

Statistical Inference Project

Matthew Henderson

14 February 2021

Overview

This project aims to display concepts from statistical inference through data analysis of simulated and real-life data. The first section uses simulated data from an exponential distribution to show evidence of well-known results in statistical theory such as the central limit theorem. In the second section there is a data analysis of a ToothGrowth dataset which constructs a hypothesis test to distinguish between different groups in the data.

Simulations

The exponential distribution

The exponential distribution models the time between counts of the Poisson distribution. The exponential density function is:

$$f(x; \lambda) = \lambda e^{-(\lambda x)}$$

where x is on the interval $[0, \infty)$. The expected value of this distribution is $E[X] = 1/\lambda$ and the variance is $Var[X] = 1/\lambda^2$.

Simulating data

In this project I will use simulated data from an exponential distribution with $\lambda = 0.2$. The PDF of this function is $0.2e^{0.2x}$ defined for values of x greater than 0. I will simulate 1000 samples of 40 random variables from this distribution.

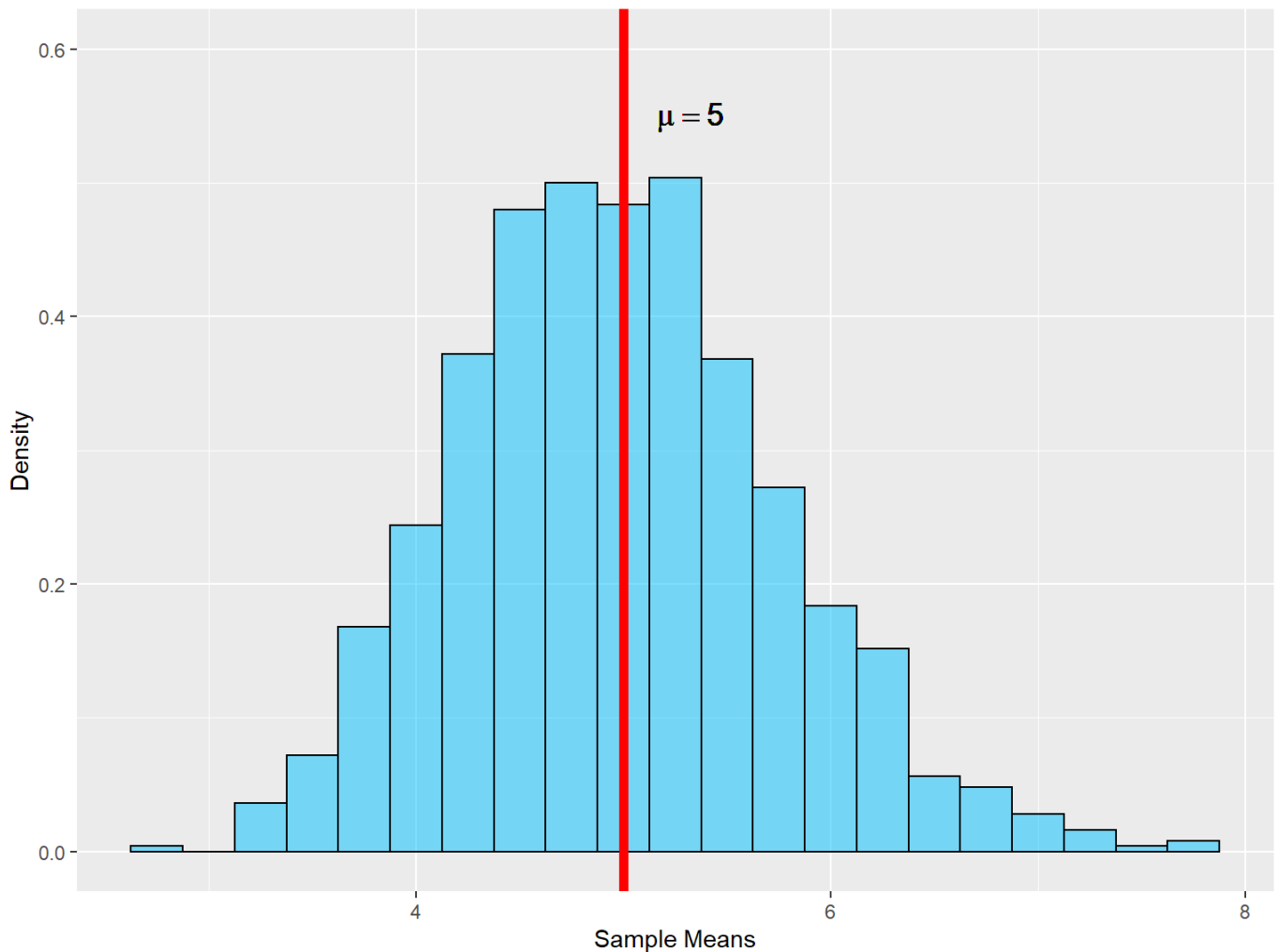
```
set.seed(1234)
lambda <- 0.2
n <- 40
nosim <- 1000
expMatrix <- matrix(rexp(n*nosim, rate = lambda), nrow = nosim)
mns <- apply(expMatrix, 1, mean)
sds <- apply(expMatrix, 1, sd)
```

