

Bios 6301: Assignment 9

Jonathan Lifferth

Due Tuesday, 28 November, 1:00 PM

$5^{n=\text{day}}$ points taken off for each day late.

40/40 Great!

40 points total.

Submit a single knitr file (named `homework9.rmd`), along with a valid PDF output file. Inside the file, clearly indicate which parts of your responses go with which problems (you may use the original homework document as a template). Add your name as `author` to the file's metadata section. Raw R code/output or word processor files are not acceptable.

Failure to name file `homework9.rmd` or include author name may result in 5 points taken off.

Question 1

15 points

Consider the following very simple genetic model (*very* simple – don't worry if you're not a geneticist!). A population consists of equal numbers of two sexes: male and female. At each generation men and women are paired at random, and each pair produces exactly two offspring, one male and one female. We are interested in the distribution of height from one generation to the next. Suppose that the height of both children is just the average of the height of their parents, how will the distribution of height change across generations?

Represent the heights of the current generation as a dataframe with two variables, `m` and `f`, for the two sexes. We can use `rnorm` to randomly generate the population at generation 1:

```
pop <- data.frame(m = rnorm(100, 160, 20), f = rnorm(100, 160, 20))
```

The following function takes the data frame `pop` and randomly permutes the ordering of the men. Men and women are then paired according to rows, and heights for the next generation are calculated by taking the mean of each row. The function returns a data frame with the same structure, giving the heights of the next generation.

```
next_gen <- function(pop) {  
  pop$m <- sample(pop$m)  
  pop$m <- rowMeans(pop)  
  pop$f <- pop$m  
  pop  
}
```

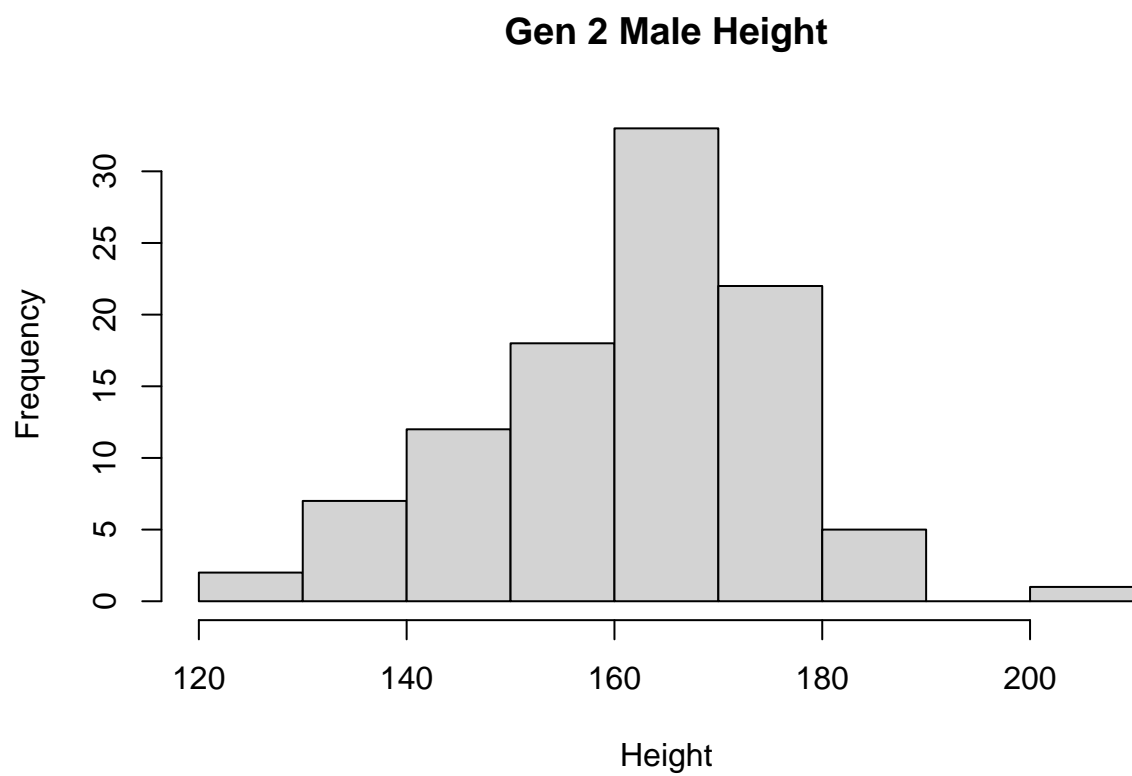
Use the function `next_gen` to generate nine generations (you already have the first), then use the function `hist` to plot the distribution of male heights in each generation (this will require multiple calls to `hist`). The phenomenon you see is called regression to the mean. Provide (at least) minimal decorations such as title and x-axis labels.

```
hist(pop$m, main = "Gen 1 Male Height", xlab = "Height")
```



```
gen2 <- next_gen(pop)  
hist(gen2$m, main = "Gen 2 Male Height", xlab = "Height")
```

