Lecture 7

Linear Models

1. The basic model

We now study a *linear statistical model*. That is, we study the models where the observations $\mathbf{Y} := (Y_1, \ldots, Y_n)'$ has the following assumed property:

$$Y = X\beta + \varepsilon,$$

where $\boldsymbol{\beta} := (\beta_0, \beta_1, \dots, \beta_{p-1})$ is a vector of p unknown parameters, and

$$X := \begin{pmatrix} x_{1,0} & \cdots & x_{1,p-1} \\ \vdots & & \vdots \\ x_{n,0} & \cdots & x_{n,p-1} \end{pmatrix}$$

is the socalled "regression matrix," or "design matrix." The elements of the $n \times p$ matrix X are assumed to be known; these are the "descriptive" or "explanatory" variables, and the randomness of the observed values is inherited from the "noise vector," $\boldsymbol{\epsilon} := (\epsilon_1, \ldots, \epsilon_n)'$, which we may think of as being "typically small." Note that we are changing our notation slightly; X is no longer assumed to be a random vector [this is done in order to conform with the historical development of the subject].

Throughout, we assume always that the ε_i 's are independent with mean zero and common variance σ^2 , where $\sigma > 0$ is possibly [in fact, typically] unknown.

In particular, it follows from this assumption that

$$\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0} \text{ and } \operatorname{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \boldsymbol{I}.$$
 (1)

41

Let us emphasize that our linear model, once written out coordinatewise, is

$$X_i = \beta_0 x_{i,0} + \dots + \beta_{p-1} x_{i,p-1} + \varepsilon_i \qquad (1 \le i \le n).$$

It is intuitively clear that unless $n \ge p$, we cannot hope to effectively use our n observed values in order to estimate the p + 1 unknowns $\sigma^2, \beta_0, \ldots, \beta_{p-1}$. Therefore, we assume always that

$$n \ge p + 1. \tag{2}$$

This condition guarantees that our linear model is not overspecified.

The best-studied linear models are the *normal models*. Those are linear models for which we assume the more stringent condition that

$$\boldsymbol{\varepsilon} \sim \mathrm{N}_n \left(\boldsymbol{0} \,, \, \sigma^2 \boldsymbol{I} \right) \,. \tag{3}$$

Example 1 (A measurement-error model). Here we study a socalled measurement-error model: Suppose the observations Y_1, \ldots, Y_n satisfy

$$Y_i = \mu + \varepsilon_i$$
 $(1 \le i \le n)$

for an unknown parameter μ . This is a simplest example of a linear model, where $\beta = \mu$ is 1×1 , and $X := \mathbf{1}_{n \times 1}$ is a vector of *n* ones.

Example 2 (Simple linear regression). In simple linear regression we assume that the observed values have the form

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \qquad (1 \le i \le n),$$

where x_i is the predictive variable the corresponds to observation *i*, and β_0, β_1 are unknown. Simple linear regression fits into our theory of linear models, once we set the design matrix as

$$X := \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}.$$

Example 3 (Polynomial regression). Consider a nonlinear regression model

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_{p-1} x_i^{p-1} + \varepsilon_i \qquad (1 \le i \le n),$$

where *p* is a known integer ≥ 1 [*p* – 1 denotes the degree of the polynomial approximation to the observed *y*'s]. Then polynomial regression models are linear models with design matrices of the form

$$X := \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{p-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^{p-1} \end{pmatrix}.$$

Example 4 (One-way layout with equal observations). In the simplest case of one-way layout, in analysis of variance, our observations are indexed by vectors themselves as follows:

$$Y_{i,j} = \mu_i + \varepsilon_{i,j}$$
 $(1 \leq i \leq I, 1 \leq j \leq J).$

For instance, suppose we are interested in the effect of *I* different fertilizers. We apply these fertilizers to *J* different blocks, independently, and "measure the effect." Then, $Y_{i,j}$ is the effect of fertilizer *i* in block *j*. The preceding model is assuming that, up to sampling error, the effect of fertilizer *i* is μ_i . This is a linear model. Indeed, we can create a new random vector **Y** of *IJ* observations by simply "vectorizing" the $Y_{i,j}$'s:

$$\mathbf{Y} := (Y_{1,1}, \dots, Y_{1,J}, Y_{2,1}, \dots, Y_{2,J}, \dots, Y_{I,1}, \dots, Y_{I,J})'$$

The vector $\boldsymbol{\beta}$ of unknowns is $\boldsymbol{\beta} := (\mu_1, \dots, \mu_I)^{\prime}$, and the design matrix is the following $IJ \times I$ matrix:

$$X := \begin{pmatrix} \mathbf{1}_{J \times 1} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{J \times 1} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{J \times 1} \end{pmatrix},$$

where $\mathbf{1}_{J \times 1} := (1, ..., 1)'$ is a *J*-vector of all ones. It is possible to show that one-way layout with unequal number of observations is also a linear model. That is the case where $Y_{i,j} = \mu_i + \varepsilon_{i,j}$, where $1 \le i \le I$ and $1 \le j \le J_i$ [the number of observed values might differ from block to block].