So now I have a 40×1000 matrix of simulated data from the $exp(\lambda = 0.2)$ distribution. A 1000-dim vector of the mean of each row and a 1000-dim vector of the standard deviation of each row.

Sample Mean vs Theoretical Mean

The theoretical mean of the exponential distribution is $E[X] = 1/\lambda$ so for this distribution the expected values is $E[X] = 1/0.2 = 5$. We can compare this theoretical value with the simulated data by creating a histogram of the sample means.

```
library(ggplot2)
## Create a density histogram for the mean data
g <- ggplot(data.frame(mns = mns), aes(x = mns)) + geom_histogram(aes(y = ..density..), colour = "black", fill = "deepskyblue1", binwidth = 0.25, alpha = 0.5) + ylim(0,0.6) + labs(x = "Sample Means", y = "Density")
## add a vertical line at the theoretical mean = 5
g <- g + geom_vline(xintercept = 5, colour = "red", size = 2) +
  geom_text(aes(x = 5, y = 0.55), label = expression(mu==5), hjust = -0.5, size = 5)
g
```

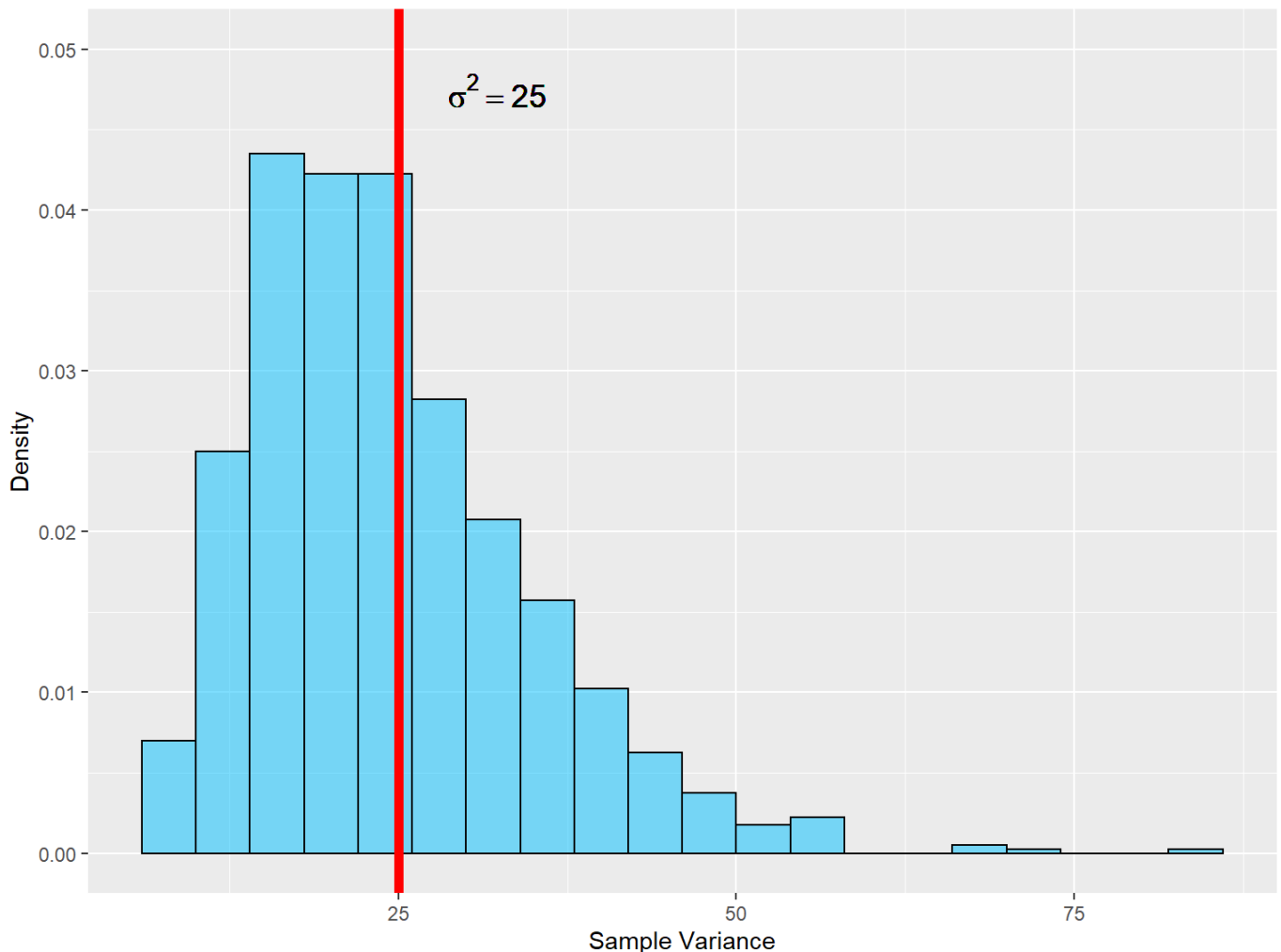


The sample means create a bell-shaped curve around the true population mean and look to be a very good estimator of the true mean.

Sample Variance vs Theoretical Variance

The theoretical variance of the exponential distribution is $Var[X] = 1/\lambda^2$ hence, for this distribution $Var[X] = 1/0.2^2 = 25$. We can compare the simulated data to the theoretical variance.

```
## create a density histogram for the variance data
g2 <- ggplot(data.frame(var = sds ^ 2), aes(x = var)) + geom_histogram(aes(y = ..density..),
  colour = "black", fill = "deepskyblue1", binwidth = 4, alpha = 0.5) + labs(x = "Sample Variance", y = "Density") + ylim(0, 0.05)
## add a vertical line at the theoretical variance = 25
g2 <- g2 + geom_vline(xintercept = 25, colour = "red", size = 2) +
  geom_text(aes(x = 25, y = 0.0475), label = expression(sigma^2==25), hjust = -0.5, size = 5)
g2
```



The sample variance is concentrated around the true population variance however, it is noticeable that the sample variances have a more skewed looking distribution than the sample means with a higher chance of more extreme values to the right side of the distribution with some more than three standard deviations away from the mean. For this distribution we probably want a higher sample size n for the sample variance to be a good approximator of the population variance.