can use ggarrange or gridextra to put all hists in one grid

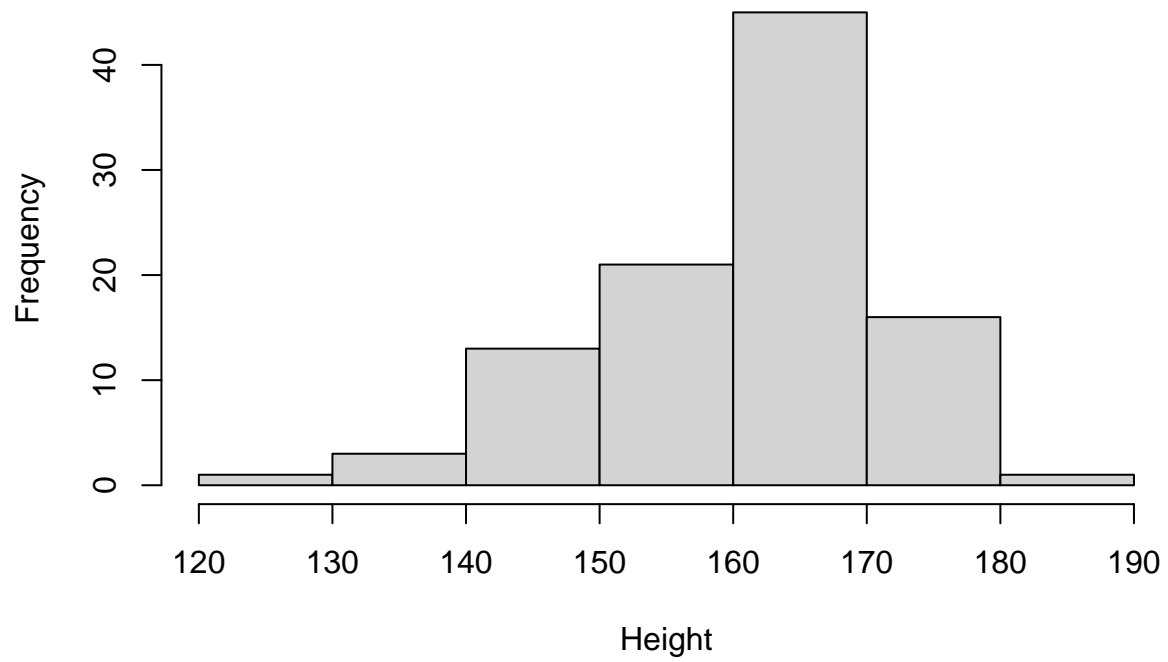


```
gen3 <- next_gen(gen2)
hist(gen3$m, main = "Gen 3 Male Height", xlab = "Height")
```



```
gen4 <- next_gen(gen2)
hist(gen3$m, main = "Gen 3 Male Height", xlab = "Height")
```

