

Lecture Notes for Theory of Linear Models

(chapter 7)

- Estimation in Linear Models
- BLUE
- Distribution of SSR, SSE, $\hat{\beta}$
- R^2 and R_a^2
- Overfitting and Underfitting

Longhai Li
Department of Mathematics and Statistics
University of Saskatchewan

Assumptions in Linear Models

Suppose that on a random sample of n units (patients, animals, trees, etc.) we observe a response variable Y and explanatory variables X_1, \dots, X_k .

Our data are then $(y_i, x_{i1}, \dots, x_{ik})$, $i = 1, \dots, n$, or, in vector/matrix form $\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_k$ where $\mathbf{y} = (y_1, \dots, y_n)$ and $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})^T$; or \mathbf{y}, \mathbf{X} where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$.

Either by design or by conditioning on their observed values, $\mathbf{x}_1, \dots, \mathbf{x}_k$ are regarded as vectors of known constants.

The linear model in its classical form makes the following assumptions:

- A1. (additive error) $\mathbf{y} = \boldsymbol{\mu} + \mathbf{e}$ where $\mathbf{e} = (e_1, \dots, e_n)^T$ is an unobserved random vector with $E(\mathbf{e}) = \mathbf{0}$. This implies that $\boldsymbol{\mu} = E(\mathbf{y})$ is the unknown mean of \mathbf{y} .
- A2. (linearity) $\boldsymbol{\mu} = \beta_1 \mathbf{x}_1 + \dots + \beta_k \mathbf{x}_k = \mathbf{X}\boldsymbol{\beta}$ where β_1, \dots, β_k are unknown parameters. This assumption says that $E(\mathbf{y}) = \boldsymbol{\mu} \in \mathcal{L}(\mathbf{x}_1, \dots, \mathbf{x}_k) = C(\mathbf{X})$ lies in the column space of \mathbf{X} ; i.e., it is a linear combination of explanatory vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$ with coefficients the unknown parameters in $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)^T$.
 - Linear in β_1, \dots, β_k not in the x 's.
- A3. (independence) e_1, \dots, e_n are independent random variables (and therefore so are y_1, \dots, y_n).
- A4. (homoscedasticity) e_1, \dots, e_n all have the same variance σ^2 ; that is, $\text{var}(e_1) = \dots = \text{var}(e_n) = \sigma^2$ which implies $\text{var}(y_1) = \dots = \text{var}(y_n) = \sigma^2$.
- A5. (normality) $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}, \quad i = 1, \dots, n,$$

or

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix}}_{=\mathbf{X}} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_{=\boldsymbol{\beta}} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

or

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad \mathbb{E}(\mathbf{e}) = \mathbf{0}, \quad \text{Var}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$$

where $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$.

or

$$y = \beta_0 \hat{j}_n + \beta_1 x_1 + \dots + \beta_k x_k + e$$

Taken together, all five assumptions can be stated more succinctly as $\mathbf{y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$.

or

$$\mathbf{y} \sim N_n(\boldsymbol{\mu}_y, \sigma^2 \mathbf{I}_n)$$

$$\boldsymbol{\mu}_y = \mathbf{X}\boldsymbol{\beta} \in c(\mathbf{X})$$

$$\mathbf{y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$$

A Note:

In addition, this effect depends upon what other explanatory variables are present in the model. For example, β_0 and β_1 in the model

$$\mathbf{y} = \beta_0 \mathbf{j}_n + \beta_1 \mathbf{x}_1 + \beta_2 \mathbf{x}_2 + \mathbf{e}$$

will typically be different than β_0^* and β_1^* in the model

$$\mathbf{y} = \beta_0^* \mathbf{j}_n + \beta_1^* \mathbf{x}_1 + \mathbf{e}^*.$$

We will first consider the
case that $\text{rank}(X) = k+1$

Example 7.3.1a. We use the data in Table 7.1 to illustrate computation of $\hat{\boldsymbol{\beta}}$ using (7.6).

$$\mathbf{y} = \begin{pmatrix} 2 \\ 3 \\ 2 \\ 7 \\ 6 \\ 8 \\ 10 \\ 7 \\ 8 \\ 12 \\ 11 \\ 14 \end{pmatrix}, \quad \mathbf{X} = \begin{matrix} \beta_0 & \beta_1 & \beta_2 \\ \begin{pmatrix} 1 & 0 & 2 \\ 1 & 2 & 6 \\ 1 & 2 & 7 \\ 1 & 2 & 5 \\ 1 & 4 & 9 \\ 1 & 4 & 8 \\ 1 & 4 & 7 \\ 1 & 6 & 10 \\ 1 & 6 & 11 \\ 1 & 6 & 9 \\ 1 & 8 & 15 \\ 1 & 8 & 13 \end{pmatrix} \end{matrix}, \quad \mathbf{X}'\mathbf{X} = \begin{pmatrix} 12 & 52 & 102 \\ 52 & 395 & 536 \\ 102 & 536 & 1004 \end{pmatrix},$$

$$\mathbf{X}'\mathbf{y} = \begin{pmatrix} 90 \\ 482 \\ 872 \end{pmatrix}, \quad (\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} .97476 & .24290 & -.22871 \\ .24290 & .16207 & -.11120 \\ -.22871 & -.11120 & .08360 \end{pmatrix},$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \begin{pmatrix} 5.3754 \\ 3.0118 \\ -1.2855 \end{pmatrix}.$$

Theorem 7.3b. If $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$, then $\hat{\boldsymbol{\beta}}$ is an unbiased estimator for $\boldsymbol{\beta}$.

PROOF

$$\begin{aligned}
 E(\hat{\boldsymbol{\beta}}) &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{y}) \quad \text{[by (3.38)]} \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \\
 &= \boldsymbol{\beta}.
 \end{aligned}$$

(7.13)

□

$$E(A\mathbf{y}) = A \cdot E(\mathbf{y})$$

Theorem 7.3c. If $\text{cov}(\mathbf{y}) = \sigma^2\mathbf{I}$, the covariance matrix for $\hat{\boldsymbol{\beta}}$ is given by $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.

PROOF

$$\begin{aligned}
 \text{cov}(\hat{\boldsymbol{\beta}}) &= \text{cov}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{cov}(\mathbf{y})[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' \quad \text{[by (3.44)]} \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\sigma^2\mathbf{I})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.
 \end{aligned}$$

(7.14)

$$\begin{aligned}
 \text{cov} &= \text{Var} \\
 \text{Var}(A\mathbf{y}) &= A \cdot \text{Var}(\mathbf{y}) \cdot A'
 \end{aligned}$$

$$V(\hat{\boldsymbol{\beta}}) = \sigma^2 \cdot (\mathbf{X}'\mathbf{X})^{-1}$$

Note: no assumption of normality.

$$\hat{A} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$$

Best Linear Unbiased Estimator BLUE

Theorem 7.3d (Gauss–Markov Theorem). If $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ and $\text{cov}(\mathbf{y}) = \sigma^2\mathbf{I}$, the least-squares estimators $\hat{\beta}_j, j = 0, 1, \dots, k$, have minimum variance among all linear unbiased estimators

$$V(\hat{\beta}_j) \geq V(\hat{\beta}_j)$$

PROOF. We consider a linear estimator $\mathbf{A}\mathbf{y}$ of $\boldsymbol{\beta}$ and seek the matrix \mathbf{A} for which $\mathbf{A}\mathbf{y}$ is a minimum variance unbiased estimator of $\boldsymbol{\beta}$. In order for $\mathbf{A}\mathbf{y}$ to be an unbiased estimator of $\boldsymbol{\beta}$, we must have $E(\mathbf{A}\mathbf{y}) = \boldsymbol{\beta}$. Using the assumption $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$, this can be expressed as

$$E(\mathbf{A}\mathbf{y}) = \mathbf{A}E(\mathbf{y}) = \mathbf{A}\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta},$$

$$\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{y}$$

which gives the unbiasedness condition

$$\mathbf{A}\mathbf{X} = \mathbf{I}_{k+1}$$

$$\text{cov}(\hat{\boldsymbol{\beta}}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

since the relationship $\mathbf{A}\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}$ must hold for any possible value of $\boldsymbol{\beta}$ [see (2.44)].

The covariance matrix for the estimator $\mathbf{A}\mathbf{y}$ is given by

$$\text{cov}(\tilde{\boldsymbol{\beta}}) = \text{cov}(\mathbf{A}\mathbf{y}) = \mathbf{A}(\sigma^2\mathbf{I})\mathbf{A}' = \sigma^2\mathbf{A}\mathbf{A}'.$$

$$\hat{\boldsymbol{\beta}} = [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}$$

$$= \hat{\mathbf{A}}\mathbf{y}$$

The variances of the $\hat{\beta}_j$'s are on the diagonal of $\sigma^2\mathbf{A}\mathbf{A}'$, and we therefore need to choose \mathbf{A} (subject to $\mathbf{A}\mathbf{X} = \mathbf{I}$) so that the diagonal elements of $\mathbf{A}\mathbf{A}'$ are minimized.

