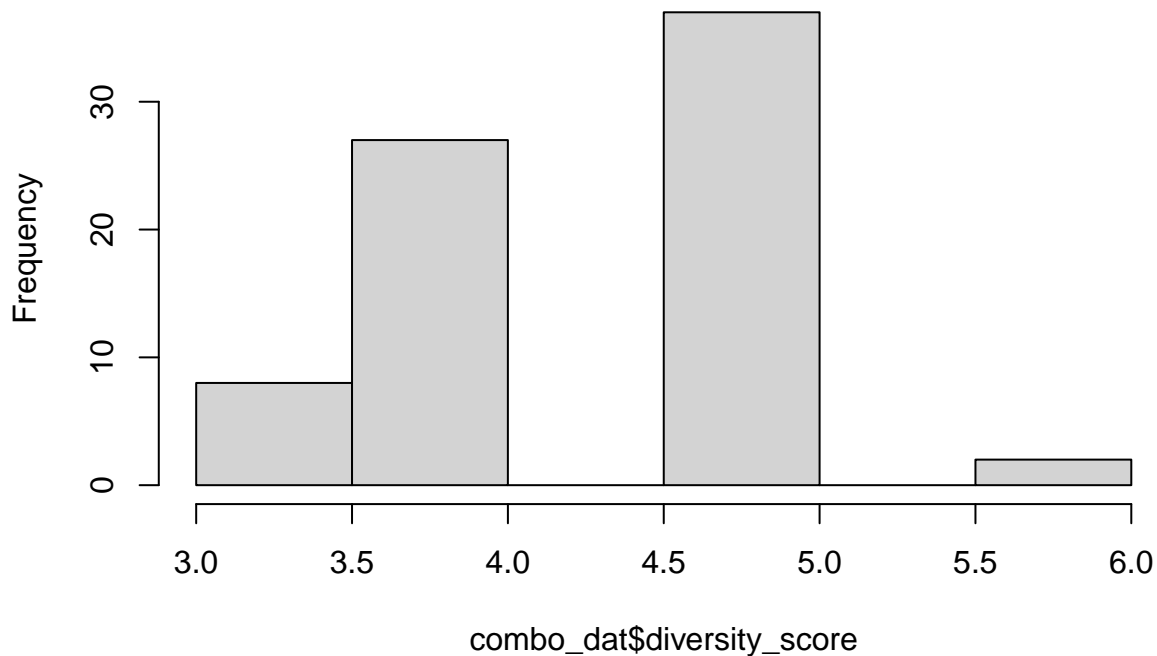
```
## Adding missing grouping variables: `team.id`
```

```
summary(diversity_vars[, -1])
```

```
##      Gender      Ethnicity      Country
## Female:133   Asian      : 61   USA      :213
## Male  :237   Black      : 17   China    : 19
##                               Hispanic : 40   India    : 16
##                               Other    : 9    Korea    : 10
##                               South Asian : 35   Argentina: 9
##                               South East Asian: 5   Canada   : 8
##                               White     :203   (Other)  : 95
```

```
hist(combo_dat$diversity_score)
```

