



Tentative Solutions to TMA4267 Linear Statistical Models 22 May 2014

Problem 1 Random vector

a)

$$\begin{aligned} E(\mathbf{Y}) &= \mathbf{A}\boldsymbol{\mu} = \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \\ \text{Cov}(\mathbf{Y}) &= \mathbf{A} \text{Cov}(\mathbf{X}) \mathbf{A}^T = \mathbf{A} \mathbf{I} \mathbf{A}^T = \mathbf{A} \mathbf{A}^T \\ &= \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{pmatrix} \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{pmatrix} \\ &= \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{pmatrix} = \mathbf{A} \end{aligned}$$

The covariance between X_1 and X_2 is zero. This doesn't imply that X_1 and X_2 are independent, unless the vector with elements X_1 and X_2 are binormal. We have no information about the distribution of \mathbf{X} , and can't assume that zero covariance implies independence.

For Y_1 and Y_2 the covariance is $-\frac{1}{3}$, and Y_1 and Y_2 is dependent.

Find the expected value of $\mathbf{X}^T \mathbf{A} \mathbf{X}$.

We may use the trace-formula:

$$E(\mathbf{X}^T \mathbf{A} \mathbf{X}) = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} = \text{tr} \mathbf{A} + \boldsymbol{\mu}^T \mathbf{0} = 3 \cdot \frac{2}{3} = 2$$

b) \mathbf{A} is clearly symmetric, which we can see by $\mathbf{A}^T = \mathbf{A}$. A projection matrix is an idempotent matrix, that is, $\mathbf{A}\mathbf{A} = \mathbf{A}$. We have already seen this in a).

The rank of a symmetric idempotent matrix equals its trace, which we found in **a)** to be $\text{tr}(\mathbf{A}) = 2$.

Derive the distribution of $\mathbf{X}^T \mathbf{A} \mathbf{X}$.

From T3.26 in BF2010 we know that if $\mathbf{X} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ and \mathbf{A} is a symmetric, idempotent matrix with rank r then $\mathbf{X}^T \mathbf{A} \mathbf{X} \sim \sigma^2 \chi_r^2$.

We have that $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{I})$, and \mathbf{A} is symmetric and idempotent with rank 2. We need to rewrite our expression so that we have a normally distributed random vector with mean zero and identity covariance matrix.

We subtract the mean and write

$$(\mathbf{X} - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{X} - \boldsymbol{\mu}) = \mathbf{X}^T \mathbf{A} \mathbf{X} - \boldsymbol{\mu}^T \mathbf{A} \mathbf{X} - \mathbf{X}^T \mathbf{A} \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} = \mathbf{X}^T \mathbf{A} \mathbf{X}$$

since $\mathbf{A} \boldsymbol{\mu} = \mathbf{0}$. Define $\mathbf{Z} = \mathbf{X} - \boldsymbol{\mu}$, where $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$. We may thus use the above theorem, to find that $\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{Z}^T \mathbf{A} \mathbf{Z} \sim \chi_2^2$, that is χ^2 with 2 degrees of freedom.

Tabeller og formler i statistikk, page 5, we see that 6 is the critical value in the χ^2 -distribution with 2 degrees of freedom and probability 0.05 (value in table is 5.991). The probability that the quadratic form is smaller than 6 is thus 95%.

Problem 2 Galapagos species

- a) The fitted regression model is:

$$\widehat{\text{Species}} = 7.07 - 0.02 \cdot \text{Area} + 0.32 \cdot \text{Elevation} + 0.009 \cdot \text{Nearest} - 0.24 \cdot \text{Scruz} - 0.75 \cdot \text{Adjacent}$$

This model explains 77% of the variability in the data. The regression is significant (the hypothesis that all regression coefficients are zero is rejected) and t -tests claim that **Elevation** and **Adjacent** are significant covariates.

The residual plots: The plot of studentized residuals vs. fitted values hints to heteroscedasticity in the errors (differing variances), and the qq-plot shows deviance from the normal distribution in the tails. The latter is also observed by looking at the Anderson-Darling normality test, which gives a p -value of 0.0002 (reject the null hypothesis that the errors are normal). The Box-Cox plot doesn't include 1 in the 95% confidence interval (dotted lines in the plot), and suggests that the cube root transform ($\lambda = 1/3$) may be suitable as a variance stabilizing transform.

