



NTNU
Norwegian University of
Science and Technology

TMA4267 Linear Statistical Models V2014 (13)
Multiple linear regression, properties of estimators [3.2-3.4]

Mette Langaas

To be lectured: February 17, 2014
wiki.math.ntnu.no/emner/tma4267/2014v/start/

Acid rain in Norwegian lakes

Measured pH in Norwegian lakes explained by content of

- x_1 : SO_4 : sulfate (the salt of sulfuric acid),
- x_2 : NO_3 : nitrate (the conjugate base of nitric acid),
- x_3 : Ca : calcium,
- x_4 : latent Al : aluminium,
- x_5 : organic substance,
- x_6 : area of lake,
- x_7 : position of lake (Telemark or Trøndelag),

TMA 4267 : Lecture 13, Ch 3.2-3.3-3.4

Multiple linear regression (MLR)

$$\begin{array}{ccccccc}
 Y & = & X & \beta & + & \varepsilon \\
 n \times 1 & & n \times p & p \times 1 & & n \times 1 \\
 | & & | & | & & | \\
 \text{response} & & \text{design} & \text{parameter} & & \text{error} \\
 & & \text{matrix} & \text{vector} & &
 \end{array}$$

i) where ε_i are i.i.d with $E(\varepsilon_i) = 0$ and $\text{Var}(\varepsilon_i) = \sigma^2$.

ii) And add normality $\varepsilon \sim N_n(0, \sigma^2 I)$.

Normal equations: $(X^T X) \hat{\beta} = X^T Y$

Estimator: $\hat{\beta} = (X^T X)^{-1} X^T Y$

Multiple linear regression model

$$Y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i, \text{ for } i = 1, \dots, n$$

Matrix formulation:

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

where

- $\boldsymbol{\beta}$ and σ^2 are unknown parameters and
- the design matrix \mathbf{X} has i th row $[x_{i1}, x_{i2}, \dots, x_{ip}]$.

Performing MLR in practice

Numerical note:

The normal equations can in practice be solved by using the QR decomposition of X . See book p 68-69 for details.

$$X = QR$$

$$(X^T X) b = X^T y \iff R b = Q^T y$$

upper triangular orthogonal

solved by
backsubstitution

R: try out the LIS.r script to analyse acid rain

Main function to use is "lm".

Finding : $\hat{\beta} = \begin{bmatrix} 5.68 \\ -0.32 \\ \vdots \\ 0.09 \end{bmatrix}$

Least squares estimation

- $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2\mathbf{I}$.
- Let \mathbf{X} has full rank $p \leq n$.
- The Least Squares Estimate (LSE) of $\boldsymbol{\beta}$ is given by

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

Properties of $\hat{\beta}$ (3.3)

First, not assume normality, only assume

$$E(\varepsilon) = 0 \text{ and } \text{Cov}(\varepsilon) = \sigma^2 \mathbf{I}.$$

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

$$E(Y) = X\beta \quad \text{Cov}(Y) = \sigma^2 \mathbf{I}$$

The LS estimator $\hat{\beta} = \underbrace{(X^T X)^{-1} X^T}_{G} Y = G Y$

1) Observe that $\hat{\beta}$ is linear in the data Y .

2) Mean: $E(\hat{\beta}) = E(GY) = G \underbrace{E(Y)}_{X\beta} = \underbrace{(X^T X)^{-1} X^T}_{\mathbf{I}} X\beta = \beta$

$\hat{\beta}$ is an unbiased estimator for β .

3) Covariance:

$$\begin{aligned} \text{Cov}(\hat{\beta}) &= \text{Cov}(GY) = G \underbrace{\text{Cov}(Y)}_{\sigma^2 \mathbf{I}} G^T \\ &= (X^T X)^{-1} X^T \sigma^2 \mathbf{I} [(X^T X)^{-1} X^T]^T \\ &= \sigma^2 \underbrace{(X^T X)^{-1} X^T X (X^T X)^{-1}}_{\mathbf{I}} = \underline{\underline{\sigma^2 (X^T X)^{-1}}} \\ &= \sigma^2 C^{-1} \quad (\text{book notation}) \end{aligned}$$

Properties of LS-estimates

$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ has

$$E(\hat{\beta}) = \beta \text{ and } \text{Cov}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

The matrix $C = X^T X$ is called the information matrix, and is SPD if X has full rank.
(symmetric, positive definite)

Note that in a designed experiment we may ourselves define \mathcal{X} . In the topic of Design of Experiments (DOE) \mathcal{X} is designed to be optimal in some sense.

