

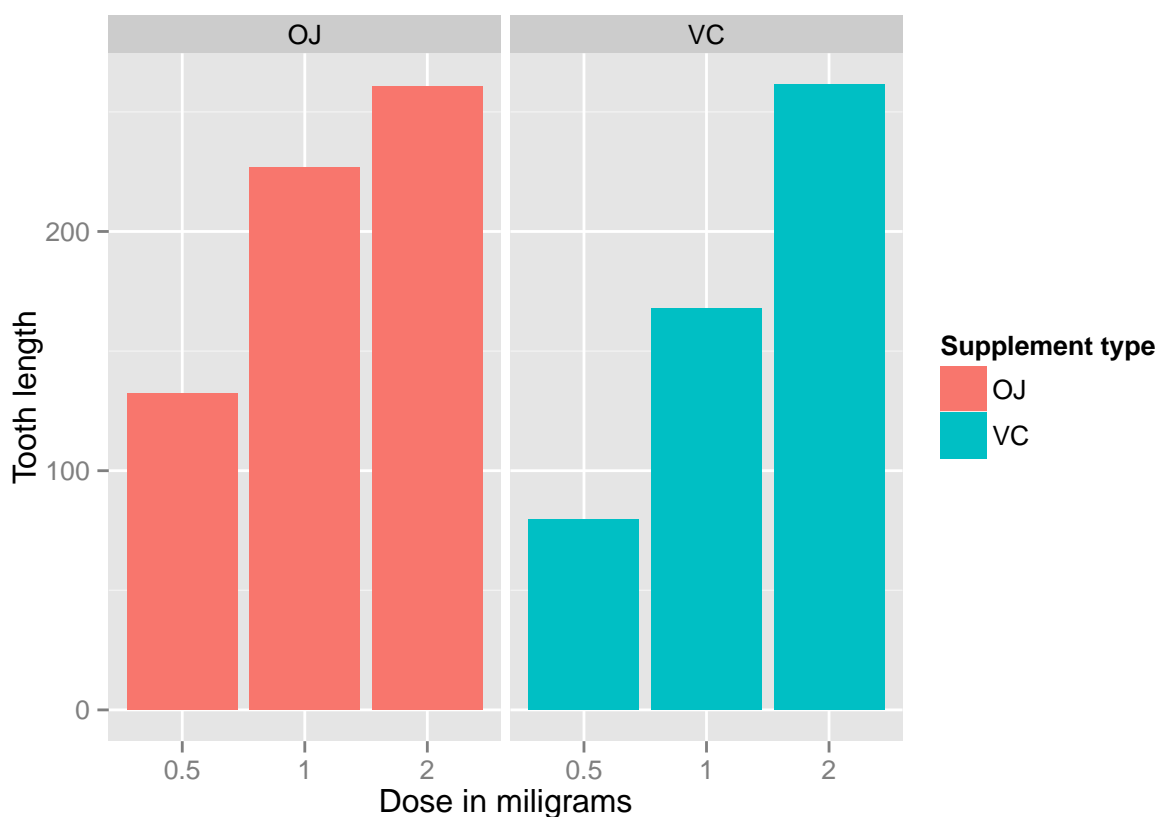
Statistical Inference Coursera Assignment: Part II

ROSTISLAV UHLIR

23 September 2015

In the second part of the project, we analyse the `ToothGrowth` data in the R `datasets` package. The data is set of 60 observations, length of odontoblasts (teeth) in each of 10 guinea pigs at each of three dose levels of Vitamin C (0.5, 1 and 2 mg) with each of two delivery methods (orange juice or ascorbic acid).

```
library(datasets)
library(ggplot2)
ggplot(data=ToothGrowth, aes(x=as.factor(dose), y=len, fill=supp)) +
  geom_bar(stat="identity",) +
  facet_grid(. ~ supp) +
  xlab("Dose in milligrams") +
  ylab("Tooth length") +
  guides(fill=guide_legend(title="Supplement type"))
```



We observe that there is a clear positive correlation between the tooth length and the dose levels of Vitamin C, for both delivery methods.

The effect of the dose can also be identified using regression analysis. One interesting question that can also be addressed is whether the supplement type (i.e. orange juice or ascorbic acid) has any effect on the tooth length. In other words, how much of the variance in tooth length, if any, can be explained by the supplement type?

```
fit <- lm(len ~ dose + supp, data=ToothGrowth)
summary(fit)
```

```
##
## Call:
## lm(formula = len ~ dose + supp, data = ToothGrowth)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.600 -3.700  0.373  2.116  8.800
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.2725     1.2824   7.231 1.31e-09 ***
## dose         9.7636     0.8768  11.135 6.31e-16 ***
## suppVC       -3.7000     1.0936  -3.383  0.0013 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.236 on 57 degrees of freedom
## Multiple R-squared:  0.7038, Adjusted R-squared:  0.6934
## F-statistic: 67.72 on 2 and 57 DF,  p-value: 8.716e-16
```

The model explains 70% of the variance in the data. The intercept is 9.2725, meaning that with no supplement of Vitamin C, the average tooth length is 9.2725 units. The coefficient of `dose` is 9.7635714. This means that increasing the delivered dose by 1 mg and keeping the rest equal (i.e. no change in the supplement type) would increase the tooth length 9.7635714 units. The last coefficient is for the supplement type. The calculated coefficient is for `suppVC` and the value is -3.7 meaning that delivering a given dose of ascorbic acid, without changing the dose, would result in 3.7 units of decrease in the tooth length. Since there are only two categories, we can also conclude that on average, delivering the dosage of orange juice would increase the tooth length by 3.7 units.

95% confidence intervals for two variables and these intercept as shown.

```
confint(fit)
```

```
##              2.5 %    97.5 %
## (Intercept)  6.704608 11.840392
## dose        8.007741 11.519402
## suppVC      -5.889905 -1.510095
```

The confidence intervals mean that if we collect a different set of data and estimate parameters of the linear model many times, 95% of the time, the coefficient estimations will be in these ranges. For each coefficient (i.e. intercept, `dose` and `suppVC`), if the coefficients are equal to zero, it means that no tooth length variation is explained by that variable (null hypothesis). All p -values are less than 0.05, rejecting this null hypothesis and suggesting that each variable explains a significant portion of variability in tooth length, assuming the significance level is 5%.