

9.3 Homoscedasticity

34

Assumed model (to be checked)

35

$$M: Y|Z \sim (X\beta, \sigma^2 I_n), \quad \varepsilon = Y - X\beta$$

$$E(\varepsilon|Z) (= E(\varepsilon)) = 0_n$$

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

Elementwise: $Y_i = X_i^T \beta + \varepsilon_i, \quad i=1, \dots, n$

$$E(\varepsilon_i | Z_i) = 0 \quad (= \text{assumption A1})$$

$$\text{var}(\varepsilon_i | Z_i) = \sigma^2 \quad (= \text{assumption A2})$$

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

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

additionally: (+normality): $\text{var}(U_i^{\text{std}} | Z) = 1$
used on residual plots

9.3.1 Tests of homoscedasticity

36

\rightarrow based on considering the following hypotheses

$$H_0: \text{var}(\varepsilon_i | Z_i) = \text{const} (= \sigma^2)$$

$$H_1: \text{var}(\varepsilon_i | Z_i) = \text{some function of } Z_i$$

chosen

- depending on some (chosen) function,
different sensitivity towards
violation of H_0 is reached

9.3.2 Score tests of homoscedasticity

Model under H_0 (\equiv full-rank normal linear model)

$$M: Y|Z \sim N_n(X\beta, \sigma^2 I_n), \text{rank}(X) = k$$

elementwise: $Y_i = X_i^T \beta + \varepsilon_i$, $\varepsilon_i | Z_i \sim N(0, \sigma^2)$
(independent)

Model under H_1 (\equiv generalization of a general normal linear model)

$$M_{\text{hetero}}: Y|Z \sim N_n(X\beta, \sigma^2 W^{-1}), \begin{matrix} \mathbb{R}^n \\ \cup \\ \mathbb{R}^k \\ \cup \\ \mathbb{R}^p \end{matrix}$$

$$W = \text{diag}(w_1, \dots, w_n), \quad w_i^{-1} = \tau(\alpha, \beta, Z_i)$$

$i=1, \dots, n$

a known function such that $\tau(0, \beta, z) = 1$

for all $\beta \in \mathbb{R}^k, z \in \mathbb{R}^p$

elementwise: $Y_i = X_i^T \beta + \varepsilon_i$,

$$\varepsilon_i | Z_i \sim N(0, \sigma^2 \underbrace{\tau(\alpha, \beta, Z_i)}_{=1 \text{ if } \alpha=0})$$

Test of homoscedasticity

$$H_0: \alpha = 0$$

$$H_1: \alpha \neq 0$$

under normality, both M and M_{hetero} are full parametric models

- ML inference is possible

$$M \subset M_{hetero}$$

$$M = M_{hetero} \text{ with } \lambda = 0$$

- Test on $\lambda \equiv$ tests of homoscedasticity in a normal LM

→ classical asymptotic tests based on ML theory (Wald, score, likelihood ratio)

will follow: some special choices of τ function ($\text{var}(\epsilon_i | Z_i) = \sigma^2 \tau(\alpha, \beta, Z_i)$) used in practice

Breusch-Pagan test

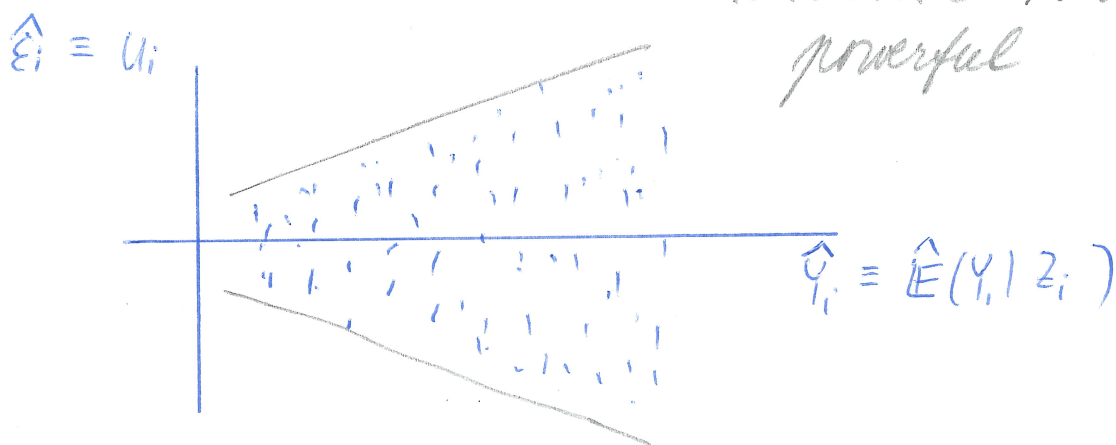
38

score test of $H_0: \alpha = 0$ in a model
with $\text{var}(\varepsilon_i | Z_i) = \sigma^2 \underbrace{\exp(\alpha X_i^T \beta)}_{\tau(\alpha, \beta, Z_i)}$