To relate $\mathbf{A}\mathbf{y}$ to $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, we add and subtract $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ to obtain

$$\mathbf{A}\mathbf{A}' = (\mathbf{A} - \hat{\mathbf{A}} + \hat{\mathbf{A}})(\mathbf{A} - \hat{\mathbf{A}} + \hat{\mathbf{A}})'$$

$$\mathbf{A}\mathbf{A}' = [\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'][\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']'$$

Expanding this in terms of $\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ and $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, we obtain four terms, two of which vanish because of the restriction $\mathbf{A}\mathbf{X} = \mathbf{I}$. The result is

$$\mathbf{A}\mathbf{A}' = [\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'][\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' + (\mathbf{X}'\mathbf{X})^{-1}. \quad (7.17)$$

The matrix $[\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'][\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']'$ on the right side of (7.17) is positive semidefinite (see Theorem 2.6d), and, by Theorem 2.6a (ii), the diagonal elements are greater than or equal to zero. These diagonal elements can be made equal to zero by choosing $\mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. (This value of \mathbf{A} also satisfies the unbiasedness condition $\mathbf{A}\mathbf{X} = \mathbf{I}$.) The resulting minimum variance estimator of $\boldsymbol{\beta}$ is

$$\mathbf{A}\mathbf{y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y},$$

which is equal to the least-squares estimator $\hat{\boldsymbol{\beta}}$. □

proof of $(A - \hat{A})\hat{A}' = 0$

$$\hat{A} = (X'X)^{-1}X', \quad \hat{\beta}_{LS} = \hat{A}y$$

$$(A - \hat{A})\hat{A}'$$

$$= A \cdot \hat{A}' - \hat{A} \cdot \hat{A}'$$

$$= A \cdot \hat{A}' - (X'X)^{-1}$$

$$= A \cdot (X \cdot (X'X)^{-1}) - (X'X)^{-1}$$

$$= \overset{\uparrow}{I_{k+1}} (X'X)^{-1} - (X'X)^{-1} = 0$$

$$AX = I$$

$$\text{Let } \hat{A} = (X'X)^{-1}X', \text{ i.e. } \hat{\beta}_{LS} = \hat{A} \cdot y$$

$$AA' = (A - \hat{A}) \cdot (A - \hat{A})' + (X'X)^{-1} + 0$$

$$\forall \alpha \in \mathbb{R}^{k+1}$$

$$\frac{1}{\sigma^2} V(\alpha' \hat{\beta}) = V(\alpha' Ay) / \sigma^2 = \alpha' AA' \alpha$$

$$= \alpha' (A - \hat{A})(A - \hat{A})' \alpha + \alpha' (X'X)^{-1} \alpha$$

$$\geq \alpha' (X'X)^{-1} \alpha = \frac{1}{\sigma^2} \text{Var}(\alpha' \hat{\beta})$$

$$\text{Let } \alpha_j = (0, \dots, 0, 1, 0, 0, \dots, 0)$$

$$V(\hat{\beta}_j) = \alpha_j' AA' \alpha_j \cdot \sigma^2$$

$$\geq \alpha_j' (X'X)^{-1} \alpha_j \cdot \sigma^2$$

$$= V(\hat{\beta}_j)$$

Notes on Gauss-Markov Thm:

1) The remarkable feature of the Gauss-Markov theorem is its distributional generality. The result holds for any distribution of \mathbf{y} ; normality is not required. The only assumptions used in the proof are $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ and $\text{cov}(\mathbf{y}) = \sigma^2\mathbf{I}$. If these assumptions do not hold, $\hat{\boldsymbol{\beta}}$ may be biased or each $\hat{\beta}_j$ may have a larger variance than that of some other estimator.

2) **Corollary 1.** If $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ and $\text{cov}(\mathbf{y}) = \sigma^2\mathbf{I}$, the best linear unbiased estimator of $\mathbf{a}'\boldsymbol{\beta}$ is $\mathbf{a}'\hat{\boldsymbol{\beta}}$, where $\hat{\boldsymbol{\beta}}$ is the least-squares estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.

3) **Theorem 7.3e.** If $\mathbf{x} = (1, x_1, \dots, x_k)'$ and $\mathbf{z} = (1, c_1x_1, \dots, c_kx_k)'$, then $\hat{y} = \hat{\boldsymbol{\beta}}'\mathbf{x} = \hat{\boldsymbol{\beta}}'_z\mathbf{z}$, where $\hat{\boldsymbol{\beta}}_z$ is the least squares estimator from the regression of y on \mathbf{z} .

PROOF. From (2.29), we can rewrite \mathbf{z} as $\mathbf{z} = \mathbf{D}\mathbf{x}$, where $\mathbf{D} = \text{diag}(1, c_1, c_2, \dots, c_k)$. Then, the \mathbf{X} matrix is transformed to $\mathbf{Z} = \mathbf{X}\mathbf{D}$ [see (2.28)]. We substitute $\mathbf{Z} = \mathbf{X}\mathbf{D}$ in the least-squares estimator $\hat{\boldsymbol{\beta}}_z = (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y}$ to obtain

$$\begin{aligned}\hat{\boldsymbol{\beta}}_z &= (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y} = [(\mathbf{X}\mathbf{D})'(\mathbf{X}\mathbf{D})]^{-1}(\mathbf{X}\mathbf{D})'\mathbf{y} \\ &= \mathbf{D}^{-1}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad [\text{by (2.49)}] \\ &= \mathbf{D}^{-1}\hat{\boldsymbol{\beta}},\end{aligned}\tag{7.18}$$

where $\hat{\boldsymbol{\beta}}$ is the usual estimator for y regressed on the x 's. Then

$$\hat{\boldsymbol{\beta}}'_z\mathbf{z} = (\mathbf{D}^{-1}\hat{\boldsymbol{\beta}})'\mathbf{D}\mathbf{x} = \hat{\boldsymbol{\beta}}'\mathbf{x}.$$

□



Estimator of σ^2 $\sigma^2 = V(\epsilon_i)$ $Y \sim N(X\beta, \sigma^2 I_n)$

We estimate σ^2 by a corresponding average from the sample

$$s^2 = \frac{1}{n-k-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (7.22)$$

where n is the sample size and k is the number of x 's. Note that, by the corollary to Theorem 7.3d, $\mathbf{x}'_i \hat{\beta}$ is the BLUE of $\mathbf{x}'_i \beta$.

Using (7.7), we can write (7.22) as

$$s^2 = \frac{1}{n-k-1} (\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})$$

$$= \frac{\mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}}{n-k-1} = \frac{\text{SSE}}{n-k-1},$$

$$\text{SSE} = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2$$

$$= \|\mathbf{y}\|^2 - \|\mathbf{X}\hat{\beta}\|^2$$

where $\text{SSE} = (\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}$. With the denominator $n-k-1$, s^2 is an unbiased estimator of σ^2 , as shown below.

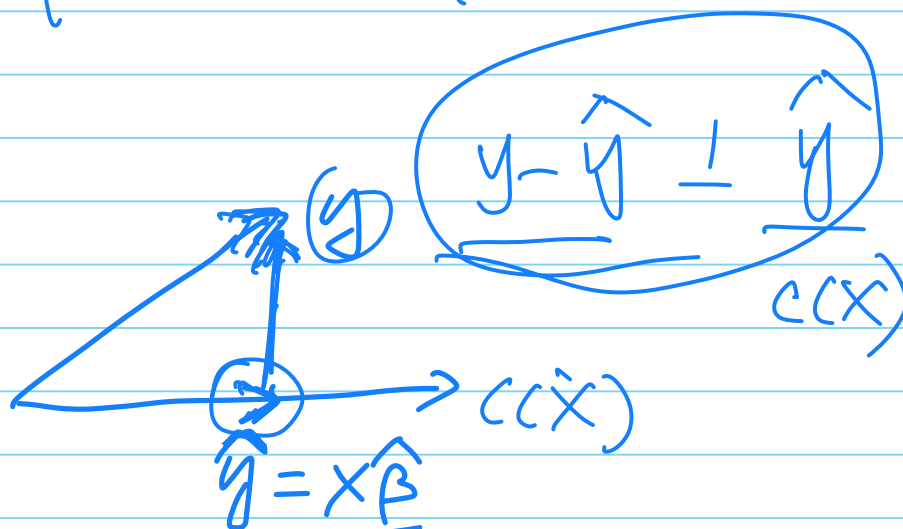
Theorem 7.3f. If s^2 is defined by (7.22), (7.23), or (7.24) and if $E(\mathbf{y}) = \mathbf{X}\beta$ and $\text{cov}(\mathbf{y}) = \sigma^2 \mathbf{I}$, then

$$E(s^2) = \sigma^2. \quad (7.25)$$

SSE, RSS

$$E\left(\frac{\text{SSE}}{n-k-1}\right) = \sigma^2$$

Various expressions of $\|\hat{y}\|^2$



$$\langle y, \hat{y} \rangle = \langle \hat{y}, \hat{y} \rangle$$

$$\langle y - \hat{y}, \hat{y} \rangle = \langle y, \hat{y} \rangle - \langle \hat{y}, \hat{y} \rangle = 0$$

$$\|\hat{y}\|^2 = \langle \hat{y}, \hat{y} \rangle = \langle y, \hat{y} \rangle$$

$$= \hat{\beta}' X' y$$

$$= \hat{\beta}' X' X \hat{\beta}$$

$$= \|X \hat{\beta}\|^2$$

Directly, $\hat{y}' y = \hat{y}' (\hat{y} + y - \hat{y})$

$$= \hat{y}' \hat{y} + 0$$

Recall:

$$E(y'Ay) = \text{tr}(A \cdot \Sigma) + u'Au \quad \text{when}$$

or $E(\|Py\|^2) = \text{rank}(P) \cdot \sigma^2 + \|Pu\|^2$; $\text{var}(y) = \sigma^2 I$

proof of Thm 7.3 f:

$$H = \underline{X(X'X)^{-1}X'} \quad (\text{proj matrix onto } C(X))$$

$$\hat{y} = H \cdot y = X\hat{\beta}$$

$$y - \hat{y} = (I_n - H)y$$

$$SSE = \|y - \hat{y}\|^2$$

$$= y' \cdot \underbrace{(I_n - H)}_{\leftarrow A} y$$

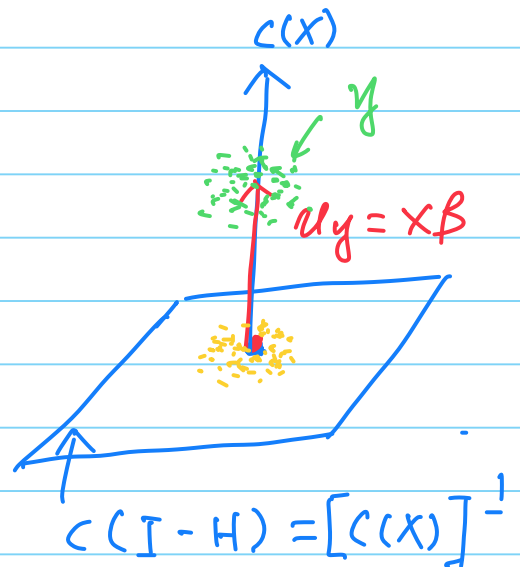
$$= \| \underline{(I_n - H)y} \|^2$$

$$E(SSE)$$

$$= \text{tr} \left(\underline{(I_n - H)} \sigma^2 I \right)$$

$$+ u_y' (I_n - H) u_y \quad \text{0}$$

$$= \underline{(n - k + 1)\sigma^2} + 0$$



$$u_y = X\beta$$

$$(I_n - H)u_y = 0$$

proof of $(I-H)u_y = 0$, where $u_y = X\beta$
 $u_y = X\beta \in (X)$

$$(I-H)u_y = 0$$

$$u_y - H u_y$$
$$= X\beta - X \cdot \cancel{(X'X)^{-1} X'} X\beta$$

$$= X\beta - X\beta = 0$$

$$\text{so } \|(I-H)u_y\|^2 = 0$$

Variance-Cov. matrix of $\hat{\beta}$

Corollary 1. An unbiased estimator of $\text{cov}(\hat{\beta})$ in (7.14) is given by

$$\widehat{\text{cov}}(\hat{\beta}) = s^2(\mathbf{X}'\mathbf{X})^{-1}. \quad (7.27)$$

□

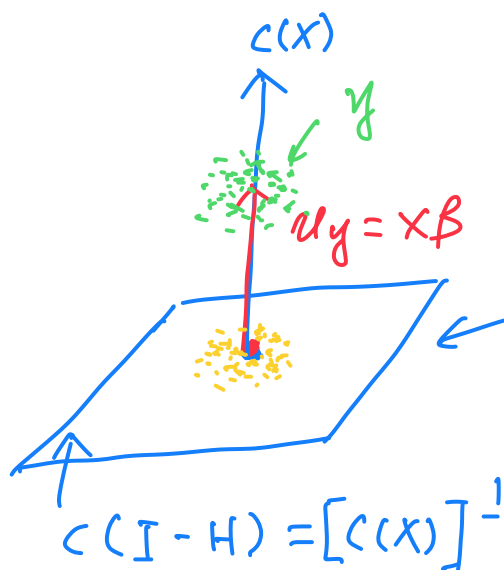
$$\text{cov}(\hat{\beta}) = \sigma^2 \cdot (\mathbf{X}'\mathbf{X})^{-1}$$

Distributions of $\hat{\beta}$ and s^2

$$y \sim N(X\beta, \sigma^2 I)$$

Theorem: Under assumptions A1-A5 of the classical linear model,

- i. $\hat{\beta} \sim N_{k+1}(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$,
- ii. $(n - k - 1)s^2/\sigma^2 \sim \chi^2(n - k - 1)$, and
- iii. $\hat{\beta}$ and s^2 are independent.



$$\begin{aligned} \text{rank}(\mathbf{I} - \mathbf{H}) \\ = \text{tr}(\mathbf{I} - \mathbf{H}) = n - k - 1 \end{aligned}$$

$$u_y = X\beta$$

$$(\mathbf{I}_n - \mathbf{H})u_y = 0$$

proof:
(i)

$$\hat{\beta} = (X'X)^{-1} X'y$$

$$y \sim N(X\beta, \sigma^2 I)$$

$$E(\hat{\beta}) = \beta$$

$$V(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

$$\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1})$$

(ii)

$$SSE = y' (I - H) y$$

$$= y' P_{(X)}^\perp y$$

$$\frac{SSE}{\sigma^2} \sim \chi^2 \left(\underset{n-k-1}{\text{rank}(I-H)}, \frac{1}{2} \cdot \frac{\| (I-H)u \|^2}{\sigma^2} \right)$$

$$(I-H)u = (I-H)X\beta = 0$$

(iii)

$$\underline{H} = P_{C(X)} = X(X'X)^{-1}X'$$

$$\left[(X'X)^{-1}X' \right] X \cdot (X'X)^{-1}X'y$$

$$= (X'X)^{-1}X'y = \hat{\beta}$$

$$\text{That is, } \hat{\beta} = (X'X)^{-1}X' \cdot \hat{y}$$

$\hat{y} = Hy$ and $(I - H)y$ are indep

because $H \cdot (I - H) = 0$

Maximum Likelihood Estimator

Theorem 7.6a. If \mathbf{y} is $N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, where \mathbf{X} is $n \times (k + 1)$ of rank $k + 1 < n$, the maximum likelihood estimators of $\boldsymbol{\beta}$ and σ^2 are

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \quad (7.48)$$

$$\hat{\sigma}^2 = \frac{1}{n}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \frac{\text{SSE}}{n} \quad (7.49)$$

PROOF. We sketch the proof. For the remaining steps, see Problem 7.21. The likelihood function (joint density of y_1, y_2, \dots, y_n) is given by the multivariate normal density (4.9)

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma^2) = f(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) &= \frac{1}{(2\pi)^{n/2} |\sigma^2\mathbf{I}|^{1/2}} e^{-(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\sigma^2\mathbf{I})^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/2} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/2\sigma^2}. \end{aligned} \quad (7.50)$$

[Since the y_i 's are independent, $L(\boldsymbol{\beta}, \sigma^2)$ can also be obtained as $\prod_{i=1}^n f(y_i; \mathbf{x}_i'\boldsymbol{\beta}, \sigma^2)$.] Then $\ln L(\boldsymbol{\beta}, \sigma^2)$ becomes

$$\ln L(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (7.51)$$

Taking the partial derivatives of $\ln L(\boldsymbol{\beta}, \sigma^2)$ with respect to $\boldsymbol{\beta}$ and σ^2 and setting the results equal to zero will produce (7.48) and (7.49). To verify that $\hat{\boldsymbol{\beta}}$ maximizes (7.50) or (7.51), see (7.10). \square

Note: $S^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n - k - 1}$

$E(\hat{\sigma}^2) = \frac{(n - k - 1)\sigma^2}{n} < \sigma^2$

For $k=1$, $\text{MSE}(\hat{\sigma}^2) < \text{MSE}(S^2)$

Linear Models in Centered Form

$$X = (\mathbf{j}_n, x_1, \dots, x_k)$$

The regression model can be written

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + e_i \\ &= \alpha + \beta_1(x_{i1} - \bar{x}_1) + \beta_2(x_{i2} - \bar{x}_2) + \dots + \beta_k(x_{ik} - \bar{x}_k) + e_i, \end{aligned}$$

for $i = 1, \dots, n$, where

$$\alpha = \beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 + \dots + \beta_k \bar{x}_k, \quad (\heartsuit)$$

and where $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$.

$$\beta_1 = (\beta_1, \beta_2, \dots, \beta_k)'$$

In matrix form, the equivalence between the original model and centered model that we've written above becomes

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{e} = (\mathbf{j}_n, \mathbf{X}_c) \begin{pmatrix} \alpha \\ \beta_1 \end{pmatrix} + \mathbf{e} = \alpha \mathbf{j}_n + \mathbf{X}_c \beta_1 + \mathbf{e}$$

where $\beta_1 = (\beta_1, \dots, \beta_k)'$, and

$$\mathbf{X}_c = \underbrace{\left(\mathbf{I} - \frac{1}{n} \mathbf{J}_{n,n} \right)}_{=\mathbf{P}_{\mathcal{L}(\mathbf{j}_n)^\perp}} \mathbf{X}_1 = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1k} - \bar{x}_k \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \dots & x_{2k} - \bar{x}_k \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \dots & x_{nk} - \bar{x}_k \end{pmatrix} = C(X) L(\mathbf{j}_n)'$$

and \mathbf{X}_1 is the matrix consisting of all but the first columns of \mathbf{X} , the original model matrix.