2. The least-squares estimator of $\theta := X\beta$

Let us return to our general linear model

$$Y = X\beta + \boldsymbol{\epsilon}.$$

Ultimately, our goal is to first find and then analyze the least-squares estimator $\hat{\beta}$ of β . But first, let us find the least-squares estimate for $\theta := X\beta$. In other words, we wish to perform the following optimization problem:

$$\min_{\boldsymbol{\beta}\in\mathbf{R}^p} \|\boldsymbol{Y}-\boldsymbol{X}\boldsymbol{\beta}\| = \min_{\boldsymbol{\theta}\in\mathcal{C}(\boldsymbol{X})} \|\boldsymbol{Y}-\boldsymbol{\theta}\|.$$

Abstractly speaking, the minimizer solves

$$\boldsymbol{\theta} = \boldsymbol{P}_{\mathcal{G}(\boldsymbol{X})}\boldsymbol{Y}.$$

But is there an optimal β ? As we shall see next, there certainly is a unique $\hat{\beta}$ when X has full rank.

3. The least-squares estimator of β

Our least-squares estimator $\hat{\beta}$ of the vector parameter β is defined via

$$\min_{\boldsymbol{\beta}\in\mathbf{R}^p} \|\boldsymbol{Y}-\boldsymbol{X}\boldsymbol{\beta}\| = \left\|\boldsymbol{Y}-\boldsymbol{X}\widehat{\boldsymbol{\beta}}\right\|$$

We aim to solve this minimization problem under natural conditions on the design matrix X. But first, let us introduce some notation. The vector

$$\widehat{\mathbf{Y}} := \widehat{\boldsymbol{\theta}} := X\widehat{\boldsymbol{\beta}}$$

is called the vector of *fitted values*, and the coordinates of

$$\mathbf{e} := \mathbf{Y} - \widehat{\mathbf{Y}}$$

are the socalled *residuals*.

Now we write our minimization problem in the following form: First find the minimizing $\widehat{\beta} \in \mathbf{R}^p$ that solves

$$\min_{\boldsymbol{z}\in \mathcal{G}(X)} \|\boldsymbol{Y}-\boldsymbol{z}\| = \left\|\boldsymbol{Y}-\boldsymbol{X}\widehat{\boldsymbol{\beta}}\right\|.$$

Now we know from Proposition 18 (p. 20) that the vector \mathbf{z} that achieves this minimum does so uniquely, and is given by $\mathbf{P}_{\mathcal{G}(X)}\mathbf{Y}$, where we recall $\mathbf{P}_{\mathcal{G}(X)} := X(X'X)^{-1}X'$ is projection onto the column space of X; this of course is valid provided that $(X'X)^{-1}$ exists. Now the matrix $\mathbf{P}_{\mathcal{G}(X)}$ plays such an important role that it has its own name: It is called the *hat matrix*, and is denoted as

$$H := P_{\mathcal{C}(X)} = X(X'X)^{-1}X'.$$

H is called the *hat matrix* because it maps the observations *Y* to the fitted values \widehat{Y} [informally, it puts a "hat" over *Y*]. More precisely, the defining feature of *H* is that

$$\widehat{Y} = HY$$
,

once again provided that X'X is nonsingular. The following gives us a natural method for checking this nonsingularity condition in terms of X directly.

Lemma 5. X'X is nonsingular if and only if rank(X) = p.

Proof. Basic linear algebra tells us that the positive semidefinite X'X is nonsingular if and only if it is positive definite; i.e., if and only if it has full rank. Since the rank of X'X is the same as the rank of X, and since n > p, X'X has full rank if and only if its rank, which is the same as rank(X), is p.

From now on we always assume [unless we state explicitly otherwise] that rank(X) = p. We have shown that, under this condition, if there is a $\hat{\beta}$, then certainly $\hat{Y} = X\hat{\beta} = HY = X(X'X)^{-1}X'Y$. In order to find $\hat{\beta}$ from this, multiply both sides by $(X'X)^{-1}X'$ to see that $\hat{\beta} = (X'X)^{-1}X'Y$.

