

Trees Data Analysis

Milica Cudina

The data set `trees` is built-in. Let's take a look at it.

```
names(trees)
```

```
## [1] "Girth" "Height" "Volume"
```

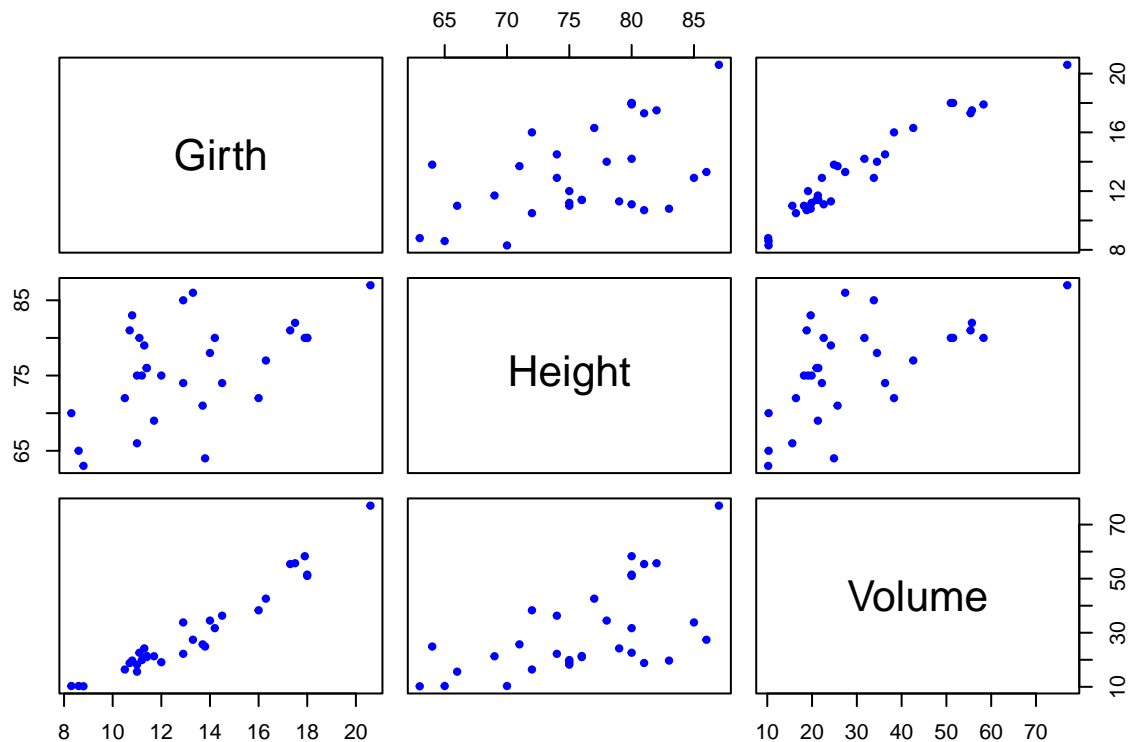
```
dim(trees)
```

```
## [1] 31 3
```

It contains measurements of 31 cherry trees, namely, their `Girth`, `Height`, and `Volume`.

Again, we undertake a rudimentary exploratory data analysis. It's natural to be interested in pairwise interactions. So, we create an array of scatterplots with which we can visually assess the shape of the dependence and the correlations of each pair of variables.

```
plot(trees,  
     col="blue", pch=20)
```



