

III. Basic Regression Diagnostics

10

DATA: $(Y_i, \mathbf{z}_i^T)^T, i=1, \dots, n$

→ regressors $X_i = t(\mathbf{z}_i)$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad Z = \begin{pmatrix} \mathbf{z}_1^T \\ \vdots \\ \mathbf{z}_n^T \end{pmatrix} \rightarrow X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix} =: t(Z)$$

Is it reasonable to assume, for chosen t :

• $E(Y|Z) = X\beta$ for some $\beta \equiv E(Y|Z) \in \mathcal{H}(X)$

• $\text{var}(Y|Z) = \sigma^2 I_n$

• $Y|Z \sim \mathcal{N}$?
(not needed for many things)

3.1 (Normal) linear model assumptions

1-2

(A1) $E(Y_i | X_i = x) = x^T \beta$ for some β and (almost all) $x \in \mathcal{X}$
 \equiv correct regression function

(A2) $\text{var}(Y_i | X_i = x) = \sigma^2$ for some $\sigma^2 > 0$ irrespective of
(almost all) values of $x \in \mathcal{X}$
 \equiv homoscedasticity

(A3) $\text{cov}(Y_i, Y_\ell | X = x) = 0, i \neq \ell$, for (almost all) $x \in \mathcal{X}^n$
 \equiv the responses are conditionally uncorrelated

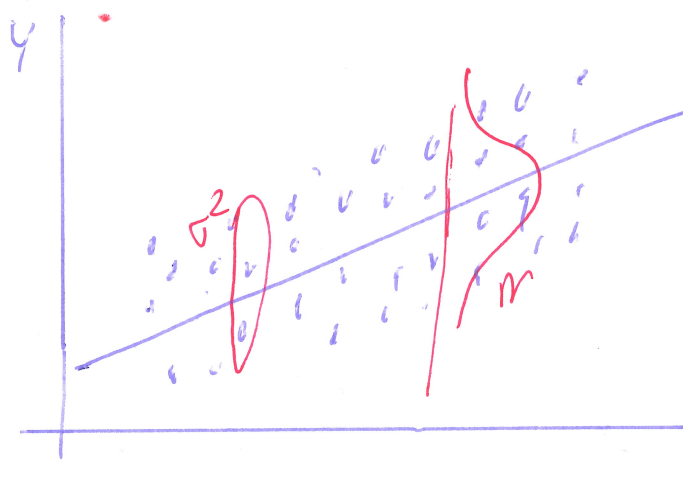
(A4) $Y_i | X_i = x \sim \mathcal{N}(x^T \beta, \sigma^2)$ for (almost all) $x \in \mathcal{X}$
 \equiv normality

Discussion of assumpt.

1

Assumptions in terms of the errors ϵ

3



$$X = (1, z)$$

$$E(Y|Z=z) = \beta_0 + \beta_1 z \quad (\text{MODEL})$$

$$\epsilon = Y - (\beta_0 + \beta_1 z)$$

(A1) $E(\epsilon_i | X_i = x) = 0$ for (almost all) $x \in X$

($\Rightarrow E\epsilon_i = 0, i = 1, \dots, n$)

\equiv the regression function of the model is correctly specified

(A2) $\text{var}(\epsilon_i | X_i = x) = \sigma^2$ for some $\sigma^2 > 0$ which is constant irrespective of (almost all) values of $x \in X$

($\Rightarrow \text{var}(\epsilon_i) = \sigma^2, i = 1, \dots, n$)

\equiv homoscedasticity of the errors

(A3) $\text{cov}(\epsilon_i, \epsilon_l | X = x) = 0, i \neq l$ for (almost all) $x \in X^n$

($\Rightarrow \text{cov}(\epsilon_i, \epsilon_l) = 0, i \neq l$)

\equiv the errors are uncorrelated

(A4) $\epsilon_i | X_i = x \sim N(0, \sigma^2)$ for (almost all) $x \in X$

($\Rightarrow \epsilon_i \sim N(0, \sigma^2), i = 1, \dots, n$)

\equiv ~~the~~ the errors are normally distributed and owing to previous assumptions,

$\epsilon_1, \dots, \epsilon_n \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$

Most checks performed through residuals

4

$$U = MY \quad , \quad M = I_n - H = I_n - X(X^T X)^{-1} X^T \\ = Y - \hat{Y}$$

Assumptions and residual properties

$$(A1) \Rightarrow E(U|X) = 0_n$$

$$(A1) \& (A2) \& (A3) \Rightarrow \text{var}(U|X) = \sigma^2 M.$$

$$(A1) \& (A2) \& (A3) \& (A4) \Rightarrow U|X \sim N_n(0, \sigma^2 M)$$

Strategy

RHS of implication checked

- if not satisfied then also LHS not valid

- if satisfied, still no guarantee that the LHS valid

one more complication

$$Y - X\hat{\beta} = U \equiv \hat{\varepsilon} \quad , \quad \varepsilon = Y - X\beta$$

$$\text{var}(U|X) = \sigma^2 M$$

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

↑
- not diagonal matrix
- not const diagonal

↑
- uncorrelated
- homoscedast.