For example X may be chosen to minimize a function of $(X^T X)^{-1}$, e.g. make the variances or covariances small.

We will look at a so-called 2^k experiment, where X is chosen with orthogonal columns so that $(X^T X)$ is a diagonal matrix, and thus all $\text{Cov}(\hat{\beta}_j, \hat{\beta}_k) = 0$.

Best linear unbiased estimator (BLUE) ↙ general matrix

Among all unbiased linear estimator $\hat{\beta} = BY$

$\hat{\beta} = (X^T X)^{-1} X^T Y$ has the minimum variance in each component

$\Rightarrow \hat{\beta}$ is BLUE.

Proof: page 73 of BF.

Gauss' LS theorem

- $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2\mathbf{I}$.
- And \mathbf{X} has full rank $p \leq n$.
- $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$.

Then, for any vector \mathbf{c} , the estimator

$$\mathbf{c}^T\hat{\boldsymbol{\beta}}$$

of $\mathbf{c}^T\boldsymbol{\beta}$ has the *smallest possible variance* among all linear estimators that are *unbiased* for $\mathbf{c}^T\boldsymbol{\beta}$.

Maximum likelihood estimator of σ

Assume Y_i independent $N(X_i\beta, \sigma^2)$
row vector of X matrix

Log-likelihood ↙ lecture 12

$$l(\beta, \sigma^2) = \ln L(\beta, \sigma)$$

$$= \frac{n}{2} \ln(2\pi) - n \log \sigma - \frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta)$$

$$\frac{\partial l}{\partial \sigma} = 0 - \frac{n}{\sigma} + \frac{1}{\sigma^3} (y - X\beta)^T (y - X\beta) = 0$$

At the likelihood maximum $\beta = \hat{\beta}$, so we insert $\hat{\beta}$ and solve for σ^2 to get $\hat{\sigma}^2$:

$$\hat{\sigma}^2 = \frac{1}{n} \underbrace{(Y - X\hat{\beta})^T (Y - X\hat{\beta})}$$

SSE \leftarrow a RV and a statistic

New notation

error

$$Y = X\beta + \underbrace{\varepsilon}_{\text{error}}$$

$$\hat{Y}_{n \times 1} = X\hat{\beta} \quad \text{is called fitted values}$$

$$e_{n \times 1} = Y - \hat{Y} = Y - X\hat{\beta} \quad \leftarrow \text{called } \underline{\text{residuals}}$$

$$\begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix} \quad SSE = (Y - X\hat{\beta})^T (Y - X\hat{\beta}) = e^T e$$

↑
sums of squared residuals

ML estimation of σ

The MLE for σ^2 is

$$\hat{\sigma}^2 = \frac{1}{n}(\mathbf{Y} - \mathbf{X}\hat{\beta})^T(\mathbf{Y} - \mathbf{X}\hat{\beta}) = \frac{1}{n} \sum_{i=1}^n (y_i - X_i\hat{\beta})^2$$

This may also be written in a slightly different way:

— First, fitted values of \mathbf{Y} :

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$

— The residuals:

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{Y} - \mathbf{X}\hat{\beta})$$

— and then the Sums-of-squares-of-error:

$$\text{SSE} = \mathbf{e}^T\mathbf{e} = (\mathbf{Y} - \mathbf{X}\hat{\beta})^T(\mathbf{Y} - \mathbf{X}\hat{\beta})$$

— with leads to

$$\hat{\sigma}^2 = \frac{1}{n}\text{SSE}$$

Projection matrices

$$\hat{Y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_H Y = HY$$

^H
putting a hat on y
(in book use \hat{y} - overboxed)

$$e = Y - \hat{Y} = Y - HY = (I - H)Y$$

Observe:

$$1) \underbrace{H}_{n \times n} H = X \underbrace{(X^T X)^{-1} X^T X}_{I} (X^T X)^{-1} X^T = X (X^T X)^{-1} X^T = H$$

$$2) \text{Homework: } (I - H)(I - H) = (I - H)$$

$$3) H \cdot (I - H) = H - HH = H - H = 0$$

\uparrow
 $A^2 = A$
idempotent

In general: A linear transformation A is a projection (into a vector space V) if $A^2 = AA = A$ (idempotent).

$$\text{Then } V = \underbrace{\text{Im}(A)} \oplus \underbrace{\text{Ker}(A)}$$

Image of A
= subspace spanned
by the columns of A

direct sum

kernel of A
= subspace spanned by
the column vectors x such that
 $Ax = 0$. Also called null space

any vector in V can be uniquely
written as a sum of a vector in $\text{Im}(A)$
and a vector in $\text{Ker}(A)$

Also: $\text{Im}(A) = \text{Ker}(I - A)$ and $\text{Ker}(A) = \text{Im}(I - A)$

In our case:

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

$$H^T = H$$

$$\text{and } (I-H)^T = (I-H)$$

If a projection matrix A is symmetric it is also orthogonal. That is, A is an orthogonal projection.

