

# Methods of Data Analysis I

Derek Li

## Contents

<b>1</b>	<b>Review</b>	<b>3</b>
1.1	Expectation . . . . .	3
1.2	Variance and Covariance . . . . .	3
1.3	Correlation . . . . .	3
1.4	Distributions . . . . .	3
1.4.1	Bivariate Normal Distribution . . . . .	4
<b>2</b>	<b>Sample Linear Regression</b>	<b>5</b>
2.1	Statistical Model . . . . .	5
2.2	Estimating $\beta_0, \beta_1$ . . . . .	5
2.2.1	Least Squares Method . . . . .	5
2.2.2	Interpretation . . . . .	6
2.2.3	Estimation in R . . . . .	6
2.3	Properties of Fitted Regression Line . . . . .	6
2.4	Assumptions . . . . .	7
2.5	Estimating the Variance of the Random Error Term . . . . .	7
2.6	Properties of Least Squares Estimators . . . . .	8
2.7	Normal Error Regression Model . . . . .	9
2.8	Inference for the Parameter . . . . .	10
2.8.1	Significance Test . . . . .	10
2.8.2	Confidence Interval . . . . .	10
2.9	The Pooled Two-Sample $t$ -Procedure . . . . .	10
2.10	Regression Analysis of Variance . . . . .	11
2.10.1	Regression ANOVA Table . . . . .	12
2.10.2	Coefficient of Determination . . . . .	12
2.10.3	Sample Correlation Coefficient . . . . .	13
2.11	Confidence Interval for the Population Regression Line . . . . .	14
2.12	Prediction Interval for Actual Value of $Y$ . . . . .	14
<b>3</b>	<b>Diagnostics and Transformations for Simple Linear Regression</b>	<b>16</b>
3.1	Phenomena and Fallacies in Regression . . . . .	16
3.2	Validity of SLR Model: Model Linearity . . . . .	16
3.2.1	Residuals . . . . .	16
3.2.2	Diagnostics . . . . .	18
3.3	Validity of SLR Model: Uncorrelated Errors . . . . .	18
3.4	Validity of SLR Model: Homoscedasticity . . . . .	18

3.4.1	Standardized Residuals . . . . .	18
3.4.2	Diagnostics . . . . .	18
3.5	Normality . . . . .	18
3.5.1	Q-Q Plots . . . . .	18
3.5.2	The Shapiro-Wilk Test . . . . .	19
3.5.3	Diagnostics . . . . .	19
3.6	Outliers, Leverage Points, and Influential Points . . . . .	20
3.6.1	Quantifying Leverage . . . . .	20
3.6.2	Cook's Distance . . . . .	21
3.6.3	Recommendations . . . . .	21
3.7	Diagnostic Plots from R . . . . .	21
3.8	Transformations . . . . .	22
3.8.1	Common Transformations . . . . .	22
3.8.2	Variance Stabilizing Transformations . . . . .	22

# 1 Review

## 1.1 Expectation

- $\mathbb{E}[a] = a, a \in \mathbb{R}$ .
- $\mathbb{E}[aY] = a\mathbb{E}[Y]$ .
- $\mathbb{E}[X \pm Y] = \mathbb{E}[X] \pm \mathbb{E}[Y]$ .
- $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$  if  $X$  and  $Y$  are independent.
- Tower rule:  $\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$ .

## 1.2 Variance and Covariance

- $\text{Var}[a] = 0, a \in \mathbb{R}$ .
- $\text{Var}[aY] = a^2\text{Var}[Y]$ .
- $\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$ .
- $\text{Cov}(Y, Y) = \text{Var}[Y]$ .
- $\text{Var}[Y] = \text{Var}[\mathbb{E}[Y|X]] + \mathbb{E}[\text{Var}[Y|X]]$ .
- $\text{Var}[X \pm Y] = \text{Var}[X] + \text{Var}[Y] \pm 2\text{Cov}(X, Y)$ .
- $\text{Cov}(X, Y) = 0$  if  $X$  and  $Y$  are independent.
- $\text{Cov}(aX + bY, cU + dW) = ac\text{Cov}(X, U) + ad\text{Cov}(X, W) + bc\text{Cov}(Y, U) + bd\text{Cov}(Y, W)$ .

