TMA4267 Linear Statistical Models Part 2: Linear regression Solutions to recommended exercise 3 - V2017

February 7, 2017

Problem 1: Simple linear regression

See the file RecEx31.R, accessible from the course website, for more details.

a) The forbes data set consists of n = 17 observations, and we are fitting a simple linear regression model which includes an intercept, so the design matrix X has 17 rows and 2 columns. Since the two columns are linearly independent, the rank of X is 2. We can confirm this using R, by calling the rankMatrix() function from the Matrix package.

library(Matrix) # Package needs to be installed, if it isn't already
rankMatrix(X)

- b) The boiling point increases with increasing pressure. When the pressure is around 0.7 bar, the boiling point is approximately 90° C, and when the pressure is around 1 bar, the boiling point is approximately 100° C. The relationship between pressure and boiling point does look linear.
- c) The estimated regression parameters are

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0\\ \hat{\beta}_1 \end{bmatrix} \approx \begin{bmatrix} 68.5\\ 31.2 \end{bmatrix},$$

and the fitted regression function has the expression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

where x is pressure, and \hat{y} is the expected boiling point in that pressure. From this expression it is clear that $\hat{\beta}_0$ is the boiling point the model would predict for x = 0 bar, i.e. in a vacuum, and $\hat{\beta}_1$ is the estimated rate of change of the boiling point with pressure; If the pressure increases by 1 bar, we expect the boiling point to increase by $\hat{\beta}_1$ degrees Celsius.

- d) Two things are apparent when looking at the plot of the raw residuals $\hat{\epsilon} = Y \hat{Y}$ versus the pressure:
 - The plot has an inverted U-shape. Going from left to right, the residuals start off negative, then gradually become positive, before decreasing and becoming negative again.

- With the exception of one data point (observation number 12), the residuals all lie inside the interval (-0.3, 0.3), and there is a roughly even split between positive and negative residuals.
- e) Is the linear model appropriate?

Yes. Although there may be a systematic relationship between the residuals and the covariate, it is clear that the relationship between the response and covariate is predominantly linear, and that the fitted model does a satisfactory job of representing that relationship.

• Are the error variances homoscedastic?

Maybe. It's a little hard to tell, because of the remaining trend in the residuals. It looks like the residual variance may vary somewhat with pressure, but even if this is the case, there is only weak heteroscedasticity, and no serious violation of this model assumption.

• Are the error terms uncorrelated?

No. The plot of residuals vs. pressure plainly shows that observations with similar pressures, have similar residuals. This is systematic variation not accounted for by the fitted model. If we want to improve the fit, we might try to use the logarithm of the pressure as a covariate, instead of the pressure.

• Does an additive model seem appropriate?

Yes. The class of additive models subsumes linear models, and we have seen that a linear model is able to achieve a reasonably good fit in this case. There is nothing about this data set, or the analysis done so far, suggesting that a non-additive model would be preferable.

Normality of residuals can be ascertained by making a normal quantile-quantile plot of $\hat{\epsilon}$, and comparing it to a straight line.

qqnorm(eps)

f) The parameter estimates, fitted response values and raw residuals all match the results from
b) and c). The residual plots (residuals vs. pressure and normal quantile-quantile) look the same as the ones in d) and e).

Problem 2: Happiness

a) The regression coefficient for work is 0.4761. Assume that we look at two individuals that have scored the same values for sex, love and money. Further assume that one of the individuals has reported work to have value 1 (seeking other employment) and the other has 2 (inbetween seeking other employment and OK). Then, on average, we would expect that the happiness for the last individual is 0.4761 higher than for the first individual. Keeping the other variables fized, the effect of work on happiness is that happiness increases on average with 0.4761 units for every one unit increase in the work variable.

We can perform a *t*-test to see if work is significant in the full model (given that the other variables are present). Test statistic:

$$t = \frac{\hat{\beta}_4}{\widehat{\operatorname{Var}}(\hat{\beta}_4)} = \frac{0.4761}{0.1994} = 2.39$$

which under the null hypothesis that $\beta_4 = 0$ (vs. the two-sided alternative that $\beta_4 \neq 0$) is referred to a *t*-distribution with n - k - 1 = 34 degrees of freedom (n = 39 is the number of observations and k = 4 is the number of explanatory variables, and 1 for the intercept). We use significance level 0.05 and the critical value at level 0.025 (since two-sided test) in the t_{34} distribution is approximately 2.03 (found for 35 in the table), which means that we reject the null hypothesis.

To test if the regression is significant, that is, not all coefficients are zero, we look at the F = 20.83 value. Using a significance level of 0.05 this is referred to a critical value for the Fisher distribution with 4 and 34 degrees of freedom: approximately 2.64 (with 35 from the table). This means that the regression is highly significant.

