

(Hoff p. 149ff.)

Regression modeling is concerned with describing how the sampling distribution of one random variable  $Y$  varies with another variable or set of variables  $\mathbf{x} = (x_1, \dots, x_p)$ . Specifically, a regression model postulates a form for  $p(y|\mathbf{x})$ , the conditional distribution of  $Y$  given  $\mathbf{x}$ . Estimation of  $p(y|\mathbf{x})$  is made using data  $y_1, \dots, y_n$  that are gathered under a variety of conditions  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

The normal linear regression model specifies that, in addition to  $E[Y|\mathbf{x}]$  being linear, the sampling variability around the mean is i.i.d. normal:

$$\begin{aligned}\epsilon_1, \dots, \epsilon_n &\stackrel{\text{i.i.d.}}{\sim} \text{normal}(0, \sigma^2) \\ Y_i &= \beta^T \mathbf{x}_i + \epsilon_i\end{aligned}$$

This model provides a complete specification of the joint probability density of observed data  $y_1, \dots, y_n$  conditional upon  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and values of  $\beta$  and  $\sigma^2$ :

$$\begin{aligned}p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \beta, \sigma^2) &= \prod_{i=1}^n p(y_i | \mathbf{x}_i, \beta, \sigma^2) \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T \mathbf{x}_i)^2 \right\}\end{aligned}\quad (1)$$

Another way to write this joint probability density is in terms of the multivariate normal distribution: Let  $\mathbf{y}$  be the  $n$ -dimensional column vector  $(y_1, \dots, y_n)^T$  and let  $\mathbf{X}$  be the  $n \times p$  matrix whose  $i$ th row is  $\mathbf{x}_i = \{x_{i,1}, x_{i,2}, \dots, x_{i,p}\}$ . Then the normal regression model is

$$\{\mathbf{y} | \mathbf{X}, \beta, \sigma^2\} \sim \text{multivariate normal}(\mathbf{X}\beta, \sigma^2 \mathbf{I}),$$

where  $\mathbf{I}$  is the  $p \times p$  identity matrix and

$$\mathbf{X}\beta = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} = \begin{pmatrix} \beta_1 x_{1,1} + \dots + \beta_p x_{1,p} \\ \vdots \\ \beta_1 x_{n,1} + \dots + \beta_p x_{n,p} \end{pmatrix} = \begin{pmatrix} E[Y_1 | \beta, \mathbf{x}_1] \\ \vdots \\ E[Y_n | \beta, \mathbf{x}_n] \end{pmatrix}$$

The density (1) depends on  $\beta$  through the residuals  $(y_i - \beta^T \mathbf{x}_i)$ . We compute the ordinary least squares estimates

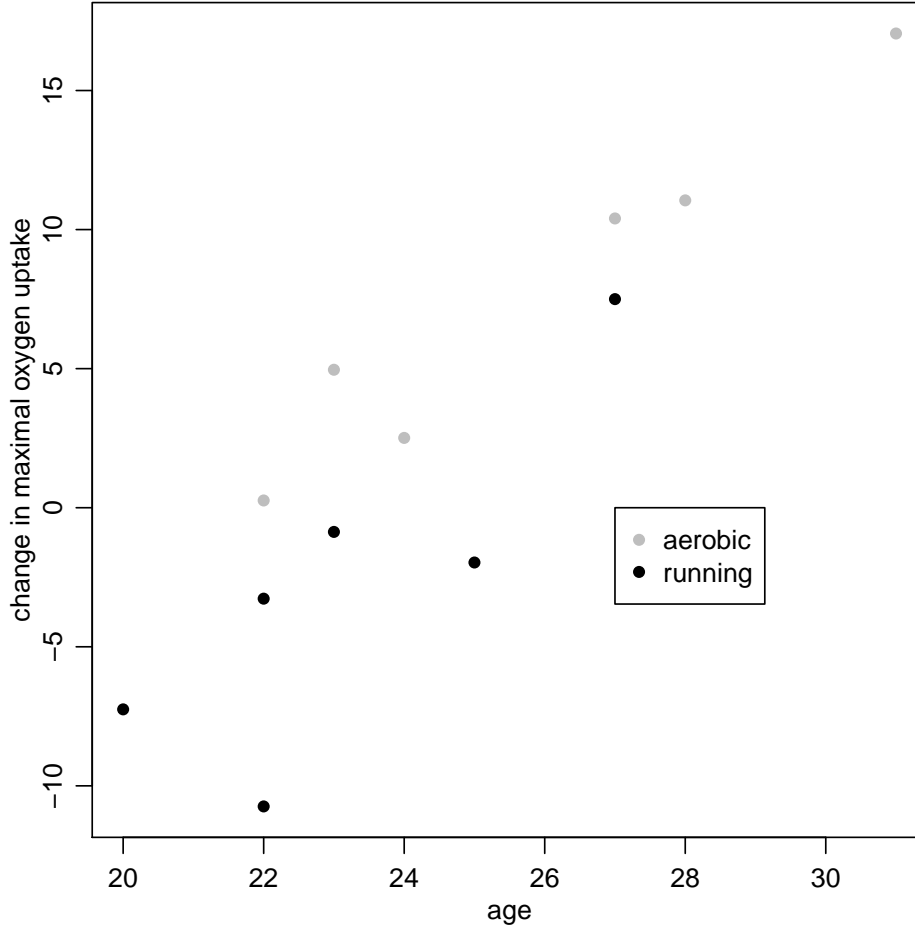
$$\hat{\beta}_{ols} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