- b) Let us assume that we have p covariates plus the intercept (notation from Part 6-7 of the course, in Parts 1-5 we have used p to include the intercept).

Estimate (Intercept): $t \text{ value} \cdot \text{Std. Error} = 7.365 \cdot 0.305 = 2.25$

Meaning: estimate for the regression coefficient. Intercept associated with first column of design matrix (first column of ones) $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$.

p -value of Area: two tails of t -distribution with 24 degrees of freedom. Can't find precise value, but from table 4 of Tabeller og formler i statistikk we see that the critical value in the t -distribution with 24 degrees of freedom is 2.064 for $\alpha = 0.025$. This means that the p -value will approximately be 0.05.

Meaning: Test the null hypothesis that $\beta_{Area} = 0$ vs. $\beta_{Area} \neq 0$, with the other four covariates present in the model, and produce a p -value of the test.

Std. Error of Nearest: $\text{estimate} / \text{tobs} = 0.012 / 0.7 = 0.017$

Meaning: the estimated standard deviation for the regression estimate. Mathematically the corresponding (Nearest) diagonal element of the square root of $(\mathbf{X}^T \mathbf{X})^{-1} s^2$, where s^2 is the estimate for the regression variance σ^2 .

Adjusted R-squared: $1 - (1 - R^2)(n - 1) / (n - p - 1) = 1 - (1 - 0.7543) \cdot 29 / 24 = 0.7032$, or in a two stage process by first observing $\text{SSE} = s^2 \cdot (n - p) = 0.9716^2 \cdot 24 = 22.65$ and then finding SST from $R^2 = 1 - \text{SSE} / \text{SST}$, $\text{SST} = \text{SSE} / (1 - R^2)$, and finally using $R_{\text{adj}}^2 = 1 - \frac{\frac{\text{SSE}}{n-p-1}}{\frac{\text{SST}}{n-1}}$.

Yes, I would prefer model B to A. The plot of standardized residuals vs fitted values shows no clear structure, and the qq-plot looks much better for B and A. The Anderson-Darling normality test doesn't reject the null hypothesis of normal data.

- c) Let SSE be the sum-of-squares of error, SSR be the regression sum-of-squares, and SST be the total som of squares. Then R^2 : coefficient of multiple determination is defined as

$$1 - \text{SSE} / \text{SST} = \text{SSR} / \text{SST}$$

and is interpreted as the amount of variability in the data that is accounted for by the regression. R^2 will increase when a regressors are added to the model, even if the new regressors are independent of the response. Why? The least squares estimator will minimize SSE and if the regression coefficient for the new regressor is estimated to be a value different from zero, this means that the SSE of this larger model will be smaller than the SSE of the smaller model.

The R_{adj}^2 is constructed to also include information about the number of parameters estimated and the number of observations in the data set. Assume we have 1 intercept and in addition p regression parameters, then

$$R_{\text{adj}}^2 = 1 - \frac{\frac{\text{SSE}}{n-p-1}}{\frac{\text{SST}}{n-1}}$$

R^2 will always increase then new covariates are added to the model, so R^2 can only be used to select the best model among models with the same number of covariates. This

is done when in best subset selection one model is reported for each total number of covariates. To choose between these models a criterion taking into account the number of covariates in the model need to be used, and one such criterion is R_{adj}^2 . We therefore use R_{adj}^2 to choose between the best models of each size.

In our example the best model is according to this strategy the model with four covariates. These are all covariates except **Nearest**. To write down the estimated regression equation we need to refit this model.

Lasso regression adds a penalty term to the least squares criterion to make the model more sparse. This may give a robust fit and avoid overfitting. The penalty term for lasso is the sum of the absolute value of the regression coefficients, and the optimization procedure is to minimize with respect to β the following quantity

$$(\mathbf{Y} - \mathbf{X}\beta)^T(\mathbf{Y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j|$$

The difficulty with using R_{adj}^2 is that the number of parameters that enter into the formula is not defined since the lasso regression shrinks parameter estimates - some by a certain amount and some all the way down to zero. The model complexity for the lasso is defined from λ and we can't put λ into the R_{adj}^2 -formula.

