

2.3 Gauss-Markov theorem

REMIND: Data: $(Y_i, X_i^T)^T, i=1, \dots, n$ not necessarily iid

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix} = (X^0, \dots, X^{k-1})$$

It is assumed $Y|X \sim (X\beta, \sigma^2 I_n)$
 $\text{rank}(X_{n \times k}) = r \leq k < n$

$$E(Y|X) = X\beta$$

$$\text{var}(Y|X) = \sigma^2 I_n$$

This is assumed

\equiv MODEL (might be wrong)

• if $r=k$, $\hat{\beta} = (X^T X)^{-1} X^T Y$ is unbiased estimator of β

• $\hat{Y} = X\hat{\beta} = HY$ is projection of Y into $\mathcal{M}(X)$

$$H = X(X^T X)^{-1} X^T \quad (\text{if } r=k)$$

$$= Q Q^T \quad (r \leq k),$$

$$= X(X^T X)^{-1} X^T, \quad Q = (q_1, \dots, q_r) \equiv \text{orthonormal basis of } \mathcal{M}(X)$$

• Statistical properties of \hat{Y} ?

\rightarrow Gauss-Markov theorem

C.F. Gauss: 1821-23

A.A. Markov: 1912

- more "elegant"

version of the proof

Theorem 2.4 Gauss-Markov

Assume a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$,
rank $(X_{n \times k}) = r \leq k < n$. The vector of fitted
values \hat{Y} is, conditionally given X ,
the best linear unbiased estimator (BLUE)
of a vector parameter $\mu := E(Y|X)$ ($= X\beta$
for some $\beta \in \mathbb{R}^k$). Further,

$$\text{var}(\hat{Y}|X) = \sigma^2 H.$$

Proof:

Linearity: $\hat{Y} = HY \equiv$ linear function of Y

Unbiasedness: $E(\hat{Y}|X) = E(HY|X) =$

$$= HE(Y|X) \stackrel{\text{MODEL}}{=} \underbrace{HX}_{\substack{\text{projection of cols. of } X \\ \text{into } \mathcal{N}(X)}} \beta = X\beta$$

Optimality: Let $\tilde{Y} = a + BY$ be some other
LUE of $\mu = X\beta$, $B = \text{fun. of } X$

That is $\forall \beta \in \mathbb{R}^k \quad E(\tilde{Y}|X) = X\beta$

$$\forall \beta \in \mathbb{R}^k \quad E(a + BY|X) = X\beta$$

$$\forall \beta \in \mathbb{R}^k \quad a + BE(Y|X) = X\beta$$

$$\forall \beta \in \mathbb{R}^k \quad a + BX\beta = X\beta$$

• Take $\beta = 0 \Rightarrow \alpha = 0$ (to maintain LUE property)

• That is, $\forall \beta \in \mathbb{R}^k \quad BX\beta = X\beta$

• Take $\beta = (0, \dots, 1, \dots, 0)^T$: $(BX)^j = X^j$
^j-th place

$\Rightarrow BX = X$ (to maintain LUE property)

That is, if $\tilde{Y} = \alpha + BX$ is LUE of $X\beta$

\Leftrightarrow *trivial* $\alpha = 0, BX = X$

$BX = X$ implies: $BX = X \quad | \cdot (X^T X)^{-1} X^T$

$$\underbrace{BX(X^T X)^{-1} X^T}_H = \underbrace{X(X^T X)^{-1} X^T}_H$$

$$BH = H \quad (= H^T)$$

$$H^T B^T = H^T \quad (= H)$$

$$H B^T = H \quad (= H^T)$$

now: $\text{var}(\tilde{Y}|X) = \text{var}(HY|X) = H \underbrace{\text{var}(Y|X)}_{\sigma^2 I_n} H^T =$
 $= \sigma^2 H H^T = \sigma^2 H$