The quantity RSS := $\mathbf{e}'\mathbf{e} = \|\mathbf{e}\|^2 := \|\mathbf{Y} - \widehat{\mathbf{Y}}\|^2$ is called the *sum* of *squared residuals*, also known as the *residual sum* of *squared errors*, and is given by $\|(\mathbf{I} - \mathbf{H})\mathbf{Y}\|^2 = \|\mathbf{P}_{\mathcal{C}(\mathbf{X})^{\perp}}\mathbf{Y}\|^2$. In particular, we have the following:

Proposition 6. If rank(X) = p, then the least-squares estimator of β is $\hat{\boldsymbol{\beta}} = (\mathbf{Y}'\mathbf{Y})^{-1}\mathbf{Y}'\mathbf{Y}$

$$\boldsymbol{\beta} := (X'X)^{-1}X'Y, \tag{4}$$

and RSS = $\|(\boldsymbol{I} - \boldsymbol{H})\boldsymbol{Y}\|^2$.

Let us make some elementary computations with the least squares estimator of β .

Lemma 7. $\hat{\boldsymbol{\beta}}$ is an unbiased estimator of $\boldsymbol{\beta}$, and $\operatorname{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (X'X)^{-1}$.

Proof. Because $E\widehat{\beta} = (X'X)^{-1}X'EY = \beta$, it follows that $\widehat{\beta}$ is unbiased. Also, $Var(\widehat{\beta}) = (X'X)^{-1}X'Var(Y)X(X'X)^{-1} = \sigma^2(X'X)^{-1}$, because of (1). \Box

4. Optimality

Theorem 8. Let $\theta := X\beta$ be estimated, via least squares, by $\hat{\theta} := HY$. Then, for all linear unbiased estimates of $\mathbf{c}'\theta$, the estimator $\mathbf{c}'\hat{\theta}$ uniquely minimizes the variance.

By a "linear unbiased estimator of $\mathbf{c}'\boldsymbol{\theta}''$ we mean an estimator of the form $\sum_{j=1}^{n} a_j Y_j$ whose expectation is $\mathbf{c}'\boldsymbol{\theta}$. In this sense, $\mathbf{c}'\hat{\boldsymbol{\theta}}$ is the best linear unbiased estimator [or "BLUE"] of $\mathbf{c}'\boldsymbol{\theta}$. The preceding can be improved upon as follows, though we will not prove it:

Theorem 9 (Rao). Under the normal model (3), $c'\hat{\theta}$ is the unique UMVUE of $c'\theta$, for every nonrandom vector $c \in \mathbf{R}^n$.

Let us consider Theorem 8 next.

Proof of Theorem 8. We saw on page 43 that $\hat{\theta} := HY$ irrespective of whether or not $(X'X)^{-1}$ exists. Therefore,

$$E(\mathbf{c}'\widehat{\boldsymbol{\theta}}) = \mathbf{c}' E\widehat{\boldsymbol{\theta}} = \mathbf{c}'\boldsymbol{\theta}, \quad \text{Var}(\mathbf{c}'\widehat{\boldsymbol{\theta}}) = \mathbf{c}' \text{Var}(\boldsymbol{H}\boldsymbol{Y})\mathbf{c} = \mathbf{c}'\boldsymbol{H}' \text{Var}(\boldsymbol{Y})\boldsymbol{H}\mathbf{c}$$
$$= \sigma^2 \|\boldsymbol{H}\mathbf{c}\|^2.$$

Any other linear estimator has the form a'Y, and satisfies

$$\mathrm{E}(\boldsymbol{a}'\boldsymbol{Y}) = \boldsymbol{a}'\mathrm{E}\boldsymbol{Y} = \boldsymbol{a}'\boldsymbol{\theta}, \quad \mathrm{Var}(\boldsymbol{a}'\boldsymbol{Y}) = \boldsymbol{a}'\mathrm{Var}(\boldsymbol{Y})\boldsymbol{a} = \sigma^2 \|\boldsymbol{a}\|^2.$$

