

IX. Checking Model Assumptions

Data $(Y_i, Z_i), i=1, \dots, n$, $Z_i \in Z \subseteq \mathbb{R}^p$

possibly two sets of regressors

$$X_i = t_x(Z_i), \quad t_x: \mathbb{R}^p \rightarrow \mathbb{R}^k$$

$$V_i = t_v(Z_i), \quad t_v: \mathbb{R}^p \rightarrow \mathbb{R}^l$$

→ two model matrices

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

$$V = \begin{pmatrix} V_1^T \\ \vdots \\ V_n^T \end{pmatrix} = (V^1, \dots, V^l)$$

Assumed linear model for $Y = (Y_1, \dots, Y_n)^T$

$$E(Y|Z) = X\beta \quad \text{for some } \beta$$

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

sometimes also assumed that $Y|Z \sim N$

Assumptions behind $Y/Z \sim N(X\beta, \sigma^2 I_n)$ 2

$$\epsilon := Y - X\beta = (Y_1 - X_1^T \beta, \dots, Y_n - X_n^T \beta)^T = (\epsilon_1, \dots, \epsilon_n)^T$$

remember: $X_i = t_x(Z_i)$

1. Correct regression function

$$E(Y_i | Z_i) = X_i^T \beta \quad \text{for some } \beta$$

$$\equiv E(\epsilon_i | Z_i) = 0 \quad \forall i$$

$\overline{\text{no dependence on } Z_i}$

2. (Conditionally) homoscedasticity of errors

$$\text{var}(Y_i | Z_i) = \sigma^2 \quad \text{for some } \sigma^2$$

$$\equiv \text{var}(\epsilon_i | Z_i) = \sigma^2 \quad \forall i$$

$\overline{\text{no dependence on } Z_i}$

3. (Conditionally) uncorrelated / independent errors

$\text{var}(Y/Z)$ is diagonal

4. (Conditionally) normal errors

$$Y_i | Z_i \sim N, \quad \epsilon_i | Z_i \sim N \quad \forall i$$

1. $E(\epsilon/Z) = 0$

2. & 3. $\text{var}(\epsilon/Z) = \sigma^2 I_n$

4. $\epsilon/Z \sim N$

also marginally, but now not that important

1. $\Rightarrow E\epsilon_i = 0 \quad \forall i$

2. $\Rightarrow \text{var}\epsilon_i = \sigma^2 \quad \forall i$

3. $\Rightarrow \text{cov}(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$

4. $\Rightarrow \epsilon_i \overset{iid}{\sim} N$

9.1 Model with added regressors

3

TECHNICAL SECTION

Let us assume two (competing) models:

$$M: Y|Z \sim (X\beta, \sigma^2 I_n), \quad X_{n \times k}$$

$$M_g: Y|Z \sim (X\beta + V\gamma, \sigma^2 I_n), \quad V_{n \times l}$$

$$G \cdot \begin{pmatrix} \beta \\ \gamma \end{pmatrix}, \quad G = (X, V)$$

Quantities derived while assuming model
 $M: Y|Z \sim (X\beta, \sigma^2 I_n)$

4

$$b = (X^T X)^{-1} X^T Y \quad = \text{any solution to normal equations } X^T X \beta = X^T Y$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad = \text{LSE of } \beta \text{ if } \text{rank}(X) = k$$

$$H = X(X^T X)^{-1} X^T = (h_{it})_{i,t=1,\dots,n} \quad = \text{hat matrix} \\ = \text{projections to } \mathcal{V}(X)$$

$$\hat{Y} = HY = (\hat{y}_1, \dots, \hat{y}_n)^T = \text{fitted values}$$

$$M = I_n - H = (m_{it})_{i,t=1,\dots,n} \quad = \text{projections to } \mathcal{V}(X)^\perp$$

$$U = Y - \hat{Y} = MY = (u_1, \dots, u_n)^T = \text{residuals}$$

$$SSE = \|U\|^2$$

3

Quantities derived while assuming model

$$M_g: Y|Z \sim (X\beta + V\gamma, \sigma^2 I_n), \quad G = \begin{pmatrix} X & V \\ \hline & \end{pmatrix} \begin{matrix} k & l \end{matrix}$$

$$(b_g^T, c_g^T)^T = (G^T G)^{-1} G^T Y \equiv \text{any solution to normal eq. } G^T G \begin{pmatrix} \beta \\ \gamma \end{pmatrix} = G^T Y$$

