

stat421_hw3

2022-10-14

Lecture 6

$$\textcircled{1} \quad \sum_i^n |y_i - \bar{y}| \Rightarrow 1 \text{ constraint : } \sum_{i=1}^n (y_i - \bar{y}) = 0 \\ \therefore df = n-1$$

1-pchisq(q=43.77, df=30)

```
## [1] 0.05003083
```

$$\textcircled{2a} \quad P\left(S^2 > 1.459\right) = P\left(\frac{(n-1)s^2}{\sigma^2} > \frac{30(1.459)}{1}\right) \quad s^2 = \frac{SS}{n-1} \Rightarrow X^2 = \frac{(n-1)s^2}{\sigma^2} \\ = P(X^2 > 43.77) \\ = 1 - P(X^2 < 43.77) \\ = 1 - pchisq(q=43.77, df=31-1) \\ = 0.05003083$$

$$\textcircled{2b} \quad s^2 = 1.6 \quad H_0: \sigma^2 \leq 1.0 \quad H_1: \sigma^2 > 1.0$$

From part a) we know that if the population variance (σ^2) is 1.0, the probability of the sample variance exceeding 1.459 is 0.05 which is very small. Therefore, with the population variance (σ^2) being 1.0, the sample variance (s^2) = 1.6 is even more unlikely than 0.05. This means that it is not plausible for the population variance (σ^2) to be 1.0.

1-pf(q=2.33, 20, 24)

```
## [1] 0.02484508
```

$$\textcircled{3} \quad V[Y|X=1] = V[Y|X=2]. \quad H_0: \sigma_1^2 = \sigma_2^2 \quad n_1 = 21 \quad n_2 = 25$$

$$\begin{aligned}
P(S_1^2 > 2.33 S_2^2) &= P\left(\frac{S_1^2}{S_2^2} > 2.33\right) \\
&= P\left(\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} > \frac{2.33}{\sigma_2^2}\right) \\
&= P(F > 2.33) \\
&= 1 - pf(q=2.33, 20, 24) \\
&= 0.02484508
\end{aligned}$$

//

(qt(0.05, df=20)*(sqrt(2/21)))+3

[1] 2.46774

4.	$H_0: \mu = 3$	$H_1: \mu < 3$	$\alpha = 0.05$	Population mean = μ
	$s^2 = 2 \Rightarrow n = 21 \Rightarrow df = 21 - 1 = 20$			

\textcircled{4a} We assume that $H_0: \mu = 3$ is true.

$$\begin{aligned}
\frac{\bar{Y}_{obs} - \mu_0}{S_{obs}/\sqrt{n}} &< t_\alpha \\
\bar{Y}_{obs} &< t_{0.05} \cdot \frac{S_{obs}}{\sqrt{n}} + \mu_0 \\
\bar{Y}_{obs} &< t_{0.05} \cdot \frac{\sqrt{2}}{\sqrt{21}} + 3 \\
\bar{Y}_{obs} &< -1.724718 \cdot \sqrt{\frac{2}{21}} + 3 \\
\bar{Y}_{obs} &< 2.46774
\end{aligned}$$

//

\textcircled{4b} Since the rejection region is $\bar{Y}_{obs} < 2.46774$ and the observed sample mean is 2.5, it is not in the rejection region. Therefore, we cannot reject H_0 .

$$\begin{aligned}
 \textcircled{qC} \quad P(\bar{Y} < \bar{y}_{obs} \mid H_0 = T) &\approx P\left(\frac{\bar{Y} - \mu}{S/\sqrt{n}} < \frac{\bar{Y}_{obs} - \mu}{S/\sqrt{n}} \mid \mu = 3\right) \\
 &= P\left(t < \frac{2.5 - 3}{\sqrt{2/21}}\right) \\
 &= 0.06043
 \end{aligned}$$

```
pt(q=(2.5-3)/(sqrt(2)/sqrt(21)),df=20)
```

```
## [1] 0.06042675
```

4d) We know that the p-value is 0.06042675 from part c, and that $\alpha = 0.05$. Since $p\text{-value} > 0.05$, we cannot reject the H_0 .

