

## 9.6 Transformation of response

50

- possible remedy measure, especially in situations when heteroscedasticity and/or non-normality detected

51

heteroscedasticity / non-normality

- serious problem especially when prediction is the main aim of modelling

often  $Y|X \sim N_n(X\beta, \sigma^2 I_n)$

but for  $Y^* = (t(Y_1), \dots, t(Y_n))^T$  model  $Y^*|X \sim N_n(X\beta, \sigma^2 I_n)$  is OK

cannot be assumed

$t: \mathbb{R} \rightarrow \mathbb{R}$ , chosen (non-linear) transformation.

→ normal linear model for transformed response

WARNING: model for transformed response is not always suitable, especially when effect of covariates on response expectation is of interest

$$E(t(Y)|X=x) \neq t(E(Y)|X=x)$$

(but see an exception of log-normal Lot)

## 9.6.1 Prediction based on a model with transformed response

52

If prediction is the main aim of modelling then no problems at all.

AIM: predict  $Y_{new}$ , given  $X_{new} = x_{new}$

ASSUME:  $t$  is strictly increasing  $t(Y_i)$

MODEL:  $M: Y^* | X \sim N_n(X\beta, \sigma^2 I)$   $Y_i^* = X_i^T \beta + \epsilon_i$   
 $Y^* = (t(Y_1), \dots, t(Y_n))^T$   $\epsilon_i | X \sim N(0, \sigma^2)$

$$Y_{new}^* = t(Y_{new})$$

1.  $\hat{Y}_{new}^*$  and  $(\hat{Y}_{new}^{*L}, \hat{Y}_{new}^{*U})$ :

point and interval prediction for  $Y_{new}^*$  based on model  $M$ , i.e.

$$1-\alpha = P\left((\hat{Y}_{new}^{*L}, \hat{Y}_{new}^{*U}) \ni Y_{new}^*; \text{param. values}\right)$$

$$= P\left((\hat{Y}_{new}^{*L}, \hat{Y}_{new}^{*U}) \ni t(Y_{new}); \text{param. values}\right)$$

$$2. = P\left(\underbrace{t^{-1}(\hat{Y}_{new}^{*L})}_{\hat{Y}_{new}^L}, \underbrace{t^{-1}(\hat{Y}_{new}^{*U})}_{\hat{Y}_{new}^U}\right) \ni Y_{new}; \text{param. values}$$

3.  $\hat{Y}_{new} = t^{-1}(\hat{Y}_{new}^*)$ : point prediction.  
(not necessarily in the middle of the prediction interval but who cares)

## 9.6.2 Log-normal model

→ reasonable interpretation of regression coefficients can still be derived

Let the following (log-normal linear) model hold:

$$\log(Y_i) = X_i^T \beta + \epsilon_i, \quad i=1, \dots, n$$

$$\epsilon_i | X \stackrel{\text{indep.}}{\sim} N(0, \sigma^2)$$

≡ normal linear model for log-transformed response

⇒ multiplicative model for the original response

$$Y_i = \exp(X_i^T \beta) \eta_i, \quad i=1, \dots, n$$

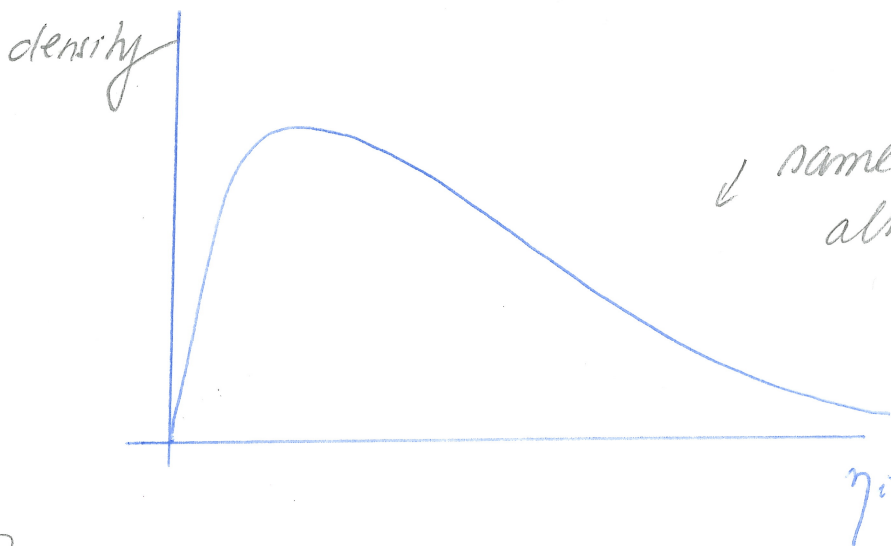
$$\eta_i = \exp(\epsilon_i), \quad \eta_i | X \stackrel{\text{indep.}}{\sim} \underline{\underline{LN}}(0, \sigma^2)$$

log-normal distrib.

