# Statistical Inference Project

*JP Dunlap*

*July 27, 2017*

## Part 2: Basic Inferential Data Analysis

In 1947 E.W. Crampton published "The growth of the odontoblast of the incisor teeth as a criterion of vitamin C intake of the guinea pig" in the Journal of Nutrition. (see Crampton, E. W. (1947) The growth of the odontoblast of the incisor teeth as a criterion of vitamin C intake of the guinea pig. The Journal of Nutrition 33(5): 491–504. http://jn.nutrition.org/content/33/5/491.full.pdf (http://jn.nutrition.org/content/33/5/491.full.pdf)). In his paper he discussed the impact of two different forms of vitamin C, administered in three does levels, on tooth growth in Guinea Pigs.

The research design utilized 60 Guinea Pigs each of which was administered Vitamin C either in the form of orange juice (coded OJ), or ascorbic acid (coded as VC). Withing the two groups the animals were further divided into three group based on dosage (0.5, 1, or 2 mg/day). The length of the odontoblasts were then measured (presumably in microns).

This analysis will attempt to determine the impact of the various delivery methods/dosages on the length of the odontoblasts.

### Research Assumptions

The analysis and conclusions for this research are based on the following assumptions:

1. All animals in the study come from the same population of Guinea Pigs
2. Animals are randomly assigned to each of the six possible treatment groups
3. The growth or odondoblasts is normally distributed with an unknown mu and sigma

### Data Mangement

As the ToothGrowth database is a part of the R package installation it needs only to be loaded in order for analysis to be conducted. A frequency graph of the odontoblast length for each possible combination is shown.

```
data(ToothGrowth)
```

### Exploratory Data Analysis - Examining Confidence Intervals

The first step is to examine the data grouped by the two effect variables, dose and delivery method. A histogram of each combination is shown below. Layering the upper and lower confidence intervals may provide some insight.

```
## Get Mean and stdev for each combination, store id dataframe ci.in

ci.mn <- summarise(group_by(ToothGrowth, dose, supp), mean(len))
ci.sd <- summarise(group_by(ToothGrowth, dose, supp), sd(len))
ci.in <- merge(ci.mn,ci.sd, by = c("dose","supp"))

colnames(ci.in) <- c("dose","supp","Mean","StDev")


## Now get the t-distribution based CI for all 6

ci.in <- mutate(ci.in, LowerBound = ci.in$Mean - (qt(.975,9) * (ci.in$StDev/sqrt(10))))
ci.in <- mutate(ci.in, UpperBound = ci.in$Mean + (qt(.975,9) * (ci.in$StDev/sqrt(10))))

ToothGrowth <- left_join(ToothGrowth, ci.in, by = c("dose","supp"))
names(ToothGrowth)[3] <- "Dose"
names(ci.in)[1] <- "Dose"
names(ToothGrowth)[2] <- "Delivery.Meth"
names(ci.in)[2] <- "Delivery.Meth"

## Repeat the plots

g <- ggplot(ToothGrowth, aes(len)) + geom_histogram(bins = 10, col = "brown", fill  = "orange")
g <- g + geom_vline(aes(xintercept = LowerBound)) + geom_vline(aes(xintercept = UpperBound))
g <- g + labs(x = "Length of Odontoblasts in Microns", title = "Length of Odontoblasts by Delivery Method and Dose", y = "Count by Length")
g <- g + facet_grid(Delivery.Meth ~ Dose) + theme_grey()
g
```
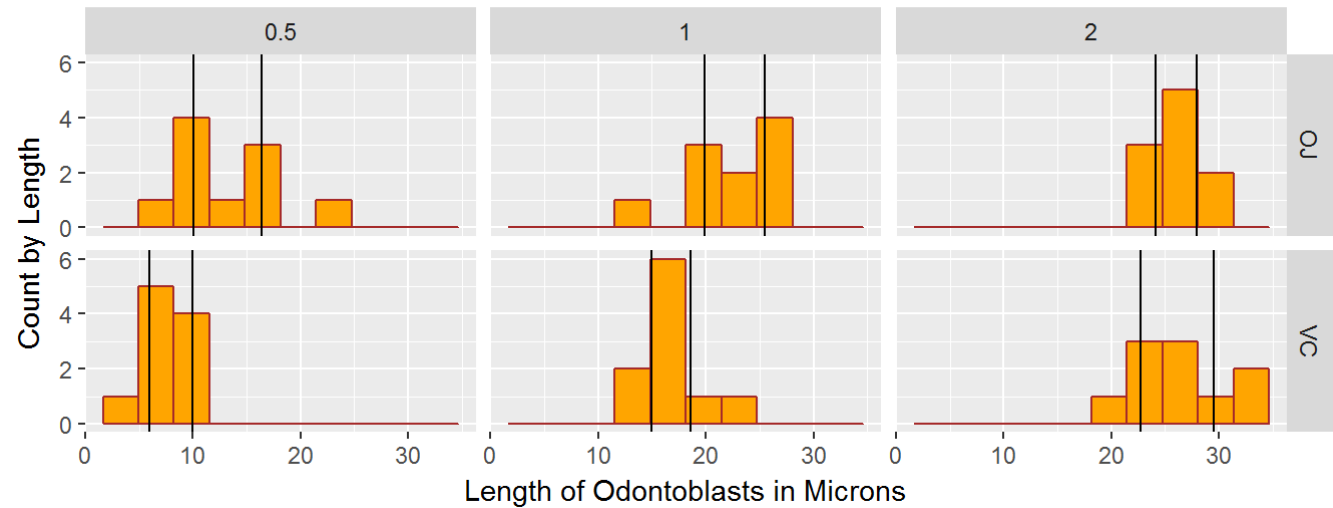
## Length of Odontoblasts by Delivery Method and Dose



