

Question 1. First, read the data with the following code:

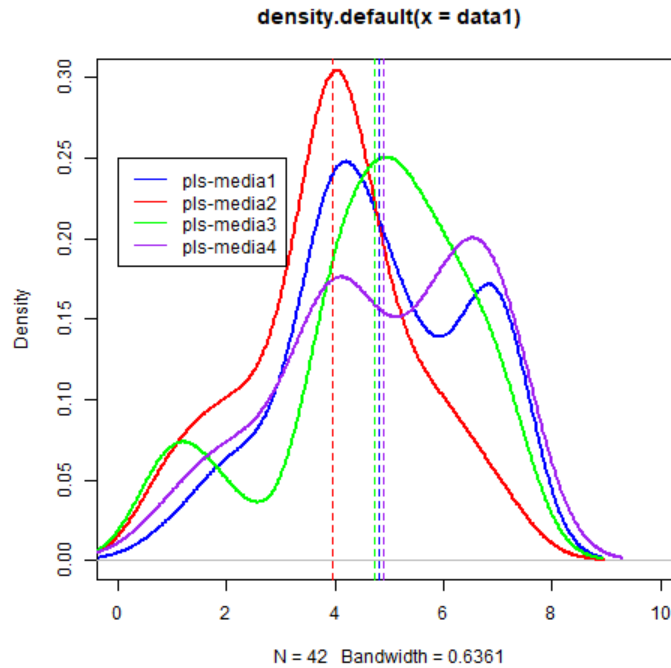
```
1 library(reshape2)
2 library(tidyverse)
3
4 # Question 1
5 data1 <- read_csv('pls-media1.csv')$INTEND.0 # read the data
6 data2 <- read_csv('pls-media2.csv')$INTEND.0
7 data3 <- read_csv('pls-media3.csv')$INTEND.0
8 data4 <- read_csv('pls-media4.csv')$INTEND.0
```

(a) Using the `mean` function directly,

```
1 # Question 1 (a)
2 data1_mean <- mean(data1)
3 data2_mean <- mean(data2)
4 data3_mean <- mean(data3)
5 data4_mean <- mean(data4)
```

(b) The density plot and the vertical lines (marking the means) with same colors are from the same type of data.

```
1 # Question 1 (b)
2 png(filename = "1b.png")
3 plot(density(data1), col="blue", lwd=2, xlim=c(0, 10), ylim=c(0,0.3))
4 lines(density(data2), col="red", lwd=2)
5 lines(density(data3), col="green", lwd=2)
6 lines(density(data4), col="purple", lwd=2)
7 abline(v=data1_mean, lty="dashed", col="blue")
8 abline(v=data2_mean, lty="dashed", col="red")
9 abline(v=data3_mean, lty="dashed", col="green")
10 abline(v=data4_mean, lty="dashed", col="purple")
11 legend(0, 0.25,
12       legend = c("pls-media1", "pls-media2", "pls-media3", "pls-media4"),
13       col = c("blue", "red", "green", "purple"),
14       lty = 1:1, cex = 1)
15 dev.off()
```



(c) Yes, the mean of data from `pls-media2.csv` (red lines) looks like having a different mean. More specifically, the mean of the other types are close. ■

Question 2. (a) Let the four types of data have means μ_1, \dots, μ_4 . The null and alternative hypotheses of ANOVA are $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$; $H_1: \mu_1 \neq \mu_2, \mu_1 \neq \mu_3, \mu_1 \neq \mu_4, \mu_2 \neq \mu_3, \mu_2 \neq \mu_4, \text{ and } \mu_3 \neq \mu_4$.

(b) I conduct the one-way ANOVA by the following hand-carved code:

```

1 # Question 2 (b)
2 # Compute MSTR
3 data <- list(data1, data2, data3, data4)
4 sstr <- sum(sapply(data, length)*(sapply(data, mean) - mean(sapply(data, mean)))^2)
5 df_mstr <- 4-1
6 mstr <- sstr/df_mstr
7 # Compute MSE
8 sse <- sum((sapply(data, length)-1)*sapply(data, var))
9 df_mse <- sum(sapply(data, length)) - 4
10 mse <- sse/df_mse
11 # Compute F-value and p-value
12 f_value <- mstr/mse # F-value
13 qf(p=0.95, df1=df_mstr, df2=df_mse) # The cutoff value of F-value
14 p_value <- pf(f_value, df_mstr, df_mse, lower.tail=FALSE) # The p-value

```

(i) MSTR= 7.53239, MSE= 2.869151, and $F = 2.625303$.