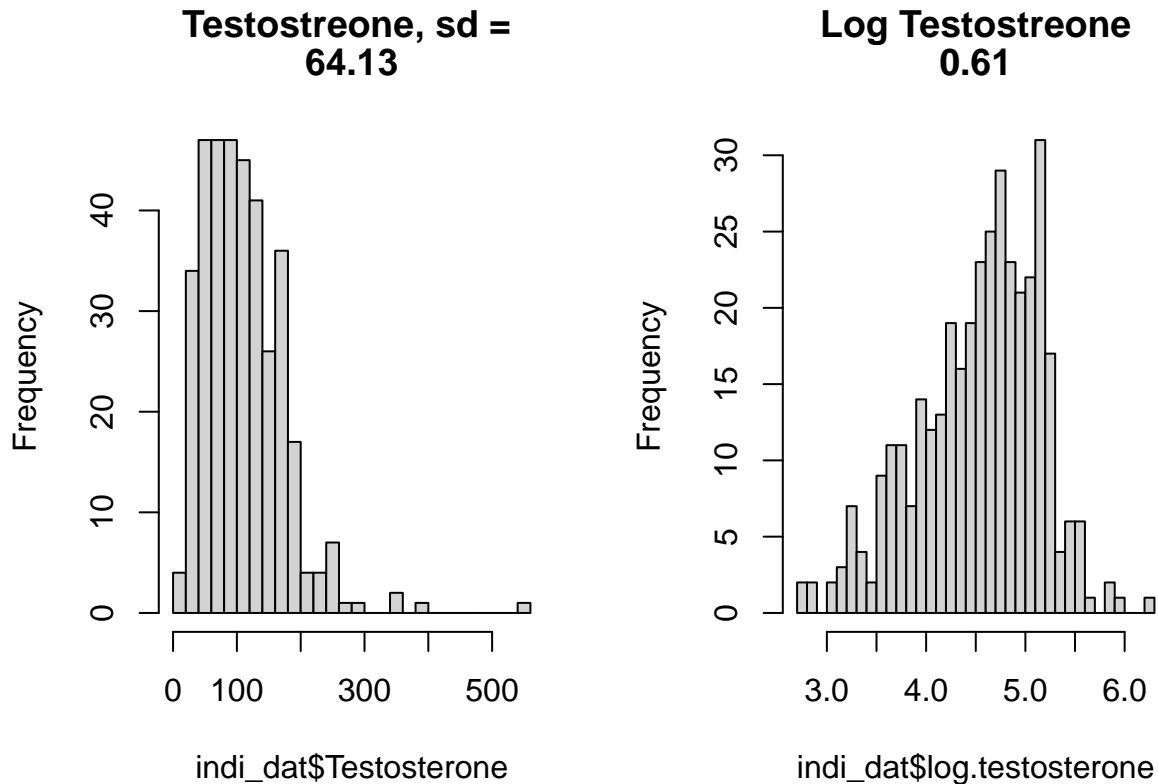
**Histogram of combo\_dat\$diversity\_score**



This dataset has a lot more men than women and mostly white Americans. There is a adequate amount of variability in the diversity scores.

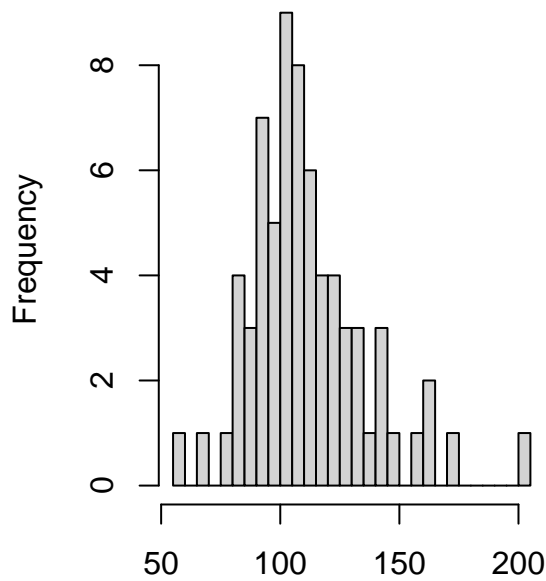
testostreone vs log testostreone

```
par(mfrow=c(1,2))
sd_testo = round(sd(indi_dat$Testosterone, na.rm = T),2)
sd_log_testo = round(sd(indi_dat$log.testosterone, na.rm = T),2)
hist(indi_dat$Testosterone, breaks = 30, main = c("Testostreone, sd = ",sd_testo))
hist(indi_dat$log.testosterone, breaks = 30, main = c("Log Testostreone", sd_log_testo))
```



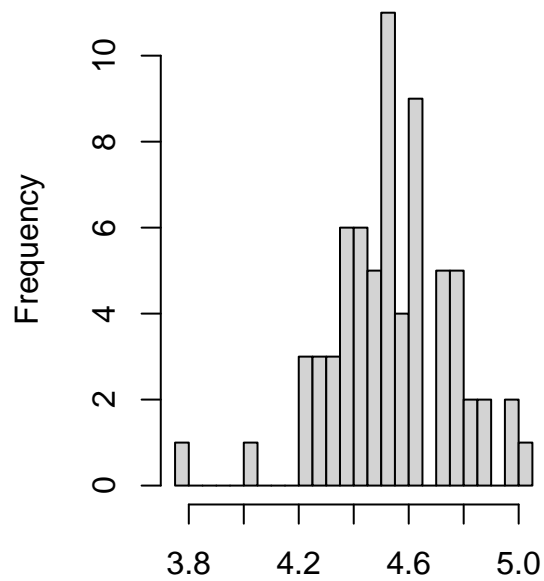
```
par(mfrow=c(1,2))
sd_testo = round(sd(agg_indi_dat$mean_testo, na.rm = T),2)
sd_log_testo = round(sd(agg_indi_dat$mean_log_testo, na.rm = T),2)
hist(agg_indi_dat$mean_testo, breaks = 30, main = c("Mean Group Testostreone, sd = ",sd_testo))
hist(agg_indi_dat$mean_log_testo, breaks = 30, main = c("Mean Group Log Testostreone", sd_log_testo))
```