Lecture 7

Lecture 7

$$\textcircled{1} \quad t = \frac{(\bar{y} - \mu) / (\sigma / \sqrt{n})}{\sqrt{x^2 / (n-1)}} \sim t_{n-1} \quad X^2 = \frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1} \quad T = \frac{\frac{\bar{y} - \mu}{\sigma / \sqrt{n}} + \Delta}{\sqrt{x^2 / (n-1)}}$$

$$-T_{\frac{\alpha}{2}, n-1, \Delta} < \frac{\frac{\bar{y} - \mu}{\sigma / \sqrt{n}} + \Delta}{\sqrt{x^2 / (n-1)}} < T_{\frac{\alpha}{2}, n-1, \Delta}$$

$$-T_{\frac{\alpha}{2}, n-1, \Delta} < \frac{\frac{\bar{y} - \mu}{\sigma / \sqrt{n}} + \Delta}{\frac{s/\sigma}{s/\sigma}} < T_{\frac{\alpha}{2}, n-1, \Delta} \quad \frac{\bar{y} - \mu + \Delta \left(\frac{\sigma}{\sqrt{n}} \right)}{\frac{\sigma}{\sqrt{n}} \cdot \frac{s}{\sigma}}$$

$$\text{let } \Delta = -z_{\alpha/2} \sqrt{n} \quad z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$$

$$-T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} < \frac{\frac{\bar{y} - \mu}{\sigma / \sqrt{n}} - z_{\alpha/2} \sqrt{n}}{s/\sigma} < T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}}$$

$$-T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} < \frac{\frac{\bar{y} - \mu - z_{\alpha/2} \sigma}{\sigma / \sqrt{n}}}{s/\sigma} < T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}}$$

$$-T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} < \frac{\bar{y} - \mu - z_{\alpha/2} \sigma}{s/\sqrt{n}} < T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}}$$

$$-T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} \cdot \frac{s}{\sqrt{n}} < \bar{y} - (\mu + z_{\alpha/2} \sigma) < T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} \cdot \frac{s}{\sqrt{n}}$$

$$\bar{y} - T_{\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} \cdot \frac{s}{\sqrt{n}} < \mu + z_{\alpha/2} \sigma < \bar{y} + T_{1-\frac{\alpha}{2}, n-1, -z_{\alpha/2} \sqrt{n}} \cdot \frac{s}{\sqrt{n}}$$

$$\textcircled{2} \quad P\left(-t_{\frac{\alpha}{2}, n-1} < t < t_{\frac{\alpha}{2}, n-1}\right) = 1-\alpha \quad \bar{y}_{\text{obs}} \pm t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{\sqrt{n}} \text{ observed}$$

$$P\left(-t_{\frac{\alpha}{2}, n-1} < \frac{\bar{y} - \mu}{s/\sqrt{n}} < t_{\frac{\alpha}{2}, n-1}\right) = 1-\alpha \quad \bar{y} \pm t_{\frac{\alpha}{2}, n-1} \cdot \frac{s}{\sqrt{n}} \text{ random}$$

$$P\left(\bar{y}_{\text{obs}} - t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{\sqrt{n}} < \bar{y} < \bar{y}_{\text{obs}} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{\sqrt{n}}\right) = 1-\alpha$$

$$P\left(\frac{\bar{y}_{\text{obs}} - t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{\sqrt{n}} - \mu}{s/\sqrt{n}} < \frac{\bar{y} - \mu}{s/\sqrt{n}} < \frac{\bar{y}_{\text{obs}} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{\sqrt{n}} - \mu}{s/\sqrt{n}}\right) = 1-\alpha$$

$$P\left(\frac{\bar{y}_{\text{obs}} - \mu}{s/\sqrt{n}} - t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{s} < \frac{\bar{y} - \mu}{s/\sqrt{n}} < \frac{\bar{y}_{\text{obs}} - \mu}{s/\sqrt{n}} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{s_{\text{obs}}}{s}\right) = 1-\alpha$$

