



Tentative Solutions to TMA4267 Linear Statistical Models August 2014

Problem 1 The Multivariate Normal Distribution

- a) Since \mathbf{Y} can be written as $\mathbf{A}\mathbf{X}$ where $\mathbf{A} = \begin{pmatrix} 3 & -2 \\ 1 & 1 \end{pmatrix}$, we see that \mathbf{Y} is a bivariate vector of linear combinations of \mathbf{X} . Since \mathbf{X} is multivariate normal, then also \mathbf{Y} will be multivariate normal. The mean and covariance of \mathbf{Y} is given as:

$$\mathbf{E}(\mathbf{Y}) = \mathbf{A}\boldsymbol{\mu} = \begin{pmatrix} 3 & -2 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} -1 \\ 3 \end{pmatrix}$$

$$\text{Cov}(\mathbf{Y}) = \mathbf{A} \text{Cov}(\mathbf{X}) \mathbf{A}^T = \begin{pmatrix} 3 & -2 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix} \begin{pmatrix} 3 & 1 \\ -2 & 1 \end{pmatrix} = \begin{pmatrix} 11 & -0.5 \\ -0.5 & 4 \end{pmatrix}$$

Since $\begin{pmatrix} Z \\ Y_2 \end{pmatrix}$ is bivariate normal (same reason as for \mathbf{Y} above), then Z and Y_2 are independent if $\text{Cov}(Z, Y_2) = 0$. Let $\begin{pmatrix} Z \\ Y_2 \end{pmatrix}$ be written as $\mathbf{B}\mathbf{X}$ where $\mathbf{B} = \begin{pmatrix} 1 & a \\ 1 & 1 \end{pmatrix}$.

$$\text{Cov} \begin{pmatrix} Z \\ Y_2 \end{pmatrix} = \mathbf{B} \text{Cov}(\mathbf{X}) \mathbf{B}^T = \begin{pmatrix} 1 & a \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ a & 1 \end{pmatrix} = \begin{pmatrix} 2a^2 + a + 1 & \frac{3}{2} + \frac{5}{2}a \\ \frac{3}{2} + \frac{5}{2}a & 4 \end{pmatrix}$$

Thus

$$\begin{aligned} \text{Cov}(Z, Y_2) &= \frac{3}{2} + \frac{5}{2}a = 0 \\ a &= -\frac{3}{5} = -0.6 \end{aligned}$$

b)

$$\begin{aligned} f(\mathbf{x}) &= ce^{-\frac{1}{2}Q(x_1, x_2)} \\ c &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \\ Q(x_1, x_2) &= (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \end{aligned}$$

where $\mu_1 = 1$, $\mu_2 = 2$, $\sigma_1 = 1$, $\sigma_2 = \sqrt{2}$ and $\rho = \text{Cov}(X_1, X_2)/(\sigma_1\sigma_2) = \frac{1}{2\sqrt{2}}$, so that $c = 0.12$.

$$\begin{aligned} f(\mathbf{x}) &= ce^{-\frac{1}{2}Q(x_1, x_2)} = d \\ Q(x_1, x_2) &= 2(\ln(c) - \ln(d)) \\ (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= 2(\ln(c) - \ln(d)) = D^2 \end{aligned}$$

This is an ellipse with principal axes along the eigenvectors of the covariance matrix $\boldsymbol{\Sigma}$, and half-lengths $D\sqrt{\lambda_i}$. See exercise E1P2 for details on this result.

Moreover, $(\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu})$ is distributed as χ_2^2 . And, the solid ellipsoid of \mathbf{x} values satisfying

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \leq \chi_2^2(\alpha)$$

has probability $1 - \alpha$. Let $\alpha = 0.05$, and we find that $\chi_2^2(0.05) = 5.99$. So, $2(\ln(c) - \ln(d)) = 5.99$. Solving for d gives $d = \exp(\ln(c) - 5.99/2) = 0.006$.

In the drawing the halflengths are given as $\sqrt{5.99\lambda_i}$, where $\lambda_1 = 2.2$ and $\lambda_2 = 0.8$ from the R print-out. Halflengths are then 3.6 and 2.2. Thus, the points on the ellipse at the principal axes are for the first principal axis $(-0.4, -1.4)$ and $(2.4, 5.4)$, and for the second principal axis $(-1.0, 2.8)$ and $(3.0, 1.2)$.

Problem 2 Predicting fat content in meat

a) The null- and alternative hypotheses are:

$$H_0 : \beta_{100} = 0 \text{ vs. } H_1 : \beta_{100} \neq 0$$

The t -statistic is related to a t -distribution with $n - p = 215 - 101 = 114$ degrees of freedom. The value of the t -statistics is -0.343. The 2.5% quantile of the t -distribution with 114 degrees of freedom is approximately 1.98 from the statistical tables. This means that the p -value of the test will be much larger than 0.05 and we will not reject the null hypothesis.

How would you judge the model fit?