(ii) The p -value of F is 0.05230686. Since $\alpha = 0.05$, The F -value is not significant. Hence we cannot reject null hypothesis H_0 .

(c) To use the built-in function in R, we need to apply the `melt` function first to transform the 4 types of data into a dataframe.

```

1 # Question 2 (c)
2 names(data) <- c("pls-media1", "pls-media2", "pls-media3", "pls-media4")
3 data_aov <- melt(data,
4                 id.vars = NULL,
5                 variable.name = "type",
6                 value.name = "intend")
7 anova_model <- aov(data_aov$intend ~ factor(data_aov$L1))
8 summary(anova_model)

```

It gives $F = 2.617$. The p -value of F is 0.0529. Since $\alpha = 0.05$, The F -value is not significant. Hence we cannot reject null hypothesis H_0 .

(d) To conduct a post-hoc Tukey test,

```
1 # Question 2 (d)
2 TukeyHSD(anova_model, conf.level = 0.05)
```

Maybe in this case, a screenshot of results is more appropriate. From the results, we may conclude that type 2 data (from `pls-media2.csv`) have a different mean from other types.

	diff	lwr	upr	p adj
pls-media2-pls-media1	-0.86215539	-1.06562977	-0.6586810	0.1085727
pls-media3-pls-media1	-0.08452381	-0.28530983	0.1162622	0.9959223
pls-media4-pls-media1	0.08178054	-0.11218249	0.2757436	0.9959032
pls-media3-pls-media2	0.77763158	0.57175512	0.9835080	0.1825044
pls-media4-pls-media2	0.94393593	0.74470805	1.1431638	0.0573229
pls-media4-pls-media3	0.16630435	-0.03017708	0.3627858	0.9687417