Examining the histograms with the 95% confidence interval shown suggests that there may be a significant effect from the delivery method at the lower two doses (0.5, and 1,0 mg/day). The method does not appear to be a significant effect at the higher does of 2.0mg/day.

Examining the actual lower and upper confidence intervals will provide additional insight.

```
byDose <- arrange(ci.in, Dose, Delivery.Meth)
byMethod <- arrange(ci.in, Delivery.Meth, Dose)

f1 <- filter (ci.in[,c(2,5,6)], ci.in$Dose == 0.5)
f2 <- filter (ci.in[,c(2,5,6)], ci.in$Dose == 1)
f3 <- filter (ci.in[,c(2,5,6)], ci.in$Dose == 2)
e2 <- filter (ci.in[,c(1,5,6)], ci.in$Delivery.Meth == "OJ")
e1 <- filter (ci.in[,c(1,5,6)], ci.in$Delivery.Meth == "VC")
```

### By Dosing Levels

Looking at a comparison of the CI bounds by dosing level may help to better understand the potential differences.

**Dose at 0.50 mg/day** The confidence interval for the two delivery methods shows no overlap between the two. This seems to suggest that at the dose level of 0.50 mg/day, the delivery method does make a difference.

**Dose at 1.0 mg/day** Similarly the confidence interval for the two delivery methods shows no overlap between the two. This seems to suggest that at the dose level of 1.0 mg/day, the delivery method does make a difference.

**Dose at 2.0 mg/day** By contrast, The confidence interval for the two delivery methods shows a complete overlap between the two. This suggest that at the dose level of 2.0 mg/day, the delivery method no longer seems make a difference.

| Delivery.Meth | LowerBound | UpperBound |
|---|---|---|
| **Dose 0.5 mg/day** | | |
| OJ | 10.039717 | 16.420283 |
| VC | 6.015176 | 9.944824 |
| **Dose 1.0 mg/day** | | |
| OJ | 19.902273 | 25.497727 |
| VC | 14.970657 | 18.569343 |
| **Dose 2.0 mg/day** | | |
| OJ | 24.160686 | 27.959314 |
| VC | 22.707910 | 29.572090 |

### By Delivery Method

Looking at a comparison of the CI bounds by dosing level may help to better understand the potential differences.

**Delivery Method - Orange Juice** The confidence interval for the three dose amounts shows no overlap between the lower two level, and a small overlap between the upper two. This seems to suggest that effect of the minimal dose of orange juice is less than the middle and higher dose, but that the effect between the middle and higher dose are less clear.