If, in addition, a'Y is unbiased, then it follows that $a'\theta = c'\theta$; i.e., a - c is orthogonal to θ . This should hold no matter what value β [and hence θ] takes. Since $\mathcal{C}(X)$ is the collection of all possible values of θ , it follows that a - c is orthogonal to $\mathcal{C}(X)$. Because H is projection onto $\mathcal{C}(X)$, it follows that H(a - c) = a; equivalently, Hc = Ha. Therefore, $\operatorname{Var}(c'\hat{\theta}) = \sigma^2 ||Ha||^2$ and

$$\operatorname{Var}\left(\boldsymbol{a}'\boldsymbol{Y}\right) - \operatorname{Var}\left(\boldsymbol{c}'\widehat{\boldsymbol{\theta}}\right) = \sigma^{2}\left\{\|\boldsymbol{a}\|^{2} - \|\boldsymbol{H}\boldsymbol{a}\|^{2}\right\} = \sigma^{2}\|(\boldsymbol{I} - \boldsymbol{H})\boldsymbol{a}\|^{2} \ge 0,$$

thanks to the Pythagorean property.

5. Regression and prediction

Now that we have the least-squares estimate for β , let us use it in order to make prediction.

Recall that our model is $Y = X\beta + \epsilon$. In applications, Y_i is the *i*th observation for the *y* variable, and the linear model is really saying that given an explanatory variable $\mathbf{x} = (x_0, \dots, x_{p-1})'$,

$$y = \beta_0 x_0 + \dots + \beta_{p-1} x_{p-1} +$$
"noise."

Therefore, our prediction, for a given x, is

[predicted value]
$$y = \hat{\beta}_0 x_0 + \dots + \hat{\beta}_{p-1} x_{p-1}$$
, (5)

where $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_{p-1})'$ is the least-squares estimate of β . We may view the right-hand side of (5) as a function of x, and call (5) the equation for the "regression line."

6. Estimation of σ^2

We wish to also estimate σ^2 . The estimator of interest to us turns out to be the following:

$$S^2 := \frac{1}{n-p} \text{RSS},\tag{6}$$

as the next lemma suggests.

Lemma 10. S^2 is an unbiased estimator of σ^2 .

Proof. Recall that RSS = $\mathbf{e}'\mathbf{e} = \|\mathbf{Y} - \mathbf{H}\mathbf{Y}\|^2$. We can write the RSS as $\|(\mathbf{I} - \mathbf{H})\mathbf{Y}\|^2 = \mathbf{Y}'(\mathbf{I} - \mathbf{H})'(\mathbf{I} - \mathbf{H})\mathbf{Y} = \mathbf{Y}'(\mathbf{I} - \mathbf{H})\mathbf{Y}$. In other words, the RSS is a random quadratic form for the matrix $\mathbf{A} := \mathbf{I} - \mathbf{H}$, and hence

$$\begin{split} \mathrm{E}(\mathrm{RSS}) &= \mathrm{tr}((\boldsymbol{I} - \boldsymbol{H})\mathrm{Var}(\boldsymbol{Y})) + (\mathrm{E}\boldsymbol{Y})'(\boldsymbol{I} - \boldsymbol{H})(\mathrm{E}\boldsymbol{Y}) \\ &= \sigma^{2}\mathrm{tr}(\boldsymbol{I} - \boldsymbol{H}) + (\boldsymbol{X}\boldsymbol{\beta})'(\boldsymbol{I} - \boldsymbol{H})\boldsymbol{X}\boldsymbol{\beta}. \end{split}$$

Because $X\beta \in \mathcal{C}(X)$, I - H projects onto the orthogonal subspace of where it is, therefore $(I - H)X\beta = 0$. And the trace of the projection matrix (I - H) is its rank, which is $n - \operatorname{tr}(H) = n - p$, since X has full rank p. It follows that $E(RSS) = \sigma^2(n - p)$, and therefore $E(S^2) = \sigma^2$. \Box

7. The normal model

In the important special case of the normal model,

$$\boldsymbol{Y} \sim N_n \left(\boldsymbol{X} \boldsymbol{\beta} \,, \sigma^2 \boldsymbol{I} \right).$$

Therefore, $\hat{\beta} \sim N_p(\beta, \sigma^2(X'X)^{-1})$. And the vector $(I - H)Y = Y - X\hat{\beta}$ of *residual errors* is also a multivariate normal:

$$(I-H)Y \sim N_{n-p}\left((I-H)X\beta, \sigma^2(I-H)(I-H)'\right) = N_{n-p}\left(0, \sigma^2(I-H)\right)$$

Therefore, in the normal model,

$$S^{2} = \frac{1}{n-p} \left\| (\boldsymbol{I} - \boldsymbol{H}) \boldsymbol{Y} \right\|^{2} \sim \sigma^{2} \frac{\chi_{n-p}^{2}}{n-p}.$$