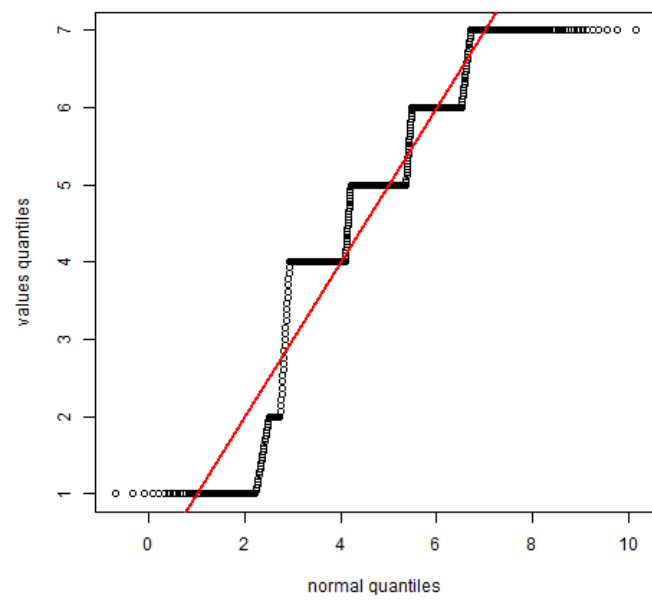
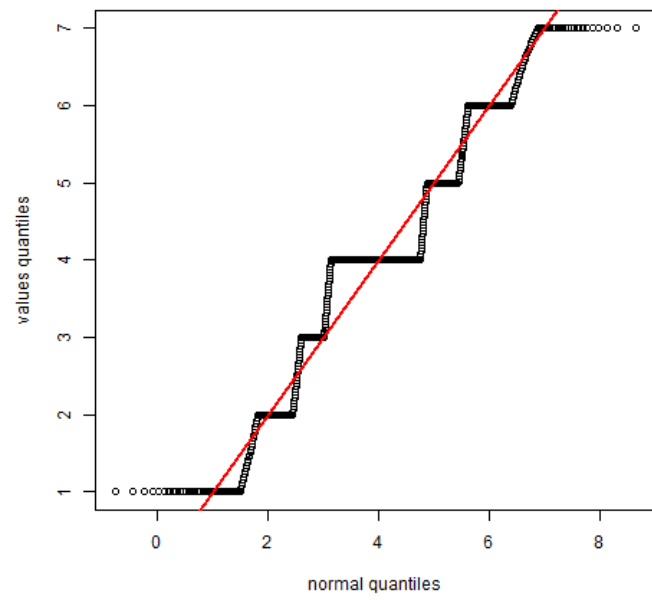
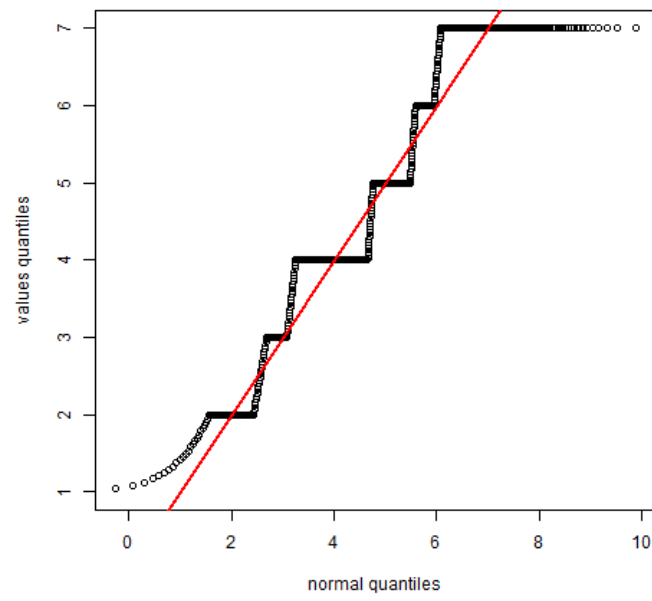
(e) Let's review the three conditions of ANOVA[1]:

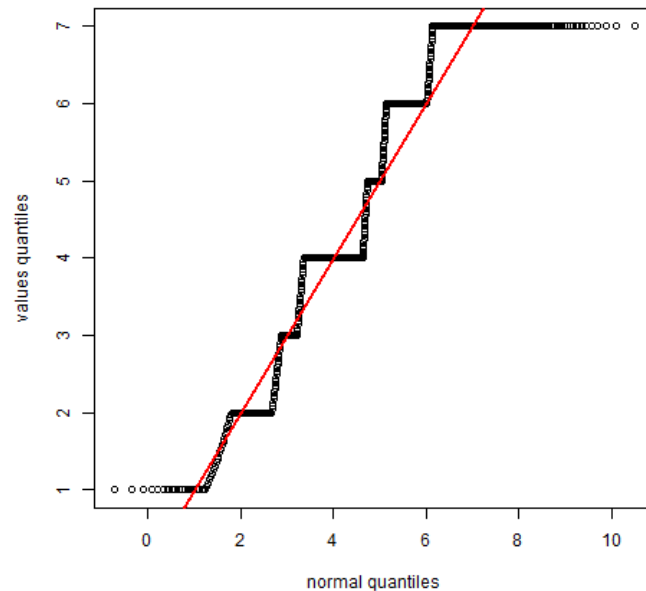
(i) **Each treatment/population's response variable is normally distributed**

The Q-Q plot based method we used in last homework[2] helps.

```
1 norm_qq_plot <- function(values)
2 {
3   probs1000 <- seq(0, 1, 0.001)
4   q_vals <- quantile(values, probs=probs1000)
5   q_norm <- qnorm(probs1000, mean=mean(values), sd=sd(values))
6   plot(q_norm, q_vals, xlab="normal quantiles", ylab="values quantiles")
7   abline(a=0, b=1, col="red", lwd=2)
8 }
9
10 png(filename = "2e-type1.png")
11 norm_qq_plot(data1)
12 dev.off()
13
14 png(filename = "2e-type2.png")
15 norm_qq_plot(data2)
16 dev.off()
17
18 png(filename = "2e-type3.png")
19 norm_qq_plot(data3)
20 dev.off()
21
22 png(filename = "2e-type4.png")
23 norm_qq_plot(data4)
24 dev.off()
```

Note that it is really unlikely that those data being normally distributed since they are discrete, whose values take ranges from 1 to 7. So we may predict that the Q-Q plot may look like step functions.