Since we know that $P\left(-t_{\frac{\alpha}{2}, n-1} < \frac{\bar{y} - \mu}{s/\sqrt{n}} < t_{\frac{\alpha}{2}, n-1}\right) = 1-\alpha$,

$$\text{Then, } \frac{\bar{y}_{\text{obs}} - \mu}{s/\sqrt{n}} = 0$$

```

# Part 3a
# H0: mu1 - mu2 = 0
# H1: mu1 - mu2 != 0
y1 = c( 65, 82, 81, 67, 57, 59, 66, 75, 82, 70 )
y2 = c( 64, 56, 71, 69, 83, 74, 59, 82, 65, 79 )
n1 = 10
n2 = 10
mean1 = mean(y1)
mean2 = mean(y2)

# REJECTION REGION
(t_alpha = qt(0.05/2, n1+n2-2))

## [1] -2.100922

(t_alpha1 = qt(1-(0.05/2), n1+n2-2))

## [1] 2.100922

(sp = sqrt((n1-1)*var(y1) + (n2-1)*var(y2))/(n1+n2-2))

## [1] 9.315459

(t_obs = (mean(y1) - mean(y2))/(sp*sqrt(1/n1 + 1/n2)))

## [1] 0.04800768

# P-VALUE
(p_val = 2*pt(-t_obs,n1+n2-2))

## [1] 0.9622388

# CONFIDENCE INTERVAL
(lower_bound = (mean1-mean2)+t_alpha*sp*(sqrt((1/n1)+(1/n2)))) 

## [1] -8.552441

(upper_bound = (mean1-mean2)-t_alpha*sp*(sqrt((1/n1)+(1/n2)))) 

## [1] 8.952441

```

The rejection region are <-2.1009 and >2.1009 . we know that the t observation is not inside the rejection region, which is why we cannot reject the H_0 that $(\mu_1-\mu_2) = 0$. Since $\alpha=0.05$ and $p\text{-value}=0.96$, we cannot reject the null hypothesis because $p\text{-value}>\alpha$. The confidence interval is $[-8.55244, 8.95244]$. Since 0 is inside the CI, we can choose to not reject the null hypothesis.

```

# Part 3b
# H0: sigma1 = sigma2
# H1: sigma1 != sigma2

# REJECTION REGION
(f_alpha = qf(0.05/2, n1-1, n2-1))

## [1] 0.2483859

(f_alpha1 = qf(1-(0.05/2), n1-1, n2-1))

## [1] 4.025994

(f_obs = var(y1)/var(y2))

## [1] 0.9782168

# P-VALUE
(p_val_f = 2*pf(f_obs,n1-1,n2-1)) # 2-sided

## [1] 0.9743665

# CONFIDENCE INTERVAL (why?)
(lower = (var(y1)/var(y2)) * qf(0.05/2, n2-1, n1-1))

## [1] 0.2429752

(upper = (var(y1)/var(y2)) * qf(1-0.05/2, n2-1, n1-1))

## [1] 3.938295

```

The rejection region are <0.248386 and >4.0259942 . We know that the f observation is not inside the rejection region, which is why we cannot reject the H0. Since alpha=0.05 and p-value=0.97, we cannot reject the null hypothesis because p-value>alpha. The confidence interval is [0.24298,3.9383]. Since the CI includes 1, we cannot reject H0.

Lecture 8

Lecture 8

$$\begin{aligned}
 \textcircled{1} \quad SS_{Treatment} &= n \sum_i^a (\bar{y}_{..} - \bar{y}_{..})^2 = \frac{1}{n} \sum_i^a y_{i..}^2 - \frac{1}{N} y_{..}^2 \\
 n \sum_i^a (\bar{y}_{i..} - \bar{y}_{..})^2 &= n \sum_i^a ((\bar{y}_{i..})^2 - 2\bar{y}_{i..}\bar{y}_{..} + (\bar{y}_{..})^2) \\
 &= n \sum_i^a \left(\left(\frac{1}{n} y_{i..} \right)^2 - \frac{2}{n} y_{..} \sum_i^a y_{i..} + \left(\frac{1}{n} y_{..} \right)^2 \right) \\
 &= \frac{1}{n} \sum_i^a (y_{i..})^2 - \frac{2}{an} y_{..} \sum_i^a y_{i..} + \frac{1}{an} (y_{..})^2 \\
 &= \frac{1}{n} \sum_i^a (y_{i..})^2 - \frac{2}{an} y_{..} \sum_i^a y_{i..} + \frac{1}{an} (y_{..})^2 \\
 &= \frac{1}{n} \sum_i^a (y_{i..})^2 - \frac{2}{N} y_{..} \sum_i^a y_{i..} + \frac{1}{N} (y_{..})^2 \\
 &= \frac{1}{n} \sum_i^a (y_{i..})^2 - \frac{1}{N} (y_{..})^2 \\
 &\quad //
 \end{aligned}$$