• $\text{var}(\tilde{Y}|X) = \text{var}(BY|X) = B \underbrace{\text{var}(Y|X)}_{\sigma^2 I_n} B^T =$
 $= \sigma^2 B B^T = \sigma^2 (H + B - H)(H + B - H)^T =$
 $= \sigma^2 \left(\underbrace{H H^T}_H + \underbrace{H(B-H)^T}_0 + \underbrace{(B-H)H^T}_0 + (B-H)(B-H)^T \right)$
 $= \sigma^2 H + \sigma^2 (B-H)(B-H)^T$

from previous page:

$$\text{var}(\hat{Y}/X) = \sigma^2 H$$

$$\text{var}(\tilde{Y}/X) = \sigma^2 H + \sigma^2 (B-H)(B-H)^T$$

$$\Rightarrow \text{var}(\tilde{Y}/X) - \text{var}(\hat{Y}/X) = \sigma^2 (B-H)(B-H)^T \geq 0$$

positive semidefinite matrix □

Gauss-Markov in short:

$$\hat{Y}/X \sim (X\beta, \sigma^2 H)$$

+ optimality

REMARK: Gauss-Markov theorem is the key tool to show ^{easily} that any linear combination of LSE $\hat{\beta}$ (if $n=k$) is BLUE of respective linear combination of β .

WILL NOW BE ASSUMED :

$$\text{rank}(X_{n \times k}) = k$$

$$l_j := (l_{j1}, \dots, l_{jk})^T \in \mathbb{R}^k, \quad l_j \neq 0$$

$$L := \begin{pmatrix} l_1^T \\ \vdots \\ l_m^T \end{pmatrix}, \quad l_j \in \mathbb{R}^k, \quad l_j \neq 0, \quad j=1, \dots, m$$

$m \leq k$, rows in L linearly independent

• Interest in estimation of

(i) ~~θ~~ $\theta := l^T \beta$

(ii) $\Theta := L\beta = \begin{pmatrix} l_1^T \beta \\ \vdots \\ l_m^T \beta \end{pmatrix}$

Theorem 2.5 Gauss-Markov for linear combinations

Assume a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$,
 $\text{rank}(X_{n \times k}) = k < n$. Then

(i) $\hat{\theta} := l^T \hat{\beta}$ is BLUE of $\theta = l^T \beta$ with

$$\text{var}(\hat{\theta}|X) = \sigma^2 l^T (X^T X)^{-1} l > 0$$

(ii) $\hat{\Theta} := L \hat{\beta}$ is BLUE of ~~$\hat{\Theta}$~~ $\Theta = L\beta$ with

$\text{var}(\hat{\Theta}|X) = \sigma^2 L (X^T X)^{-1} L$ which
is a positive definite matrix.

Proof: **KEY STEP** (the rest is more or less trivial):

$\mathcal{N}(X^T) =$ linear span of ROWS of X ,
each row is $\in \mathbb{R}^k$, number of
rows is $n > k$, $\text{rank}(X_{n \times k}) = k$
 $= \mathbb{R}^k$

(i) That is, $\forall l \in \mathbb{R}^k$, $l \in \mathcal{N}(X^T)$.

Hence $\forall l \in \mathbb{R}^k \exists a \in \mathbb{R}^n$ $l = X^T a$

(linear comb. of
 n rows from X)

$$\begin{aligned} \text{Hence } l^T \hat{\beta} &= a^T X \hat{\beta} \quad \text{for some } a \in \mathbb{R}^n \\ &= a^T \hat{Y} \quad \text{--- " ---} \\ &= a^T H Y \end{aligned}$$

That is, $l^T \hat{\beta}$ is a linear function of \hat{Y}
this is BLUE of $X\beta$.

$$\Rightarrow l^T \hat{\beta} \text{ is BLUE of } a^T X \beta = l^T \beta$$

REMARK: Unbiasedness now holds also
unconditionally:

$$\begin{aligned} E(l^T \hat{\beta}) &= E(E(l^T \hat{\beta} | X)) = E(l^T \beta) = \\ &= l^T \beta \quad \text{for any } \beta \in \mathbb{R}^k \end{aligned}$$

$$\begin{aligned}
 \bullet \text{var}(l^T \hat{\beta} | X) &= \text{var}(a^T \hat{Y} | X) = a^T \text{var}(\hat{Y} | X) a = \\
 &= a^T (\sigma^2 H) a = \sigma^2 \underbrace{a^T X^T}_{l^T} (X^T X)^{-1} \underbrace{X a}_{l} = \\
 &= \sigma^2 l^T (X^T X)^{-1} l
 \end{aligned}$$

OR CALCULATED DIRECTLY:

$$\begin{aligned}
 \text{var}(l^T \hat{\beta} | X) &= l^T \text{var}(\hat{\beta} | X) l = \\
 &\stackrel{\text{Lemma 2.2}}{=} l^T (\sigma^2 (X^T X)^{-1}) l = \sigma^2 l^T (X^T X)^{-1} l
 \end{aligned}$$

$\text{var}(l^T \hat{\beta} | X) > 0$ since $l \neq 0$ and $(X^T X)^{-1}$ is positive definite matrix

(ii) The BLUE property of $\hat{\theta} = L \hat{\beta}$ as an estimator of $\theta = L \beta$ follows from (i).