**Mean Group Testosterone, sd = 24.93**



agg\_indi\_dat\$mean\_testo

**Mean Group Log Testosterone 0.22**



agg\_indi\_dat\$mean\_log\_testo

The above plots shows the histogram of the testosterone level compared to the log of the testosterone level, for overall as well as group mean. The log transformation helps with lessen the impact of outliers on our analysis. We don't want to let the group that contain the individual with really high testosterone level impact our analysis disproportionately since mean is very sensitive to outliers. The standard deviation is drastically reduced in the log transformed variable in both cases and the histogram approximate a normal distribution.

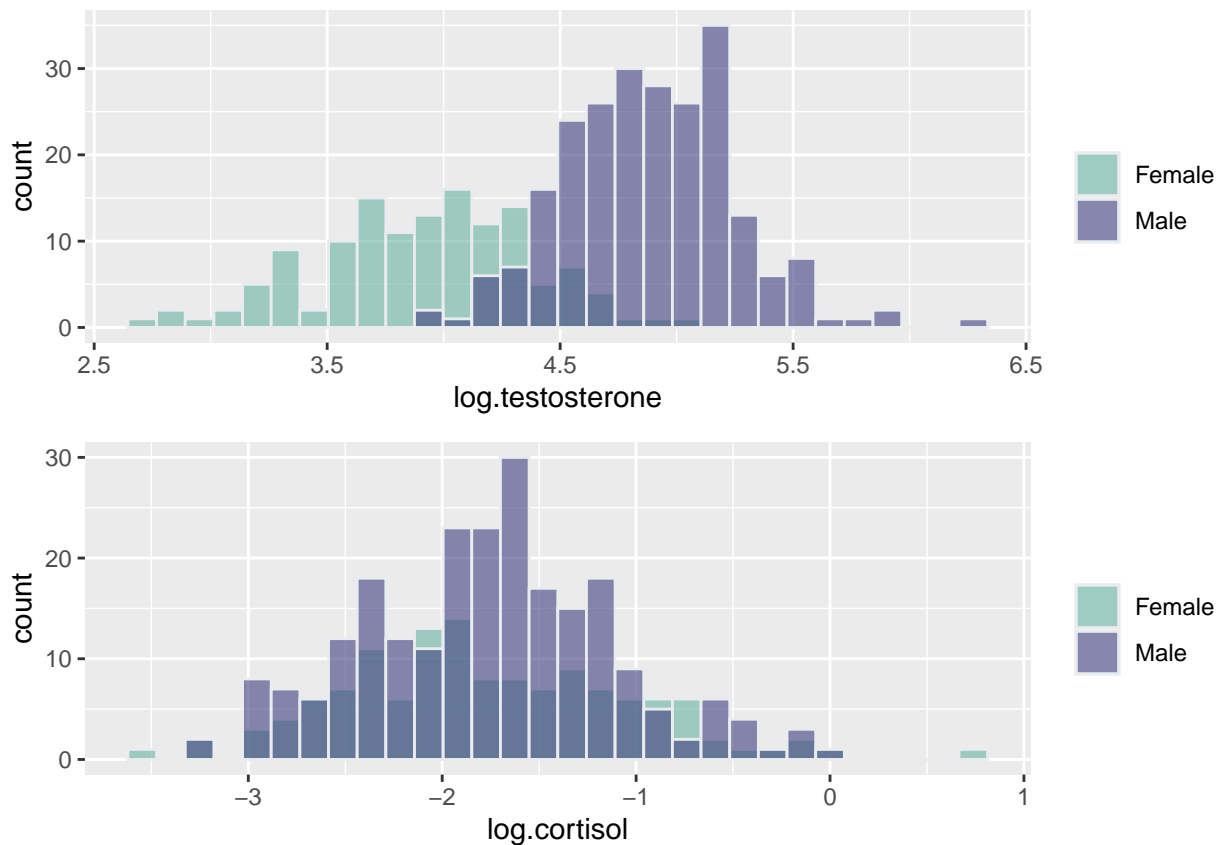
#### testosterone, cortisol by gender

```
testo <- indi_dat %>%