**Delivery Method - Ascorbic Acid** For ascorbic acid, the confidence interval for the three dose level shows no overlap at all. This seems to suggest that for ascorbic acid the size of the dose has a significant effect on growth.

| Dose | LowerBound | UpperBound |
|---|---|---|
| **Orange Juice** | | |
| 0.5 | 10.039717 | 16.420283 |
| 1.0 | 19.902273 | 25.497727 |
| 2.0 | 24.160686 | 27.959314 |
| **Ascorbic Acid** | | |
| 0.5 | 6.015176 | 9.944824 |
| 1.0 | 14.970657 | 18.569343 |
| 2.0 | 22.707910 | 29.572090 |

### Conducting Hypotheses Testing Using Students T-Test

T-tests are a traditional small sample statistic to examine the difference between two samples, in this case, taken from the same population. In the case of this data set, the comparison is between the various combinations of dose rate and delivery method. Note that it would also be appropriate to use other statistical tests such as one-way and two-way analysis of variance to explore these relationships, these techniques are outside of the scope of this project.

```
## Code to create all of the possible t.tests. This is stored in the source file, t-test.R and is called here. A copy of this
  code is in the appendix.

fileLoc  <- "C:/Users/Jeff/Google Drive/coursera/Data Science Course/Course Materials/6-Statistical Inference/Week 4/Exponenti
al-Distribution-Simulation/t-test.R"
source(fileLoc )
```

Hypotheses Test Group 1 - Comparing Delivery Method within Dose Amounts

H_0: mean_OJ EQ mean_VC - for each dosing amount

H-a: mean_OJ NE mean_vc - for each dosing amount

The first set of analyses include the comparison of delivery method by dosing amount. The mean length of Odontoblasts are compared with dose held constant and delivery method compared. The following table summarizes these results.

From these data it appears that the two different delivery methods due vary in effectiveness at lower and medium doses (0.5 and 1.0 mg/day) with Orange Juice appearing more effective. At the higher dose of 2.0 mg/day, the effect is not found.

Both the Bonferroni correction to correct for the Family-Wise Error Rate, and the Benjamini-Hochberg Correction to correct for False Discovery Rate, were applied to the results. While both correction do change the p-value, there are no changes to the impact of the results. Given that there are only three tests conducted, this is not surprising.

| Test_Conducted | tvalue | pvalue | Significant | Bonferroni | BH |
|---|---|---|---|---|---|
| H_0: mean_OJ = mean_VC - Dose = 0.5 mg/day | 3.170 | 0.0053 | TRUE | 0.015900 | 0.007950 |
| H_0: mean_OJ = mean_VC - Dose = 1.0 mg/day | 4.033 | 0.000781 | TRUE | 0.002343 | 0.002343 |
| H_0: mean_OJ = mean_VC - Dose = 2.0 mg/day | -0.0461 | 0.964 | FALSE | 1.000000 | 0.964000 |

As a result we are able to reject the null hypothesis at dosing levels 0.5 and 1.0 mg/day. The null hypothesis is not rejected for the dosing level at 2.0 mg/day.

Hypotheses Test Group 2 - Comparing Dosing Combinations within Delivery Method

H_0: mean_dosea EQ mean_doseb - for each possbile pair of dosing amount holding delivery method constant

H-a: mean_dosea NE mean_doseb - for each possbile pair of dosing amount holding delivery method constant

The second set of hypotheses include the comparison of dosage pairs by delivery method. The mean length of Odontoblasts are compared with delivery method held constant and pairs of dosage amounts compared (i.e., 0.5 to 1.0, 0.5 to 2.0, and 1.0 to 2.0). In all there are six possible sets of hypotheses.

The following table summarizes these results.

From these data it appears that the two different delivery methods due vary in effectiveness at lower and medium doses (0.5 and 1.0 mg/day) with Orange Juice appearing more effective. At the higher dose of 2.0 mg/day, the effect is not found.

Again, the Bonferroni correction to correct for the Family-Wise Error Rate, and the Benjamini-Hochberg Correction to correct for False Discovery Rate, were applied to the results. Examining the False Discovery Rate correction, there are no important changes. However, when correcting the Family-Wise Error Rate, one of the six tests, Orange Juice compared at 1.0 to 2.0 did fall above the alpha = 0.05 cutoff level.

It is reasonable to reject the null hypothesis in each of the six possible pairings, based on the individual t-tests, suggesting that dose rates do make a difference. However when applying the correction for Family-Wise Error Rate, it would be reasonable to fail to reject the null hypothesis comparing 1.0 to 2.0 mg/day dosing via Orange Juice delivery.