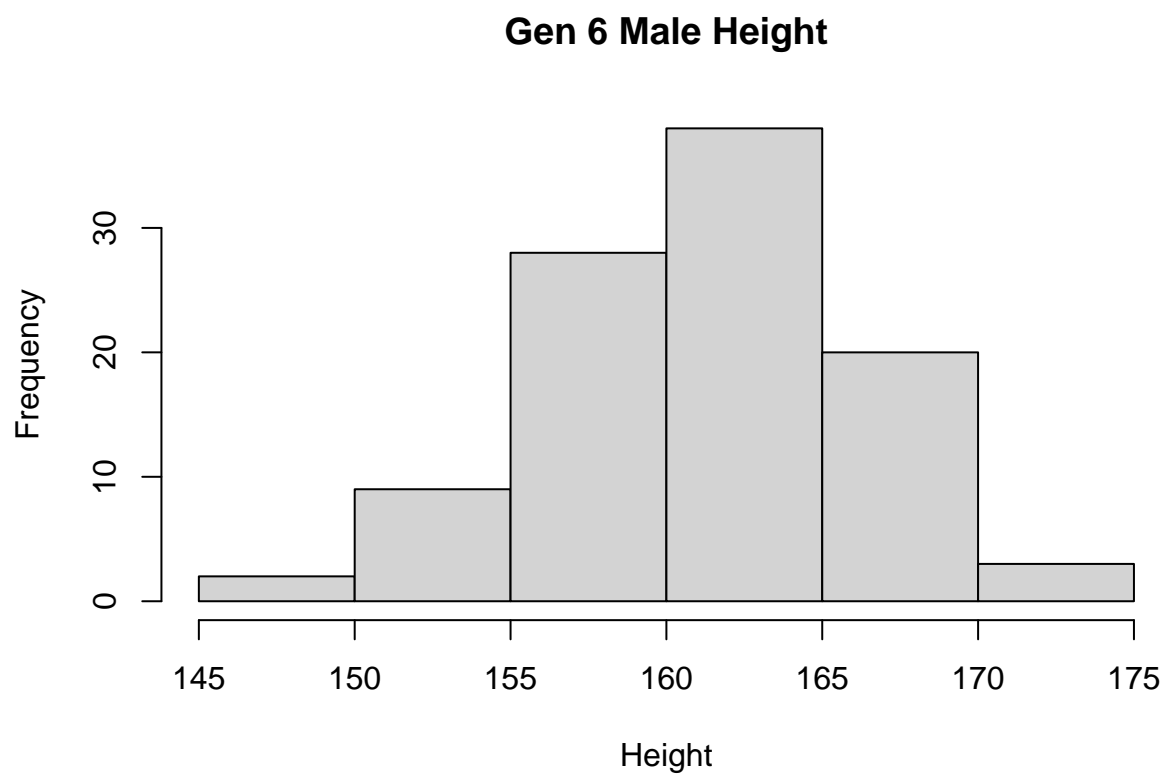
Gen 3 Male Height



```
gen5 <- next_gen(gen4)
hist(gen5$m, main = "Gen 5 Male Height", xlab = "Height")
```



```
gen6 <- next_gen(gen5)
hist(gen6$m, main = "Gen 6 Male Height", xlab = "Height")
```



```
gen7 <- next_gen(gen6)
hist(gen7$m, main = "Gen 7 Male Height", xlab = "Height")
```

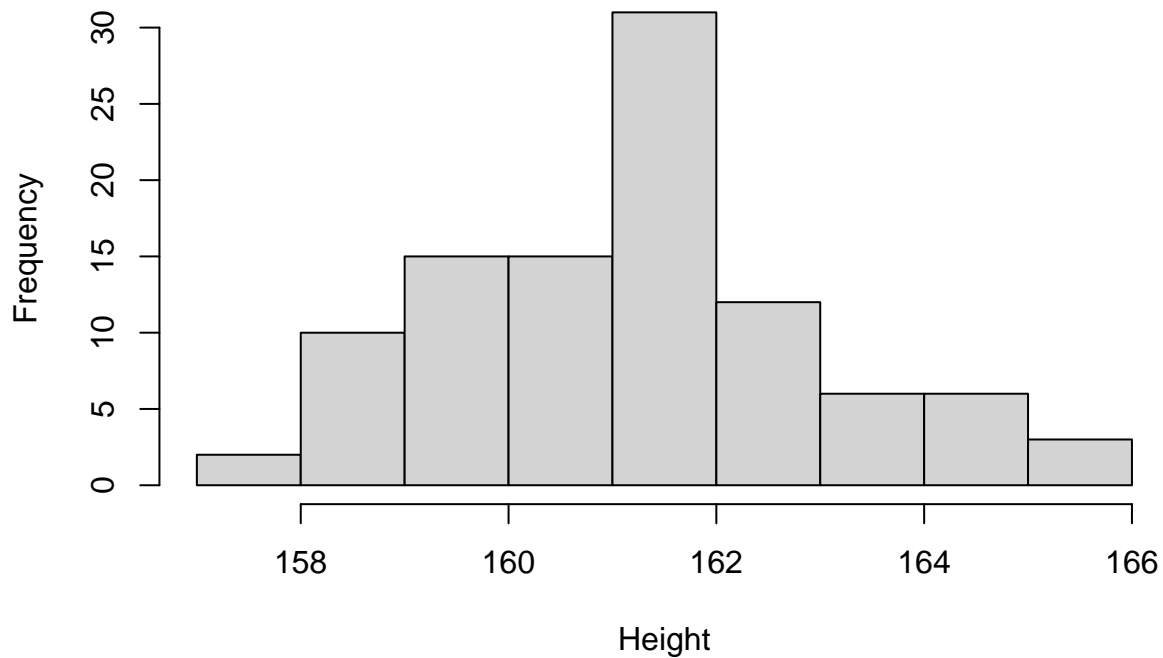


```
gen8 <- next_gen(gen7)
hist(gen8$m, main = "Gen 8 Male Height", xlab = "Height")
```




```
gen9 <- next_gen(gen8)
hist(gen9$m, main = "Gen 9 Male Height", xlab = "Height")
```