$$(\hat{\beta}_g^T, \hat{\gamma}_g^T)^T = (G^T G)^{-1} G^T Y \equiv \text{LSE of } \begin{pmatrix} \beta \\ \gamma \end{pmatrix} \text{ if rank}(G) = k+l$$

$$H_g = G (G^T G)^{-1} G^T = (h_{g,ij})_{i,j=1,\dots,n}$$

$$\hat{Y}_g = H_g Y = (\hat{y}_{g,1}, \dots, \hat{y}_{g,n})^T$$

$$M_g = I_n - H_g = (m_{g,ij})_{i,j=1,\dots,n}$$

$$U_g = Y - \hat{Y}_g = M_g Y = (u_{g,1}, \dots, u_{g,n})^T$$

$$SS_{g,e} = \|U_g\|^2$$

Lemma 9.1 Model with added regressor

6

Quantities derived while assuming model $M: \mathcal{Y} | \mathcal{Z} \sim (X\beta, \sigma^2 I_n)$ and quantities derived while assuming model $M_g: \mathcal{Y} | \mathcal{Z} \sim (X\beta + V\gamma, \sigma^2 I_n)$ are mutually in the following relationship.

$$\begin{aligned}\hat{\mathcal{Y}}_g &= \hat{\mathcal{Y}} + MV(V^T M V)^{-1} V^T U \\ &= Xb_g + Vc_g \quad \text{for some } b_g \in \mathbb{R}^k, c_g \in \mathbb{R}^l\end{aligned}$$

Vectors b_g and c_g such that $\hat{\mathcal{Y}}_g = Xb_g + Vc_g$ satisfy

$$c_g = (V^T M V)^{-1} V^T U$$
$$b_g = b - (X^T X)^{-1} X^T V c_g$$

for some $b = (X^T X)^{-1} X^T \mathcal{Y}$.

Finally $SSE - SSE_g = \|MVc_g\|^2$.

Proof: • $\hat{\mathcal{Y}}_g$ is projection of \mathcal{Y} into $\mathcal{M}(X, V) = \mathcal{M}(X, \underbrace{MV}_I)$

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

$$MV = V - X(X^T X)^{-1} X^T V$$

→ $\mathcal{M}(X, MV)$ is still a vector space being generated by columns from X, V .

• Use " $H = X(X^T X)^{-1} X^T$ " to calculate Hg :

$$Hg = (X, MV) \begin{pmatrix} X^T X & X^T MV \\ \underbrace{V^T M X}_0 & \underbrace{V^T M V}_{M \cdot M} \end{pmatrix}^{-1} \begin{pmatrix} X^T \\ V^T M \end{pmatrix} =$$

$$= (X, MV) \begin{pmatrix} (X^T X)^{-1} & 0 \\ 0 & (V^T M V)^{-1} \end{pmatrix} \begin{pmatrix} X^T \\ V^T M \end{pmatrix} = \underbrace{X(X^T X)^{-1} X^T}_H + MV(V^T M V)^{-1} V^T M$$

5

So that $\hat{y} = Hg = HY + MV(V^T M V)^{-1} V^T U =$
 $= \hat{y} + MV(V^T M V)^{-1} V^T U$

- The fitted values \hat{y} must lie in the corresponding regression space $\mathcal{M}(X, V)$
 → it must exist $b_g \in \mathbb{R}^k, c_g \in \mathbb{R}^l$ such that
 $\hat{y} = X b_g + V c_g$

At the same time, $(b_g^T, c_g^T)^T$ must minimize the sum of squares of model Mg .

Proof of Lemma 2.1: vector $(b_g^T, c_g^T)^T$ minimizes the sum of squares \Leftrightarrow
 \Leftrightarrow it solves corresponding normal eqs.

Try to rewrite \hat{y} to see what b_g and c_g could be.
 Same arguments: $\hat{y} = X b \Leftrightarrow b = (X^T X)^{-1} X^T y$
 (b solves normal eqs.)

$\hat{y} = X b + \underbrace{(I_n - X(X^T X)^{-1} X^T)}_M V (V^T M V)^{-1} V^T U =$

we need $X \cdot \underbrace{b - (X^T X)^{-1} X^T V (V^T M V)^{-1} V^T U}_{b_g} + V \cdot \underbrace{(V^T M V)^{-1} V^T U}_{c_g}$

Finally: $SSE - SSE_{reg} = \|\hat{y} - y\|^2 =$
 $= \|MV(V^T M V)^{-1} V^T U\|^2 = \|MV c_g\|^2.$