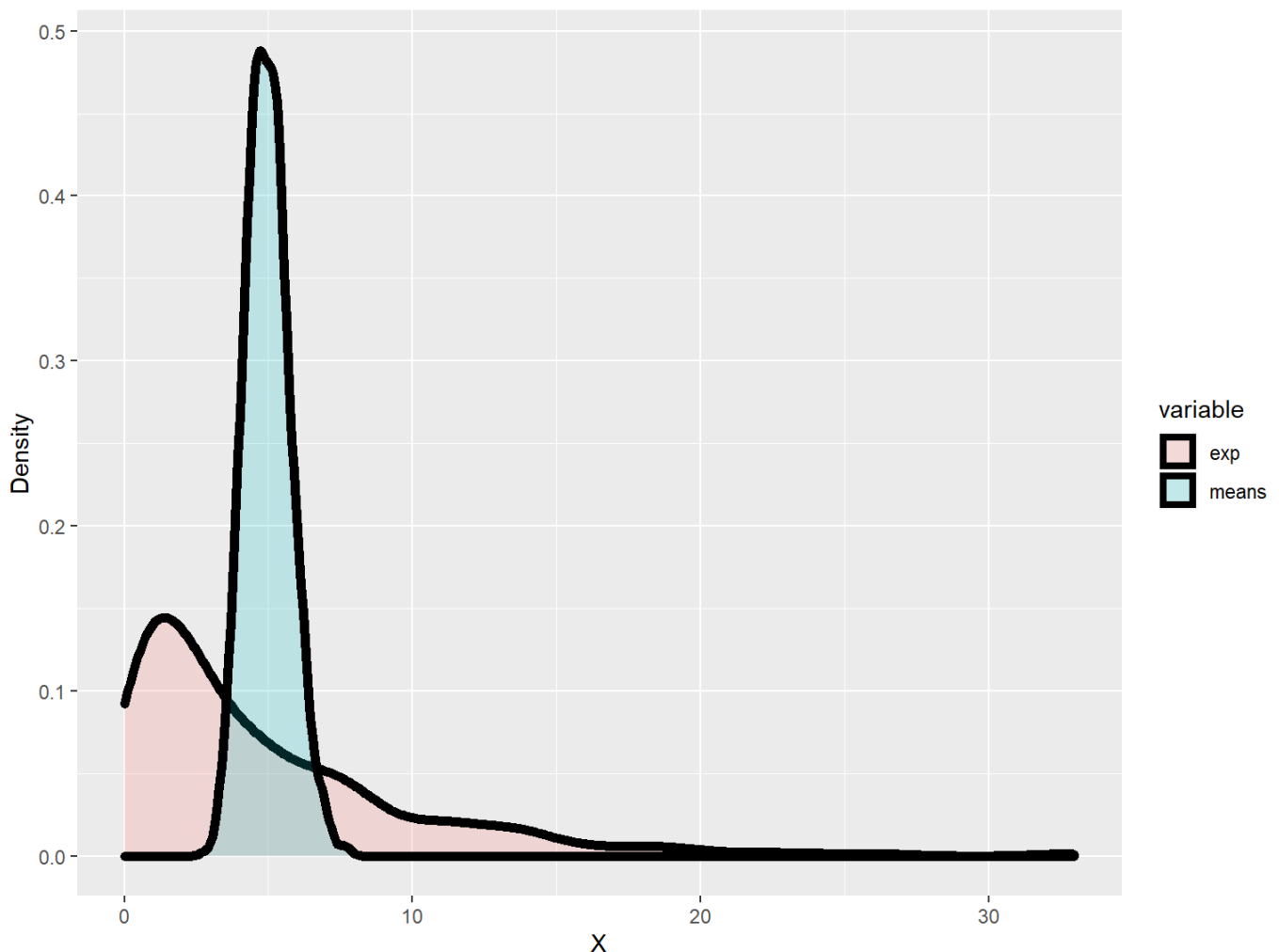
Sample Means are approximately Normal

Plotting a density curve of a large number of random variables from an exponential distribution should be an exponential looking curve with the center of mass around its true mean. If we take the average of 40 exponential random variables a large number of times we see that the distribution of averages has a gaussian looking density which is far more concentrated around the population mean.

```
library(reshape2)
df <- data.frame(exp = rexp(nosim, lambda), means = mns)
df <- melt(df)
```

```
## No id variables; using all as measure variables
```

```
g3 <- ggplot(df, aes(x = value, fill = variable)) + geom_density(alpha = 0.2, size = 2) + labs(x = "X", y = "Density")
g3
```



This results can be explained by the Central Limit Theorem which tells us that the distribution of averages of any independent and identically distributed (iid) random variable becomes increasingly more normal as the sample size increases.