Residual plots: We see no clear trend in the plot of residuals vs. fitted values, which is good. The quantile-quantile plot shows no clear deviation from normality, but at least one outlier is identified.

b) In a multiple regression the least squares estimates for the regression coefficients are found by solving a set of equations. When the explanatory variables are not orthogonally selected the value for one explanatory variable will influence the estimate of the regression coefficient for the other. In a design of experiments where explanatory variables are chosen so that they are independent of each other (orthogonal columns) the normal equations will become uncoupled and the regression coefficient estimate for the explanatory variables will not influence each other.

Problem 3: Results on $\hat{\beta}$ and SSE in multiple linear regression

(Exam K2014, Problem 4) Define the matrix $\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T$.

a) H is a orthogonal projection matrix, since it is symmetric and idempotent, $H^T = H$ and HH = H. For a symmetric and idempotent matrix the rank is equal to the trace. The rank of H is p. See proof below.

$$\begin{split} \boldsymbol{H} &= \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \\ \boldsymbol{H}^T &= (\boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T)^T = \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T = \boldsymbol{H} \\ \boldsymbol{H}^2 &= (\boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T) \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T = \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T = \boldsymbol{H} \\ \operatorname{tr}(\boldsymbol{H}) &= \operatorname{tr}(\boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T) = \operatorname{tr}(\boldsymbol{X}^T \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1}) = \operatorname{tr}(\boldsymbol{I}_{p \times p}) = p \end{split}$$

Graphically: The vector HY is a projection of the vector Y onto the space spanned by the columns of X.

The matrix I - H is also symmetric and idempotent, and thus a symmetric projection matrix. The rank of I - H is n - p. See proof below.

$$(\mathbf{I} - \mathbf{H})^T = \mathbf{I} - \mathbf{H}^T = \mathbf{I} - \mathbf{H}$$
$$(\mathbf{I} - \mathbf{H})^2 = (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) = \mathbf{I} - 2\mathbf{H} + \mathbf{H}^2 = \mathbf{I} - \mathbf{H}$$
$$\operatorname{tr}(\mathbf{I} - \mathbf{H}) = \operatorname{tr}(\mathbf{I}_{n \times n}) - \operatorname{tr}(\mathbf{H}) = n - p$$

Graphically: The vector (I - H)Y is a projection of the vector Y onto the space orthogonal to the space spanned by the columns of X.

b) Let $SSE = \mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y}$. Derive the distribution of SSE.

One of the key theorems of this course (point 2 of Theorem B.8 on p. 651 of Fahrmeir et al.) states that if D is a symmetric and idempotent matrix with rank r and $Z \sim N_n(0, I)$, then

$$\boldsymbol{Z}^T \boldsymbol{D} \boldsymbol{Z} \sim \chi_r^2,$$

which means that if $\boldsymbol{Z} \sim N_n(\boldsymbol{0}, \sigma^2 \boldsymbol{I})$, then

$$\frac{\boldsymbol{Z}^T \boldsymbol{D} \boldsymbol{Z}}{\sigma^2} \sim \chi_r^2$$

In this case we have $\boldsymbol{D} = (\boldsymbol{I} - \boldsymbol{H})$ symmetric and idempotent with rank n - p, and $\boldsymbol{Y} \sim N_n(\boldsymbol{X}\boldsymbol{\beta}, \sigma^2 \boldsymbol{I})$. To use the theorem we need to look at $\boldsymbol{Y}^* = \boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta} \sim N_n(\boldsymbol{0}, \sigma^2 \boldsymbol{I})$.

$$(I - H)Y^* = (I - H)(Y - X\beta)$$

= $(I - H)Y - (I - H)X\beta = (I - H)Y - (X\beta - HX\beta) =$
= $(I - H)Y - (X\beta - X\beta) = (I - H)Y$

since $HX = X(X^TX)^{-1}X^TX = X$. Projecting X onto the space spanned by the columns of X gives X.

Thus, we have shown that $\mathbf{Y}^T(\mathbf{I} - \mathbf{H})\mathbf{Y} = \mathbf{Y}^{*T}(\mathbf{I} - \mathbf{H})\mathbf{Y}^*$, and we may use the theorem to conclude that $\text{SSE}/\sigma^2 \sim \chi^2_{n-p}$.