## 1.3 Correlation

If  $X$  and  $Y$  are random variables, a symmetric measure of the direction and strength of the linear dependence between them is their correlation

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}[X]\text{Var}[Y]}}.$$

## 1.4 Distributions

- If  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then  $Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1)$ .
- Let  $U = Z^2$ , then  $U \sim \chi_{(1)}^2$ .
- If  $Z$  and  $X \sim \chi_{(m)}^2$  are independent, then  $\frac{Z}{\sqrt{X/m}} \sim t_{(m)}$ .
- If  $X \sim \chi_{(m)}^2, Y \sim \chi_{(n)}^2$  are independent, then  $\frac{X/m}{Y/n} \sim F_{(m,n)}$ .
- $t_{(m)} \xrightarrow{D} Z$ , as  $m \rightarrow \infty$ .

### 1.4.1 Bivariate Normal Distribution

$X$  and  $Y$  are jointly normally distributed is their joint density function is

$$f(x, y) = \frac{e^{-\frac{Q}{2}}}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}},$$

where

$$Q = \frac{1}{1-\rho^2} \left[ \frac{(x-\mu_x)^2}{\sigma_x^2} - 2\rho \frac{(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} + \frac{(y-\mu_y)^2}{\sigma_y^2} \right].$$

Two marginal distributions are

$$X \sim \mathcal{N}(\mu_x, \sigma_x^2) \text{ and } Y \sim \mathcal{N}(\mu_y, \sigma_y^2).$$

The conditional distribution of  $Y$  given  $X = x$  is

$$Y|X = x \sim \mathcal{N}\left(\mu_y + \rho\sigma_y\left(\frac{x-\mu_x}{\sigma_x}\right), (1-\rho^2)\sigma_y^2\right).$$

**Theorem 1.1.** If  $X$  and  $Y$  are jointly normally distributed, then a zero covariance between  $X$  and  $Y$  implies that they are statistically independent.

## 2 Sample Linear Regression

### 2.1 Statistical Model

$$Y = \beta_0 + \beta_1 X + e,$$

where  $Y$  is dependent or response variable,  $X$  is independent or explanatory variable,  $\beta_0$  is intercept parameter,  $\beta_1$  is slope parameter, and  $e$  is random error or noise (variation in measures that we cannot account for).

Given a specific value of  $X = x$ , we want to find the expected value of  $Y$

$$\mathbb{E}[Y|X = x].$$

### 2.2 Estimating $\beta_0, \beta_1$

Given  $n$  pairs bivariate data  $(x_1, y_1), \dots, (x_n, y_n)$ , we want to use  $\hat{\beta}_0$  and  $\hat{\beta}_1$  to estimate  $\beta_0$  and  $\beta_1$ .

Consider the residual sum of squares

$$RSS = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \left[ y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right]^2,$$

we can use least squares method that minimizes the criterion RSS to find the estimators.

#### 2.2.1 Least Squares Method

Least squares method makes no statistical assumptions. We have

$$\frac{\partial RSS}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad \text{and} \quad \frac{\partial RSS}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i.$$

Let  $\frac{\partial RSS}{\partial \hat{\beta}_0}$  and  $\frac{\partial RSS}{\partial \hat{\beta}_1}$  be 0, we get the normal equations

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad \text{and} \quad \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0.$$

Therefore, we have

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_0 - \sum_{i=1}^n \hat{\beta}_1 x_i = n\bar{y} - n\hat{\beta}_0 - n\hat{\beta}_1 \bar{x} = 0 \Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Besides,

$$\begin{aligned} \sum_{i=1}^n x_i y_i - \sum_{i=1}^n \hat{\beta}_0 x_i - \sum_{i=1}^n \hat{\beta}_1 x_i^2 &= \sum_{i=1}^n x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 \\ &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} + n\hat{\beta}_1 \bar{x}^2 - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0, \end{aligned}$$

i.e.,

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} := \frac{S_{XY}}{S_{XX}}.$$

### 2.2.2 Interpretation

$\hat{\beta}_0$  : The expected value of  $y$  when  $x = 0$ . No practical interpretation unless 0 is within the range of the predictor values.

$\hat{\beta}_1$  : When  $x$  changes by 1 unit, the corresponding average change in  $y$  is the slope.