and

$$\hat{\sigma}_{ols}^2 = \frac{SSR(\hat{\beta}_{ols})}{(n-p)} = \frac{\sum (y_i - \hat{\beta}_{ols}^T \mathbf{x}_i)^2}{(n-p)}.$$

*Example: Oxygen uptake (from Kuehl (2000), Hoff p. 149ff)*

Twelve healthy men who did not exercise regularly were recruited to take part in a study of the effects of two different exercise regimens on oxygen uptake. Six of the twelve men were randomly assigned to a 12-week flat-terrain running program, and the remaining six were assigned to a 12-week step aerobics program. The maximum oxygen uptake of each subject was measured (in liters per minute) while running on an inclined treadmill, both before and after the 12-week program. Of interest is how a subject's change in maximal oxygen uptake may depend on which program they were assigned to. However, other factors, such as age, are expected to affect the change in maximal uptake as well. The results are shown here:



Hoff's regression model:

$$Y_i = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \beta_3 x_{i,3} + \beta_4 x_{i,4} + \epsilon_i, \text{ where} \quad (2)$$

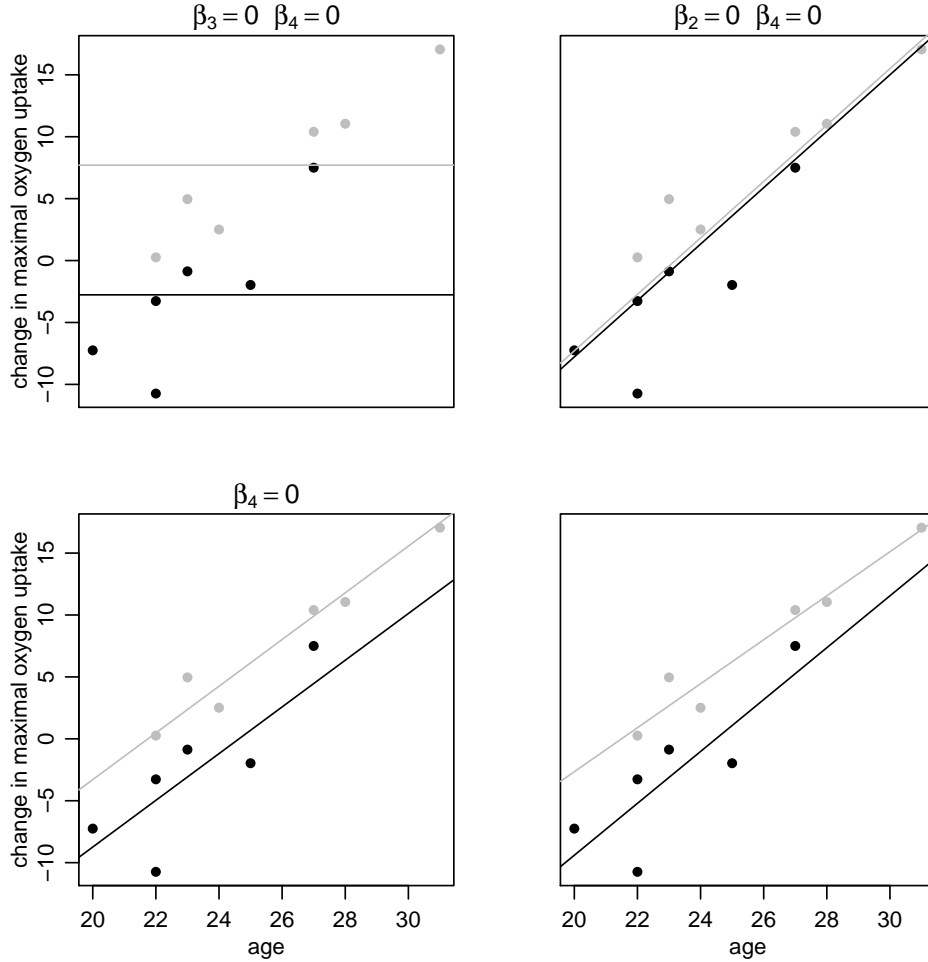
$x_{i,1} = 1$  for each subject  $i$   
 $x_{i,2} = 0$  if subject  $j$  is on the running program, 1 if on aerobic  
 $x_{i,3} = \text{age of subject } i$   
 $x_{i,4} = x_{i,2} \times x_{i,3}$