Finally:

$$\begin{aligned}
 \text{var}(L \hat{\beta} | X) &= L \text{var}(\hat{\beta} | X) L^T = \\
 &= \sigma^2 L (X^T X)^{-1} L^T
 \end{aligned}$$

IS THIS MATRIX POSITIVE DEFINITE?

one more observation:

$$\mathcal{N}(\underbrace{X^T X}_{k \times k \text{ regular matrix}}) = \mathbb{R}^k = \mathcal{N}(X^T)$$

Hence, for each row from L : $l_j \in \mathcal{N}(X^T X)$,

$$\text{i.e. } \mathcal{N}(L^T) \subset \mathcal{N}(X^T X)$$

$k \times m$

\Rightarrow There exist $B_{k \times m}$ with linearly indep. columns (since L has lin. indep. rows)

such that

$$L^T = X^T X B$$

$\hat{=} A_{n \times m}$ (also lin. indep. columns)

$$\begin{aligned} \Rightarrow L(X^T X)^{-1} L^T &= B^T X^T X (X^T X)^{-1} X^T X B = \\ &= B^T X^T X B = \underbrace{A^T A}_{m \times m \text{ matrix}} \end{aligned}$$

A is has linear indep. columns,

hence $A^T A$ is regular

and positive definite matrix. \square

In particular (take $L = I_k$):

$$\hat{\beta} = (X^T X)^{-1} X^T Y \text{ is BLUE of } \beta.$$

2.4 Residuals, properties

REMNDR: Data $(y_i, X_i^T)^T, i=1, \dots, n$ not neces, iid

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix} = (X^0, \dots, X^{k-1})$$

It is assumed

$$Y|X \sim (X\beta, \sigma^2 I_n)$$

$$\text{rank}(X_{n \times k}) = r \leq k < n$$

MODEL $\left\{ \begin{array}{l} E(Y|X) = X\beta \\ \text{var}(Y|X) = \sigma^2 I_n \end{array} \right. \longrightarrow$ its BLUE: $\hat{Y} = H \cdot Y \in \mathcal{M}(X)$

$$= Q Q^T Y$$

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

how to estimate this?

Residuals $U = Y - \hat{Y} = M Y \in \mathcal{M}(X)^\perp$

$$= N N^T Y$$

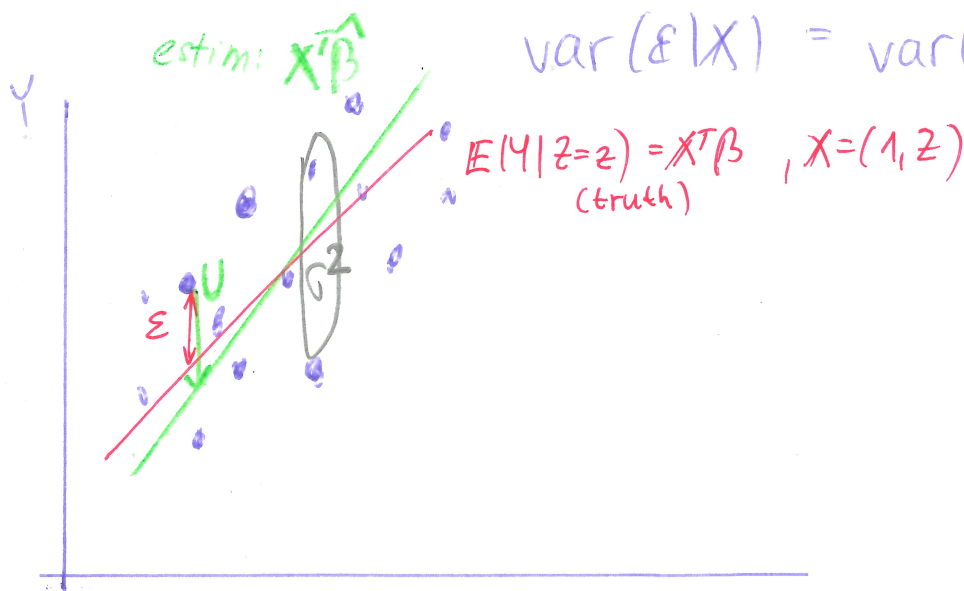
$$= (I_n - H) Y$$

We also have:

Error terms $\epsilon = Y - X\beta = Y - E(Y|X)$

already know: $E(\epsilon|X) = E(\epsilon) = 0$

$$\text{var}(\epsilon|X) = \text{var}(\epsilon) = \sigma^2 I_n$$



Def 2.5 Residual sum of squares

Consider a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$,
 $\text{rank}(X_{n \times k}) = r \leq k < n$. The quantity

$$SSE = \|U\|^2 \quad (= \sum_{i=1}^n U_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \|Y - \hat{Y}\|^2)$$

will be called the residual sum of squares.