Gen 9 Male Height



Question 2

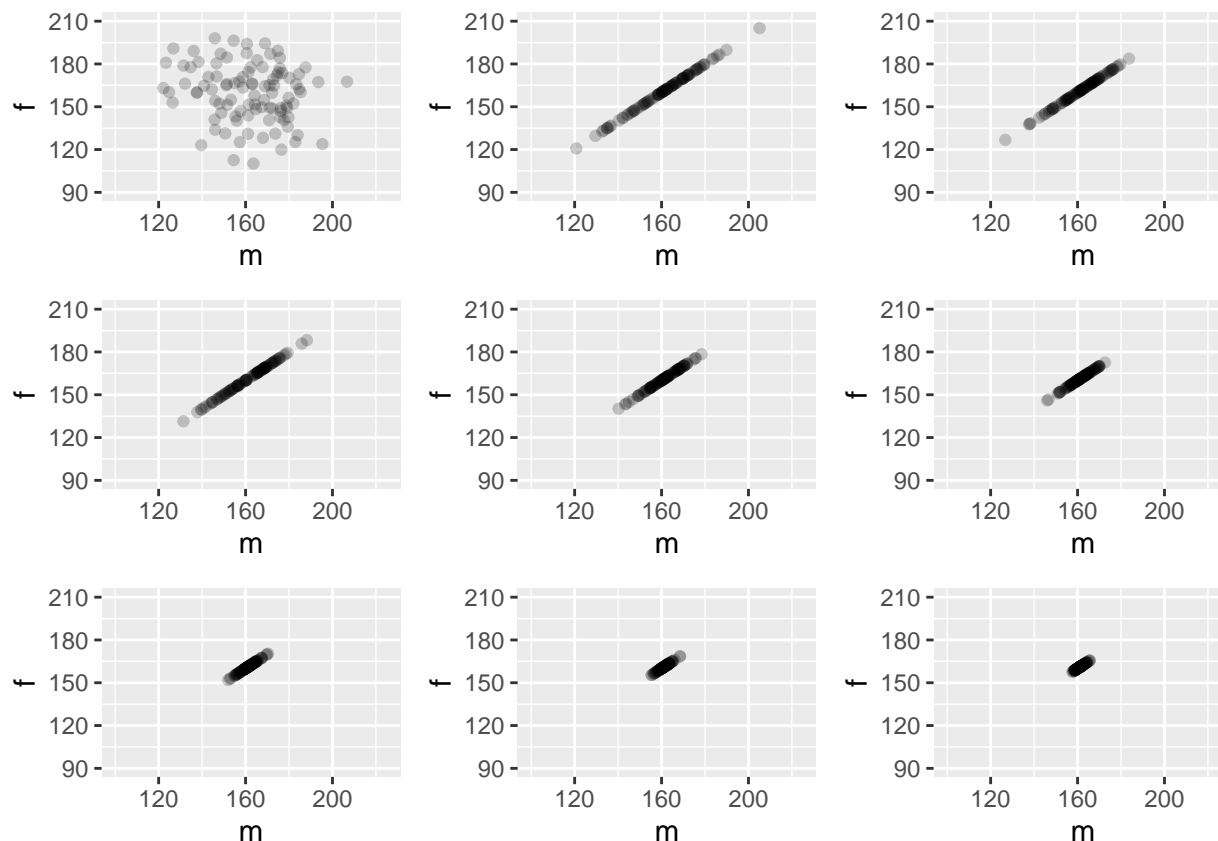
10 points

Use the simulated results from question 1 to reproduce (as closely as possible) the following plot in ggplot2.

```
library(ggplot2)
library(gridExtra)
p1 <- ggplot(pop, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p2 <- ggplot(gen2, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p3 <- ggplot(gen3, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p4 <- ggplot(gen4, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p5 <- ggplot(gen5, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p6 <- ggplot(gen6, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p7 <- ggplot(gen7, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p8 <- ggplot(gen8, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)
p9 <- ggplot(gen9, aes(x=m, y=f)) + geom_point(alpha = 1/5) + xlim(100, 225) + ylim(90, 210)

grid.arrange(p1, p2, p3, p4, p5, p6, p7, p8, p9, nrow = 3)
```

```
## Warning: Removed 2 rows containing missing values ('geom_point()').
```



Question 3

15 points

You calculated the power of a study design in question #1 of assignment 3. The study has two variables, treatment group and outcome. There are two treatment groups (0, 1) and they should be assigned randomly with equal probability. The outcome should be a random normal variable with a mean of 60 and standard deviation of 20. If a patient is in the treatment group, add 5 to the outcome.

Starting with a sample size of 250, create a 95% bootstrap percentile interval for the mean of each group. Then create a new bootstrap interval by increasing the sample size by 250 until the sample is 2500. Thus you will create a total of 10 bootstrap intervals. Each bootstrap should create 1000 bootstrap samples. (9 points)

```
first_sample_size <- 250
num_bootstrap_intervals <- 10
num_bootstrap_samples <- 1000

quantiles_0 <- data.frame(matrix(ncol=4,nrow=0, dimnames=list(NULL, c("mean", "upper", "lower", "size")))
quantiles_1 <- data.frame(matrix(ncol=4,nrow=0, dimnames=list(NULL, c("mean", "upper", "lower", "size")))