The regression is significant, and explains 95% of the variability in the data. A number of the 100 coefficients looks to be significant (from Figure 4). The residual plots shows no obvious trend, but the tail behaviour of the qq-plot looks a bit off the normal line. However, the Andreson-Darling test will not reject the null hypothesis of normal data. So, the model fit seems ok wrt residual analysis.

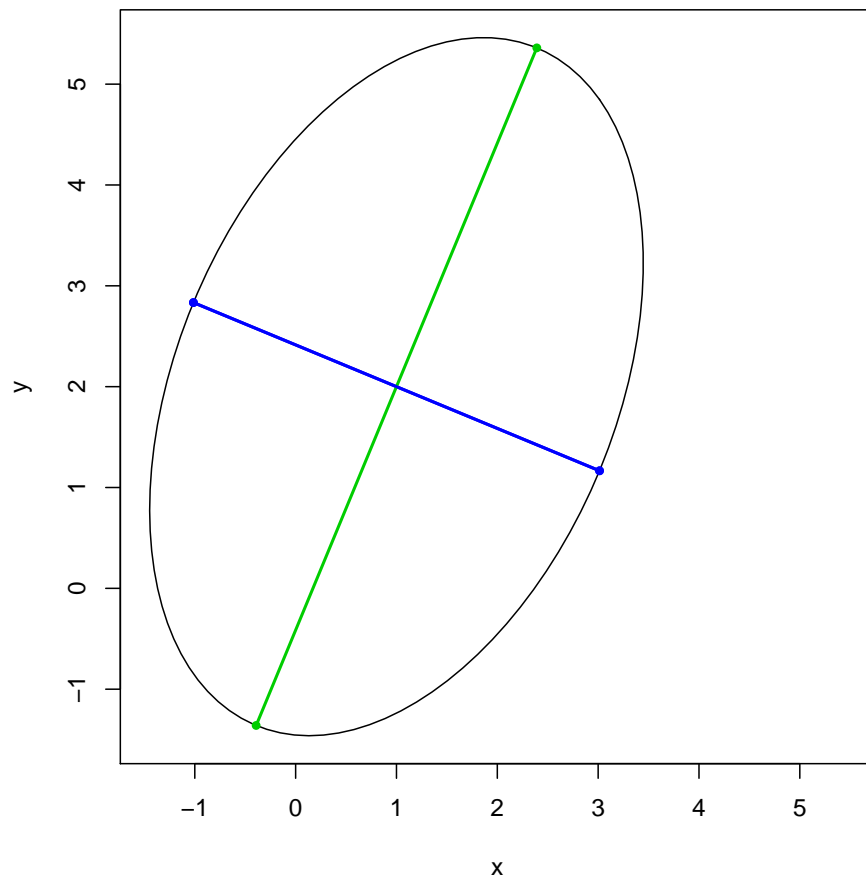


Figure 1: Contour in Problem 1b

What is overfitting? When there are many covariates as compared to observations (yes, 100 covariates is much compared to 215 observations) we have the potential problem of not just fitting the signal in the data, but also the noise. Fitting the data noise is called overfitting. Might that be a problem here? Yes, possibly. From the plot of the estimated coefficients (Figure 4) we see that some of the estimated coefficients are very large, which might point towards overfitting.

- b) What is the mathematical definition of the principal component loadings and scores?

Let \mathbf{X} be the random variable under study and let Σ be the population covariance matrix. Further, let $(\lambda_i, \mathbf{e}_i)$ be eigenvalue/vector pairs from a spectral decomposition of the covariance matrix. The pairs are ordered so that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. The linear combination $\mathbf{e}_i^T \mathbf{X}$ is the i th principal component, and the entries in \mathbf{e}_i are called the loadings and the values $\mathbf{e}_i^T \mathbf{X}$ (with the data here) are called the scores.

Alternatively: The principal components are uncorrelated linear combinations of the original variables, found to maximize the variance of the components. The PCA loadings are the weights in the linear combinations, and are given as the eigenvectors of the covariance or correlation matrix of the population (or data) under study. The PCA scores are the numerical values of the linear combinations when applied to the samples. If the number of principal components are p there are p scores for each individual sample.

In our data the first principal component have nearly equal loadings for all original covariates, and can be seen as an average effect. The second principal component gives a high value for the first absorbances, then low, then high and then low again. The third principal component gives positive values for the first absorbances, and then negative values for the last absorbances. Observe the high degree of smoothness in the loadings.

The percentage of total variance explained by the first three principal components are 99.875%.

The PC scores can be used as covariates into a regression, taking the place of the original covariates. If all scores are used the regression in the PCs will be mathematically the same as the regression in the original variables. If less than p scores are used, then this will result in a shrinkage effect on the coefficients, similar to performing a ridge regression analysis. Using a few PCs as covariates might help towards possible overfitting (and also towards multicollinearity).

Problem 3 Design of experiments

- a) Design matrix:

	A	B	C	D	AB
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

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 $D=AB$ (which is seen from the table above after the AB column is added). The defining relation is then $I=ABD$.

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.