### 2.2.3 Estimation in R

---

```
model = lm(y ~ x)
summary(model)
```

---

## 2.3 Properties of Fitted Regression Line

**Property 2.1.**

$$\sum_{i=1}^n \hat{e}_i = 0.$$

*Proof.* By definition,

$$\begin{aligned} \sum_{i=1}^n \hat{e}_i &= \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \sum_{i=1}^n (y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i) \\ &= n\bar{y} - n\bar{y} + n\hat{\beta}_1 \bar{x} - n\hat{\beta}_1 \bar{x} = 0. \end{aligned}$$

□

**Property 2.2.** The sum of squares of residuals is not 0 unless the fit to the data is perfect.

**Property 2.3.**

$$\sum_{i=1}^n \hat{e}_i x_i = 0.$$

*Proof.* By definition,

$$\begin{aligned} \sum_{i=1}^n \hat{e}_i x_i &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = \sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i + \hat{\beta}_1 \bar{x} \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 \\ &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} - \hat{\beta}_1 \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = 0. \end{aligned}$$

□

**Property 2.4.**

$$\sum_{i=1}^n \hat{e}_i \hat{y}_i = 0.$$

*Proof.* By definition,

$$\sum_{i=1}^n \hat{e}_i \hat{y}_i = \sum_{i=1}^n \hat{e}_i (\hat{\beta}_0 + \hat{\beta}_1 x_i) = \hat{\beta}_0 \sum_{i=1}^n \hat{e}_i + \hat{\beta}_1 \sum_{i=1}^n \hat{e}_i x_i = 0 + 0 = 0.$$

□

**Property 2.5.**

$$\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n y_i.$$

*Proof.* We have

$$\sum_{i=1}^n \hat{e}_i = 0 = \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{y}_i \Rightarrow \sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n y_i.$$

□

## 2.4 Assumptions

The Gauss-Markov conditions are:

1.  $\mathbb{E}[e_i] = 0$ .
2.  $\text{Var}[e_i] = \sigma^2$ , i.e., homoscedastic.
3. The errors are uncorrelated or  $\text{Cov}(e_i, e_j) = \rho(e_i, e_j) = 0$ .

**Theorem 2.1** (Gauss-Markov Theorem). Under the conditions of the simple linear regression model, the least-squares parameter estimators are best linear unbiased estimators.

We assume that  $Y$  is related to  $x$  by the simple linear regression model

$$Y_i = \beta_0 + \beta_1 x_i + e_i, i = 1, \dots, n.$$

Under the conditions we have

$$\mathbb{E}[Y|X = x_i] = \beta_0 + \beta_1 x_i$$

and

$$\text{Var}[Y|X = x_i] = \text{Var}[\beta_0 + \beta_1 x_i + e_i|X = x_i] = \text{Var}[e_i] = \sigma^2.$$

## 2.5 Estimating the Variance of the Random Error Term

The variance  $\sigma^2$  is another parameter of the SLR model and we want to estimate  $\sigma^2$  to measure the variability of our estimates of  $Y$ , and carry out inference on the model.

An unbiased estimate of  $\sigma^2$  is

$$S^2 = \frac{\sum_{i=1}^n \hat{e}_i^2}{n-2} = \frac{RSS}{n-2}.$$

## 2.6 Properties of Least Squares Estimators

Since  $\sum_{i=1}^n (x_i - \bar{x}) = 0$ ,

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i - \bar{y} \sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x})y_i.$$

Let  $c_i = \frac{x_i - \bar{x}}{SXX}$ , we can rewrite  $\hat{\beta}_1$  as

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i,$$

which is a linear combination of  $y_i$ .

We have

$$\begin{aligned} \mathbb{E}[\hat{\beta}_1|X] &= \mathbb{E}\left[\sum_{i=1}^n c_i y_i | X = x_i\right] = \sum_{i=1}^n c_i \mathbb{E}[y_i | X = x_i] \\ &= \sum_{i=1}^n c_i \mathbb{E}[\beta_0 + \beta_1 x_i] = \beta_0 \sum_{i=1}^n c_i + \beta_1 \sum_{i=1}^n c_i x_i \\ &= \frac{\beta_0}{SXX} \sum_{i=1}^n (x_i - \bar{x}) + \beta_1 \sum_{i=1}^n \frac{(x_i - \bar{x})x_i}{SXX} \\ &= \beta_1 \frac{\sum_{i=1}^n x_i^2 - n\bar{x}^2}{SXX} = \beta_1. \end{aligned}$$