In our case with have taken the average of 40 exponentials ($\lambda = 0.2$) 1000 times, hence each average is independant of the last with the exact same probability distribution. The central limit theorem states:

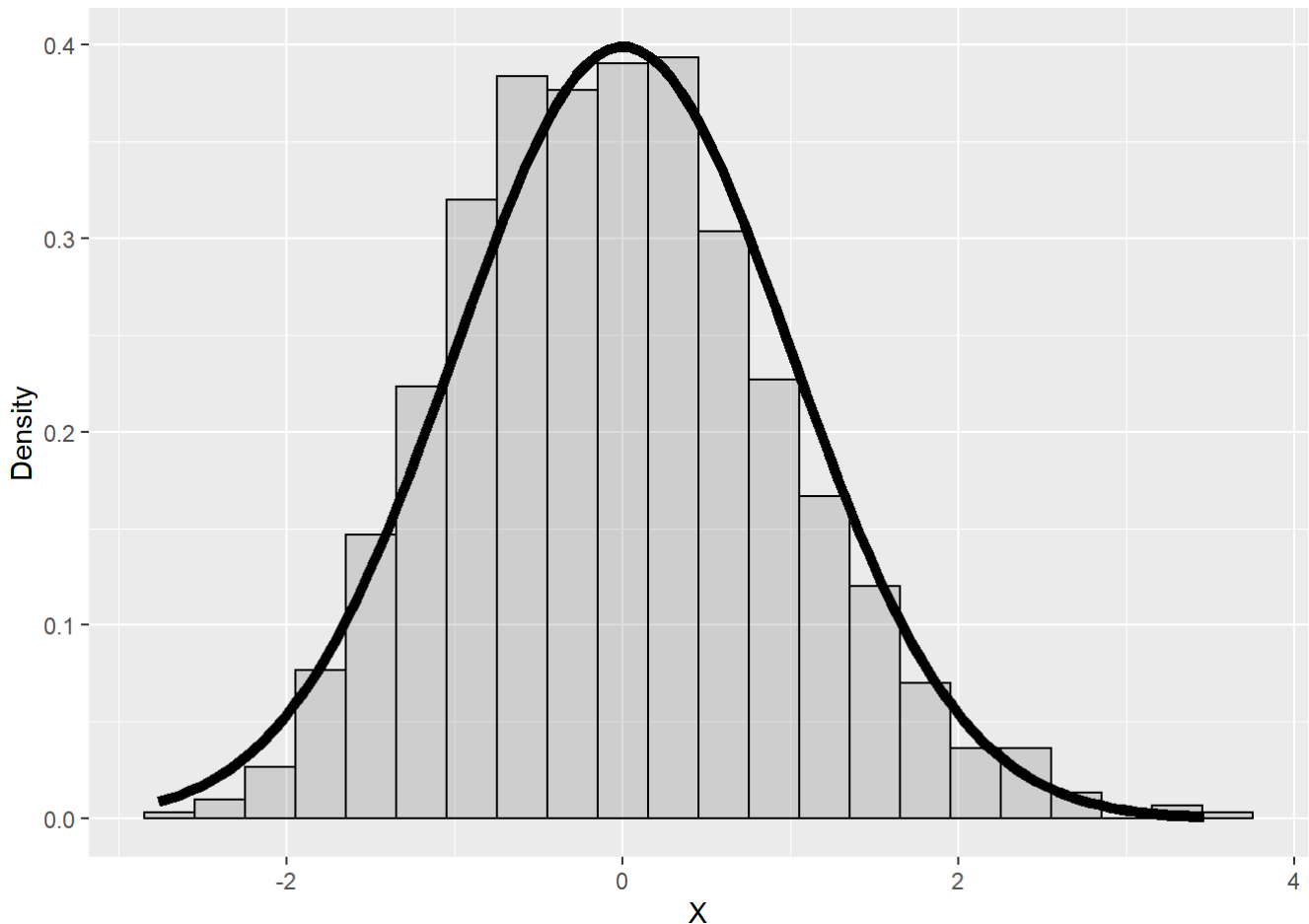
$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\text{sigma}} = \frac{\text{Estimate} - \text{Mean of Estimate}}{\text{Std. Err. of estimate}}$$

follows a distribution like that of a standard normal. We can apply the CLT to our simulated data by applying the transformation:

$$\frac{\sqrt{40}(\bar{X}_n - 5)}{5}$$

since $\mu = 5$, $\sigma = 5$ and the sample size $n = 40$.

```
# normalise the data using variables from the population
cltFunc = function(x, n) sqrt(n) * (mean(x) - 5) / 5
df2 <- data.frame(x = apply(expMatrix, 1, cltFunc, 40))
g4 <- ggplot(df2, aes(x = x)) + geom_histogram(aes(y = ..density..), alpha = 0.2, binwidth =
0.3, colour = "black") + labs(x = "X", y = "Density")
g4 <- g4 + stat_function(fun = dnorm, size = 2)
g4
```



The approximation of the standard normal distribution from the distribution of averages of 40 random variables from the $\text{exp}(\lambda = 0.2)$ is a very good fit.