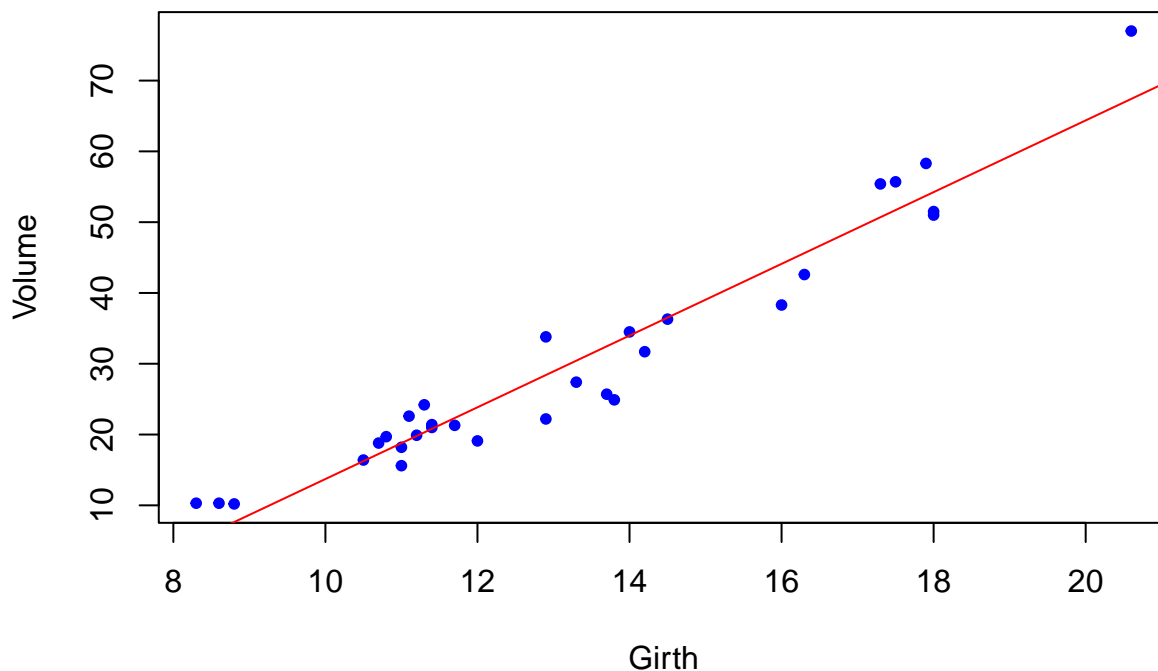
We might be interested in looking at, say, `Girth` as the explanatory and `Volume` as the response. This would be a simple linear regression.

```
lm.fit<-lm(Volume ~ Girth, data=trees)  
summary(lm.fit)
```

```
##
## Call:
## lm(formula = Volume ~ Girth, data = trees)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.065 -3.107  0.152  3.495  9.587
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept) -36.9435     3.3651  -10.98 0.000000000000762 ***
## Girth         5.0659     0.2474   20.48    < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.252 on 29 degrees of freedom
## Multiple R-squared:  0.9353, Adjusted R-squared:  0.9331
## F-statistic: 419.4 on 1 and 29 DF,  p-value: < 2.2e-16

attach(trees)
plot(Girth, Volume,
     pch=20, col="blue",
     main="Dependence of Volume on Girth")
abline(lm.fit, col="red")
```

Dependence of Volume on Girth



Now, let's see what happens when we add Height as an additional explanatory variable, thus creating a multiple linear regression.

```
lm.fit.m<-lm(Volume ~ Height + Girth)
summary(lm.fit.m)
```

```
##
## Call:
## lm(formula = Volume ~ Height + Girth)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.4065 -2.6493 -0.2876  2.2003  8.4847
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -57.9877      8.6382  -6.713 0.000000275 ***
## Height       0.3393      0.1302   2.607  0.0145 *
## Girth        4.7082      0.2643  17.816 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.882 on 28 degrees of freedom
## Multiple R-squared:  0.948, Adjusted R-squared:  0.9442
## F-statistic: 255 on 2 and 28 DF, p-value: < 2.2e-16
```

Let's compare the **coefficient of determination** R^2 for the above two fits.

For anyone who has ever seen trees, it's natural to suspect that there is a correlation between **Height** and **Girth**. Let's check

```
cor(Height, Girth)
```

```
## [1] 0.5192801
```

So, it might be a good idea to introduce the interaction term in our multiple linear regression.

```
lm.fit.mi<-lm(Volume ~ Girth*Height)
summary(lm.fit.mi)
```

```
##
## Call:
## lm(formula = Volume ~ Girth * Height)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.5821 -1.0673  0.3026  1.5641  4.6649
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  69.39632   23.83575   2.911  0.00713 **
## Girth       -5.85585    1.92134  -3.048  0.00511 **
## Height      -1.29708    0.30984  -4.186  0.00027 ***
## Girth:Height  0.13465    0.02438   5.524 0.00000748 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.709 on 27 degrees of freedom
## Multiple R-squared:  0.9756, Adjusted R-squared:  0.9728
## F-statistic: 359.3 on 3 and 27 DF, p-value: < 2.2e-16
```

We should take note, again, of any changes (improvements?) in the R^2 and/or the p -values.