Therefore,  $\hat{\beta}_1$  is unbiased for  $\beta_1$ . Besides,

$$\begin{aligned} \text{Var}[\hat{\beta}_1|X] &= \text{Var}\left[\sum_{i=1}^n c_i y_i | X\right] = \sum_{i=1}^n c_i^2 \text{Var}[y_i | X = x_i] \\ &= \sigma^2 \sum_{i=1}^n c_i^2 = \sigma^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{SXX^2} = \frac{\sigma^2}{SXX}. \end{aligned}$$

We have

$$\begin{aligned} \mathbb{E}[\hat{\beta}_0|X] &= \mathbb{E}[\bar{y} - \hat{\beta}_1 \bar{x} | X = x_i] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n y_i - \hat{\beta}_1 \bar{x} | X = x_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\beta_0 + \beta_1 x_i + e_i | X = x_i] - \bar{x} \mathbb{E}[\hat{\beta}_1 | X = x_i] \\ &= \frac{1}{n} n \beta_0 + \frac{1}{n} n \beta_1 \bar{x} - \bar{x} \beta_1 = \beta_0. \end{aligned}$$

Therefore,  $\hat{\beta}_0$  is unbiased for  $\beta_0$ . Besides,

$$\begin{aligned} \text{Var}[\hat{\beta}_0|X] &= \text{Var}[\bar{y} - \hat{\beta}_1 \bar{x} | X = x_i] \\ &= \text{Var}[\bar{y} | X = x_i] + \text{Var}[\hat{\beta}_1 \bar{x} | X = x_i] - 2\text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x} | X = x_i) \\ &= \frac{\sigma^2}{n} + \frac{\bar{x}^2 \sigma^2}{SXX} - 0 = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right). \end{aligned}$$



Note that  $\text{Cov}\left(\bar{y}, \hat{\beta}_1 \bar{x} | X = x_u\right) = \frac{\bar{x}\sigma^2}{n} \sum_{i=1}^n c_i = 0$ .

## 2.7 Normal Error Regression Model

Given distributional assumption:

$$e_i \sim \mathcal{N}(0, \sigma^2),$$

we know:

- (1) the errors are independent since  $\rho = 0$ ;
- (2) since  $y_i = \beta_0 + \beta_1 x_i + e_i$ , then  $Y_i | X \sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$ ;
- (3) the least squares estimates of  $\beta_0, \beta_1$  are equivalent to their maximum likelihood estimators.
- (4) since  $\hat{\beta}_1 = \sum_{i=1}^n c_i y_i$  is a linear combination of the  $y_i$ 's,  $\hat{\beta}_1 | X \sim \mathcal{N}\left(\beta_1, \frac{\sigma^2}{S_{XX}}\right)$ ; since  $\bar{y}$  is normally distributed,  $\hat{\beta}_0 | X \sim \mathcal{N}\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}}\right)\right)$ .

**Property 2.6.** Under the normal error SLR model, where

$$e_i \sim \mathcal{N}(0, \sigma^2) \text{ and } S^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{e}_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2,$$

we have

$$\frac{(n-2)S^2}{\sigma^2} = \sum_{i=1}^n \left( \frac{Y_i - \hat{Y}_i}{\sigma} \right)^2 \sim \chi_{(n-2)}^2.$$

**Property 2.7.** Under the normal error SLR model,

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{S^2}{S_{XX}}}} \sim t_{(n-2)}.$$

*Proof.* We have  $\hat{\beta}_1 | X = x_i \sim \mathcal{N}\left(\beta_1, \frac{\sigma^2}{S_{XX}}\right)$ , and thus

$$\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{S_{XX}}} \sim \mathcal{N}(0, 1).$$

Wherefore

$$\frac{\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{S_{XX}}}}{\sqrt{(n-2)S^2/\sigma^2/(n-2)}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{S^2}{S_{XX}}}} \sim t_{(n-2)}.$$

□

## 2.8 Inference for the Parameter

### 2.8.1 Significance Test

- Step 1:  $H_0 : \beta_1 = \beta_1^0$  against  $H_a : \beta_1 \neq \beta_1^0$ .
- Step 2: Test statistic  $t = \frac{\hat{\beta}_1 - \beta_1^0}{\sqrt{S^2/SXX}}$ , and under  $H_0, t \sim t_{(n-2)}$ .
- Step 3:  $p\text{-value} = 2P(t_{(n-2)} \geq |t|)$ .
- Step 4: The smaller the  $p$ -value, the greater the evidence against  $H_0$  and the larger  $p$ -value indicate that the data is consistent with  $H_0$ .

