The simple linear regression model for a paired data set  $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$ 

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

If we assume the  $\varepsilon_i$  are a *random sample* from a mean 0 normal distribution with variance  $\sigma^2$ , then there are 3 parameters to estiamte from the data:  $\beta_0, \beta_1$ , and  $\sigma$ . I'll use hats to indicate the estimates based on the sample, e.g.  $\hat{\beta}_1$  is the estimate for the unknown value of  $\beta_1$ .

For the  $\beta$ 's we have the following fact about the sampling distribution

$$T = \frac{\mathrm{obs} - \mathrm{exp}}{\mathrm{exp}} = \frac{\hat{\beta} - \beta}{\mathsf{SE}}$$

has a t distribution with n-2 degrees of freedom. This fact allows us to do

1. Confidence intervals for either  $\beta_0$  or  $\beta_1$ :

$$\hat{\beta}_i \pm t^* \mathsf{SE}(\hat{\beta}_i)$$

2. Significance tests about either use T as a test statistic.

Now, how to find the estimates and the standard errors using the computer? Let's look at the data stored in the **fat** data set. First download it:

## > source("http://wiener.math.csi.cuny.edu/st/R/fat.R")

Then attach the variables  $^{1}$ 

## > attach(fat)

The data set has many variables, including measurements of neck and wrist for different people. The names() functions lists them all

## > names(fat)

[1]	"case"	"body.fat"	"body.fat.siri"	"density"
[5]	"age"	"weight"	"height"	"BMI"
[9]	"ffweight"	"neck"	"chest"	"abdomen"
[13]	"hip"	"thigh"	"knee"	"ankle"
[17]	"bicep"	"forearm"	"wrist"	

<sup>1</sup>The warning about density is due to the data set containing a variable named density that overrides the name for the density() function.

To regress the neck size on the wrist size means to fit a model with neck size as a *response* variable and wrist size as a *predictor* variable. The lm() function does the work: (Just put the response on the left of the tilde and the predictor on the right.)

The estimate for  $\beta_0$  (labeled (Intercept)) and  $\beta_1$  (labeled wrist) are printed. To get more information we need to ask for it. The summary command is used to do so:

```
> summary(res)
```

```
Call:
lm(formula = neck ~ wrist)
Residuals:
     Min
                    Median
               1Q
                                 ЗQ
                                          Max
-7.99799 -1.08890 -0.01920
                           1.11489
                                     7.05953
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)
              2.6370
                         2.0058
                                  1.315
                                             0.19
              1.9394
                         0.1099 17.649
wrist
                                           <2e-16 ***
___
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 1.625 on 250 degrees of freedom
Multiple R-Squared: 0.5548,
                                   Adjusted R-squared: 0.553
F-statistic: 311.5 on 1 and 250 DF, p-value: < 2.2e-16
```

The output contains

- 1. The "call" which just echoes your command
- 2. A summary of the Residuals. These have a median of -0.01920. The mean is 0, why?

- 3. The important summary of the coefficients. This is where the action is. More below
- 4. A summary of the standard error. The value 1.625 is found from

$$\sqrt{\frac{1}{n-2}\sum e_i^2}$$

Or in R

> sqrt(sum(resid(res)^2)/(252 - 2))

[1] 1.625288

That value of R-squared might sound familiar, especially if you type cor(neck, wrist)2

The Coefficients have not only the estimates, but also the standard errors computed. Additionally a two-sided significance test for whether  $\beta$  is 0 is performed and summarized with an observed value (t value) and a *p*-value, Pr(>|t|).

From the output we see that the estimate  $\hat{\beta}_0 = 2.6370$  is *not* statistically significant from 0. (Why?) Yet the smaller value  $\hat{\beta}_1 = 1.9364$  is statistically significant from 0.

What about the natural question prompted by the following passage from *Gulliver's Travels* by Jonathan Swift (in the giant's voice)

Then they measured my right Thumb, and desired no more; for by a mathematical Computation, that twice round the Thumb is once round the Wrist, and so on to the Neck and the Waist, and by the help of my old Shirt, which I displayed on the Ground before them for a Pattern, they fitted me exactly.

That is, in the language of MTH 214, is the data consistent with

$$H_0: \beta_1 = 2, \quad H_A: \beta_1 \neq 2?$$

Here we can use T as the test statistic. The observed value is

$$Tobs = (1.9364 - 2)/0.1099Tobs$$

The *p*-value is computed with the *t*-distribution and n-2 degrees of freedom:

pt(Tobs, df = 250)

That is the difference between the 1.9364 and the hypothesized 2 is not statistically significant.

Okay, your turn.

- 1. Use T to find a 95% confidence interval for  $\beta_0$  based on  $\hat{\beta}_0$ . You can use  $t^* = 1.96 \text{ (why?)}$  or see what it is yourself with qt(.975, df = 250).
- 2. Is twice around the neck once around the waist? Let's model by the hip values using the waist values as a predictor.

Do two-sided tests of the following

$$H_0: \beta_0 = 0$$
$$H_0: \beta_1 = 2$$

What are your *p*-values?

3. The BMI is a simple measurement of fitness: weight divided by height squared. Larger BMIs mean a person is heavier than taller. The body fat is a much more precise measurement of fitness, but is much harder to measure. Does the BMI predict the body fat well? and if so, how does one convert?

Model the variable body.fat using BMI as a predictor.

Find a 95% CI for  $\beta_0$ .

Perform a two-sided significance test of

$$H_0: \beta_1 = 1.5$$

Comment on this proposed relationship between BMI and body fat:

To compute body fat from BMI do the following triple the BMI, divide by 2 and subtract 20.

- 4. Does body fat depend on height? Model body fat using height as a predictor and perform a one-sided significance test of  $H_0: \beta_1 = 0$  against  $H_A: \beta_1 < 0$ . Is the difference significant at the  $\alpha = 0.05$  level?
- 5. For a model of ankle size modeled by wrist size find 95% CIs for  $\beta_0$  and  $\beta_1$ . Is 0 in the first one? Is 1 in the latter? If so, comment on the summary that wrist and ankle sizes are the same *on average*.