$$(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

$$\begin{aligned} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \end{pmatrix} &= [(\mathbf{j}_n, \mathbf{X}_c)'(\mathbf{j}_n, \mathbf{X}_c)]^{-1} (\mathbf{j}_n, \mathbf{X}_c)' \mathbf{y} = \begin{pmatrix} n & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_c' \mathbf{X}_c \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{j}_n' \\ \mathbf{X}_c' \end{pmatrix} \mathbf{y} \\ &= \begin{pmatrix} n^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}_c' \mathbf{X}_c)^{-1} \end{pmatrix} \begin{pmatrix} n\bar{y} \\ \mathbf{X}_c' \mathbf{y} \end{pmatrix} = \begin{pmatrix} \bar{y} \\ (\mathbf{X}_c' \mathbf{X}_c)^{-1} \mathbf{X}_c' \mathbf{y} \end{pmatrix}, \end{aligned}$$

$\hat{\alpha} = \bar{y}$, and

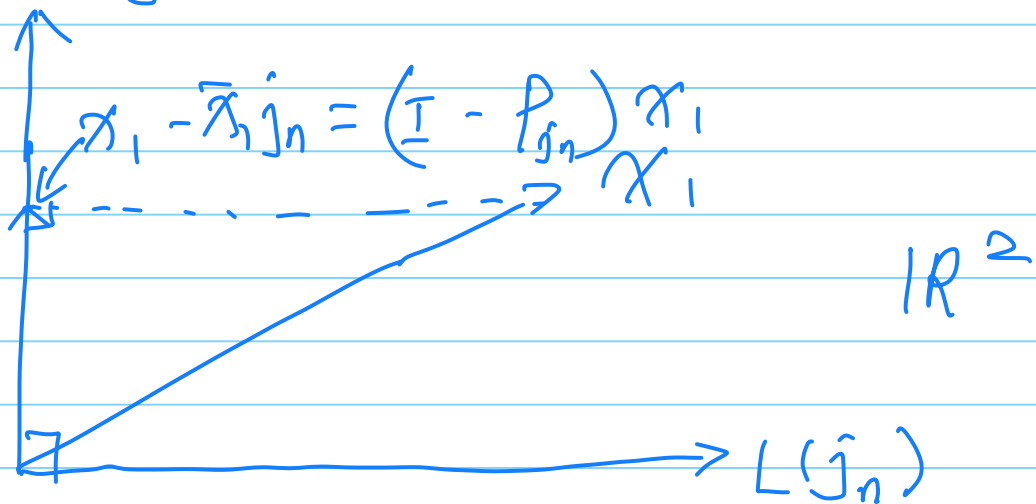
$$\hat{\beta}_1 = (\mathbf{X}_c' \mathbf{X}_c)^{-1} \mathbf{X}_c' \mathbf{y} = S_{xx}^{-1} S_{xy}$$

$$\hat{\beta}_0 = \hat{\alpha} - \hat{\beta}_1 \bar{x}_1 - \hat{\beta}_2 \bar{x}_2 - \dots - \hat{\beta}_k \bar{x}_k = \bar{y} - \hat{\beta}_1' \bar{\mathbf{x}} = \bar{y} - \bar{\mathbf{x}}' \hat{\beta}_1$$

$$\mathbf{P}_{\mathbf{j}_n} \mathbf{y} = \hat{\alpha} \mathbf{j}_n = \bar{y} \mathbf{j}_n$$

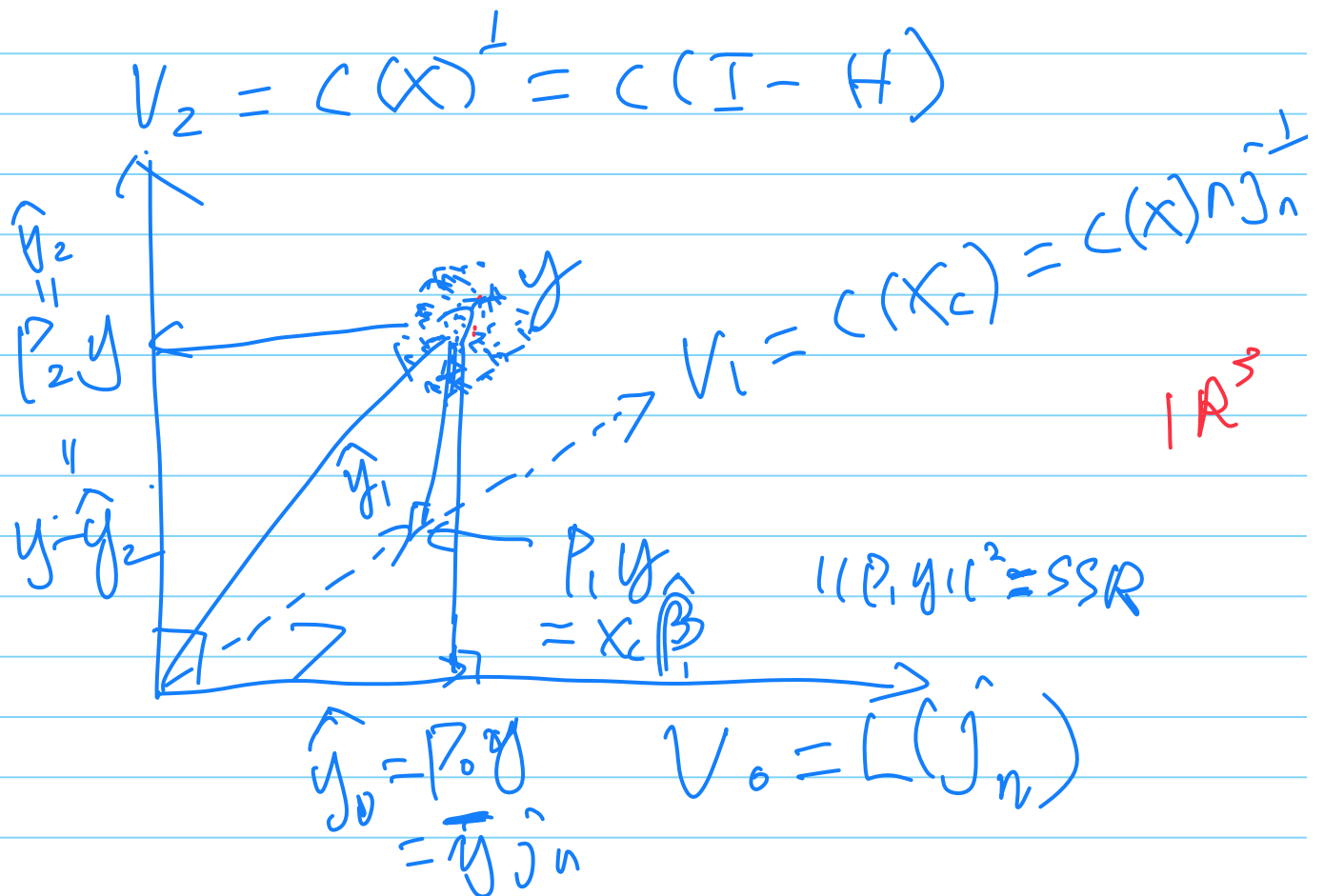
$$\begin{cases} S_{xx} = \frac{\mathbf{X}_c' \mathbf{X}_c}{n-1} \\ S_{xy} = \frac{\mathbf{X}_c' \mathbf{y}}{n-1} \end{cases}$$

$$L(\hat{j}_n)^\perp \cdot L(x_c) = c(x_{1,c})$$



$$\begin{aligned} c(x_c) &= c(x) \cap [L(\hat{j}_n)]^\perp \\ &= c(P_{c(x)} - P_{L(\hat{j}_n)}) \\ &= c(P_{L(\hat{j}_n)^\perp} \cdot x) \\ &= c(x - P_{L(\hat{j}_n)} x) \end{aligned}$$

$$c(x) = c(x_c) \oplus L(\hat{j}_n)$$



$$H = X \cdot (X'X)^{-1} X'$$

$$C(X) = C([L(j_n), X_c])$$

$$P_1 y = X_c \hat{\beta}_1 = X_c (X_c' X_c)^{-1} X_c' y$$

$$SSR = \|P_1 y\|^2 = \hat{\beta}_1' X_c' X_c \hat{\beta}_1$$

Notations of S_{xx} & S_{xy}

$$y_c = y - \bar{y} \hat{j}_n = (I - P_{\hat{j}_n}) y$$

$$X_c' y_c = X_c' (y - \bar{y} \hat{j}_n)$$

$$= X_c' y - \bar{y} X_c' \hat{j}_n$$

$$= X_c' y$$

$$\frac{X_c' y}{n-1}$$

$$= \widehat{\text{Cov}}(x, y) \equiv S_{xy}$$

$$\frac{X_c' X_c}{n-1}$$

$$= \widehat{\text{Cov}}(x, x) \equiv S_{xx}$$

Example 7.3.1a. We use the data in Table 7.1 to illustrate computation of $\hat{\beta}$ using (7.6).

