TMA4267 Linear Statistical Models V2017 (L9) Part 2: Linear regression: Parameter estimation [F:3.2]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 14, 2017

The classical linear model

$$\begin{array}{ll} \mathbf{Y} &= & \mathbf{X} \stackrel{\boldsymbol{\beta}}{}_{(n \times 1)} + \stackrel{\boldsymbol{\varepsilon}}{}_{(n \times 1)} \\ E(\boldsymbol{\varepsilon}) = & \mathbf{0} \\ (n \times 1) \end{array} \text{ and } \quad Cov(\boldsymbol{\varepsilon}) = \stackrel{\sigma^{2}}{}_{(n \times n)} \mathbf{I} \end{array}$$

where

• $\boldsymbol{\beta}$ and σ^2 are unknown parameters and

• the design matrix \boldsymbol{X} has full rank, with *i*th row $[x_{i1}x_{i2}\cdots x_{ip}]$. Today

- 1. find estimator for β ,
- 2. find estimator for σ^2 , and
- 3. look at two idempotent matrices H and I H to arrive at
- 4. geometric interpretation.

Rules for derivatives with respect to a vector

- Let β be a *p*-dimensional column vector of interest,
- and let ∂/∂β denote the p-dimensional vector with partial derivatives wrt the p elements of β.
- ▶ Let *d* be a *p*-dimensional column vector of constants and
- **D** be a $p \times p$ symmetric matrix of constants.

Rule 1:

$$rac{\partial}{\partialoldsymbol{eta}}(oldsymbol{d}^{ op}oldsymbol{eta}) = rac{\partial}{\partialoldsymbol{eta}}(\sum_{j=1}^p d_jeta_j) = oldsymbol{d}$$

Rule 2:

$$\frac{\partial}{\partial \boldsymbol{\beta}}(\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{D}\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}}(\sum_{j=1}^{p}\sum_{k=1}^{p}\beta_{j}d_{jk}\beta_{k}) = 2\boldsymbol{D}\boldsymbol{\beta}$$

See Härdle and Simes (2015), page 65, Equation (2.23) and (2.24).

Two questions

Have found least squares and maximum likelihood estimator for β :

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{\, au} oldsymbol{X})^{-1} oldsymbol{X}^{\, au} oldsymbol{Y}$$

and we have assumed that the rank(\boldsymbol{X}) = p for $n \times p$ design matrix (where n > p).

- Q1: What can we say about $X^T X$?
- Q2: Why is the following wrong?

Using $(AB)^{-1} = B^{-1}A^{-1}$,

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^Toldsymbol{X})^{-1}oldsymbol{X}^Toldsymbol{Y} = oldsymbol{X}^{-1}(oldsymbol{X}^T)^{-1}oldsymbol{X}^Toldsymbol{Y} = oldsymbol{X}^{-1}oldsymbol{Y}$$

The classical linear model

The model

$$m{Y} = m{X}m{eta} + m{arepsilon}$$

is called a classical linear model if the following is true:

1.
$$E(\varepsilon) = 0$$
.

2.
$$\operatorname{Cov}(\varepsilon) = \operatorname{E}(\varepsilon \varepsilon^{T}) = \sigma^{2} I.$$

3. The design matrix has full rank $rank(\mathbf{X}) = k + 1 = p$. The classical *normal* linear regression model is obtained if additionally

4. $\boldsymbol{\varepsilon} \sim N_n(\boldsymbol{0}, \sigma^2 \boldsymbol{I})$

holds. For random covariates these assumptions are to be understood conditionally on \boldsymbol{X} .

Acid rain

occurs when emissions of sulfur dioxide (SO2) and oxides of nitrogen (NOx) react in the atmosphere with water, oxygen, and oxidants to form various acidic compounds. These compounds then fall to the earth in either dry form (such as gas and particles) or wet form (such as rain, snow, and fog).