The mean of a χ^2 -distributed variable equals the number of degrees of freedom, so

$$E\left(\frac{SSE}{\sigma^2}\right) = n - p$$
$$E(SSE) = (n - p)\sigma^2$$
$$E\left(\frac{SSE}{n - p}\right) = \sigma^2$$

Thus, $\hat{\sigma}^2 = \frac{\text{SSE}}{n-p}$ will be an unbiased estimator for σ^2 . The variance of a χ^2 -distributed random variable is twice the number of degrees of freedom, so the variance of $\hat{\sigma}^2$ is

$$\operatorname{Var}(\hat{\sigma}^2) = \operatorname{Var}\left(\frac{\operatorname{SSE}}{n-p}\right) = \frac{1}{(n-p)^2} \operatorname{Var}(\operatorname{SSE})$$
$$= \frac{1}{(n-p)^2} \operatorname{Var}\left(\sigma^2 \frac{\operatorname{SSE}}{\sigma^2}\right) = \frac{\sigma^4}{(n-p)^2} \operatorname{Var}\left(\frac{\operatorname{SSE}}{\sigma^2}\right)$$
$$= \frac{1}{(n-p)^2} 2(n-p)\sigma^4 = \frac{2\sigma^4}{n-p}.$$

c) We consider the two matrices

$$\boldsymbol{A} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T$$
 and $\boldsymbol{B} = \boldsymbol{I} - \boldsymbol{H} = \boldsymbol{I} - \boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T$

and the corresponding transformed random variables $Z_1 = AY$ and $Z_2 = BY$. Here A is a $p \times n$ matrix (since X is $n \times p$), and B is the same dimension as H, that is, $n \times n$, and is symmetric and idempotent (found previously).

Since $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ then $\mathbf{A}\mathbf{Y}$ and $\mathbf{B}\mathbf{Y}$ are independent random variables if $\sigma^2 \mathbf{A}\mathbf{B}^T = \mathbf{0}$.

$$\begin{split} \boldsymbol{A}\boldsymbol{B}^T &= (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T(\boldsymbol{I} - \boldsymbol{X}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T) \\ &= (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T - (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{X}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T \\ &= (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T - (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T = 0 \end{split}$$

We have proven that $Z_1 = AY$ and $Z_2 = BY$ are independent random variables. Then it follows that Z_1 and $Z_2^T Z_2$ are also independent random variables. Since $Z_1 = \hat{\beta}$ and $Z_2^T Z_2 = Y^T BY = SSE$, we have proven that $\hat{\beta}$ and SSE are independent random variables.

Application to multiple linear regression: The independence of $\hat{\beta}$ and SSE is used in the construction of a *t*-test for hypotheses about β .

Problem 4: Weighted linear regression

(Exam V2014 Problem 4)

a) Let (λ_i, e_i) , i = 1, ..., p be the eigenvalues and eigenvectors of V. Let P be the $(p \times p)$ matrix of eigenvectors,

$$oldsymbol{P} = [oldsymbol{e}_1 oldsymbol{e}_2 \cdots oldsymbol{e}_p]$$

and Λ be a diagonal matrix with the eigenvalues $\lambda_1, \lambda_1, ..., \lambda_p$ on the diagonal. Then $V^{-\frac{1}{2}}$ is defined as

$$oldsymbol{V}^{-rac{1}{2}} = oldsymbol{P} oldsymbol{\Lambda}^{-rac{1}{2}} oldsymbol{P}^T$$

Observe that $V^{-\frac{1}{2}}$ is symmetric, and that $V^{-\frac{1}{2}}V^{-\frac{1}{2}} = V^{-1}$.

$$egin{aligned} &oldsymbol{Y} = oldsymbol{X}eta + arepsilon \ &oldsymbol{V}^{-rac{1}{2}}oldsymbol{Y} = oldsymbol{V}^{-rac{1}{2}}oldsymbol{X}eta + oldsymbol{V}^{-rac{1}{2}}arepsilon \ &oldsymbol{Y}^* = oldsymbol{X}^*eta + arepsilon^* \end{aligned}$$

where $\boldsymbol{\varepsilon}^* \sim N(\boldsymbol{0}, \sigma^2 \boldsymbol{I})$. To see this calculate $\operatorname{Cov}(\boldsymbol{\varepsilon}^*) = \boldsymbol{V}^{-\frac{1}{2}} \operatorname{Cov}(\boldsymbol{\varepsilon}) \boldsymbol{V}^{-\frac{1}{2}} = \boldsymbol{V}^{-\frac{1}{2}} \sigma^2 \boldsymbol{V} \boldsymbol{V}^{-\frac{1}{2}} = \sigma^2 \boldsymbol{I}$.

