

---

---

---

---

---



# Transformations on X's or Y to meet model assumptions

$$\mathbb{E}Y = X\beta = \begin{pmatrix} x_1^T \beta \\ \vdots \\ x_n^T \beta \end{pmatrix}$$

- $\mathbb{E}e = 0 \quad \text{cov } Y = \sigma^2 I$

- Consider  $\mathbb{E}Y = \alpha + \beta^T x$

Transform  $x$  by  $g(x)$

$$\mathbb{E}Y = \alpha + \beta g(x)$$

Consider simple transformations like

$$g_\lambda(x) = \begin{cases} \ln x, \lambda=0 \\ \frac{x^\lambda - 1}{\lambda}, \lambda \neq 0 \end{cases}$$

- $\lim_{\lambda \rightarrow 0} \frac{x^\lambda - 1}{\lambda} = \ln x.$

In practice, run the regression with

$$\lambda = \{-2, -\frac{1}{2}, \frac{1}{3}, 0, \frac{1}{3}, \frac{1}{2}, 1, 2, -1\}$$

Look at  $SSE_\lambda$  when we

- reg  $y$  on  $g_\lambda(x)$  for each  $\lambda$ . Pick the  $\lambda$  with smallest  $SSE$ .

- Do same procedure on the response ( $y > 0$ ) & ad-hoc.
- Alternatively, use variance stabilizing transformation

so that  $V(y)$  is a constant

$$\stackrel{\text{EY}}{g(Y)} = g(\mu) + (Y - \mu)g'(\mu) + \sigma_p^2(Y - \mu)$$

- $\text{Var}(g(Y)) \approx \text{Var}(Y)[g'(\mu)]^2 \equiv \sigma^2(g(\mu))^2$

$$\Rightarrow g'(\mu) = \sqrt{\frac{\sigma^2}{\text{Var}(Y)}}$$

$$g(\mu) = \int \frac{\sigma}{\sqrt{\text{Var}(Y)}} d\mu$$

if Poisson:  $\mathbb{E}Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$

- In ANOVA,  $k$  groups may have  $\sigma_i^2, i=1, 2, \dots, k$ .

$$H_0: \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2 \quad \text{vs.}$$

$$H: \sigma_i^2 \neq \sigma_j^2, i \neq j$$

Let  $\hat{\sigma}_i^2$  be the unbiased estimate of  $\sigma_i^2$ ,  $i=1, 2, \dots, k$ . The ratio

Hartley's test  $\frac{\max_{1 \leq i \leq k} \hat{\sigma}_i^2}{\min_{1 \leq i \leq k} \hat{\sigma}_i^2}$

is close to 1 if  $H_0$  holds, otherwise reject.

Levene's test

$$Y_{11}, Y_{12}, \dots, Y_{1n_1}$$

$$Y_{21}, Y_{22}, \dots, Y_{2n_2}$$

$$\vdots \quad \vdots \quad \vdots$$

$$Y_{k1}, Y_{k2}, \dots, Y_{kn_k}$$

$$Z_{ij} = |Y_{ij} - \bar{Y}_{i\cdot}|$$

$$Z_{..} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} Z_{ij}$$

$$Z_{i\cdot} = \frac{1}{n_i} \sum_{j=1}^{n_i} Z_{ij}$$

use  $\Sigma_i n_i (Z_{i\cdot} - Z_{..})^2 / (k-1)$

$$W = \frac{\Sigma_i \Sigma_j (Z_{ij} - Z_{i\cdot})^2 / (n-k)}{\Sigma_i \Sigma_j (Z_{ij} - Z_{i\cdot})^2 / (n-k)}$$

If  $\bar{Y}_{i\cdot}$  is replaced by  $\tilde{Y}_{i\cdot}$ : median in  $i^{\text{th}}$  group, then it is

Brown-Forsythe test

to test  $H_0$   
if  $W > F_{k-1, nk}$

# Collinearity

Reg  $Y$  on  $X_1 \dots X_p$

$$\hat{Y} = X\hat{\beta} = (X_1 \dots X_p) \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix}$$

$$\text{Var}(\hat{\beta}_k) = \frac{6^2}{\|X_k^T X\|_2^2} = 6^2 \|X_k - P(X_k | V_{(k)})\|_2^{-2}$$

$V_{(k)} = \{x_1, \dots, x_p \text{ but without } x_k\}$

If  $x_k \approx \text{linear comb. of the other } x_i$ 's.  
then  $\hat{\beta}_k$  will be estimated unreliable.