$$\textcircled{2a} \quad \sum_i^a \sum_j^b (y_{ij} - \mu)^2 \Rightarrow 0 \text{ constraints} \\
 \text{Therefore, } df = an$$

$$\textcircled{2b} \quad \sum_i^a \sum_j^b (y_{ij} - \bar{y}_{..})^2 \Rightarrow 1 \text{ constraint: } \sum_i^a \sum_j^b (y_{ij} - \bar{y}_{..}) = 0 \\
 \text{Therefore, } df = an - 1$$

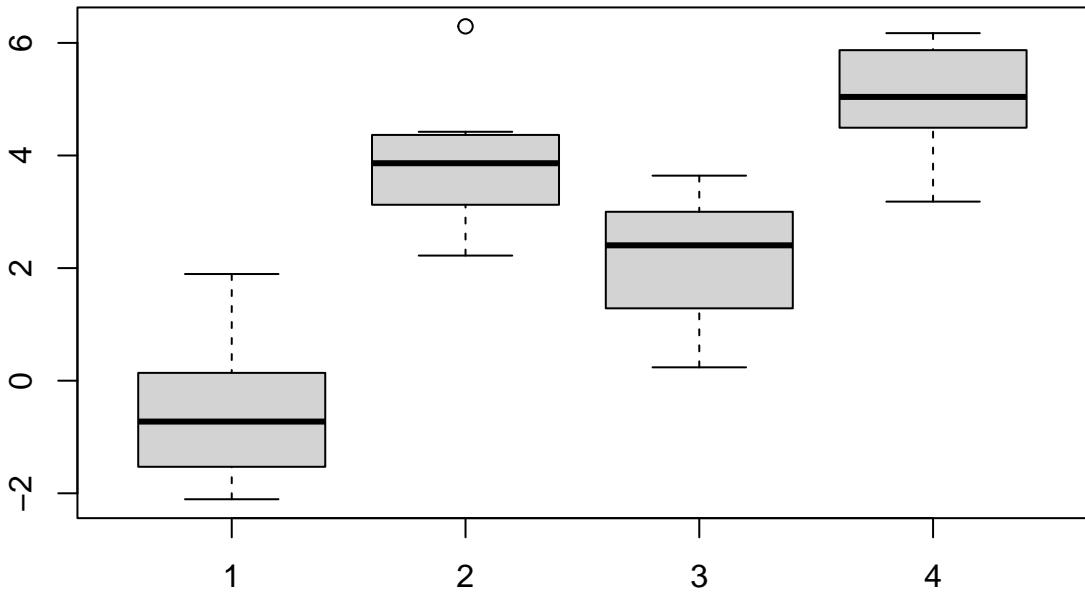
$$\textcircled{2c} \quad \sum_i^a (\bar{y}_{i..} - \bar{y}_{..})^2 \Rightarrow 1 \text{ constraint: } \sum_i^a (\bar{y}_{i..} - \bar{y}_{..}) = 0 \\
 \text{Therefore, } df = a - 1$$

```

y = matrix(nrow=4, ncol=10)
y[1,] = c(-2.10552316, 1.89491371, -1.52919682, -0.99265143, -0.45911960, 1.09271028, -1.54680778, 0.11
y[2,] = c(4.25667943, 4.36518096, 4.42108835, 3.77229146, 2.22264903, 3.95354759, 6.29377745, 3.5850108
y[3,] = c(3.64209745, 2.76932242, 1.46001019, 0.23739519, 0.27629510, 2.83897173, 2.99999590, 3.5465783
y[4,] = c(3.18088593, 5.44976665, 5.87116946, 4.01275036, 6.00826692, 5.19220036, 6.17338313, 4.8884607

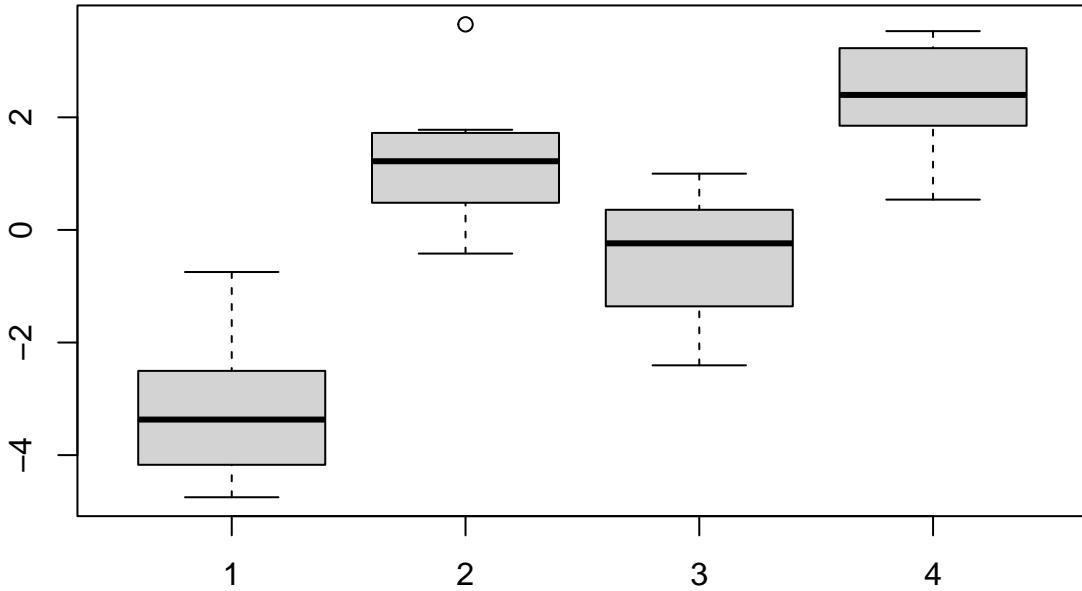
# part a
boxplot(y[1,], y[2,], y[3,], y[4,])