Source: http://myecoproject.org/get-involved/pollution/acid-rain/



http://www.eoearth.org/view/article/149814/

Acid rain in Norwegian lakes

Measured pH in Norwegian lakes explained by content of

- ▶ x1: SO₄: sulfate (the salt of sulfuric acid),
- ▶ x2: N0₃: nitrate (the conjugate base of nitric acid),
- x3: Ca: calsium,
- x4: latent AI: aluminium,
- x5: organic substance,
- x6: area of lake,
- ▶ x7: position of lake (Telemark or Trøndelag),

pH is a measure of the acidity of alkalinity of water, expressed in terms of its concentration of hydrogen ions. The pH scale ranges from 0 to 14. A pH of 7 is considered to be neutral. Substances with pH of less that 7 are acidic; substances with pH greater than 7 are basic.



http://www.eoearth.org/view/article/149814/





0=Telemark, 1=Trondelag

Acid rain data



Output from fitting the full model in R

```
> fit=lm(y~.,data=ds)
```

```
> summary(fit)
```

```
Coefficients:
```

	E	Stimate	Std.	Error	t valu	le P	r(> t	1)					
(Intercep	t) 5.	6764334	0.13	389162	40.86	52	< 2e-1	16	**	*			
x1	-0.	3150444	0.0	587512	-5.36	52 4	.27e-0)5	**	*			
x2	-0.	0018533	0.0	012587	-1.47	2	0.15	58					
x3	0.	9751745	0.14	449075	6.73	30 2	.62e-0	06	**	*			
x4	-0.	0002268	0.0	010038	-0.22	26	0.82	24					
x5	-0.	0334242	0.0	225009	-1.48	35	0.15	55					
x6	-0.	0039399	0.0	724339	-0.05	54	0.95	57					
x7	0.	0888722	0.10	025724	0.86	6	0.39	98					
Signif. c	odes:	0 ****	, 0.00	01 '**	0.01	, _* ,	0.05	'.	,	0.1	,	,	1

Residual standard error: 0.1165 on 18 degrees of freedom Multiple R-squared: 0.93,Adjusted R-squared: 0.9027 F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

Question: explain how to interpret $\hat{\beta}_0$ and $\hat{\beta}_3$.



3.10 Asymptotic Properties of the Least Squares Estimator

- 1. The least squares estimator $\hat{\beta}_n$ for β and the ML or REML estimator $\hat{\sigma}_n^2$ for the variance σ^2 are consistent.
- 2. The least squares estimator asymptotically follows a normal distribution, specifically

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \xrightarrow{d} \mathrm{N}(\boldsymbol{0}, \sigma^2 V^{-1}).$$

That is the difference $\hat{\beta}_n - \beta$ normalized with \sqrt{n} converges in distribution to the normal distribution on the right-hand side.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.120)

Projection matrix: definition and properties

- A matrix **A** is a *projection matrix* if it is idempotent, $\mathbf{A}^2 = \mathbf{A}$.
- An idempotent matrix is an orthogonal projection matrix if, in the decomposition of a vector, v = Av + (v − Av), Av and v − Av = (I − A)v are always orthogonal, that is, (Av)^T(v − Av) = 0.
- A symmetric projection matrix is orthogonal.
- The eigenvalues of a projection matrix are 0 and 1.
- If a (n × n) symmetric projection matrix A has rank r then r eigenvalues are 1 and n − r are 0.
- The trace and rank of a symmetric projection matrix are equal: tr(A) = rank(A).

Results so far

Least squares and maximum likelihood estimator for β:

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{T} \boldsymbol{X})^{-1} \boldsymbol{X}^{T} \boldsymbol{Y}$$

• Restricted maximum likelihood estimator for σ^2 :

$$\hat{\sigma^2} = \frac{1}{n-p} (\boldsymbol{Y} - \boldsymbol{X}\hat{\beta})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\beta}) = \frac{\mathsf{SSE}}{n-p}$$

Projection matrices: idempotent, symmetric/orthogonal:

$$H = X(X^T X)^{-1} X^T$$
$$I - H = I - X(X^T X)^{-1} X^T$$

with important connection:

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$

 $\hat{\mathbf{\varepsilon}} = \mathbf{I} - \mathbf{H}\mathbf{Y}$

Results from Mathematics 3

Best approximation theorem

The vector $\hat{\mathbf{Y}}$ in the column space of \mathbf{X} that makes $||\mathbf{Y} - \hat{\mathbf{Y}}||$ as small as possible, is the orthogonal projection of \mathbf{Y} on the column space of \mathbf{X} .