Inferential Data Analysis of ToothGrowth Dataset

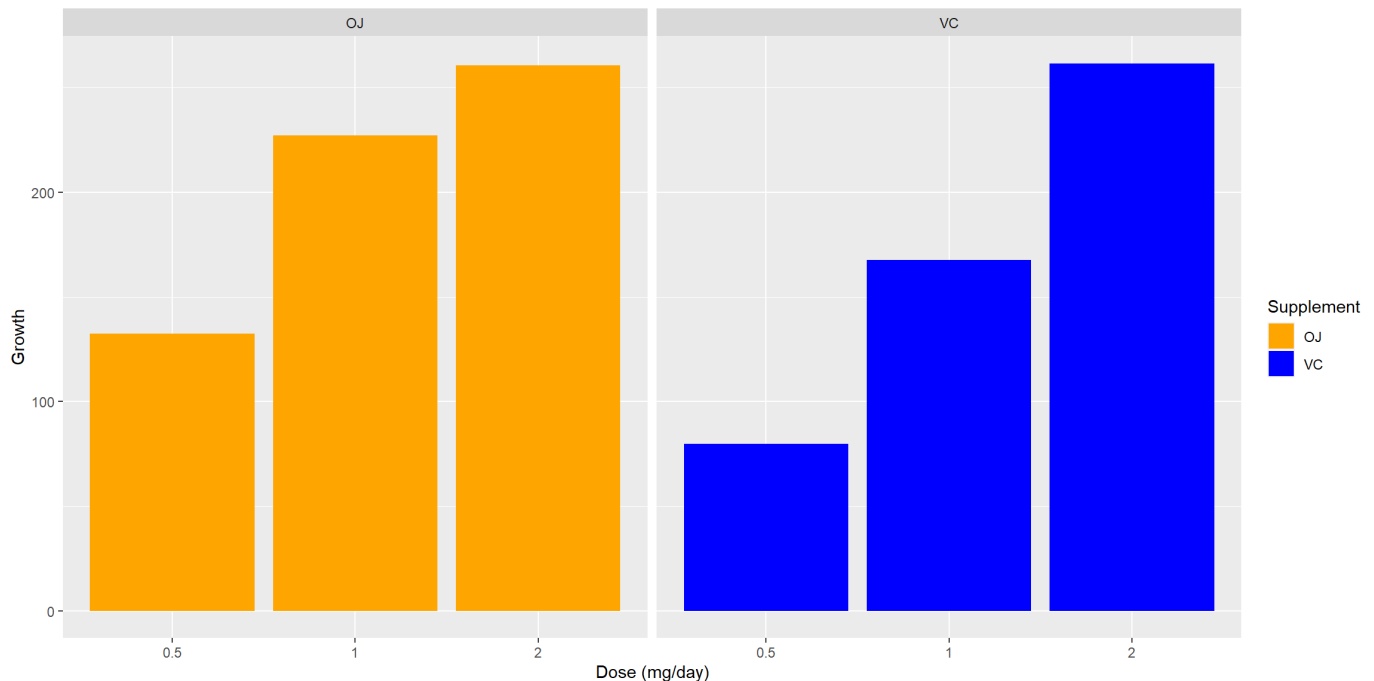
The ToothGrowth dataset from the datasets library contains data on the effect of Vitamin C on tooth growth in Guinea Pigs. The output is the length of odontoblasts (cells responsible for tooth growth) in 60 guinea pigs. There were two different delivery methods for the vitamin C, orange juice or ascorbic acid and three different doses (0.5, 1 or 2mg/day). Each Guinea pig recieved vitamin C by one of the two delivery methods and one of the three doses.

```
library(datasets)
data("ToothGrowth")
str(ToothGrowth)
```

```
## 'data.frame':   60 obs. of  3 variables:
## $ len : num  4.2 11.5 7.3 5.8 6.4 10 11.2 11.2 5.2 7 ...
## $ supp: Factor w/ 2 levels "OJ","VC": 2 2 2 2 2 2 2 2 2 2 ...
## $ dose: num  0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5 ...
```