$$\mathbf{y} = \begin{pmatrix} 2 \\ 3 \\ 2 \\ 7 \\ 6 \\ 8 \\ 10 \\ 7 \\ 8 \\ 12 \\ 11 \\ 14 \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & 0 & 2 \\ 1 & 2 & 6 \\ 1 & 2 & 7 \\ 1 & 2 & 5 \\ 1 & 4 & 9 \\ 1 & 4 & 8 \\ 1 & 4 & 7 \\ 1 & 6 & 10 \\ 1 & 6 & 11 \\ 1 & 6 & 9 \\ 1 & 8 & 15 \\ 1 & 8 & 13 \end{pmatrix}, \quad \mathbf{X}'\mathbf{X} = \begin{pmatrix} 12 & 52 & 102 \\ 52 & 395 & 536 \\ 102 & 536 & 1004 \end{pmatrix},$$

$$\mathbf{X}'\mathbf{y} = \begin{pmatrix} 90 \\ 482 \\ 872 \end{pmatrix}, \quad (\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} .97476 & .24290 & -.22871 \\ .24290 & .16207 & -.11120 \\ -.22871 & -.11120 & .08360 \end{pmatrix},$$

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \begin{pmatrix} 5.3754 \\ 3.0118 \\ -1.2855 \end{pmatrix}.$$

$$\hat{\alpha} = \bar{y}$$

$$\mathbf{x}_{1,c} = (0 - \bar{x}_1, 2 - \bar{x}_1, \dots, 8 - \bar{x}_1)'$$

$$\mathbf{x}_{2,c} = (2 - \bar{x}_2, 6 - \bar{x}_2, \dots, 13 - \bar{x}_2)'$$

$$\mathbf{X}_c = (\mathbf{x}_{1,c}, \mathbf{x}_{2,c}) : 12 \times 2$$

$$\hat{\beta}_1 = (\mathbf{X}_c' \mathbf{X}_c)^{-1} \mathbf{X}_c' \mathbf{y}$$

Example 7.5. For the data in Table 7.1, we calculate $\hat{\beta}_1$ and $\hat{\beta}_0$ using (7.46) and (7.47).

$$\hat{\beta}_1 = \mathbf{S}_{xx}^{-1} \mathbf{s}_{yx} = \begin{pmatrix} 6.4242 & 8.5455 \\ 8.5455 & 12.4545 \end{pmatrix}^{-1} \begin{pmatrix} 8.3636 \\ 9.7273 \end{pmatrix}$$

$\frac{X_c' y}{n-1} = S_{xy}$

$$= \begin{pmatrix} 3.0118 \\ -1.2855 \end{pmatrix},$$

$\frac{X_c' X_c}{n-1} = S_{xx}$

$$\hat{\beta}_0 = \bar{y} - \mathbf{s}'_{yx} \mathbf{S}_{xx}^{-1} \bar{\mathbf{x}}$$

$\leftarrow \hat{\beta}_1$

$$= 7.5000 - (3.0118, -1.2855) \begin{pmatrix} 4.3333 \\ 8.5000 \end{pmatrix}$$

$$= 7.500 - 2.1246 = 5.3754.$$

These values are the same as those obtained in Example 7.3.1a. □

Sum Squares in Centered

$$SST = \|y - \bar{y} \mathbf{j}_n\|^2 = \left\| \left(\frac{\mathbf{1}}{n} - \mathbf{P}_{\mathbf{j}_n} \right) y \right\|^2$$

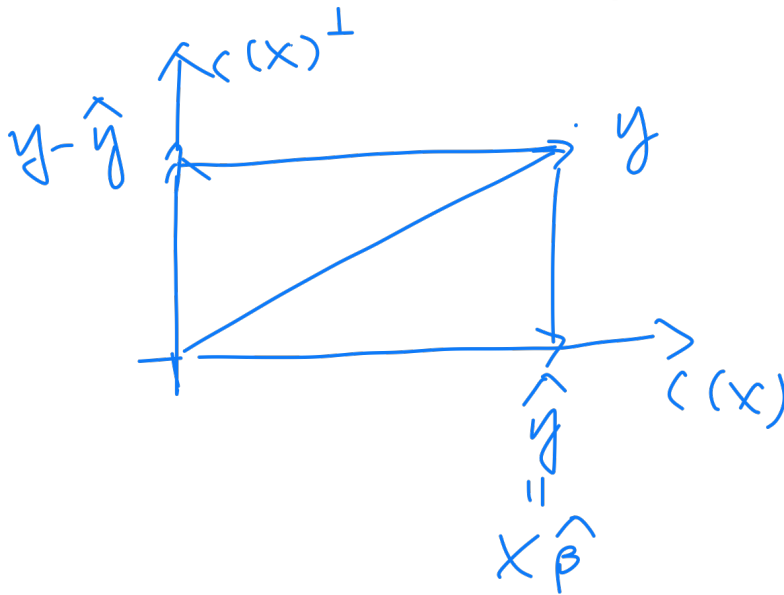
$$\begin{aligned} SSE &= (y - \mathbf{X}\hat{\beta})^T (y - \mathbf{X}\hat{\beta}) = (y - \mathbf{P}_{C(\mathbf{X})}y)^T (y - \mathbf{P}_{C(\mathbf{X})}y) \\ &= y^T y - y^T \mathbf{P}_{C(\mathbf{X})}y - y^T \mathbf{P}_{C(\mathbf{X})}y + y^T \mathbf{P}_{C(\mathbf{X})}y \\ &= y^T y - y^T \mathbf{P}_{C(\mathbf{X})}y = y^T y - \hat{\beta}^T \mathbf{X}^T y. \end{aligned}$$

$$\begin{aligned} \hat{y}_0 &= \mathbf{P}_{\mathbf{j}_n} y \\ &= \bar{y} \mathbf{j}_n \end{aligned}$$

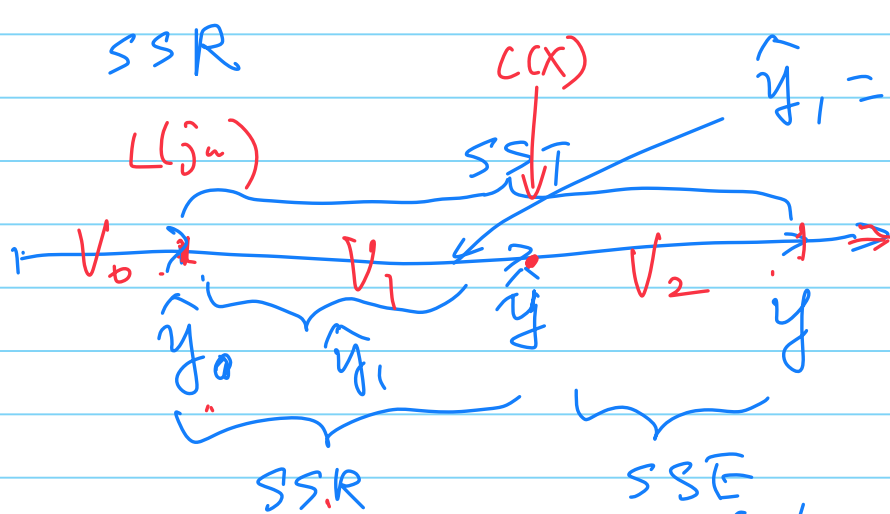
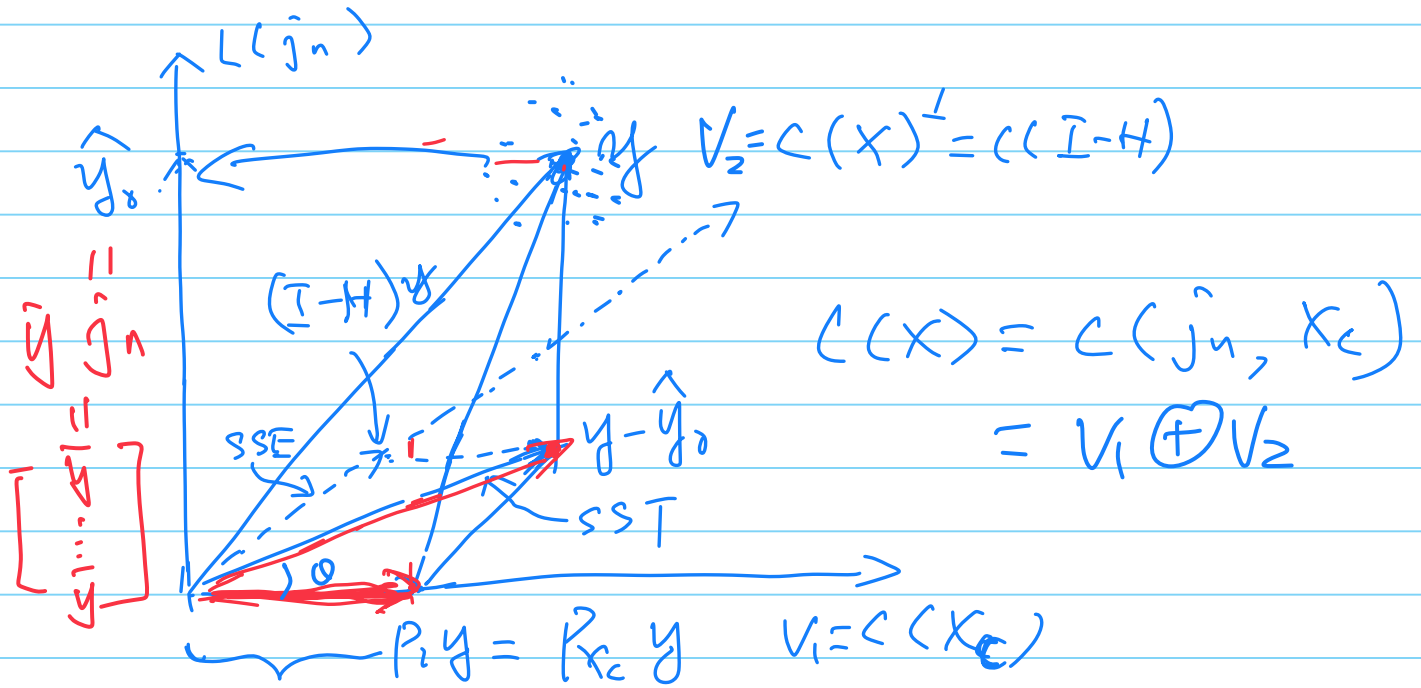
$$\begin{aligned} SSE &= y^T y - (\hat{\alpha}, \hat{\beta}_1^T) \begin{pmatrix} \mathbf{j}_n^T \\ \mathbf{X}_c^T \end{pmatrix} y \\ &= y^T y - \bar{y} \mathbf{j}_n^T y - \hat{\beta}_1^T \mathbf{X}_c^T y \\ &= (y - \bar{y} \mathbf{j}_n)^T y - \hat{\beta}_1^T \mathbf{X}_c^T y \\ &= (y - \bar{y} \mathbf{j}_n)^T (y - \bar{y} \mathbf{j}_n) - \hat{\beta}_1^T \mathbf{X}_c^T y \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 - \hat{\beta}_1^T \mathbf{X}_c^T y \end{aligned}$$