| Test_Conducted | tvalue | pvalue | Significant | Bonferroni | BH |
|---|---|---|---|---|---|
| H_0: mean_dosea = mean_doseb - 0.5 to 1.0 mg/day, with Orange Juice delivery | -5.049 | 0.0000836 | TRUE | 0.0005016000 | 0.0001003200 |
| H_0: mean_dosea = mean_doseb - 0.5 to 2.0 mg/day, with Orange Juice delivery | -7.817 | 0.00000034 | TRUE | 0.0000020400 | 0.0000010200 |
| H_0: mean_dosea = mean_doseb - 1.0 to 2.0 mg/day, with Orange Juice delivery | -2.248 | 0.0374 | TRUE | 0.2244000000 | 0.0374000000 |
| H_0: mean_dosea = mean_doseb - 0.5 to 1.0 mg/day, with Ascorbic Acid delivery | -7.463 | 0.000000649 | TRUE | 0.0000038940 | 0.0000012980 |
| H_0: mean_dosea = mean_doseb - 0.5 to 2.0 mg/day, with Ascorbic Acid delivery | -10.388 | 0.00000000496 | TRUE | 0.0000000298 | 0.0000000298 |
| H_0: mean_dosea = mean_doseb - 1.0 to 2.0 mg/day, with Ascorbic Acid delivery | -5.470 | 0.000034 | TRUE | 0.0002040000 | 0.0000510000 |

## Part 2 Conclusion

The of the confidence interval analysis and the hypotheses testing (t-test) analyses are essentially the same. As a result, it is reasonable to conclude that increasing dosage of both ascorbic acid and orange juice above a minimum amount is effective in increasing odontoblast length. At the highest dose, there is no indicated difference in the delivery method.

When comparing the various dosage pairs, all appear to have a significant impact, regardless of delivery method. However the Bonferroni correction for Family-Wise Error Rate does call one effect into question.

