

24 September, 2020

```
## 0 healthy, 1 diseased
X <- "HHHDDDDHHDDHHHHHHDDHHHHHHHHHHDDHHHHDDDDHHHHDDHHHHDDHH" %>%
  stringr::str_split("") %>% unlist() %>% `==`("D") %>% as.integer()

n <- length(X)
n_0 <- n - sum(X)
n_1 <- sum(X)
```

```

R <- 1 + sum(X[2:n] != X[1:n-1])
R_est <- 1 + 2*n_0*n_1/n
R_var <- (2*n_0*n_1*(2*n_0*n_1 - n))/(n^2*(n - 1))
Z <- (R - R_est)/sqrt(R_var)
p <- pnorm(Z)
cbind(R, R_est, R_var, Z, p)

```

```

##      R      R_est      R_var      Z      p
## [1,] 15 20.65957 7.974767 -2.004125 0.02252834

```

With p-value of 0.0225 we reject H_0 at the 5% level. We conclude that values are not randomly ordered.

d)

```

set.seed(101)

calc_R <- function(input)
  1 + sum(input[2:length(input)] != input[1:length(input)-1])

sample_R <- function(iter, input)
  input %>% sample() %>% calc_R()

obs_R <- calc_R(X)

N <- 1e4
perm_R <- 1:N %>% sapply(sample_R, input = X)

p <- sum(perm_R < obs_R)/length(perm_R)
p

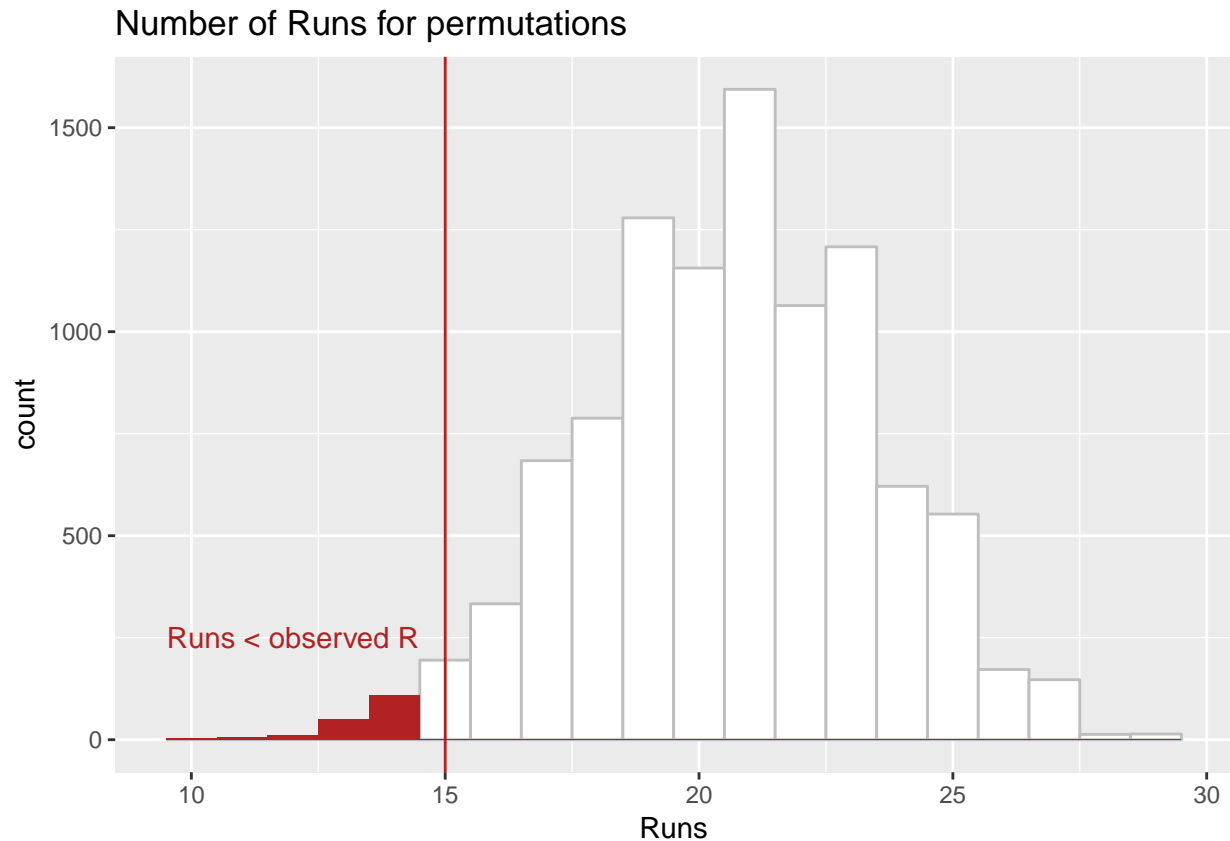
```

```
## [1] 0.0179
```

```

ggplot() +
  geom_histogram(data = data.frame(x = perm_R[perm_R >= obs_R]),
    aes(x = x), binwidth = 1, fill = "white", colour = "grey") +
  geom_histogram(data = data.frame(x = perm_R[perm_R < obs_R]),
    aes(x = x), binwidth = 1, fill = "firebrick") +
  geom_vline(xintercept = obs_R, colour = "firebrick") +
  ggtitle("Number of Runs for permutations") +
  xlab("Runs") +
  geom_text(aes(x = 12, y = 250), label = "Runs < observed R", colour = "firebrick")

```



Because the number of runs is not a continuous variable, there is some ambiguity around whether the p-value should be calculated comparing values $<$ or \leq or even using some half point. I have decided to use $<$ as it gives the test the highest power.

