

Copyright Notice

These notes are copyright @**Liqun Diao** at the University of Waterloo, Ontario, Canada.

Under no circumstances shall these notes be published, distributed, or made accessible on the Internet or any public platform without the prior explicit written consent of Professor Diao. Professor Diao retains the sole discretion to define the form, scope, and manner of consent required for such publication or distribution.

CONFIDENCE INTERVAL FOR $\mu(x) = \alpha + \beta x$

$$\hat{\mu}(x) = \hat{\alpha} + \hat{\beta}x \quad (\text{The invariance property of MLE})$$

$$= \bar{y} - \hat{\beta} \bar{x} + \hat{\beta}x$$

$$= \frac{1}{n} \sum_{i=1}^n y_i + \hat{\beta} (x - \bar{x})$$

$$= \frac{1}{n} \sum_{i=1}^n y_i + \frac{\sum_{i=1}^n (x_i - \bar{x})}{S_{xx}} y_i (x - \bar{x})$$

$$= \sum_{i=1}^n \left\{ \frac{1}{n} + \frac{(x_i - \bar{x})(x - \bar{x})}{S_{xx}} \right\} y_i$$

CONFIDENCE INTERVAL FOR $\mu(x) = \alpha + \beta x$

$$\Rightarrow \tilde{N}(x) = \sum_{i=1}^n b_i Y_i$$

$\tilde{N}(x)$ is a combination of independent Gaussian random variables. So, $\tilde{N}(x)$ is also Gaussian distributed.

CONFIDENCE INTERVAL FOR $\mu(x) = \alpha + \beta x$

$$\mathbb{E}(\hat{\mu}(x)) = \mathbb{E}\left(\sum_{i=1}^n b_i Y_i\right)$$

$$= \sum_{i=1}^n b_i \mathbb{E}(Y_i)$$

$$= \sum_{i=1}^n b_i (\alpha + \beta x_i)$$

$$= \alpha \sum_{i=1}^n b_i + \beta \sum_{i=1}^n b_i x_i$$

$$= \alpha + \beta x \quad (\text{unbiased})$$

Recall $Y_i = \alpha + \beta x_i + \epsilon_i$,

$\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma)$

Exercise:

Show that

$$(1) \sum_{i=1}^n b_i = 1$$

$$(2) \sum_{i=1}^n b_i x_i = x$$

CONFIDENCE INTERVAL FOR $\mu(x) = \alpha + \beta x$

$$\text{Var}(\tilde{\mu}(x)) = \text{Var}\left[\sum_{i=1}^n b_i Y_i\right]$$

$$= \sum_{i=1}^n b_i^2 \text{Var}(Y_i) + \sum_{i \neq j} b_i b_j \text{Cov}(Y_i, Y_j) \rightarrow 0$$

$$= \sum_{i=1}^n b_i^2 \sigma^2$$

$$= \sigma^2 \sum_{i=1}^n b_i^2$$

$$= \sigma^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}} \right]$$

Exercise:

Show that

$$\sum_{i=1}^n b_i^2 = \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}$$

CONFIDENCE INTERVAL FOR $\mu(x) = \alpha + \beta x$

In Summary:

$$\tilde{N}(x) \sim G\left(\overbrace{\alpha + \beta x}^{N(x)}, \sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}\right).$$

(Exact Distribution)

Pivotal Quantity:

$$\frac{\tilde{N}(x) - N(x)}{se \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \sim t(n-2).$$

$$\frac{\tilde{N}(x) - \mu(x)}{\sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \sim N(0, 1)$$

$$Se^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \tilde{N}_i)^2$$

$$\frac{(n-2)Se^2}{\sigma^2} \sim \chi^2(n-2)$$

Then,

Using theorem from Ch. 4

$$T = \frac{\tilde{N}(x) - \mu(x)}{\sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \bigg/ \sqrt{\frac{Se^2}{\sigma^2}} \sim t(n-2).$$

CONFIDENCE INTERVAL FOR $\mu(x) = \alpha + \beta x$

A $100p\%$ CI for $\mu(x)$ is

$$\hat{\mu}(x) \pm \alpha \cdot \text{Se} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$

where $P(T \leq \alpha) = \frac{1+p}{2}$, $T \sim t(n-2)$.

PREDICTION INTERVAL FOR FUTURE RESPONSE

To predict a future response Y that is independent of our sample (X_i, Y_i) , $i = 1, \dots, n$.

Model: $Y = \alpha + \beta X + R$, $R \stackrel{\text{ind}}{\sim} N(0, \sigma^2)$.

$$Y \perp Y_i, \quad i = 1, \dots, n.$$

We predict Y using $\hat{N}(x)$.

PREDICTION INTERVAL FOR FUTURE RESPONSE

We consider $Y - \tilde{N}(x)$

$$\mathbb{E}(Y - \tilde{N}(x)) = \mathbb{E}(Y) - \mathbb{E}(\tilde{N}(x))$$

$$= \alpha + \beta x - (\alpha + \beta x)$$

$$= 0$$

$$\begin{aligned} \text{Var}(Y - \tilde{N}(x)) &= \text{Var}(Y) + \text{Var}(\tilde{N}(x)) - 2\text{Cov}(Y, \tilde{N}(x)) \\ &= \sigma^2 + \sigma^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}} \right]. \end{aligned}$$

PREDICTION INTERVAL FOR FUTURE RESPONSE

Note $Y - \hat{\mu}(x)$ is Gaussian distributed.

$$\frac{Y - \hat{\mu}(x)}{\sigma \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \sim N(0, 1)$$

$$\frac{Y - \hat{\mu}(x)}{se \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \sim t(n-2).$$

PREDICTION INTERVAL FOR FUTURE RESPONSE

A 100% Prediction Interval for Y

$$\hat{\mu}(x) \pm a \cdot se \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}},$$

where $P(T \leq a) = \frac{1+P}{2}$, $T \sim t(n-2)$.

6.3 MODEL CHECKING - ASSUMPTIONS TO BE CHECKED

Homoscedasticity

There are two main assumptions for Gaussian linear response models:

(1) Y_i (given covariates x_i) is Gaussian with standard deviation σ which does not depend on the covariates. $Y_i = \alpha + \beta x_i + R_i, R_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma)$.

(2) $E(Y_i) = \mu(x_i)$ is a linear combination of observed covariates with unknown coefficients. $E(Y_i) = \mu(x_i) = \alpha + \beta x_i, i = 1, \dots, n$.

MODEL ASSUMPTIONS SHOULD ALWAYS BE CHECKED!!!

We will examine three graphical methods to check these assumptions.

METHOD I - SCATTERPLOT OF DATA AND FITTED REGRESSION LINE

In simple linear regression, a scatterplot of the data with the fitted line $y = \hat{\alpha} + \hat{\beta}x$ superimposed shows how well the model fits. If there are any obvious departures from the fitted line then these departures might suggest a model which would fit the data better.