$p$ -value	Evidence against $H_0$
$< 0.001$	Very strong
$(0.001, 0.01)$	Strong
$(0.01, 0.05)$	Moderate
$(0.05, 0.1)$	Weak
$> 0.1$	None

Note that the test statistic for  $\hat{\beta}_0$  is  $t = \frac{\hat{\beta}_0 - \beta_0^0}{\sqrt{S^2\left(\frac{1}{n} + \frac{\bar{x}^2}{SXX}\right)}}$

### 2.8.2 Confidence Interval

The CI is

$$\text{Estimate} \pm 100 \left(1 - \frac{\alpha}{2}\right) \text{th quantile} \times \text{Standard Error (Estimate)},$$

where  $\alpha$  is the critical value.

For  $\beta_1$ , the CI is

$$\left[ \hat{\beta}_1 \pm t_{\frac{\alpha}{2}(n-2)} \sqrt{\frac{S^2}{SXX}} \right].$$

For  $\beta_0$ , the CI is

$$\left[ \hat{\beta}_0 \pm t_{\frac{\alpha}{2}(n-2)} \sqrt{S^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)} \right].$$

Note that a  $100(1 - \alpha)\%$  CI for  $\theta$  consists of all those values of  $\theta_0$  for which  $H_0 : \theta = \theta_0$  will not be rejected at level  $\alpha$ . In other words, we do not reject  $H_0$  if  $\theta_0$  lies within the CI, and we reject  $H_0$  if the CI does not include  $\theta_0$ .

## 2.9 The Pooled Two-Sample $t$ -Procedure

We want to test  $H_0 : \mu_x = \mu_y$ , where

$$X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_x, \sigma_x^2) \text{ and } Y_1, \dots, Y_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_y, \sigma_y^2).$$

Suppose two samples are independent and  $\sigma_x^2 = \sigma_y^2 = \sigma^2$ , then we have

$$t = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{s_p \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} \sim t_{(n_x + n_y - 2)},$$

where  $s_p^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{n_x + n_y - 2}$ .

## 2.10 Regression Analysis of Variance

Notice that  $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$ . We have

$$\begin{aligned} TSS &= \sum_i^n (y_i - \bar{y})^2, \\ RSS &= \sum_i^n (y_i - \hat{y}_i)^2 = \sum_i^n \hat{e}_i^2, \\ RegSS &= \sum_i^n (\hat{y}_i - \bar{y})^2. \end{aligned}$$

$RSS$ , residual SS, is the least square criterion, representing the unexplained variation in  $y$ 's.  $RegSS$ , regression SS, is the amount of variation in  $y$ 's explained by regression line.

**Property 2.8.**  $RegSS = \hat{\beta}_1^2 SXX$ .

*Proof.* We have

$$\begin{aligned} RegSS &= \sum_i^n (\hat{y}_i - \bar{y})^2 = \sum_i^n (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2 \\ &= \sum_i^n (\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i - \bar{y})^2 = \hat{\beta}_1^2 \sum_i^n (x_i - \bar{x})^2 = \hat{\beta}_1^2 SXX. \end{aligned}$$

□

**Property 2.9.**  $TSS = RSS + RegSS$ .

*Proof.* We have

$$\begin{aligned} \sum_i^n (y_i - \bar{y})^2 &= \sum_i^n ((y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}))^2 \\ &= \sum_i^n (y_i - \hat{y}_i)^2 + \sum_i^n (\hat{y}_i - \bar{y})^2 + 2 \sum_i^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ &= RSS + RegSS + 2 \sum_i^n \hat{e}_i(\hat{y}_i - \bar{y}) \\ &= RSS + RegSS + 2 \sum_i^n \hat{e}_i \hat{y}_i - 2\bar{y} \sum_i^n \hat{e}_i \\ &= RSS + RegSS. \end{aligned}$$