3.2 Standardized residuals

5

We know: $\text{var}(U_i | X_i) = \sigma^2 m_{ii}$, $i=1, \dots, n$

$$\Rightarrow \text{var} \left(\frac{U_i}{\sqrt{\sigma^2 m_{ii}}} \mid X_i \right) = 1$$

LATER: $m_{ii} = 0 \Leftrightarrow \text{rank}(X_{(-i)}) = \text{rank}(X) - 1$
matrix X without row i

$$\boxed{? \text{ var} \left(\frac{U_i}{\sqrt{\text{MSE}_{m_{ii}}}} \mid X_i \right) = 1 ?}$$

Def 3.1 Standardized residuals

6

The standardized residuals or the vector of standardized residuals of the model is the vector $U^{\text{std}} = (U_1^{\text{std}}, \dots, U_n^{\text{std}})^T$, where

$$U_i^{\text{std}} = \begin{cases} \frac{U_i}{\sqrt{\text{MSE}_{m_{ii}}}} & , m_{ii} > 0, \\ \text{undefined} & , m_{ii} = 0, \end{cases} \quad i=1, \dots, n.$$

Lemma 3.1 Moments of standardized residuals under normality

Let $Y|X \sim N_n(X\beta, \sigma^2 I_n)$ and let for chosen $i \in \{1, \dots, n\}$, $m_{ii} > 0$. Then

$$E(U_i^{std} | X) = 0, \quad \text{var}(U_i^{std} | X) = 1.$$

Proof: See the notes (not for exam).

First Lemma B.2 $Z \sim N_m(0, \sigma^2 I_m)$

$T: \mathbb{R}^m \rightarrow \mathbb{R}$ measurable fun.,
scale transformations invariant: $\forall c > 0 \quad \forall z \in \mathbb{R}^m$
 $T(cz) = T(z) \Rightarrow T(z) \propto \|z\|.$

$$U_i^{std} = \frac{U_i}{\sqrt{MSE_{m_{ii}}}} = \frac{U_i}{\sqrt{SSE}} \sqrt{\frac{n-k}{m_{ii}}} = \frac{U_i}{\|U\|} \underbrace{\sqrt{\frac{n-k}{m_{ii}}}}_{= \text{const given } X}$$

$U = M \cdot Y = N \cdot N^T Y$, $N =$ orthonormal basis of $\mathcal{N}(X)^\perp$

given X : $Y \sim N(X\beta, \sigma^2 I_n) \Rightarrow$

$$N^T Y \sim N(\underbrace{N^T X}_{0} \beta, \sigma^2 \underbrace{N^T N}_{I_{n-k}})$$

$n-k$, $k = \text{rank}(X)$
columns of N ,
scalar products with cols of X ,
cols of $N =$ basis of $\mathcal{N}(X)^\perp$

That is, $N^T Y | X \sim N(0, \sigma^2 I_{n-r})$

Z from Lemma

$$T(Z) := \frac{j_i^T (NZ)}{\|NZ\|} = \frac{NN^T Y}{\|NN^T Y\|} = U \quad j_i = (0, \dots, 1, \dots, 0)^T$$

place i

$$T(cZ) = \frac{j_i^T NcZ}{\|cNZ\|} = T(Z)$$

Lemma

$$\Rightarrow T(Z) = \frac{U_i}{\|U\|} \quad \text{and} \quad \|Z\| = \|N^T Y\| = \sqrt{Y^T N N^T Y} = \|U\|$$

are independent (given X) $M=MM^T$

All expectations given X :

$$E\left(\frac{U_i}{\|U\|} \cdot \|U\|\right) = E\left(\frac{U_i}{\|U\|}\right) \cdot E\|U\|$$

$E U_i = 0$
we know

since $\|U\| \geq 0$
and $\frac{1}{\sigma^2} \|U\|^2 \sim \chi_{n-r}^2$

hence it cannot be $\|U\| = 0$ a.s.

$$\text{That is } 0 = E\left(\frac{U_i}{\|U\|} | X\right)$$

$$\Rightarrow E(U_i^{\text{std}} | X) = 0$$

$$E\left(\frac{U_i^2}{\|U\|^2} \cdot \|U\|^2\right) = E\left(\frac{U_i^2}{\|U\|^2}\right) \cdot E\|U\|^2$$

$E U_i^2 = \sigma^2 m_{ii}$
we know

$$\text{That is, } E\left(\frac{U_i^2}{\|U\|^2} | X\right) = \text{var}\left(\frac{U_i^2}{\|U\|^2} | X\right)$$

$$= \frac{\sigma^2 m_{ii}}{\sigma^2 (n-r)} = \frac{m_{ii}}{n-r}$$

$$\Rightarrow \text{var}(U_i^{\text{std}} | X) = 1$$

3.3 Graphical tools of regression diagnostics

8

Model matrix $X = (X^0, \dots, X^{k-1}) = t(Z)$

mostly $X^0 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$

often: additional regressors are available

$V := (V^1, \dots, V^m) = t_v(Z)$

(A1) $E(Y|Z) = X\beta$ for some $\beta \in \mathbb{R}^k$?

• perhaps $= X\beta + V\gamma$ ($\gamma \in \mathbb{R}^m$)

(A2) $\text{var}(Y|Z) = \sigma^2 I_n$

• perhaps $\sigma^2 = \sigma^2(X), \sigma^2(V), \sigma^2(X, V)$

(A1) plots

$E(U|Z) = 0$

9-10

↑
• does not depend on X, V
any function of them
 $(\hat{Y}), \dots$

(A2) homoscedasticity - plots + ...

11-13

Problems when evaluating plots:

$\text{var}(U_i | Z) = \sigma^2 m_{ii}$, raw residuals are NOT homoscedast. even if (A2) OK

$\text{var}(U_i^{\text{std}} | Z) = 1$ (only under N)

scale - location plot: $\sqrt{|U_i^{\text{std}}|}$ on y-axis

(A3) uncorrelated errors

14-15

(A4) normality

16

• if N holds $U|Z \sim N(0, \sigma^2 M)$

- not diagonal matrix
- diagonal not constant (heteroscedasticity)

$U^{\text{std}} | Z \sim \boxed{?} (0, \begin{pmatrix} 1 & * \\ * & 1 \end{pmatrix})$

• not necessarily normal

not diagonal matrix

• overview of three the most important plots

17

• Examples

18-23