How to detect collinearity?

$R_j^2 \equiv R^2$  obtained when reg.  $x_j$  on  $X_{-j}$ .

$VIF_j \equiv \frac{1}{1-R_j^2}$  variance inflation factor

If  $VIF_j$  is large, say  $> 10$ , then  $x_j$  is almost linearly related with the rest of  $x$ 's.

What to do?

① Omit the variable or variables with  $VIF > 10$ .

②  $E Y = X\beta$

$$= X R^T \beta, R \text{ non-singular}$$

where  $R$  is from Householder matrix:

the QR decomp. of  $X$ :

$$X = QR \quad Hx = \begin{pmatrix} \|x\|_2 \\ 0 \end{pmatrix}$$

$$(x_1 \ x_2 \ \dots \ x_p) = (\underbrace{w_1 \ w_2 \ \dots \ w_p}_{Q}) (\underbrace{\sigma_1 \ \sigma_2 \ \dots \ \sigma_p}_{R})$$

$$w_i^T w_j = 0 \quad \& \quad w_i^T w_i = 1.$$

$$w_1 = \frac{x_1}{\|x_1\|_2}, \quad c_2 w_1 + c_2 w_2 = x_2.$$

$$\gamma \equiv R\beta.$$

$$\hat{Y} = \widehat{R\beta} = (Q^T Q)^{-1} Q^T Y = Q^T Y = R^{-1} X^T Y$$

$$SSR_{\text{reg}} = Y^T P_{\text{col}} Y = Y Q Q^T Y$$

$$= \sum_{i=1}^p (g_i^T Y)^2 \quad Q = [g_1 \ \dots \ g_p]$$

$$HW: EY = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

I take obs at  $x=1, 2, \dots, 5$ .

$$X = \begin{pmatrix} 1 & 1 & 1 \\ \frac{1}{2} & \frac{4}{3} & \frac{9}{4} \\ \frac{1}{4} & \frac{16}{9} & \frac{27}{16} \\ \frac{1}{8} & \frac{64}{27} & \frac{125}{64} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

$$EY = X\beta = (X R^{-1})(R\beta) = (w_1 \ \dots \ w_4) \gamma$$

$$\Rightarrow O = \sum_{x=1}^5 w_i(x) w_j(x)$$

$$H_0: \beta_p = 0 \quad \text{iff} \quad H_0: \gamma_p = 0.$$

# Heteroscedasticity

(Cook's paper)

$$Y_i = X_i^T \beta + \epsilon_i, \quad E\epsilon_i = 0, \quad \text{Var}\epsilon_i = \sigma_i^2$$

$$EY = X\beta, \quad P(X) = P$$

Assume  $\sigma_i^2$  depends on covariate

$$\sigma_i^2 = f(\text{a } q\text{-vector covariate}) \quad \begin{matrix} \downarrow \\ \text{can be one} \end{matrix} \quad \text{of the } x\text{'s.}$$

$$= \exp\{z_i^T \lambda\} \quad \&$$

1st column of  $X$  is 1. Then  $z = (1, z_1, \dots, z_k)$ .

If  $\lambda = \lambda_0 = (1, 0, 0, \dots, 0)$ . Then testing

$$H_0: \lambda = \lambda_0, \quad \sigma_i^2 = e^{\lambda_0} = \text{const}$$

$$L(Y_i | \beta, \sigma_i^2) \rightarrow L = \prod_{i=1}^n L_i \rightarrow \left( \frac{\partial \log L}{\partial \beta} \right) \left( \frac{\partial \log L}{\partial \sigma_i^2} \right)^T$$