# create each interval population
for(i in 1:num_bootstrap_intervals) {
  print(i)
  sample_size <- first_sample_size * i
```

```

# quant0_bootstraps <- data.frame(matrix(ncol=3,nrow=0, dimnames=list(NULL, c("mean", "upper", "lower"),
# quant1_bootstraps <- data.frame(matrix(ncol=3,nrow=0, dimnames=list(NULL, c("mean", "upper", "lower"),

grp <- sample(0:1, sample_size, replace=TRUE)
values <- rnorm(sample_size, 60, 20)
values[grp == 1] <- values[grp == 1] + 5
interval_population <- data.frame(treatment=grp, value=values)

means_0 <- list()
means_1 <- list()

# perform bootstrap, sample from current interval population
for(j in 1:num_bootstrap_samples) {
  bootstrap_sample <- interval_population[sample(nrow(interval_population), sample_size, replace=TRUE),

  # group 0 stats
  mean0 <- mean(subset(bootstrap_sample, treatment == 0)$value)
  means_0[j] <- mean0

  # group 1 stats
  mean1 <- mean(subset(bootstrap_sample, treatment == 1)$value)
  means_1[j] <- mean1
}

# after performing bootstraps, add means to quantile dfs
means_0 <- unlist(means_0)
interval_mean_0 <- mean(means_0)
interval_upper_0 <- quantile(means_0, c(.025, .975))[2]
interval_lower_0 <- quantile(means_0, c(.025, .975))[1]
quantiles_0[i,] <- data.frame(mean=interval_mean_0, upper=interval_upper_0,
                             lower=interval_lower_0, size=sample_size)

means_1 <- unlist(means_1)
interval_mean_1 <- mean(means_1)
interval_upper_1 <- quantile(means_1, c(.025, .975))[2]
interval_lower_1 <- quantile(means_1, c(.025, .975))[1]
quantiles_1[i,] <- data.frame(mean=interval_mean_1, upper=interval_upper_1,
                             lower=interval_lower_1, size=sample_size)
}

```

```

## [1] 1
## [1] 2
## [1] 3
## [1] 4
## [1] 5
## [1] 6
## [1] 7
## [1] 8
## [1] 9
## [1] 10

```

Produce a line chart that includes the bootstrapped mean and lower and upper percentile intervals for each group. Add appropriate labels and a legend. (6 points)

You may use base graphics or ggplot2. It should look similar to this (in base).