$X_i = t_x(Z_i) \equiv$ regressors used also in
model $E(Y_i | Z_i)$ as $X_i^T \beta$

Hetero: $\text{var}(\varepsilon_i | Z_i) = \text{var}(Y_i | Z_i)$
 $= \sigma^2 \exp(\alpha \cdot E(Y_i | Z_i))$

\equiv increase / decrease of
residual variability with $E(Y|Z)$
 \rightarrow situation when BP test
powerful



+ additional comments on slide

Linear dependence on the regressors

$$M_{hetero}: \text{var}(\epsilon_i | Z_i) = \text{var}(Y_i | Z_i) = \sigma^2 \exp(\alpha^T W_i)$$

$W_i = \text{tw}(Z_i) \equiv$ some function (parameterization) of covariates

written differently:

$$M_{hetero}: \log(\text{var}(\epsilon_i | Z_i)) = \underbrace{\log \sigma^2}_{\alpha_0} + \alpha^T W_i$$

\equiv linear model for log-residual variance

most commonly used:

$X_i = (X_{i0}, X_{i1}, \dots, X_{ik-1})^T$ regressors used in the model

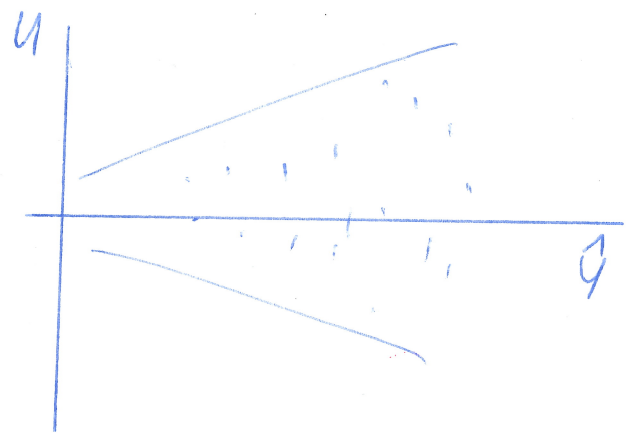
$$M_{hetero}: \text{var}(\epsilon_i | Z_i) = \sigma^2 \exp(\alpha X_{ij})$$

$$\log(\text{var}(\epsilon_i | Z_i)) = \underbrace{\log \sigma^2}_{\alpha_0} + \alpha X_{ij}$$

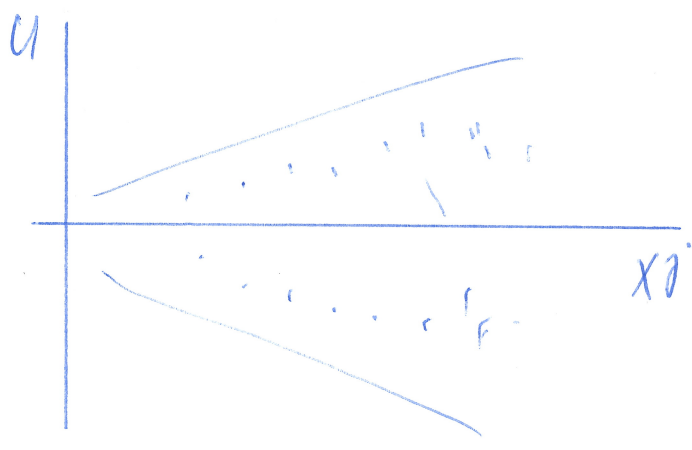
\equiv linear dependence of $\log(\text{var}(\epsilon_i | Z_i))$ on the j th regressor

- which Metrics to consider in practice?
- which type of heteroscedasticity is hypothesized?