□

### 2.10.1 Regression ANOVA Table

Source	SS	df	Mean SS
Regression Line	$RegSS = \hat{\beta}_1^2 SXX$	1	$\hat{\beta}_1^2 SXX$
Error	$RSS = \sum_{i=1}^n \hat{e}_i^2$	$n - 2$	$\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = S^2$
Total	$TSS = \sum_i^n (y_i - \bar{y})^2$		

**Property 2.10.** Let

$$F = \frac{MRegSS}{MRSS} = \frac{RegSS/1}{RSS/(n-2)}.$$

If  $\beta_1 = 0$ , then

$$F \sim F_{(1, n-2)}.$$

*Proof.* If  $\beta_1 = 0$ , then  $\hat{\beta}_1 \sim \mathcal{N}\left(0, \frac{\sigma^2}{SXX}\right)$ , i.e.,

$$\frac{\hat{\beta}_1}{\sqrt{\sigma^2/SXX}} \sim \mathcal{N}(0, 1) \Rightarrow \frac{\hat{\beta}_1^2}{\sigma^2/SXX} \sim \chi_{(1)}^2.$$

Besides,  $\frac{(n-2)S^2}{\sigma^2} \sim \chi_{(n-2)}^2$ , and we have

$$\frac{\frac{\hat{\beta}_1^2}{\sigma^2/SXX}}{\frac{(n-2)S^2}{\sigma^2}/(n-2)} = \frac{\hat{\beta}_1^2 SXX}{S^2} = F \sim F_{(1, n-2)}.$$

□

Note that  $F$  is another test of  $H_0 : \beta_1 = 0$ , and in R, we have:

---

```
anova(model)
```

---

### 2.10.2 Coefficient of Determination

Let

$$R^2 = \frac{RegSS}{TSS} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}.$$

Here are some comments about  $R^2$ :

- $R^2 \in [0, 1]$ .
- $R^2$  gives percentage of variation in  $y$ 's explained by regression line.
- $R^2$  is not resistant to outliers.
- A high  $R^2$  does not indicate that the estimated regression line is a good fit since:
  - \* we do not have absolute rules about how large it should be;
  - \*  $R^2$  can get very high by overfitting.

- It is not meaningful for models without intercept.
- To compare 2 models,  $R^2$  is only useful:
  - \* same observations,  $y$ 's in original units (not transformed);
  - \* one set of predictor variables is a subset of the other.

### 2.10.3 Sample Correlation Coefficient

The estimate of the population correlation is Pearson's Product-Moment Correlation Coefficient

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{SXY}{\sqrt{SXX \cdot SY Y}},$$

which is the MLE of  $\rho$ .  $r$  is distribution free and is always a number between -1 and 1.

**Theorem 2.2.**  $R^2 = r^2$ .

*Proof.* We have

$$R^2 = \frac{RegSS}{TSS} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\hat{\beta}_1^2 SXX}{SY Y} = \frac{\frac{SXY^2}{SXX^2} \cdot SXX}{SY Y} = \frac{SXY^2}{SXX \cdot SY Y} = r^2.$$

□

**Property 2.11.** If  $\rho = 0$ ,

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{\hat{\beta}_1}{\sqrt{S^2/SXX}} \sim t_{(n-2)},$$

where  $\hat{\beta}_1$  is the slope estimate for the normal error SLR model.

*Proof.* Since  $r^2 = R^2$ , then

$$\frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{\frac{\hat{\beta}_1 \sqrt{SXX}}{\sqrt{SXY}} \sqrt{n-2}}{\sqrt{(n-2)S^2/SXY}} = \frac{\hat{\beta}_1}{\sqrt{S^2/SXX}}.$$

If  $\rho = 0$ , then  $\beta_1 = 0$ , i.e.,

$$\frac{\hat{\beta}_1}{\sqrt{S^2/SXX}} \sim t_{(n-2)}.$$

□

## 2.11 Confidence Interval for the Population Regression Line

We want to find a CI for the unknown population regression line at a given value of  $X$ , denoted by  $x^*$ , i.e.,

$$\mathbb{E}[Y|X = x^*] = \beta_0 + \beta_1 x^*.$$

The point estimate for  $\mathbb{E}[Y|X = x^*]$  is

$$\hat{y}^* = \hat{\beta}_0 + \hat{\beta}_1 x^*.$$

We have

$$\mathbb{E}[\hat{y}^*] = \mathbb{E}[\hat{y}|X = x^*] = \beta_0 + \beta_1 x^*,$$

i.e.,  $\hat{y}^*$  is unbiased for  $\mathbb{E}[Y|X = x^*]$ .