Orthogonal decomposition

We want $\hat{\boldsymbol{\beta}}$ to minimize $|| \boldsymbol{Y} - \hat{\boldsymbol{Y}} || = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$ (least squares principle).

The column space of X consists of vectors of the form $X\hat{\beta}$, so $X\hat{\beta}$ is the orthogonal projection of Y onto the column space of $X \cdot \hat{Y} = HY$, and $H = X(X^TX)^{-1}X^T$ projects onto the column space of X. Observe: HX = X.

This is equivalent to observing that $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is in the orthogonal complement of the column space of \mathbf{X} .

 $\hat{\varepsilon} = Y - HY = (I - H)Y$, and I - H projects onto the space orthogonal to the column space of X. Observe: (I-H)X=0

That is, $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is orthogonal to all columns of \mathbf{X} , so $\mathbf{X}^{T}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$ and $\mathbf{X}^{T}\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^{T}\mathbf{Y}$.



Putanen, Styan and Isotalo: Matrix Tricks for Linear Statistical Models: Our Personal Top Twenty, Figure 8.3.

3.7 Geometric Properties of the Least Squares Estimator

The method of least squares has the following geometric properties:

- 1. The predicted values \hat{y} are orthogonal to the residuals $\hat{\varepsilon}$, i.e., $\hat{y}'\hat{\varepsilon} = 0$.
- 2. The columns \mathbf{x}^{j} of X are orthogonal to the residuals $\hat{\mathbf{e}}$, i.e., $(\mathbf{x}^{j})'\hat{\mathbf{e}} = 0$ or $X'\hat{\mathbf{e}} = \mathbf{0}$.
- 3. The average of the residuals is zero, i.e.,

$$\sum_{i=1}^{n} \hat{\varepsilon}_i = 0 \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^{n} \hat{\varepsilon}_i = 0.$$

4. The average of the predicted values \hat{y}_i is equal to the average of the observed response y_i , i.e.,

$$\frac{1}{n}\sum_{i=1}^n \hat{y}_i = \bar{y}.$$

5. The regression hyperplane runs through the average of the data, i.e.,

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \dots + \hat{\beta}_k \bar{x}_k.$$

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.112) Alternative summery of Geometry of Least Squares

• Mean response vector: $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$

- As β varies, Xβ spans the model plane of all linear combinations. I.e. the space spanned by the columns of X: the column-space of X.
- Due to random error (and unobserved covariates), Y is not exactly a linear combination of the columns of X.
- LS-estimation chooses β̂ such that Xβ̂ is the point in the column-space of X that is closes to Y.
- ► The residual vector $\hat{\varepsilon} = \mathbf{Y} \hat{\mathbf{Y}} = (\mathbf{I} \mathbf{H})\mathbf{Y}$ is perpendicular to the column-space of \mathbf{X} .
- ► Multiplication by H = X(X^TX)⁻¹X^T projects a vector onto the column-space of X.
- Multiplication by I H = I X(X^TX)⁻¹X^T projects a vector onto the space perpendicular to the column-space of X.

Today

Least squares and maximum likelihood estimator for β:

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{Y}$$

has mean $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ and $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2 (\boldsymbol{X}^T \boldsymbol{X})^{-1}$.

- For the normal model: $\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T \boldsymbol{X})^{-1}).$
- Asymptotic properties of the least squares estimator: normality.
- Orthogonal projection matrices *H* and *I H* with geometric interpretation.

Next time: properties of residuals and $\hat{\sigma}^2$, confidence intervals and hypothesis testing for regression coefficients.