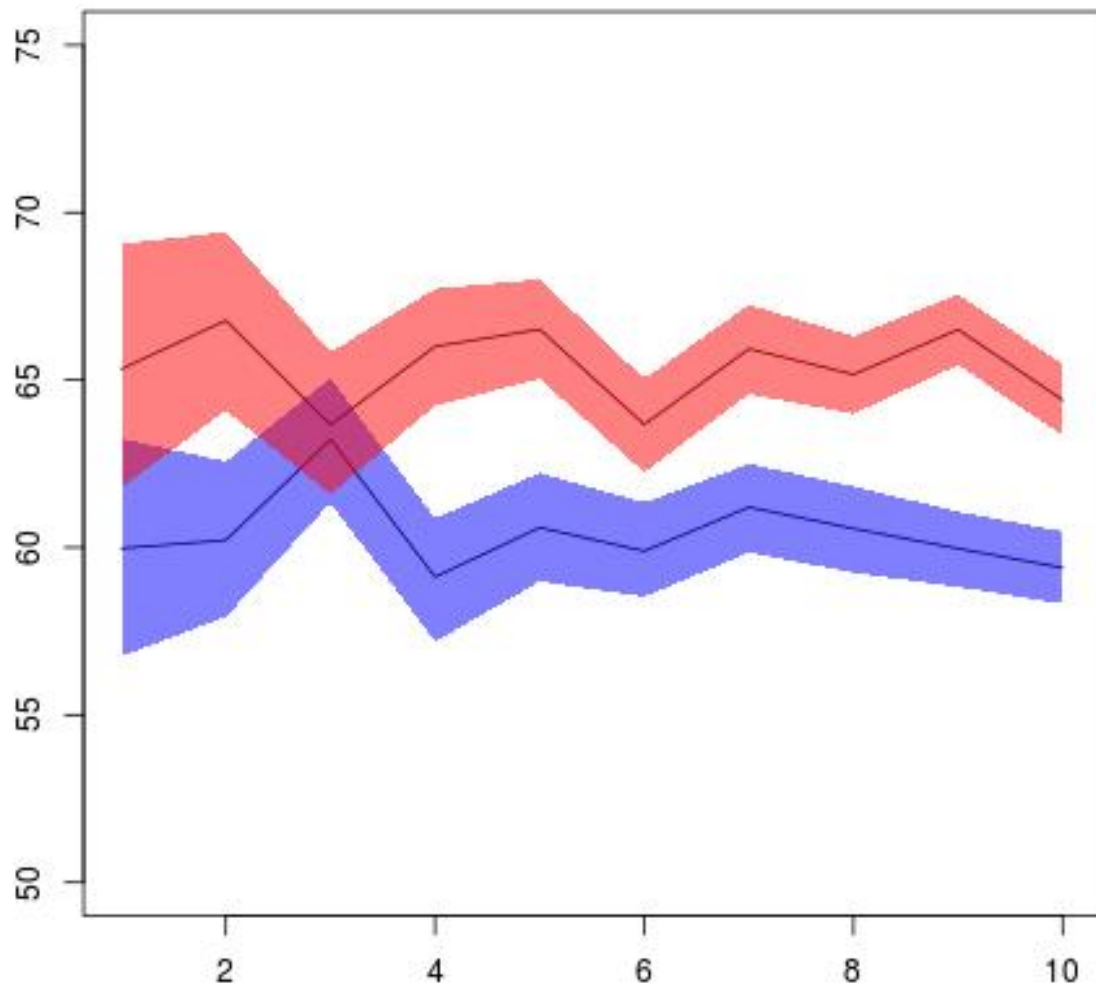


Figure 1: bp interval plot

```
library(ggplot2)
```

```
groups <- c("Group 0", "Group 1")  
colors <- c("Group 0" = "blue", "Group 1" = "red")
```

```
ggplot() +  
  geom_line(aes(x=quantiles_0$size,y=quantiles_0$mean, color='Group 0')) +  
  geom_line(aes(x=quantiles_1$size,y=quantiles_1$mean, color='Group 1')) +  
  ylab('Value') +
```

```

xlab('Interval') +
labs(x = "Sample Size",
     y = "Mean Value",
     color = "Treatment") +
scale_color_manual(values = colors) +
geom_ribbon(aes(x=quantiles_0$size, ymin=quantiles_0$lower, ymax=quantiles_0$upper), linetype=2, alpha=0.5) +
geom_ribbon(aes(x=quantiles_1$size, ymin=quantiles_1$lower, ymax=quantiles_1$upper), linetype=2, alpha=0.5) +
ggtitle("Bootstrap means by treatment group") +
geom_line()

```