```



We can see from the boxplots that x has some effect on y. We can see that y4 has a higher distribution than the other y's, while y1 has the smallest values. We can also see that there is a difference between the y2 and y3.

```
# part b
grand_mean = mean(y)
y_measure = y-grand_mean
boxplot(y_measure[1,],y_measure[2,],y_measure[3,],y_measure[4,])
```



Subtracting the mean does not change anything, but the overall boxplot are just shifted downwards. It still looks like that x has an effect on y.

```
# part c
y_bar = c(mean(y[1,]),mean(y[2,]),mean(y[3,]),mean(y[4,]))
effects = c(y_bar[1]-grand_mean, y_bar[2]-grand_mean, y_bar[3]-grand_mean, y_bar[4]-grand_mean)

## [1] -3.0961875  1.2613624 -0.5329403  2.3677655

# part d
(se = c(sd(y[1,])/sqrt(10),sd(y[2,])/sqrt(10),sd(y[3,])/sqrt(10),sd(y[4,])/sqrt(10)))

## [1] 0.3986627 0.3433781 0.3964595 0.2976147
```

8-3e) -3.096 ± 0.399 1.261 ± 0.343 -0.533 ± 0.396 2.368 ± 0.298 it looks like the third effect might be zero because zero is inside the interval of mean+se

```
# part f
n = 10
a = 4
var = c(var(y[1,]),var(y[2,]),var(y[3,]),var(y[4,]))
SS_treatment = n*sum(effects^2)
SS_error = (n-1)*(sum(var))
MS_treatment = SS_treatment/(a-1)
MS_error = SS_error/((n*a)-a)
(F_ratio = MS_treatment / MS_error)
```

```

## [1] 43.54614

(F_alpha = qf(0.05,a-1,a*n-a,lower.tail=FALSE))

## [1] 2.866266