Finally, we note that $t'\widehat{\beta} + s'(I - H)Y$ is a matrix times Y for all $t \in \mathbb{R}^p$ and $s \in \mathbb{R}^{n-p}$. Therefore, $(\widehat{\beta}, (I - H)Y)$ is also a multivariate normal. But

$$\operatorname{Cov}\left(\widehat{\boldsymbol{\beta}}, (\boldsymbol{I}-\boldsymbol{H})\boldsymbol{Y}\right) = (X'X)^{-1}X'\operatorname{Var}(\boldsymbol{Y})(\boldsymbol{I}-\boldsymbol{H})' = \sigma^2(X'X)^{-1}X'(\boldsymbol{I}-\boldsymbol{H}) = \boldsymbol{0},$$

since the columns of X are obviously orthogonal to every element in $\mathcal{C}(X)^{\perp}$ and $I - H = P_{\mathcal{C}(X)^{\perp}}$. This shows that $\hat{\beta}$ and (I - H)Y are independent, and hence $\hat{\beta}$ and S^2 are also independent. Thus, we summarize our findings.

Theorem 11. The least-squares estimator $\hat{\beta}$ of β is given by (4); it is always unbiased. Moreover, S^2 is an unbiased estimator of σ^2 . Under the normal model, S and $\hat{\beta}$ are independent, and

$$\widehat{oldsymbol{eta}}\sim \mathrm{N}_p\left(oldsymbol{eta}$$
 , $\sigma^2(X'X)^{-1}
ight)$, $S^2\sim\sigma^2rac{\chi^2_{n-p}}{n-p}.$

Recall that \boldsymbol{Y} has the nondegenerate multivariate normal distribution $N_n(\boldsymbol{X}\boldsymbol{\beta},\sigma^2\boldsymbol{I})$. Therefore, its pdf is

$$f_{\boldsymbol{Y}}(\boldsymbol{y}) = rac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-rac{1}{2\sigma^2}\|\boldsymbol{y}-\boldsymbol{X}\boldsymbol{\beta}\|^2
ight).$$

This shows readily the following.

Lemma 12. In the normal model, $\hat{\beta}$ is the MLE of β and $\left(\frac{n-p}{n}\right)S^2$ is the MLE for σ^2 .

Proof. Clearly, maximizing the likelihood function, over all $\boldsymbol{\beta}$, is the same as minimizing $\|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta}\|^2$. Therefore, MLE = least squares for $\boldsymbol{\beta}$. As for σ^2 , we write the log likelihood function:

$$L(\sigma) = -\frac{n}{2}\ln(2\pi) - n\ln\sigma - \frac{1}{2\sigma^2}\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

Then,

$$L'(\sigma) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \| \boldsymbol{Y} - \boldsymbol{X} \boldsymbol{\beta} \|^2.$$

Set $L'(\sigma) = 0$ and solve to see that the MLE of σ^2 is $\frac{1}{n} || \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} ||^2 = (\frac{n-p}{n}) S^2$, thanks to the MLE principle and the already-proven fact that the MLE of $\boldsymbol{\beta}$ is $\hat{\boldsymbol{\beta}}$.

8. Some examples

1. A measurement-error model. Recall the measurement-error model

$$Y_i = \mu + \varepsilon_i$$
 $(1 \le i \le n).$

We have seen that this is a linear model with p = 1, $\beta = \mu$, and $X := \mathbf{1}_{n \times 1}$. Since X'X = n and $X'Y = \sum_{i=1}^{n} Y_i$, we have

$$\widehat{\boldsymbol{\beta}} = \widehat{\mu} := \overline{Y}_{\mu}$$

and

$$\frac{1}{n-1}S^2 = \frac{1}{n-1} \left\| \mathbf{Y} - \bar{\mathbf{Y}}\mathbf{1}_{n \times 1} \right\|^2 = \frac{1}{n-1}\sum_{i=1}^n (Y_i - \bar{Y})^2 := S^2.$$

These are unbiased estimators for μ and σ^2 , respectively. Under the normal model, \bar{Y} and S^2 are independent, $\bar{Y} \sim N(\mu, \sigma^2)$ and $(n-1)S^2 \sim \sigma^2 \chi^2_{n-1}$. These are some of the highlights of Math. 5080.