Lemma 2.6 Alternative expressions of residuals and residual sum of squares

Let $Y|X \sim (X\beta, \sigma^2 I_n)$, $\text{rank}(X_{n \times k}) = r \leq k < n$.
The following then holds:

(i) $U = M\varepsilon$, where $\varepsilon = Y - X\beta$.

(ii) $SSE = Y^T M Y = \varepsilon^T M \varepsilon$.

Proof: (i) $U = M Y = M (X\beta + \varepsilon) =$
 $= \underbrace{M X}_{0_{n \times k}} \beta + M \varepsilon = M \varepsilon$
(projections of cols. of X into $N(X)^\perp$)

(ii) $SSE = \|U\|^2 = U^T U = (M Y)^T M Y =$
 $= Y^T \underbrace{M^T M}_M Y = Y^T M Y$

(ii) $= \varepsilon^T M^T M \varepsilon = \varepsilon^T M \varepsilon$

□

Lemma 2.7 Moments of residuals and resid. sum of sq.

Let $Y|X \sim (X\beta, \sigma^2 I_n)$, $\text{rank}(X_{n \times k}) = k \leq k < n$.

Then (i) $E(U|X) = 0_n$, $\text{var}(U|X) = \sigma^2 M$.

(ii) $E(\text{SSE}|X) = E(\text{SSE}) = (n-k)\sigma^2$.

Proof: (i) $E(U|X) = E(MY|X) =$

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

$$= M E(Y|X) = M X \beta = 0$$

$0_{n \times k}$ (projections of cols. of X into $\mathcal{N}(X)^\perp$)

$$\begin{aligned} \text{var}(U|X) &= \text{var}(MY|X) = M \underbrace{\text{var}(Y|X)}_{\sigma^2 I_n} M^T = \\ &= \sigma^2 M M^T = \sigma^2 M. \end{aligned}$$

$$(ii) E(\text{SSE}|X) \stackrel{\text{Lemma 2.6}}{=} E(\varepsilon^T M \varepsilon | X) =$$

$$= E(\text{tr}(\varepsilon^T M \varepsilon) | X) = E(\text{tr}(M \varepsilon \varepsilon^T) | X) =$$

$$= \text{tr} \{ E(M \varepsilon \varepsilon^T | X) \} = \text{tr} \{ M \underbrace{E(\varepsilon \varepsilon^T | X)}_{\text{var}(\varepsilon|X) = \sigma^2 I_n} \} =$$

$$= \text{tr}(M \cdot \sigma^2 I_n) = \sigma^2 \text{tr}(M) = \sigma^2 \text{tr}(N N^T) =$$

$$= \sigma^2 \text{tr}(N^T N) = \sigma^2 \text{tr}(I_{n-k}) =$$

$$= \sigma^2 (n-k)$$

orth. basis of $\mathcal{N}(X)^\perp$
in cols., vec-dim = $n-k$

$$E(\text{SSE}) = E(E(\text{SSE}|X)) = E(\sigma^2 (n-k)) = \sigma^2 (n-k)$$

Lemma 2.7 (i) in short: $U|X \sim (0, \sigma^2 M)$

Consequence of Lemma 2.7 (ii):

Quantity $\frac{SSE}{n-r}$ is unbiased (both conditionally given X and unconditionally) estimator of σ^2 .

Def. 2.6 Residual mean square
and residual degrees of freedom

Consider a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$,
 $\text{rank}(X_{n \times k}) = r \leq k < n$.

(i) The residual mean square of the model is the quantity $\frac{SSE}{n-r}$ and will be denoted as MSE , i.e. $MSE = \frac{SSE}{n-r}$.

(ii) The residual degrees of freedom of the model is the vector dimension of the residual space $\mathcal{N}(X)^\perp$ and will be denoted as ν_e , i.e. $\nu_e = n-r$.