We have now the ordinary least squares problem in the new quantities Y^* , X^* and ε^* , and know that the least squares solution is

$$\begin{split} \tilde{\boldsymbol{\beta}} &= (\boldsymbol{X}^{*T} \boldsymbol{X}^*)^{-1} \boldsymbol{X}^{*T} \boldsymbol{Y}^* \\ &= (\boldsymbol{X}^T \boldsymbol{V}^{-\frac{1}{2}} \boldsymbol{V}^{-\frac{1}{2}} \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{V}^{-\frac{1}{2}} \boldsymbol{V}^{-\frac{1}{2}} \boldsymbol{Y} \\ &= (\boldsymbol{X}^T \boldsymbol{V}^{-1} \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{V}^{-1} \boldsymbol{Y} \end{split}$$

Mean:

$$\begin{split} \mathbf{E}(\hat{\boldsymbol{\beta}}) &= \mathbf{E}((\boldsymbol{X}^{*T}\boldsymbol{X}^{*})^{-1}\boldsymbol{X}^{*T}\boldsymbol{Y}^{*}) = (\boldsymbol{X}^{*T}\boldsymbol{X}^{*})^{-1}\boldsymbol{X}^{*T}\mathbf{E}(\boldsymbol{Y}^{*}) \\ &= (\boldsymbol{X}^{*T}\boldsymbol{X}^{*})^{-1}\boldsymbol{X}^{*T}\boldsymbol{X}^{*}\boldsymbol{\beta} = \boldsymbol{\beta} \end{split}$$

since $E(\mathbf{Y}^*) = \mathbf{X}^* \boldsymbol{\beta}$.

The ordinary least square estimator $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ is unbiased in this model since the mean of \mathbf{Y} doesn't depend on \mathbf{V} .

$$\begin{split} \mathbf{E}(\hat{\beta}) &= \mathbf{E}((\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{Y}) = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \, \mathbf{E}(\boldsymbol{Y}) \\ &= (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{\beta} = \boldsymbol{\beta} \end{split}$$

If we just look at unbiasedness it may appear that the two estimators are equally good. However, since $\tilde{\boldsymbol{\beta}}$ is the least squares estimator (from looking at transformed quantities) we may conclude using the Gauss-Markov Theorem (p 181 in Fahrmeir et al (2013)) that $\tilde{\boldsymbol{\beta}}$ has the minimum variance in each component among all the unbiased estimators, BLUE. If we had calculated the covariance matrices of $\tilde{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{\beta}}$, we should see this. Thus, $\tilde{\boldsymbol{\beta}}$ should be preferred. Another issue is the fact that \boldsymbol{V} seldom is known, and need to be estimated. The concept of BLUE is handled in detail in our Statistical Inference course.

b) Find the expected value and covariance matrix of $\hat{\alpha}_1$ under the true model

$$\begin{split} \mathbf{E}(\hat{\boldsymbol{\alpha}_{1}}) &= \mathbf{E}((\boldsymbol{X}_{1}^{T}\boldsymbol{X}_{1})^{-1}\boldsymbol{X}_{1}^{T}\boldsymbol{Y}) \\ &= (\boldsymbol{X}_{1}^{T}\boldsymbol{X}_{1})^{-1}\boldsymbol{X}_{1}^{T}\mathbf{E}(\boldsymbol{Y}) == (\boldsymbol{X}_{1}^{T}\boldsymbol{X}_{1})^{-1}\boldsymbol{X}_{1}^{T}(\boldsymbol{X}_{1}\boldsymbol{\beta}_{1} + \boldsymbol{X}_{2}\boldsymbol{\beta}_{2}) \\ &= \boldsymbol{\beta}_{1} + (\boldsymbol{X}_{1}^{T}\boldsymbol{X}_{1})^{-1}\boldsymbol{X}_{1}^{T}\boldsymbol{X}_{2}\boldsymbol{\beta}_{2} \end{split}$$

Thus, $\hat{\boldsymbol{\alpha}}_1$ is a biased estimator for $\boldsymbol{\beta}_1$.

$$Cov(\hat{\boldsymbol{\alpha}}_1) = Cov((\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \boldsymbol{Y})$$

= $(\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T Cov(\boldsymbol{Y}) \boldsymbol{X}_1 (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1}$
= $(\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \sigma^2 \boldsymbol{I} \boldsymbol{X}_1 (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1}$
= $\sigma^2 (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1}$

Observe, $\operatorname{Cov}(\hat{\boldsymbol{\alpha}}_1)$ is not dependent on $\boldsymbol{\beta}_2$.

We see that the bias term for $\hat{\boldsymbol{\alpha}}_1$ is $(\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \boldsymbol{X}_2 \boldsymbol{\beta}_2$. When is the bias term equal to zero? When $\boldsymbol{\beta}_2 = \boldsymbol{0}$ there is no bias (but that is not so exciting). The bias is also zero when $\boldsymbol{X}_1^T \boldsymbol{X}_2 = \boldsymbol{0}$. This will happen if the two matrices are orthogonal. In Part 4: DOE this will be useful to know.