Under this model the conditional expectations of  $Y$  for the two different levels of  $x_{i,1}$  are

$$E[Y|\mathbf{x}] = \beta_1 + \beta_3 \times \text{age if } x_1 = 0, \text{ and}$$

$$E[Y|\mathbf{x}] = (\beta_1 + \beta_2) + (\beta_3 + \beta_4) \times \text{age if } x_1 = 1$$

In other words, the model assumes that the relationship is linear in age for both exercise groups, with the difference in intercepts given by  $\beta_2$  and the difference in slopes given by  $\beta_4$ . If we assumed that  $\beta_2 = \beta_4 = 0$ , then we would have identical lines for both groups. If we assumed  $\beta_4 = 0$  then we would have a different line for each group but they would be parallel. Allowing all coefficients to be non-zero gives us two unrelated lines. Some different possibilities are depicted graphically below:



Let's find the least squares regression estimates for the model in 2, and use the results to evaluate the differences between the two exercise groups. The ages of the 12 subjects, along with their observed changes in maximal oxygen uptake, are

$$\mathbf{x}_3 = (23, 22, 22, 25, 27, 20, 31, 23, 27, 28, 22, 24)$$

$$\mathbf{y} = (-0.87, -10.74, -3.27, -1.97, 7.50, -7.25, 17.05, 4.96, 10.40, 11.05, 0.26, 2.51),$$

with the first six elements of each vector corresponding to the subjects in the running group and the latter six corresponding to subjects in the aerobics group. After constructing

the 12 matrix  $\mathbf{X}$  out of the vectors  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4$  defined as in (2), the matrices  $\mathbf{X}^T\mathbf{X}$  and  $\mathbf{X}^T\mathbf{y}$  can be computed, from which we get  $\hat{\beta}_{ols} = (-51.29, 13.11, 2.09, -0.32)^T$ :

```
> n <- length(y)
> x1 <- rep(1,n)
> x4 <- x2*x3
> X = cbind(x1,x2,x3,x4)
> p = ncol(X)
> beta.ols<- solve(t(X)%*%X)%*%t(X)%*%y
> sig2.ols = (t(y-X%*%beta.ols)%*%(y-X%*%beta.ols)/(n-p))[1]
> sig2.ols

[1] 8.542477

> #sampling variance-covariance matrix of beta.ols:
> SIG2.ols = solve(t(X)%*%X)*sig2.ols
> #standard errors for the components of beta.ols:
> SE.ols = sqrt(diag(SIG2.ols))
> betadata = cbind(beta.ols,SE.ols)
> colnames(betadata) = c("beta.ols","SE.ols")
> betadata
```

|    | beta.ols    | SE.ols     |
|----|-------------|------------|
| x1 | -51.2939459 | 12.2522126 |
| x2 | 13.1070904  | 15.7619762 |
| x3 | 2.0947027   | 0.5263585  |
| x4 | -0.3182438  | 0.6498086  |

This means that the estimated linear relationship between uptake and age has an intercept and slope of -51.29 and 2.09 for the running group, and  $-51.29 + 13.11 = -38.18$  and  $2.09 - 0.32 = 1.77$  for the aerobics group. These two lines are plotted in the fourth panel of Figure XX. We obtain unbiased estimate  $\sigma^2 = SSR(\hat{\beta}_{ols})/(n - p) = 8.54$ , and use this to compute the standard error of the components of  $\hat{\beta}_{ols}$ , which are 12.25, 15.76, 0.53, and 0.65, respectively. comparing the values of  $\hat{\beta}_{ols}$  to their standard errors suggests that the evidence for differences between the two exercise regimens is not very strong.