Given this we find a p-value of 0.0179. We reject H_0 at the 5% level. We conclude that values are not randomly ordered.

Question 2

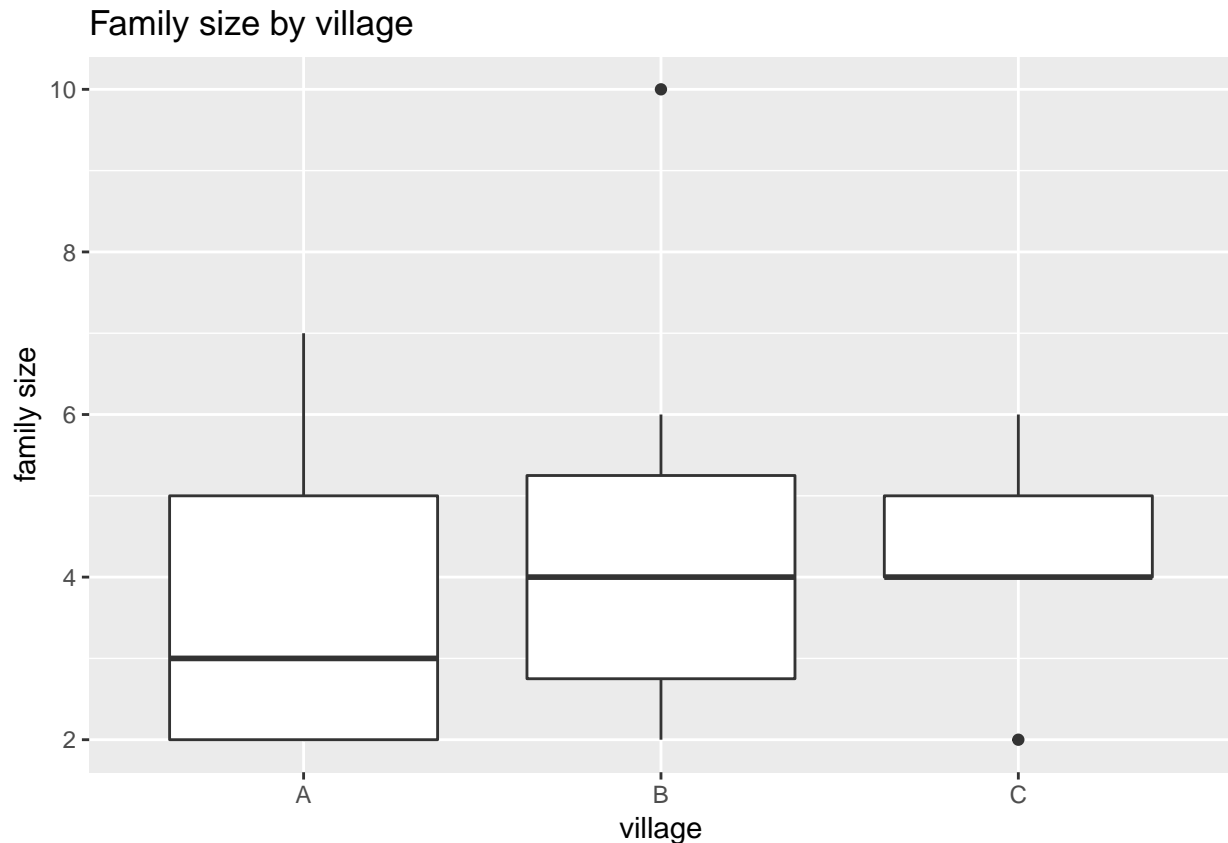
a)

We could use a Poisson distribution

b)

```
dataset <- data.frame(
  size = c(2, 3, 2, 7, 5, 5, 3, 2, 6, 10, 3, 2, 2, 5, 6, 4, 4, 5, 4, 4, 6, 5, 4, 2),
  village = rep(c("A", "B", "C"), c(9, 8, 7))
)

ggplot(dataset, aes(x = village, y = size)) +
  geom_boxplot() +
  ggtitle("Family size by village") +
  ylab("family size")
```



```
## village means
tapply(dataset$size, dataset$village, mean) %>% round(digits = 4)
```

```
##      A      B      C
## 3.8889 4.5000 4.2857
```

c)

```
set.seed(101)

calc_F <- function(input)
  lm(size ~ village, data = input) %>% anova() %>% .[["F value"]] %>% .[1]

sample_F <- function(iter, input)
  input %>% dplyr::mutate(village = sample(village)) %>% calc_F()

obs_F <- calc_F(dataset) %>% round(digits = 4)

N <- 1e4
perm_F <- 1:N %>% sapply(sample_F, input = dataset)

p <- sum(perm_F > obs_F)/length(perm_F)
p
```

```
## [1] 0.8275
```

```
ggplot() +
  geom_histogram(data = data.frame(x = perm_F),
    aes(x = x), bins = 50, fill = "white", colour = "grey") +
```

```
geom_vline(xintercept = obs_F, colour = "firebrick", linetype = "dashed") +
ggtitle("Histogram of the F Statistic") +
xlab("F statistic") +
geom_text(aes(x = 1, y = 150), label = "observed F", colour = "firebrick")
```