  ggplot( aes(x=log.testosterone, fill=Gender)) +
  geom_histogram( color="#e9ecef", alpha=0.6, position = 'identity') +
  scale_fill_manual(values=c("#69b3a2", "#404080")) +
  labs(fill="")

cortisol <- indi_dat %>%
  ggplot( aes(x=log.cortisol, fill=Gender)) +
  geom_histogram( color="#e9ecef", alpha=0.6, position = 'identity') +
  scale_fill_manual(values=c("#69b3a2", "#404080")) +
  labs(fill="")

grid.arrange(testo, cortisol)
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

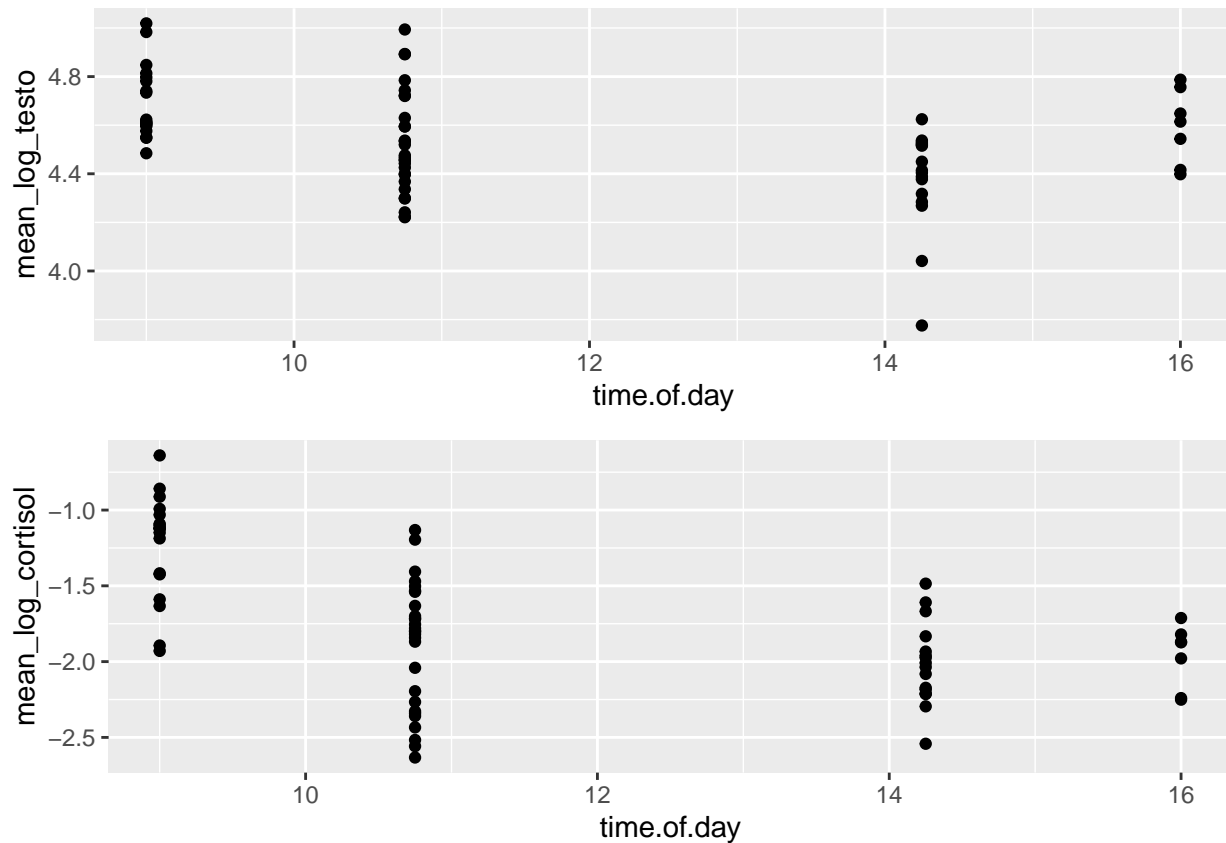


It makes sense that gender is correlated to the level of testosterone so special attention, since women tend to have lower testosterone level than men. Gender seem to have no impact on cortisol level. ### testosterone, cortisol by time of the day