2. Simple linear regression. Recall that simple linear regression is our linear model in the special case that p = 2, $\beta = (\beta_0, \beta_1)'$, and

$$X := \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}.$$

We have

$$X'X = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix}, \quad X'Y = \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n x_i Y_i \end{pmatrix},$$

and

$$(X'X)^{-1} = \frac{1}{n\sum_{i=1}^{n} x_i^2 - (\sum_{i=1}^{n} x_i)^2} \begin{pmatrix} \sum_{i=1}^{n} x_i^2 & -\sum_{i=1}^{n} x_i \\ -\sum_{i=1}^{n} x_i & n \end{pmatrix}$$
$$= \frac{1}{\sum_{i=1}^{n} (x_i - \bar{x})^2} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^{n} x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix}.$$

Therefore,

$$egin{pmatrix} \widehat{eta}_0 \ \widehat{eta}_1 \end{pmatrix} = (X'X)^{-1}X'Y,$$

which leads to

$$\widehat{\beta}_{1} = \frac{\sum_{i=1}^{n} (x_{i} - \bar{x}) Y_{i}}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} = \frac{\sum_{i=1}^{n} (x_{i} - \bar{x}) (Y_{i} - \bar{Y})}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} \text{ and } \widehat{\beta}_{0} = \bar{Y} - \widehat{\beta}_{1} \bar{x}.$$

We have derived these formulas by direct computation already. In this way we find that the fitted values are

$$\widehat{\mathbf{Y}} = \mathbf{X}\widehat{\boldsymbol{\beta}} = \begin{pmatrix} \widehat{\beta}_0 + \widehat{\beta}_1 x_1 \\ \vdots \\ \widehat{\beta}_0 + \widehat{\beta}_1 x_n \end{pmatrix}.$$

Also,

$$S^{2} = \frac{1}{n-2} \left\| \boldsymbol{Y} - \widehat{\boldsymbol{Y}} \right\|^{2} = \frac{1}{n-2} \sum_{i=1}^{n} \left(Y_{i} - \widehat{\beta}_{0} - \widehat{\beta}_{1} x_{i} \right)^{2},$$

and this is independent of $(\widehat{\beta}_0, \widehat{\beta}_1)$ under the normal model.

Recall that our linear model is, at the moment, the simple regression model,

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

Perhaps the most important first question that we can ask in this context is $\beta_1 = 0$; that is, we wish to know whether the *x* variables are [linearly] independent of the *Y*'s. Let us try to find a confidence interval for β_1 in

order to answer this question. From now on, we work under the normal model. Recall that under the normal model,

$$\widehat{\boldsymbol{\beta}} \sim \mathrm{N}_2\left(\boldsymbol{\beta}, \sigma^2(X'X)^{-1}\right) \Rightarrow \widehat{\boldsymbol{\beta}}_1 \sim \mathrm{N}\left(\boldsymbol{\beta}_1, \sigma^2\left[(X'X)^{-1}\right]_{2,2}\right) = \mathrm{N}\left(\boldsymbol{\beta}_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right).$$
Equivalently

Equivalently,

$$Z := \frac{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2}}{\sigma} \left(\widehat{\beta}_1 - \beta_1\right) \sim \mathcal{N}(0, 1)$$

Now, S^2/σ^2 is independent of Z and is distributed as $\chi^2_{n-2}/(n-2)$. Therefore,

$$\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \left(\frac{\widehat{\beta}_1 - \beta_1}{S}\right)} = \frac{Z}{\sqrt{S^2/\sigma^2}} \sim t_{n-2}.$$

Therefore, a $(1 - \alpha) \times 100\%$ confidence interval for β_1 is

$$\widehat{\beta}_1 \pm \frac{S}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} t_{n-2}^{\alpha/2},$$

where t_{ν}^{r} is the point whose right area, under the t_{ν} pdf, is r. If zero is not in this confidence interval, then our statistical prediction is that β_{1} is not zero [at the confidence level α].

3. A remark on polynomial regression. Recall that, in polynomial regression, we have $Y = X\beta + \epsilon$, where X is the following design matrix:

$$X := \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{p-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^{p-1} \end{pmatrix}.$$

If p = 2, then this is simple linear regression. Next consider the "quadratic regression" model where p = 3. That is,

$$X := \begin{pmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix} \quad \Rightarrow \quad X'X = \begin{pmatrix} n & \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i^2 \\ \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i^3 & \sum_{i=1}^n x_i^4 \end{pmatrix}.$$

Because $n \ge 4$, X is nonsingular if and only if x is not a vector of constants [a natural condition]. But you see that already it is painful to invert X'X. This example shows the importance of using a computer in linear models: Even fairly simple models are hard to work with, using only direct calculations.