This means that $\text{Im}(A)$ and $\text{ker}(A)$ are orthogonal vector spaces. And, thus are $\text{Im}(A)$ and $\text{Im}(I-A)$ orthogonal vector spaces.

Back to $H, (I-H)$, MLR:

$\hat{Y} = HY$ is the projection of Y onto the space spanned by the columns of H . And, $HX = X(X^T X)^{-1} X^T X = X$ so the column space of H and X are the same.

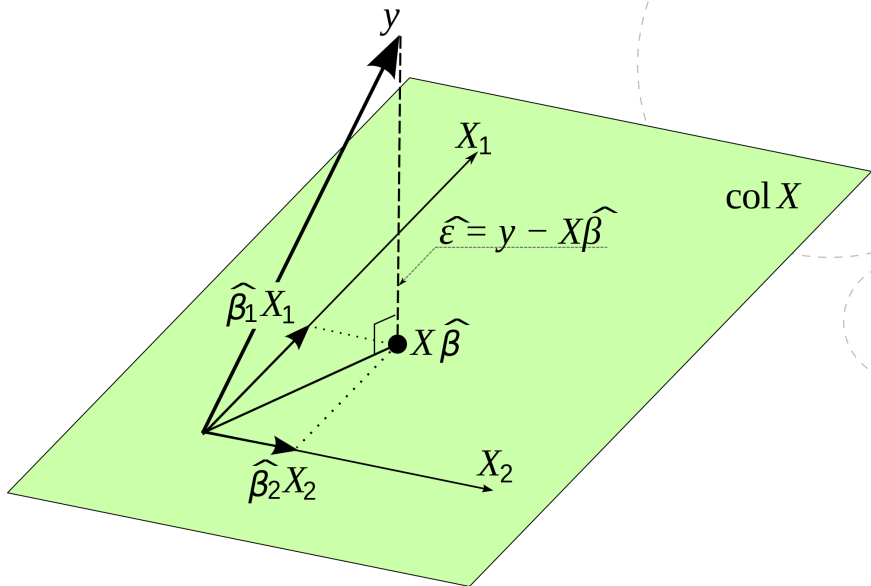
Further: $e = (I-H)Y$ is the projection of Y onto the space orthogonal to the space spanned by the columns of H .

Geometry of Least Squares

- Mean response vector: $E(\mathbf{Y}) = \mathbf{X}\beta$
- As β varies, $\mathbf{X}\beta$ spans the model plane of all linear combinations. I.e. the space spanned by the columns of \mathbf{X} : the column-space of \mathbf{X} .
- Due to random error (and unobserved covariates), \mathbf{y} is not exactly a linear combination of the columns of \mathbf{X} .
- LS-estimation chooses $\hat{\beta}$ such that $\mathbf{X}\hat{\beta}$ is the point in the column-space of \mathbf{X} that is closest to \mathbf{y} .

Geometry of Least Squares (cont.)

- The residual vector $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$ is perpendicular to the column-space of \mathbf{X} .
- Multiplication by $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ projects a vector onto the column-space of \mathbf{X} .
- Multiplication by $\mathbf{I} - \mathbf{H} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ projects a vector onto the space perpendicular to the column-space of \mathbf{X} .



http://en.wikipedia.org/wiki/File:OLS_geometric_interpretation.svg

Properties of symmetric projection matrices

A projection \mathbf{A} matrix is idempotent, $\mathbf{A}^2 = \mathbf{A}$. A symmetric projection matrix is orthogonal.

1. The eigenvalues of a projection matrix are 0 and 1.
2. The rank of a symmetric matrix (actually: a diagonalizable quadratic matrix) equals the number of nonzero eigenvalues of the matrix. Should be known from previous courses.
3. (Combining 1+2). If a $(n \times n)$ symmetric projection matrix \mathbf{A} has rank r then r eigenvalues are 1 and $n - r$ are 0.
4. The trace and rank of a symmetric projection matrix are equal: $tr(\mathbf{A}) = \text{rank}(\mathbf{A})$.