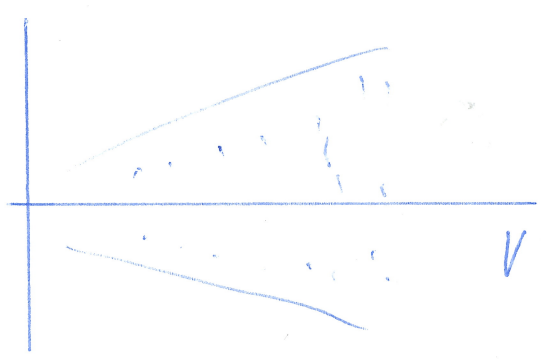
→ residual plots help



$$\tau(\alpha, \beta, z) = \exp(\alpha X^T \beta)$$



$$\tau(\alpha, \beta, z) = \exp(\alpha x_j)$$



$$\tau(\alpha, \beta, z) = \exp(\alpha V)$$

→ covariate not included in the model

R software: see slide

9.3.3 Some other tests of homoscedasticity

- for info only
- see slides

9.4 Normality

Assumed model (to be checked)

$$M: Y|Z \sim N_n(X\beta, \sigma^2 I_n), \text{rank}(X_{n \times k}) = r \leq k$$

$$\Rightarrow \varepsilon_i = Y_i - X_i^T \beta \text{ satisfy } \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad i=1, \dots, n$$

(A4) normality

$$\varepsilon_i | Z \stackrel{indep.}{\sim} N, \quad \varepsilon_i \stackrel{iid}{\sim} N$$

(A1) & (A2) & (A3)

$$\& (A4) \quad \Rightarrow U|Z \sim N_n(0, \sigma^2 M)$$

$$\Rightarrow U_i^{std} | Z \sim (0, 1) \quad i=1, \dots, n$$

Problems: Even with normality

given Z : U_1, \dots, U_n are normal
but neither independent
nor homoscedastic

$U_1^{std}, \dots, U_n^{std}$ are homoscedastic
but not normal

slide: reminder of notation
+ expression of important
quantities

41

42

43

normality of errors $\epsilon_1, \dots, \epsilon_n$

$$\Rightarrow U|Z \sim N_n(0_n, \sigma^2 M)$$

$$\Rightarrow U_i^{std}|Z \sim (0, 1), i=1, \dots, n$$

approximate approaches to test

H_0 : distribution of $\epsilon_1, \dots, \epsilon_n$ is normal

→ classical tests of normality

(Shapiro-Wilk, Lilliefors, Anderson-Darling, ...)

applied to (i) raw residuals U_1, \dots, U_n

(ii) std. residuals $U_1^{std}, \dots, U_n^{std}$

empirical studies show • significance level is maintained

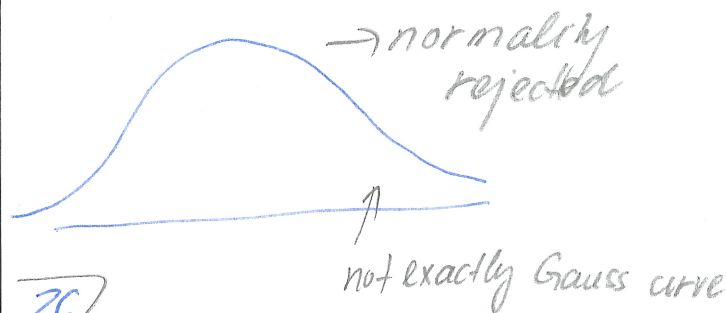
• raw residuals more recommended

standard problem on why tests of normality quite useless

(a) small sample size \equiv small power



(b) large sample size \equiv large power to reject normality even if only negligibly violated



BUT, for many things, we have asymptotics and do not need normality!

9.5 Uncorrelated errors

47

Typical situations when uncorrelated errors cannot be taken for granted

(i) Time series

Y_1 Y_2 Y_3 Y_4 Y_5 ...



$T: t_1$ t_2 t_3 t_4 t_5 ... \leftarrow time

Z_1 Z_2 Z_3 Z_4 Z_5 ... \leftarrow additional covariates

interest to model

$$E(Y_i | T_i, Z_i)$$

model: $Y_i = \boxed{m(T_i, Z_i)} + \varepsilon_i$

trend which may depend on covariates, possibly linear model \downarrow errors

\rightarrow see NMJA 409 Stochastic Processes 2

NMST 537 Time Series

(ii) Repeated measurements / longitudinal data

Data: $(Y_{i1}, \dots, Y_{im_i})^T, i=1, \dots, N$
 Y_i

Response vector $Y = (Y_1^T, \dots, Y_N^T)^T$

Covariates: Z_{i1}, \dots, Z_{im_i}

Primary interest (again) to model $E(Y_{ij} | Z_{ij})$

Example: $i \equiv$ herd, $j \equiv$ cow

herd 1,, herd N

Y_{11}

Y_{N1}

\vdots

\vdots

Y_{1,m_1}

Y_{N,m_N}

Y_{ij} = laboratory measurement of response to some vaccine

$Z_{ij} \equiv$ type of vaccine (A, B) / age of cow

cow

Possible model: $Y_{ij} = m(Z_{ij}) + \epsilon_{ij}$

regression function, perhaps linear model

$\epsilon_i = (\epsilon_{i1}, \dots, \epsilon_{i,m_i})^T$

$\epsilon_N = (\epsilon_{N1}, \dots, \epsilon_{N,m_N})^T$

|| \leftarrow
(uncorrelated)

not necessarily uncorrelated.