Moments of the log-normal distributions

$$M_i = E(\eta_i) = E(\eta_i | X) = \exp\left(\frac{\sigma^2}{2}\right) > 1 \quad (\text{if } \sigma^2 > 0)$$

$$V_i = \text{var}(\eta_i) = \text{var}(\eta_i | X) = (\exp(\sigma^2) - 1) \cdot \exp(\sigma^2)$$



↓ same shape has also assumed distrib. of  $Y_i$

MODEL:  $Y_i = \exp(X_i^T \beta) \cdot \eta_i$



$$\begin{aligned} E(Y|X=x) &= \exp(x^T \beta) \cdot E(\eta) = \\ &= \exp(x^T \beta) \cdot \underbrace{\exp\left(\frac{\sigma^2}{2}\right)}_{M > 1} \end{aligned}$$

$$\begin{aligned} \text{var}(Y|X=x) &= \exp(2x^T \beta) \cdot V = \\ &= V \cdot \left( \frac{E(Y|X=x)}{M} \right)^2 \end{aligned}$$

var(Y|X=x) increases with |E(Y|X=x)|

log-normal linear model captures situations when

1. response (conditional) distribution is skewed
2. response (conditional) variance increases with the (absolute value) of expectation

Interpretation of regression coefficients  
in case of log-normal linear model

MODEL:  $\log(Y) = X^T \beta + \epsilon$ ,  $\epsilon \sim N(0, \sigma^2)$

$$E(Y|X=x) = \exp(x^T \beta) \cdot M, \quad M = \exp\left(\frac{\sigma^2}{2}\right)$$

$$X = (x_0, \dots, x_j, \dots, x_{k-1})^T \in \mathcal{X}$$

$$X^{(j+1)} = (x_0, \dots, x_{j+1}, \dots, x_{k-1})^T \in \mathcal{X}$$

$$\beta = (\beta_0, \dots, \beta_j, \dots, \beta_{k-1})^T$$

$$\Rightarrow \frac{E(Y|X=X^{(j+1)})}{E(Y|X=X)} = \frac{M \exp(X^{(j+1)T} \beta)}{M \exp(X^T \beta)} = e^{\beta_j}$$

→ in output, report  $e^{\hat{\beta}_j}$  rather than  $\hat{\beta}_j$

→ analogously, if  $(\hat{\beta}_j^L, \hat{\beta}_j^U)$  is confidence interval for  $\beta_j$  then  $(e^{\hat{\beta}_j^L}, e^{\hat{\beta}_j^U})$  is confidence interval for  $e^{\beta_j}$  (and should be reported)

Special case, one-way classification

& log-normal linear model

- one categorical covariate parameterized by chosen (pseudo) contrasts  $C = \begin{pmatrix} C_1^T \\ \vdots \\ C_g^T \end{pmatrix}$

MODEL:

$$E(\log(Y) | Z=z) = \beta_0 + C_g^T \beta^z$$

$$\begin{aligned} \Rightarrow \frac{E(Y | Z=g)}{E(Y | Z=h)} &= \frac{M \cdot \exp(\beta_0 + C_g^T \beta^z)}{M \cdot \exp(\beta_0 + C_h^T \beta^z)} = \\ &= \exp\{(C_g^T - C_h^T) \beta^z\} \end{aligned}$$

For example: reference group (pseudo) contrasts (dummy variables)

$$C = \begin{pmatrix} 0 & \dots & 0 \\ 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{pmatrix}$$

$$E(\log(Y) | Z=g) = \beta_0 + \beta_g^z, \quad \beta_1^z = 0$$

$$\Rightarrow \frac{E(Y | Z=g)}{E(Y | Z=1)_{\text{reference}}} = \exp(\beta_g^z), \quad g=2, \dots, G$$

$$\frac{E(Y | Z=g)}{E(Y | Z=h)} = \exp(\beta_g^z - \beta_h^z), \quad g, h = 2, \dots, G$$