\uparrow
 $\| \hat{y} \|^2$

$$\begin{aligned} &= \sum (y_i - \bar{y})^2 - y^T \mathbf{X}_c (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T y \\ &= SST - \| \mathbf{P}_{\mathbf{X}_c} y \|^2 \end{aligned}$$



$$u_y = u \cdot \hat{j}_n$$



$$SSR = \| P_{X_c} y \|^2 = X_c (X_c' X_c)^{-1} X_c' y$$

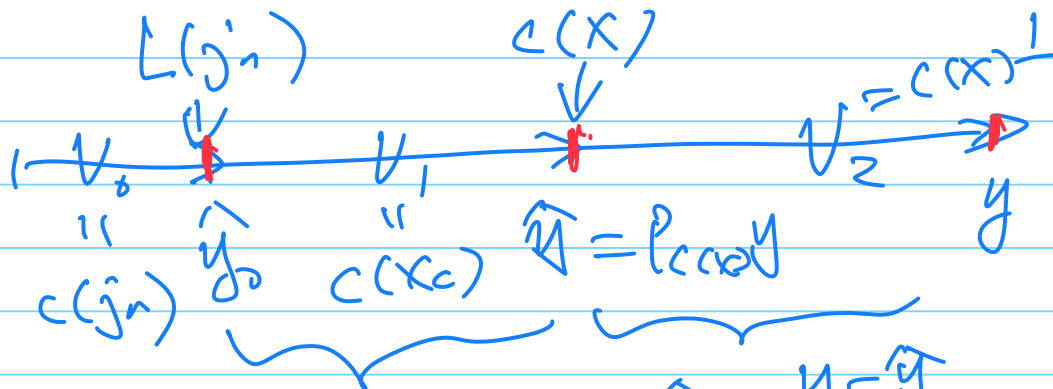
$$SSE = \| P_{C(X)}^\perp y \|^2 = \| (I-H) y \|^2$$

$$SST = SSR + SSE = \| y - \bar{y} \hat{j}_n \|^2$$

$$= \sum_{i=1}^n (y_i - \bar{y})^2$$

$$H = X(X'X)^{-1}X' = P$$

$$L(\hat{j}_n) \subseteq C(X)$$



$$P_{X_c} y = \hat{y} - \hat{y}_0$$

$$| \text{---} SSR \text{---} | \quad | \text{---} SSE \text{---} |$$

$$| \text{---} SST \text{---} |$$

$$SST = SSR + SSE$$

$$C(X) = C(\hat{j}_n, X_c)$$

$$= L(\hat{j}_n) \oplus C(X_c)$$

$$\hat{j}_n \perp X_c$$

$$SSR = \| \hat{y} - \hat{y}_0 \|^2 = \| \hat{y} \|^2 - \| \hat{y}_0 \|^2$$

$$SSE = \| y - \hat{y} \|^2 = \| y \|^2 - \| \hat{y} \|^2$$

$$SST = \| y - \bar{y} \hat{j}_n \|^2 = \| y \|^2 - n \cdot \bar{y}^2$$

R^2 , the Estimated Coefficient of Determination

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \hat{\beta}_1^T \mathbf{X}_c^T \mathbf{y} + \text{SSE}$$

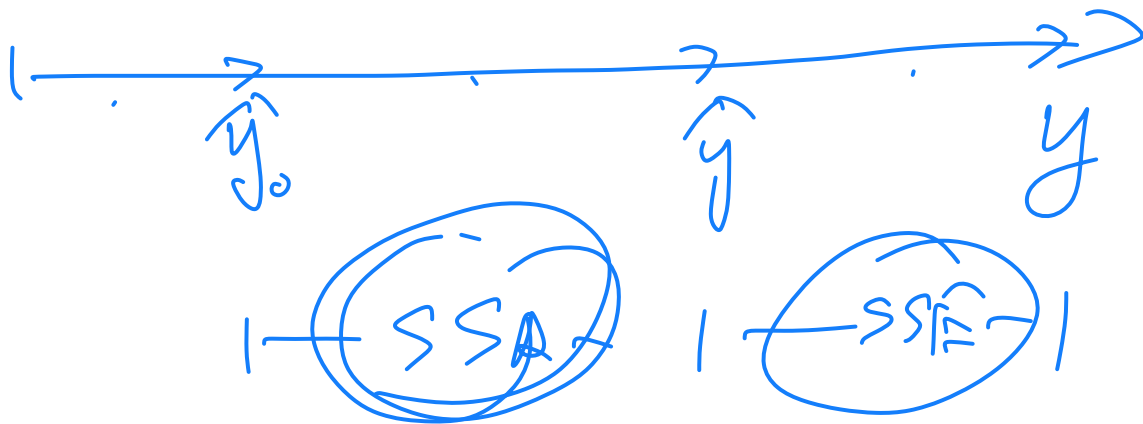
or $\text{SST} = \text{SSR} + \text{SSE}$

$\text{SSR} \perp \text{SSE}$

$$\text{SSR} = \hat{\beta}_1^T \mathbf{X}_c^T \mathbf{y} = \hat{\beta}_1^T \mathbf{X}_c^T \mathbf{X}_c (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T \mathbf{y} = (\mathbf{X}_c \hat{\beta}_1)^T (\mathbf{X}_c \hat{\beta}_1).$$

$$R^2 = \frac{\text{SSR}}{\text{SST}} = \frac{\hat{\beta}_1^T \mathbf{X}_c^T \mathbf{X}_c \hat{\beta}_1}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\hat{\beta}_1^T \mathbf{X}_c^T \mathbf{y} - n\bar{y}^2}{\mathbf{y}^T \mathbf{y} - n\bar{y}^2}.$$

$$= 1 - \frac{\text{SSE}}{\text{SST}}$$



$$0 \leq \frac{\text{SSR}}{\text{SST}} \leq 1$$

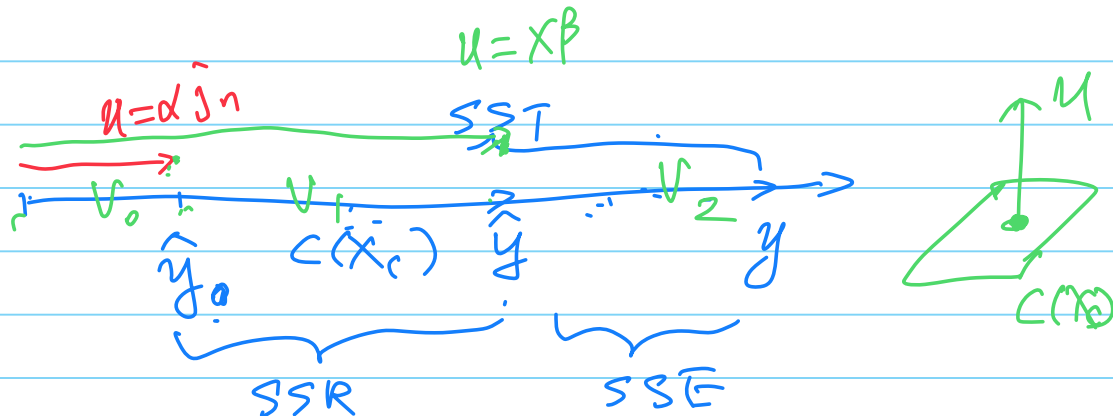
Adjusted R^2

$E(y) = \mu, V(y) = \sigma^2 I$

P is a proj matrix with $\text{rank}(P) = r$, then

$E(\|Py\|^2) = E(y'Py) = \text{tr}(P \cdot \sigma^2 I) + \mu'P\mu$

$= \text{tr}(P) \cdot \sigma^2 + \|P\mu\|^2 = \text{rank}(P) \cdot \sigma^2 + \|P\mu\|^2$