(or even independent)

= same environment, . . .

\rightarrow NMST 432 Advanced regression models

will follow: test for uncorrelated errors

47

but only for situations where ordering of observations has a practical meaning and may induce dependence between $\epsilon_1, \dots, \epsilon_n$
(e.g. time series like data)

9.5.1 Durbin-Watson test

48

H_0 : errors are uncorrelated

H_1 : are not

- test based on considering the following models

$$H_0 \equiv M: Y_i = X_i^T \beta + \epsilon_i, \quad i=1, \dots, n$$

$$\mathbb{E}(\epsilon_i | X_i) = 0, \quad \text{var}(\epsilon_i | X_i) = \sigma^2$$

$$\text{cor}(\epsilon_i, \epsilon_\ell | X) = 0, \quad i \neq \ell$$

$$\epsilon_i | X \sim \text{WN}(0, \sigma^2)$$

white noise

$$H_1 \equiv \text{MAR}: Y_i = X_i^T \beta + \epsilon_i, \quad i=1, \dots, n$$

$$\epsilon_1 = \eta_1, \quad \epsilon_i = \rho \epsilon_{i-1} + \eta_i, \quad i=2, \dots, n$$

$$\mathbb{E}(\eta_i | X) = 0, \quad \text{var}(\eta_i | X) = \sigma^2$$

$$\text{cor}(\eta_i, \eta_\ell | X) = 0, \quad i \neq \ell$$

$$\epsilon_i | X \sim \text{AR}_1(\rho)$$

$-1 < \rho < 1$: additional unknown parameter of the model

Properties of $AR_1(\rho)$: $\text{cov}(\epsilon_{i+m}, \epsilon_i | X) = \rho^m$
 $m \geq 1$

48

→ Test for uncorrelated errors $H_0: \rho = 0$

$H_1: \rho \neq 0$

(or one-sided
 $\rho > 0$)

Durbin-Watson test statistic

49

$U = (U_1, \dots, U_n)^T$: residuals from the (linear) model M

$$DW = \frac{\sum_{i=2}^n (U_i - U_{i-1})^2}{\sum_{i=1}^n U_i^2}$$

IDEA: under H_0 (and model M): $E(U|X) = 0$,

also under H_1 (and model M_{AR}): $E(U|X) = 0$

(only $E(Y|X) \in \mathcal{M}(X)$ is needed to show that $E(U|X) = 0$)

$U_1, \dots, U_n \equiv$ predictions of $\epsilon_1, \dots, \epsilon_n$

→ suitable estimator of $\text{cov}(\epsilon_t, \epsilon_{t-1} | X)$ is

$$\hat{\sigma}_{12} = \frac{1}{\underbrace{n-1}_{\text{or } n-2}} \sum_{t=2}^n U_t U_{t-1} \quad (\text{sample covariance})$$

→ suitable estimators of $\text{var}(\epsilon_t | X)$:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{t=1}^{n-1} U_t^2 \quad \text{or} \quad \frac{1}{n-1} \sum_{t=2}^n U_t^2 \quad \text{or} \quad \frac{1}{n-1} \sum_{t=1}^n U_t^2$$

$$DW = \frac{\sum_{i=2}^n (U_i - U_{i-1})^2}{\sum_{i=1}^n U_i^2} = \frac{\sum_{i=2}^n U_i^2 + \sum_{i=1}^{n-1} U_i^2 - 2 \sum_{i=2}^n U_i U_{i-1}}{\sum_{i=1}^n U_i^2}$$

$$\approx \frac{\hat{\sigma}^2 + \hat{\sigma}^2 - 2\hat{\sigma}_{12}}{\hat{\sigma}^2} = 2 \left(1 - \frac{\hat{\sigma}_{12}}{\hat{\sigma}^2} \right)$$

↙
↘

That is, $\hat{\rho} = 1 - \frac{1}{2} DW$

$$\underline{DW = 2(1 - \hat{\rho})}$$

under H_0 : $DW \approx 2$

additional comments → see slide

bootstrap → NMST 434 Modern Statistical Methods