```

The rejection region is >2.866 . Since, the F_ratio is inside the rejection region we choose to reject H0 in favor of H1. This means that x has an effect on y

```

# part g
x_vec=rep(1:4,10)
y_vec = as.numeric(y) # visually confirm y.vector is the correct y in vector form.
summary.aov(lm(y_vec ~ as.factor(x_vec)))

##           Df Sum Sq Mean Sq F value    Pr(>F)
## as.factor(x_vec)  3 170.68   56.89   43.55 4.58e-12 ***
## Residuals        36  47.03    1.31
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

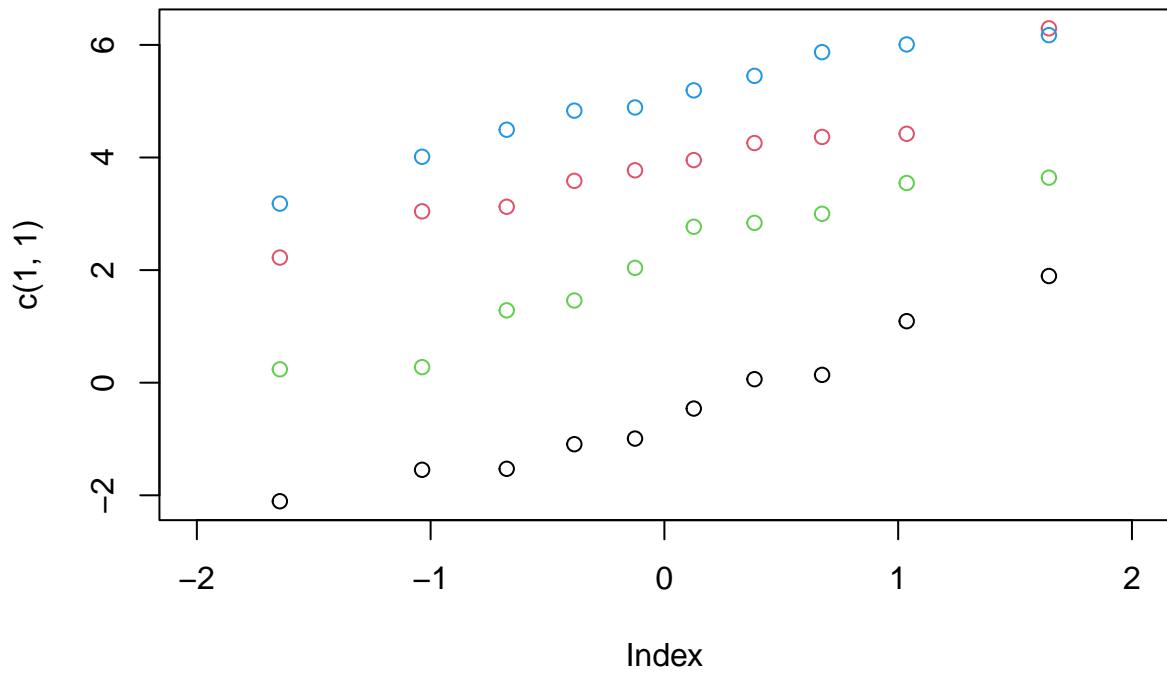
```

The p-value is 4.58e-12. Since 4.58e-12 is smaller than alpha, then we choose to reject the H0 in favor of H1

```

# part h
plot(c(1,1),cex=0, xlim=c(-2,2), ylim=range(y))
for (i in 1:a) {
  x = y[i,]
  n = length(x)
  probs = seq(0.5/n, 1-0.5/n, length = n)
  q = qnorm(probs,0,1)
  points(q,sort(x),col=i)
}

```



We can see that the qq plots are pretty much straight, and this means that they most likely follow the normal distribution. Their slopes are also pretty similar which means that their variances are also similar.