$SSR = \|P_{X_c} y\|^2$

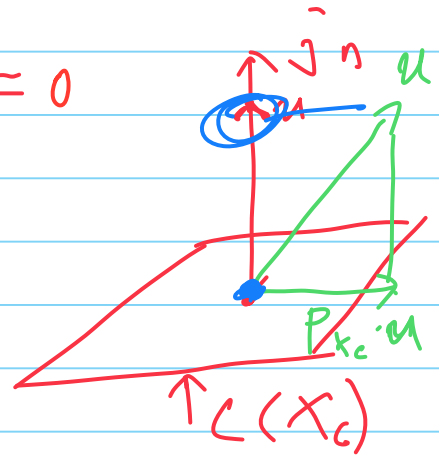
$SSE = \|P_{C(X_c)^\perp} y\|^2 = \|(I - P_x) y\|^2$

$\mu = E(y) = X\beta = \alpha j_n + X_c \beta$

$E(SSR) = \text{rank}(P_{X_c}) \sigma^2 + \|P_{X_c} \mu\|^2$ $j_n \perp X_c$
 $= k \sigma^2 + \|X_c \beta\|^2$

$E(SSE) = (n - k - 1) \sigma^2 + \|P_{C(X_c)^\perp} X\beta\|^2$
 $= (n - k - 1) \sigma^2$

$E(R^2)$ when $\beta_1 = \dots = \beta_k = 0$
 $[\beta_1 = 0]$



$$u \in C(j_n) \perp C(X_c)$$

$$P_{X_c} u = X_c \beta_1 = 0$$

$$E(SSR) = \text{rank}(X_c) \sigma^2 = k \cdot \sigma^2$$

$$E(SSE) = \text{rank}(I - P_X) \sigma^2 + \|(I - P_X) u\|^2$$

$$= (n - k + 1) \sigma^2 + 0$$

Since $u = X\beta \in C(X) \perp C(I - P_X)$.

Also, SSR indep SSE.

$$E(SST) = (n - 1) \sigma^2 = \text{rank}(\hat{j}_n^\perp) \cdot \sigma^2$$

Simply, $E(\|Py\|^2) = \text{Dim}(P) \sigma^2$

$$R^2 = \frac{SSR}{SSR + SSE} \quad \xrightarrow{E} (n-k-1)\sigma^2$$

$$\frac{1}{R^2} = 1 + \frac{SSE}{SSR} \quad \text{note } SSE \xrightarrow{E} k\sigma^2 \text{ indep } SSR$$

$$E\left(\frac{1}{R^2}\right) = 1 + \frac{n-k-1}{k} \quad \left(\begin{array}{l} \text{Dim}(\bar{y}) \\ \text{Dim}(X_c) \end{array} \right)$$

$$= \frac{n-1}{k} \quad \left(\text{Dim}(X_c) \right)$$

$$E(R^2) \approx \frac{k}{n-1} \quad \text{as } k \uparrow$$

We expect $E(R^2) = 0$ when

$$\beta_1 = \beta_2 = \dots = \beta_k = 0$$

Indeed SSE always decreases as $k \uparrow$.

$$\text{Let } dfe = n - k - 1 = \text{rank}(C(X)^T)$$

$$df_{\text{reg}} = k = \text{rank}(X_C)$$

$$dfe + df_{\text{reg}} = n - 1$$

Adjusted R^2 :

$$R_a^2 = 1 - \frac{SSE / dfe}{SST / (dfe + df_{\text{reg}})}$$

$$= 1 - \frac{SSE}{SST} \cdot \frac{n-1}{n-k-1}$$

$$E(R_a^2) \approx 0$$

when $\beta_1 = \beta_2 = \dots = \beta_k = 0$

$$E\left(\frac{SSE}{SST}\right) \approx \frac{n-k-1}{n-1} \quad L(j_n) \downarrow$$

$$E(R^2) \approx 1 - 1 = 0$$

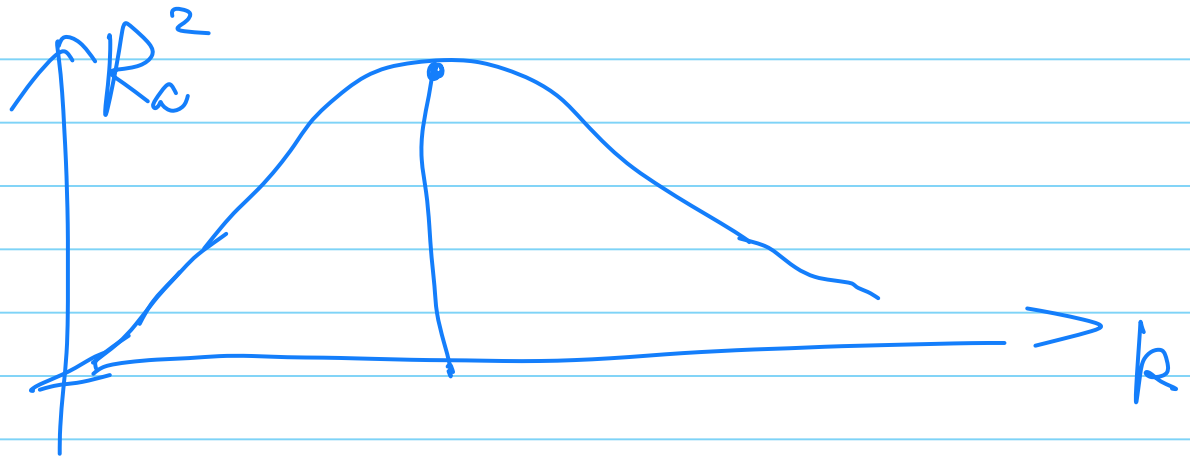
This is desired as

$$\beta_1 = \beta_2 = \dots = \beta_k = 0$$

Large R^2 but low R_a^2 is

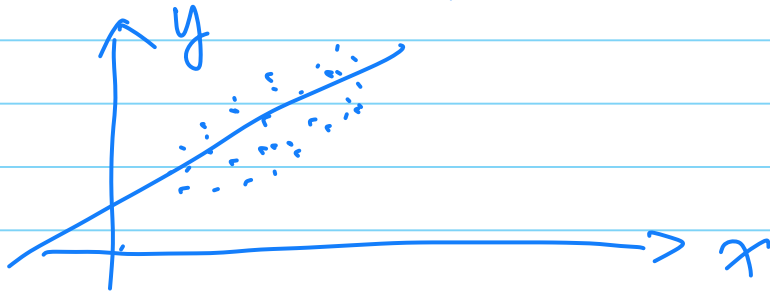
a sign of overfitting

R_a^2 is a criterion for selecting k :



R_a^2 as an estimate of $\rho^2 = 1 - \frac{E(V(y|x))}{V(y)}$

$$V(y) = E(V(y|x)) + V(E(y|x))$$



suppose $y = x\beta + e$, $\text{Cov}(x, e) = 0$

$$\text{SSE} / (n-k-1) = \widehat{\text{Var}(e)}$$

$$\text{SST} / (n-1) = \widehat{\text{Var}(y)} \quad [x \text{ is a r.v. too}]$$

$$R_a^2 = 1 - \frac{\widehat{\text{Var}(e)}}{\widehat{\text{Var}(y)}} \quad \left(= \frac{\text{Var}(e)}{\text{Var}(y)} \right)$$

$$= 1 - \frac{\text{MSE}}{\text{MST}}$$

Underfitting and Overfitting

Suppose that the true model is $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where we return to the spherical errors case: $\text{var}(\mathbf{e}) = \sigma^2\mathbf{I}$. We want to consider what happens when we omit some explanatory variable in \mathbf{X} and when we include too many x 's. So, let's partition our model as

$$\begin{aligned}\mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} = (\mathbf{X}_1, \mathbf{X}_2) \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix} + \mathbf{e} \\ &= \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{e}.\end{aligned}\tag{†}$$

Full

→

- If we leave out $\mathbf{X}_2\boldsymbol{\beta}_2$ when it should be included (when $\boldsymbol{\beta}_2 \neq \mathbf{0}$) then we are **underfitting**.
- If we include $\mathbf{X}_2\boldsymbol{\beta}_2$ when it doesn't belong in the true model (when $\boldsymbol{\beta}_2 = \mathbf{0}$) then we are **overfitting**.
- We will consider the effects of both overfitting and underfitting on the bias and variance of $\hat{\boldsymbol{\beta}}$. The book also consider effects on predicted values and on the MSE s^2 .

Underfitting:

Suppose model (†) holds, but we fit the model

$$\mathbf{y} = \mathbf{X}_1 \boldsymbol{\beta}_1^* + \mathbf{e}^*, \quad \text{var}(\mathbf{e}^*) = \sigma^2 \mathbf{I}. \quad (\clubsuit)$$

reduced

The following theorem gives the bias and var-cov matrix of $\hat{\boldsymbol{\beta}}_1^*$ the OLS estimator from \clubsuit .

Theorem: If we fit model \clubsuit when model (†) is the true model, then the mean and var-cov matrix of the OLS estimator $\hat{\boldsymbol{\beta}}_1^* = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y}$ are as follows:

(i) $E(\hat{\boldsymbol{\beta}}_1^*) = \boldsymbol{\beta}_1 + \mathbf{A} \boldsymbol{\beta}_2$, where $\mathbf{A} = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2$.

(ii) $\text{var}(\hat{\boldsymbol{\beta}}_1^*) = \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1}$.