Instead cross-validation is used for model choice.

The following was not asked for: SST is the same for all models, and thus R^2 is proportional to SSE. SSE can therefore be used as a model selection criterion in the cross-validation. In short - this involves: 10-fold, training set, test set, fit, predict, SSE on the test set, and sum to produce SSE on the whole set. Choose the penalty parameter with the smallest SSE.

The fitted lasso regression model is:

$$\widehat{\text{Species}}^{1/3} = 3.54 + 0.00028 \cdot \text{Elevation}$$

Problem 3 Design of experiments

A column for BC is added to the design,

	A	B	C	D	BC
1	-1	-1	-1	1	1
2	1	-1	-1	1	1
3	-1	1	-1	-1	-1
4	1	1	-1	-1	-1
5	-1	-1	1	-1	-1
6	1	-1	1	-1	-1
7	-1	1	1	1	1
8	1	1	1	1	1

a) What type of experiment is this?

We see that we have a full factorial design in the factors A, B, C, but there is a fourth factor D added. This is a half fraction of a 2^4 design, also called a 2^{4-1} -design.

What is the generator and the defining relation for the experiment?

The generator for the design is $D=BC$ (which is seen from the table above after the BC column is added). The defining relation is then $I=BCD$.

What is the resolution of the experiment?

The resolution of the design equals the number of letters in the defining relation, thus the resolution is III.

Write down the alias structure of the experiment.

$A=ABCD$, $B=CD$, $C=BD$, $D=BC$.

$AB=ACD$, $AC=ABD$, $AD=ABC$

$I=BCD$

Problem 4 Multiple Linear Regression

a) Let $(\lambda_i, \mathbf{e}_i)$, $i = 1, \dots, p$ be the eigenvalues and eigenvectors of \mathbf{V} . Let \mathbf{P} be the $(p \times p)$ matrix of eigenvectors,

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

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

$$\mathbf{V}^{-\frac{1}{2}} = \mathbf{P} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{P}^T$$

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

$$\begin{aligned}\mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \mathbf{V}^{-\frac{1}{2}}\mathbf{Y} &= \mathbf{V}^{-\frac{1}{2}}\mathbf{X}\boldsymbol{\beta} + \mathbf{V}^{-\frac{1}{2}}\boldsymbol{\varepsilon} \\ \mathbf{Y}^* &= \mathbf{X}^*\boldsymbol{\beta} + \boldsymbol{\varepsilon}^*\end{aligned}$$

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

We have now the ordinary least squares problem in the new quantities \mathbf{Y}^* , \mathbf{X}^* and $\boldsymbol{\varepsilon}^*$, and know that the least squares solution is

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

Mean:

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

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

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

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

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 (T3.13 in Bingham and Fry 2010) 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 \mathbf{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

$$\begin{aligned} E(\hat{\alpha}_1) &= E((\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{Y}) \\ &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T E(\mathbf{Y}) = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T (\mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2) \\ &= \boldsymbol{\beta}_1 + (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \boldsymbol{\beta}_2 \end{aligned}$$

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

$$\begin{aligned} \text{Cov}(\hat{\alpha}_1) &= \text{Cov}((\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{Y}) \\ &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \text{Cov}(\mathbf{Y}) \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \\ &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \sigma^2 \mathbf{I} \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \\ &= \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \end{aligned}$$

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

We see that the bias term for $\hat{\alpha}_1$ is $(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \boldsymbol{\beta}_2$. When is the bias term equal to zero?

When $\boldsymbol{\beta}_2 = \mathbf{0}$ there is no bias (but that is not so exciting). The bias is also zero when $\mathbf{X}_1^T \mathbf{X}_2 = \mathbf{0}$. This will happen if the two matrices are orthogonal. If we think back to Part 6: DOE this is true for DOE, and is why the interpretation of the DOE coefficients is easy and the model chosen doesn't influence the unbiasedness of the coefficients - but influences the variance thereof.