(ii) The variance (s^2) of the response variables is the same for all treatments/populations

```

1 # Question 2 (e)
2 data1_sd <- sd(data1)
3 data2_sd <- sd(data2)
4 data3_sd <- sd(data3)
5 data4_sd <- sd(data4)
6 data_sd <- sd(data_aov$intend)

```

Compute the s^2 using `sd` functions, I found that the s^2 for all four types are 1.641506, 1.52364, 1.753933, 1.816324, and s^2 among all data is 1.718571. We may run a F-test (or something similar) to test whether they are significantly different.

(iii) **The observations are independent: the response variables are not related between groups**

This is hard to test, I assume (with a faith of belief that) this is true. ■

Question 3. (a) Let \bar{R}_i be the mean of sum of ranks, $i = 1, 2, 3, 4$. The null and alternative hypotheses are $H_0: \bar{R}_1 = \bar{R}_2 = \bar{R}_3 = \bar{R}_4$; $H_1: \bar{R}_i \neq \bar{R}_j, \forall i \neq j, i, j = 1, 2, 3, 4$

(b) I compute Kruskal Wallis H by the following hand-carved code:

```

1 # Question 3 (b)
2 intend_rank <- rank(data_aov$intend)
3 data_aov$rank <- intend_rank # add as a new column to data_aov
4 group_ranks <- split(data_aov, data_aov$L1)
5 rank_sum_1 <- sum(group_ranks$`pls-media1`$rank)
6 rank_sum_2 <- sum(group_ranks$`pls-media2`$rank)
7 rank_sum_3 <- sum(group_ranks$`pls-media3`$rank)
8 rank_sum_4 <- sum(group_ranks$`pls-media4`$rank)
9 R <- c(rank_sum_1^2/length(data1),
10        rank_sum_2^2/length(data2),
11        rank_sum_3^2/length(data3),
12        rank_sum_4^2/length(data4))
13 H = (12/(length(data)*(length(data)+1))) * sum(R) - 3*(length(data)+1)
14 kw_p <- 1 - pchisq(H, df=4-1)

```

(i) The code gives $H = 706140.1$.

- (ii) The p -value is 0, we can reject the null hypothesis H_0 . Hence The rank sum of all types are all different.
 (c) The built in function of Kruskal Wallis test:

```
1 # Question 3 (c)
2 kruskal.test(intend~L1, data=data_aov)
```

It gives p -value= 0.03166. Since this $p < 0.05$, we can reject the null hypothesis H_0 and having the same conclusion of the hypotheses as in (b).

- (d) The built in function of the post-hoc Dunn test:

```
1 # Question 3 (d)
2 dunnTest(intend~L1, data=data_aov, method = "bonferroni")
```

	Comparison	Z	P.unadj	P.adj
1	pls-media1 - pls-media2	2.30087819	0.021398517	0.12839110
2	pls-media1 - pls-media3	-0.09233644	0.926430736	1.00000000
3	pls-media2 - pls-media3	-2.36408588	0.018074622	0.10844773
4	pls-media1 - pls-media4	-0.31452459	0.753122646	1.00000000
5	pls-media2 - pls-media4	-2.65613380	0.007904225	0.04742535
6	pls-media3 - pls-media4	-0.21613379	0.828883460	1.00000000

We may conclude that type 2 data (from `pls-media2.csv`) have a different rank sum from other types.



References

- [1] Handouts of week 7.
- [2] A Q-Q Plot Dissection Kit, <https://seankross.com/2016/02/29/A-Q-Q-Plot-Dissection-Kit.html>