The question is does the method of delivery or dose of Vitamin C effect the length of growth?

```
library(dplyr)
sumTG <- ToothGrowth %>% group_by(supp, dose) %>% summarise(Growth = sum(len))
g5 <- ggplot(sumTG, aes(x = as.factor(dose), y = Growth, fill = supp)) + geom_bar(stat = "identity")
g5 <- g5 + facet_grid(. ~ supp) +
  scale_fill_manual(values = c("orange", "blue")) +
  labs(x = "Dose (mg/day)", y = "Growth", fill = "Supplement")
g5
```



It looks like the Orange juice is a more effective method of delivery for doses of 0.5 and 1mg/day however there does not seem to be much of a difference for 2mg/day.

Hypothesis Testing

We want to see if there is any statistical evidence of a difference in lengths between the two delivery methods for Vitamin C. To do this we split the data into two groups for the guinea pigs that had orange juice and those that had ascorbic acid. The null hypothesis will be that the mean difference between these two groups is 0 and the alternative hypothesis is that the mean difference is not equal to zero. Hence:

$$H_0 : \mu = 0$$

versus:

$$H_a : \mu > 0$$

where μ is the mean difference between the two groups. We will assume this sample mean follows a t-distribution with $n_1 + n_2 - 2 = 58$ degrees of freedom ($\alpha = 0.05$). The two groups are independent of each other so this is not a paired t-test.

```
g1 <- filter(ToothGrowth, supp == "OJ")$len
g2 <- filter(ToothGrowth, supp == "VC")$len
t.test(g1, g2, paired = FALSE, var.equal = FALSE)$p.value
```

```
## [1] 0.06063451
```

```
t.test(g1, g2, paired = FALSE, var.equal = FALSE)$conf.int
```

```
## [1] -0.1710156  7.5710156
## attr(,"conf.level")
## [1] 0.95
```

The output of the test tells us that we cannot reject the null hypothesis with 95% confidence since the 95% confidence value for the mean contains 0.

Our exploratory analysis, seems to show that the 2mg/day dose seems to be an outlier in the dataset so we might be interested to see what happens if we split the 2mg/day group from the 0.5mg and 1mg/day groups.

```
# split the 0.5mg and 1mg doses from the 2mg dose
lowDose1 <- g1[1:20]
lowDose2 <- g2[1:20]
t.test(lowDose1, lowDose2, paired = FALSE, var.equal = FALSE)$p.value
```

```
## [1] 0.00423861
```

```
t.test(lowDose1, lowDose2, paired = FALSE, var.equal = FALSE)$conf
```

```
## [1] 1.875234 9.304766
## attr(,"conf.level")
## [1] 0.95
```

```
# test the 2mg dose group
highDose1 <- g1[21:30]
highDose2 <- g2[21:30]
t.test(highDose1, highDose2, paired = FALSE, var.equal = FALSE)$p.value
```

```
## [1] 0.9638516
```

```
t.test(highDose1, highDose2, paired = FALSE, var.equal = FALSE)$conf
```

```
## [1] -3.79807  3.63807
## attr(,"conf.level")
## [1] 0.95
```

With a p-value of ($p = 0.01$) we can reject the null hypothesis ($H_0 : \mu = 0$) and accept the alternative hypothesis ($H_a : \mu > 0$) that the orange juice delivery method creates a larger tooth growth for doses 0.5mg and 1mg/day. However, we cannot reject the null hypothesis when the dose is 2mg/day.

Conclusion

The tests show that we can be confident that the orange juice delivery method is a more effective supplement of Vitamin C than the ascorbic acid for smaller doses of 0.5mg and 1mg/day. However, when the dose is increased to 2mg/day there is no statistical evidence of a difference between the two delivery methods.