# Appendix to Report 2

```
## This lengthy block of code calculates all of the t-test required for this analysis. It is omitted from the
## report file because of its length. It is called by Part 2.Rmd.


ii = 0
for (i in c(0.5, 1.0, 2.0)){
        ii = ii+1
        jj = 0
        for (j in c("OJ", "VC")){
                jj = jj + 1
                assign(paste0("df",ii,jj), filter(ToothGrowth, Delivery.Meth == j & Dose == i))


        }
}


## create 3 possible combinations of dose by treatment
df1. <- rbind(df11,df12) ## treatment by dose at 0.5
df2. <- rbind(df21,df22) ## treatment by dose at 1.0
df3. <- rbind(df31,df32) ## treatment by dose at 2.0

## create 6 possible combinations of treatment by dose
df.1.1 <- rbind(df11,df21) ## OJ at .05 and 1.0
df.1.2 <- rbind(df11,df31) ## OJ at .05 and 2.0
df.1.3 <- rbind(df21,df31) ## OJ at 1.0 and 2.0


df.2.1 <- rbind(df12,df22) ## VC at .05 and 1.0
df.2.2 <- rbind(df12,df32) ## VC at .05 and 2.0
df.2.3 <- rbind(df22,df32) ## VC at 1.0 and 2.0


## run first three t.tests on possible combinations of dose by treatment
tt1. <- t.test(len ~ Delivery.Meth, data = df1., alternative = "t", paired = F, var.equal = T)
tt2. <- t.test(len ~ Delivery.Meth, data = df2., alternative = "t", paired = F, var.equal = T)
tt3. <- t.test(len ~ Delivery.Meth, data = df3., alternative = "t", paired = F, var.equal = T)


ttt1 <- data.frame(Test_Conducted=as.character(),tvalue=as.double(),pvalue=as.double(),
                  Significant = as.logical(), Bonferroni = as.double(), BH = as.double(), stringsAsFactors = FALSE)
ttt1[1,1] <- "H_0: mean_OJ = mean_VC - Dose = 0.5 mg/day"
ttt1[1,2] <- format(tt1.$statistic, digits = 3, nsmall = 3)
ttt1[1,3] <- format(tt1.$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt1.$p.value <= 0.05) {ttt1[1,4] <- TRUE} else {ttt1[1,4] <- FALSE}


ttt1[2,1] <- "H_0: mean_OJ = mean_VC - Dose = 1.0 mg/day"
ttt1[2,2] <- format(tt2.$statistic, digits = 3, nsmall = 3)
ttt1[2,3] <- format(tt2.$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt2.$p.value <= 0.05) {ttt1[2,4] <- TRUE} else {ttt1[2,4] <- FALSE}


ttt1[3,1] <- "H_0: mean_OJ = mean_VC - Dose = 2.0 mg/day"
ttt1[3,2] <- format(tt3.$statistic, digits = 3, nsmall = 3)
ttt1[3,3] <- format(tt3.$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt3.$p.value <= 0.05) {ttt1[3,4] <- TRUE} else {ttt1[3,4] <- FALSE}


## Conduct Bonferroni and Benjamini-Hochberg corrections


ttt1[,5] <- format(p.adjust(ttt1[,3], method = "bonferroni"),scientific = FALSE, digits = 5, nsmall = 5)
ttt1[,6] <- format(p.adjust(ttt1[,3], method = "BH"),scientific = FALSE, digits = 5, nsmall = 5)


## run three t.tests on possible combinations of OJ by dose
tt.1.1 <- t.test(len ~ Dose, data = df.1.1, alternative = "t", paired = F, var.equal = T)
tt.1.2 <- t.test(len ~ Dose, data = df.1.2, alternative = "t", paired = F, var.equal = T)
tt.1.3 <- t.test(len ~ Dose, data = df.1.3, alternative = "t", paired = F, var.equal = T)


ttt2 <- data.frame(Test_Conducted=as.character(), tvalue=as.double(),pvalue=as.double(),
                  Significant = as.logical(), stringsAsFactors = FALSE, Bonferroni = as.double(), BH = as.double())
ttt2[1,1] <- "H_0: mean_dosea = mean_doseb - 0.5 to 1.0 mg/day, with Orange Juice delivery"
ttt2[1,2] <- format(tt.1.1$statistic, digits = 3, nsmall = 3)
ttt2[1,3] <- format(tt.1.1$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt.1.1$p.value <= 0.05) {ttt2[1,4] <- TRUE} else {ttt2[1,4] <- FALSE}


ttt2[2,1] <- "H_0: mean_dosea = mean_doseb - 0.5 to 2.0 mg/day, with Orange Juice delivery"
ttt2[2,2] <- format(tt.1.2$statistic, digits = 3, nsmall = 3)
ttt2[2,3] <- format(tt.1.2$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt.1.2$p.value <= 0.05) {ttt2[2,4] <- TRUE} else {ttt2[2,4] <- FALSE}


ttt2[3,1] <- "H_0: mean_dosea = mean_doseb - 1.0 to 2.0 mg/day, with Orange Juice delivery"
ttt2[3,2] <- format(tt.1.3$statistic, digits = 3, nsmall = 3)
ttt2[3,3] <- format(tt.1.3$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt.1.3$p.value <= 0.05) {ttt2[3,4] <- TRUE} else {ttt2[3,4] <- FALSE}
```

```
## run three t.tests on possible combinations of VC by dose
tt.2.1 <- t.test(len ~ Dose, data = df.2.1, alternative = "t", paired = F, var.equal = T)
tt.2.2 <- t.test(len ~ Dose, data = df.2.2, alternative = "t", paired = F, var.equal = T)
tt.2.3 <- t.test(len ~ Dose, data = df.2.3, alternative = "t", paired = F, var.equal = T)


ttt2[4,1] <- "H_0: mean_dosea = mean_doseb - 0.5 to 1.0 mg/day, with Ascorbic Acid delivery"
ttt2[4,2] <- format(tt.2.1$statistic, digits = 3, nsmall = 3)
ttt2[4,3] <- format(tt.2.1$p.value,scientific = FALSE, digits = 3, nsmall = 3)


if (tt.2.1$p.value <= 0.05) {ttt2[4,4] <- TRUE} else {ttt2[4,4] <- FALSE}


ttt2[5,1] <- "H_0: mean_dosea = mean_doseb - 0.5 to 2.0 mg/day, with Ascorbic Acid delivery"
ttt2[5,2] <- format(tt.2.2$statistic, digits = 3, nsmall = 3)
ttt2[5,3] <- format(tt.2.2$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt.2.2$p.value <= 0.05) {ttt2[5,4] <- TRUE} else {ttt2[5,4] <- FALSE}


ttt2[6,1] <- "H_0: mean_dosea = mean_doseb - 1.0 to 2.0 mg/day, with Ascorbic Acid delivery"
ttt2[6,2] <- format(tt.2.3$statistic, digits = 3, nsmall = 3)
ttt2[6,3] <- format(tt.2.3$p.value,scientific = FALSE, digits = 3, nsmall = 3)
if (tt.2.3$p.value <= 0.05) {ttt2[6,4] <- TRUE} else {ttt2[6,4] <- FALSE}


## Conduct Bonferroni and Benjamini-Hochberg corrections
ttt2[,5] <- format(p.adjust(ttt2[,3], method = "bonferroni"), scientific = FALSE, digits = 3, nsmall = 3)
ttt2[,6] <- format(p.adjust(ttt2[,3], method = "BH"), scientific = FALSE, digits = 3, nsmall = 3)
```