Proof:

(i)

$$\begin{aligned} E(\hat{\boldsymbol{\beta}}_1^*) &= E[(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y}] = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T E(\mathbf{y}) \\ &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T (\mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2) \\ &= \boldsymbol{\beta}_1 + \mathbf{A} \boldsymbol{\beta}_2. \quad \neq \boldsymbol{\beta}_1 \end{aligned}$$

(ii)

$$\begin{aligned} \text{var}(\hat{\boldsymbol{\beta}}_1^*) &= \text{var}[(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y}] \\ &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T (\sigma^2 \mathbf{I}) \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \\ &= \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1}. \end{aligned}$$

■

- This result says that when underfitting, $\hat{\boldsymbol{\beta}}_1^*$ is biased by an amount that depends upon both the omitted and included explanatory variables.

Corollary If $\mathbf{X}_1^T \mathbf{X}_2 = \mathbf{0}$, i.e.. if the columns of \mathbf{X}_1 are orthogonal to the columns of \mathbf{X}_2 , then $\hat{\boldsymbol{\beta}}_1^*$ is unbiased.

Remarks:

- 1) orthogonal design
- 2) randomization

Over fitting

True model: (Reduced)

$$y = X_1 \beta_1^* + \epsilon$$

Fitted Model: (Full)

$$y = X_1 \beta_1 + X_2 \beta_2 + \epsilon$$

Note that the fitted model

is the true model for

$y | X_1, X_2$ with the true parameter

$$\begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \beta_1^* \\ 0 \end{bmatrix}.$$

$$E \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \beta_1^* \\ 0 \end{pmatrix} \Rightarrow E(\hat{\beta}_1) = \beta_1^*$$

$\hat{\beta}_1$ (from the Full model) is unbiased.

Note that in the above theorem the var-cov matrix of $\hat{\beta}_1^*$, $\sigma^2(\mathbf{X}_1^T \mathbf{X}_1)^{-1}$ is not the same as the var-cov matrix of $\hat{\beta}_1$, the corresponding portion of the OLS estimator $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ from the full model. How these var-cov matrices differ is established in the following theorem:

*n. n. d.
not
negative
definite*

Theorem: Let $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ from the full model (†) be partitioned as

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

and let $\hat{\beta}_1^* = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y}$ be the estimator from the reduced model ♣. Then

$$\text{var}(\hat{\beta}_1) - \text{var}(\hat{\beta}_1^*) = \mathbf{A} \mathbf{B}^{-1} \mathbf{A}^T$$

a n.n.d. matrix. Here, $\mathbf{A} = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2$ and $\mathbf{B} = \mathbf{X}_2^T \mathbf{X}_2 - \mathbf{X}_2^T \mathbf{X}_1 \mathbf{A}$.

- Thus $\text{var}(\hat{\beta}_j) \geq \text{var}(\hat{\beta}_j^*)$, meaning that underfitting results in smaller variances of the $\hat{\beta}_j$'s and overfitting results in larger variances of the $\hat{\beta}_j$'s.

Proof: Partitioning $\mathbf{X}^T \mathbf{X}$ to conform to the partitioning of \mathbf{X} and β , we have

$$\begin{aligned} \text{var}(\hat{\beta}) &= \text{var} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 \begin{pmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 \end{pmatrix}^{-1} \\ &= \sigma^2 \begin{pmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{H}_{21} & \mathbf{H}_{22} \end{pmatrix}^{-1} = \sigma^2 \begin{pmatrix} \mathbf{H}^{11} & \mathbf{H}^{12} \\ \mathbf{H}^{21} & \mathbf{H}^{22} \end{pmatrix}, \end{aligned}$$

(X1', X2')

where $\mathbf{H}_{ij} = \mathbf{X}_i^T \mathbf{X}_j$ and \mathbf{H}^{ij} is the corresponding block of the inverse matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ (see p. 54).

So, $\text{var}(\hat{\beta}_1) = \sigma^2 \mathbf{H}^{11}$. Using the formulas for inverses of partitioned matrices,

$$\mathbf{H}^{11} = \mathbf{H}_{11}^{-1} + \mathbf{H}_{11}^{-1} \mathbf{H}_{12} \mathbf{B}^{-1} \mathbf{H}_{21} \mathbf{H}_{11}^{-1},$$

where

$$\mathbf{B} = \mathbf{H}_{22} - \mathbf{H}_{21} \mathbf{H}_{11}^{-1} \mathbf{H}_{12}.$$

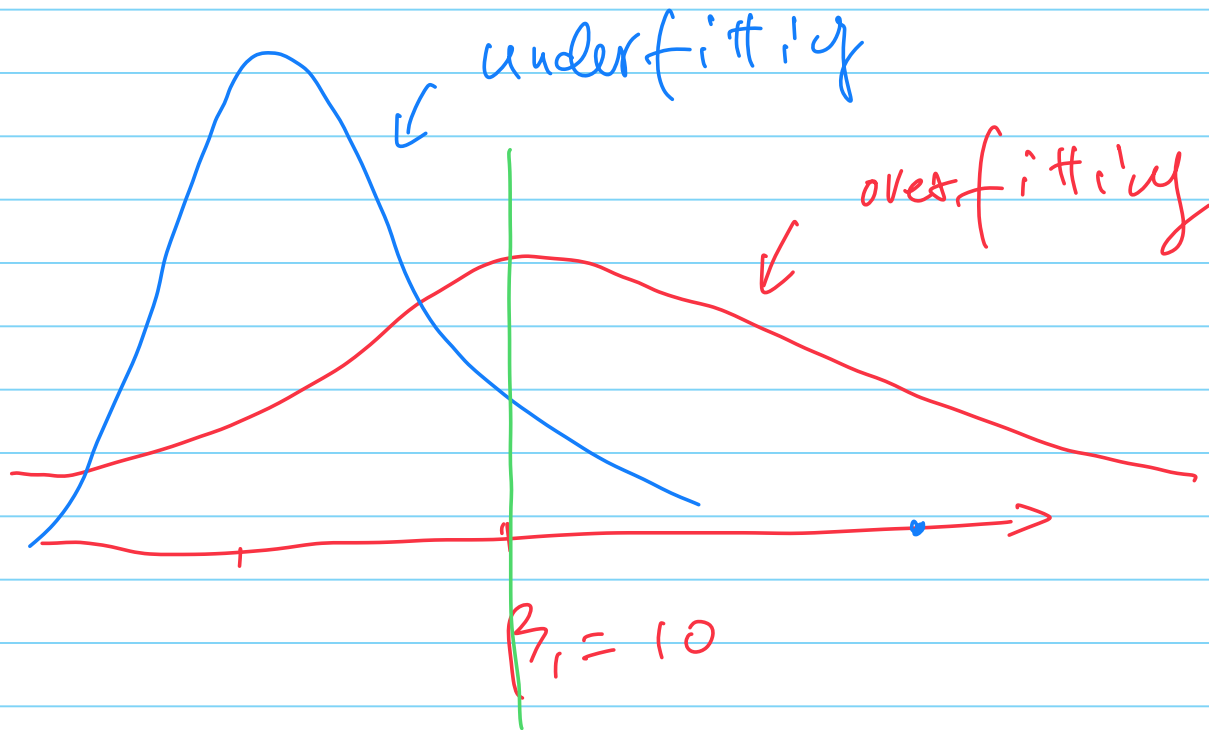
B is n.n.d. b/c X'X is n.n.d.

In the previous theorem, we showed that $\text{var}(\hat{\beta}_1^*) = \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} = \sigma^2 \mathbf{H}_{11}^{-1}$. Hence,

$$\begin{aligned} \text{var}(\hat{\beta}_1) - \text{var}(\hat{\beta}_1^*) &= \sigma^2 (\mathbf{H}^{11} - \mathbf{H}_{11}^{-1}) \\ &= \sigma^2 (\mathbf{H}_{11}^{-1} + \mathbf{H}_{11}^{-1} \mathbf{H}_{12} \mathbf{B}^{-1} \mathbf{H}_{21} \mathbf{H}_{11}^{-1} - \mathbf{H}_{11}^{-1}) \\ &= \sigma^2 (\mathbf{H}_{11}^{-1} \mathbf{H}_{12} \mathbf{B}^{-1} \mathbf{H}_{21} \mathbf{H}_{11}^{-1}) \\ &= \sigma^2 [(\mathbf{X}_1^T \mathbf{X}_1)^{-1} (\mathbf{X}_1^T \mathbf{X}_2) \mathbf{B}^{-1} (\mathbf{X}_2^T \mathbf{X}_1) (\mathbf{X}_1^T \mathbf{X}_1)^{-1}] \\ &= \sigma^2 \mathbf{A} \mathbf{B}^{-1} \mathbf{A}^T. \end{aligned}$$

$$\text{var}(\alpha' \hat{\beta}_1) - \text{var}(\alpha' \hat{\beta}_1^*) = \sigma^2 \alpha' \mathbf{A} \mathbf{B}^{-1} \mathbf{A}' \alpha$$

where $\alpha' = (\alpha', 0, \dots, 0)$
↑
for X_2



true value

To summarize, we've seen that underfitting reduces the variances of regression parameter estimators, but introduces bias. On the other hand, overfitting produces unbiased estimators with increased variances. Thus it is the task of a regression model builder to find an optimum set of explanatory variables to balance between a biased model and one with large variances.

Occam's razor, Ockham's razor, Ocham's razor (**Latin: *novacula Occami***), also known as the principle of parsimony or the law of parsimony (**Latin: *lex parsimoniae***), is the problem-solving **principle** that "entities should not be multiplied beyond necessity".