```
testo <- combo_dat %>%
  ggplot(aes(x=time.of.day, y=mean_log_testo)) +
  geom_point()

cortisol <- combo_dat %>%
  ggplot(aes(x=time.of.day, y=mean_log_cortisol)) +
  geom_point()

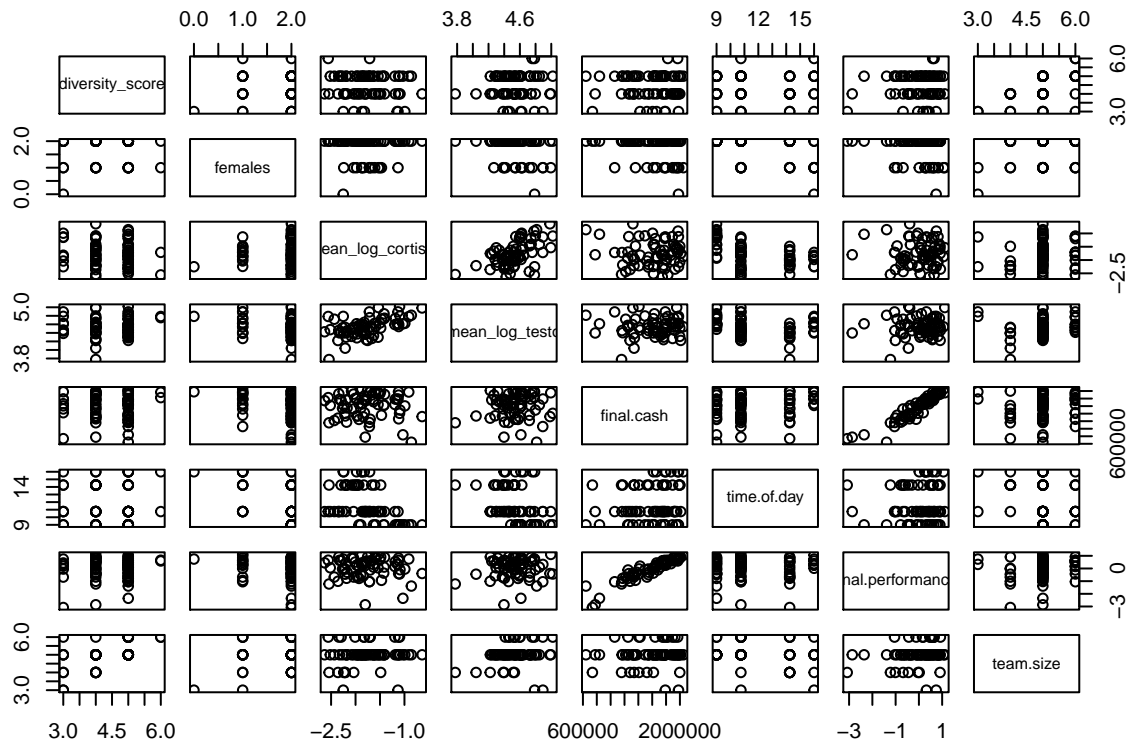
grid.arrange(testo, cortisol)
```



Time of the day could also be a confounder since the testosterone and cortisol level could change base on what time of the day you measure it.

### relationships between a subset of variables

```
combo_dat_subset <- combo_dat %>% select(diversity_score,females, mean_log_cortisol, mean_log_testo, f
pairs(combo_dat_subset) # not including team id
```



- Worth noting that there are exactly 2 females in all the groups, if we want to control for gender it would make more sense to use percentage.
- Cortisol seems to be correlated with Testosterone levels, since Cortisol is related to stress maybe it will help test/explain why there could be an interaction effect between diversity and testosterone (through causing stress on the team).