$A=BD$, $B=AD$, $C=ABCD$, $D=AB$

$AC=BCD$, $BC=ACD$, $CD=ABC$

$I=ABD$

Problem 4 Multiple linear regression

Define the matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

- a) \mathbf{H} is a symmetric projection matrix, since $\mathbf{H}^T = \mathbf{H}$ and $\mathbf{H}\mathbf{H} = \mathbf{H}$. For a symmetric and idempotent matrix the rank is equal to the trace. The rank of \mathbf{H} is p . See proof below.

$$\begin{aligned}
 \mathbf{H} &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\
 \mathbf{H}^T &= (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)^T = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H} \\
 \mathbf{H}^2 &= (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H} \\
 \text{tr}(\mathbf{H}) &= \text{tr}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \text{tr}(\mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}) = \text{tr}(\mathbf{I}_{p \times p}) = p
 \end{aligned}$$

Graphically: The vector $\mathbf{H}\mathbf{Y}$ is a projection of the vector \mathbf{Y} onto the space spanned by the columns of \mathbf{X} .

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

$$\begin{aligned}(\mathbf{I} - \mathbf{H})^T &= \mathbf{I} - \mathbf{H}^T = \mathbf{I} - \mathbf{H} \\(\mathbf{I} - \mathbf{H})^2 &= (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) = \mathbf{I} - 2\mathbf{H} + \mathbf{H}^2 = \mathbf{I} - \mathbf{H} \\ \text{tr}(\mathbf{I} - \mathbf{H}) &= \text{tr}(\mathbf{I}_{n \times n}) - \text{tr}(\mathbf{H}) = n - p\end{aligned}$$

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

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

One of the key theorems of this course (theorem 3.26 in Bingham&Fry), state that if \mathbf{D} is a symmetric and idempotent matrix with rank r and $\mathbf{Z} \sim N_n(0, \sigma^2 \mathbf{I})$, then $\mathbf{Z}^T \mathbf{D} \mathbf{Z} \sim \sigma^2 \chi_r^2$.

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

$$\begin{aligned}(\mathbf{I} - \mathbf{H})\mathbf{Y}^* &= (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ &= (\mathbf{I} - \mathbf{H})\mathbf{Y} - (\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = (\mathbf{I} - \mathbf{H})\mathbf{Y} - (\mathbf{X}\boldsymbol{\beta} - \mathbf{H}\mathbf{X}\boldsymbol{\beta}) = \\ &= (\mathbf{I} - \mathbf{H})\mathbf{Y} - (\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{I} - \mathbf{H})\mathbf{Y}\end{aligned}$$

since $\mathbf{H}\mathbf{X} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} = \mathbf{X}$. Projecting \mathbf{X} onto the space spanned by the columns of \mathbf{X} gives \mathbf{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 $\text{SSE} \sim \sigma^2 \chi_{n-p}^2$.

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

$$\begin{aligned}\text{E}\left(\frac{\text{SSE}}{\sigma^2}\right) &= n - p \\ \text{E}(\text{SSE}) &= (n - p)\sigma^2 \\ \text{E}\left(\frac{\text{SSE}}{n - p}\right) &= \sigma^2\end{aligned}$$

Thus, $\hat{\sigma}^2 = \frac{\text{SSE}}{n-p}$ will be an unbiased estimator for σ^2 .

Variance:

$$\begin{aligned}\text{Var}\left(\frac{\text{SSE}}{n-p}\right) &= \frac{1}{(n-p)^2} \text{Var}(\text{SSE}) = \frac{1}{(n-p)^2} \text{Var}\left(\sigma^2 \frac{\text{SSE}}{\sigma^2}\right) \\ &= \frac{1}{(n-p)^2} 2(n-p)\sigma^4 = \frac{2\sigma^4}{n-p}\end{aligned}$$

c)

$$\begin{aligned}\mathbf{A} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{B} &= \mathbf{I} - \mathbf{H} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\end{aligned}$$

Here \mathbf{A} is a $p \times n$ matrix (since \mathbf{X} is $n \times p$), and \mathbf{B} is the same dimension as \mathbf{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 $\sigma^2 \mathbf{A}\mathbf{B}^T = \mathbf{0}$.

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

We have proven that $\mathbf{Z}_1 = \mathbf{A}\mathbf{Y}$ and $\mathbf{Z}_2 = \mathbf{B}\mathbf{Y}$ are independent random variables. Then it follows that \mathbf{Z}_1 and $\mathbf{Z}_2^T \mathbf{Z}_2$ are also independent random variables. Since $\mathbf{Z}_1 = \hat{\boldsymbol{\beta}}$ and $\mathbf{Z}_2^T \mathbf{Z}_2 = \mathbf{Y}^T \mathbf{B}\mathbf{Y} = \text{SSE}$, we have proven that $\hat{\boldsymbol{\beta}}$ and SSE are independent random variables.

Use in MLR: The independence of $\hat{\boldsymbol{\beta}}$ and SSE is used in the construction of a t -test for hypothesis about $\boldsymbol{\beta}$.