Recall that  $\text{Var}[\hat{\beta}_0|X] = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)$ ,  $\text{Var}[\hat{\beta}_1|X] = \frac{\sigma^2}{SXX}$ , then

$$\text{Cov}[\hat{\beta}_0, \hat{\beta}_1|X] = \text{Cov}[\bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1|X] = -\bar{x} \text{Var}[\hat{\beta}_1|X] = -\frac{\bar{x} \sigma^2}{SXX}.$$

Wherefore

$$\begin{aligned} \text{Var}[\hat{y}^*] &= \text{Var}[\hat{y}|X = x^*] = \text{Var}[\hat{\beta}_0 + \hat{\beta}_1 x^*|X = x^*] \\ &= \text{Var}[\hat{\beta}_0|X = x^*] + (x^*)^2 \text{Var}[\hat{\beta}_1|X = x^*] + 2x^* \text{Cov}[\hat{\beta}_0, \hat{\beta}_1|X = x^*] \\ &= \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right) + (x^*)^2 \frac{\sigma^2}{SXX} - \frac{2x^* \bar{x} \sigma^2}{SXX} = \sigma^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right). \end{aligned}$$

Hence, as  $n \uparrow$ ,  $\text{Var}[\hat{y}^*] \downarrow$ ; as  $x^*$  closer to  $\bar{x}$ ,  $\text{Var}[\hat{y}^*] \downarrow$ .

Using  $S^2 = MRSS$ , we get the standard error of the estimate of  $\mathbb{E}[Y|X = x^*]$ ,

$$\sqrt{S^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)}.$$

Hence, a  $100(1 - \alpha)\%$  CI for  $\mathbb{E}[Y|X = x^*]$ , the mean response for all the elements in the population with  $X = x^*$  is

$$\left[ \hat{y}^* \pm t_{\frac{\alpha}{2}(n-2)} \sqrt{S^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)} \right].$$

Notice that it is only valid for  $x^*$  in the range of the original data values of  $X$  but not for extrapolation.

## 2.12 Prediction Interval for Actual Value of $Y$

A confidence interval is always reported for a parameter while a prediction interval is reported for the value of a random variable. We want to find a PI for the actual value of  $Y$  at  $X = x^*$ , i.e.,  $Y^* = Y|X = x^*$ .

The point estimate for  $Y^*$  is

$$\hat{y}^* = \hat{\beta}_0 + \hat{\beta}_1 x^*.$$

The error in our prediction is

$$\varepsilon^* = Y^* - \hat{y}^*.$$

The predicted value  $\hat{y}^*$  has two sources of variability:

- Since the regression line is estimated at  $\hat{\beta}_0 + \hat{\beta}_1 X$ ;
- due to  $\varepsilon^*$ , some points do not fall exactly on the line.

We have

$$\begin{aligned} \text{Var}[Y^* - \hat{y}^*] &= \text{Var}[Y - \hat{y}|X = x^*] \\ &= \text{Var}[Y|X = x^*] + \text{Var}[\hat{y}|X = x^*] - 2\text{Cov}(Y, \hat{y}|X = x^*) \\ &= \sigma^2 + \sigma^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right) - 0 = \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right). \end{aligned}$$

Notice that  $\text{Cov}(Y, \hat{y}|X = x^*) = 0$  since  $Y^*$  is a new observation.

Hence, a  $100(1 - \alpha)\%$  PI for  $Y|X = x^*$  is

$$\left[ \hat{y}^* \pm t_{\frac{\alpha}{2}(n-2)} \sqrt{S^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)} \right].$$

PIs for  $Y^*$  have the same center but are wider than CIs for  $\mathbb{E}[Y|X = x^*]$ .

## 3 Diagnostics and Transformations for Simple Linear Regression

### 3.1 Phenomena and Fallacies in Regression

- The regression effect: Regression to the mean - more values near the average than away from it; unusually large or small measurements tend to be followed by measurements that are closer to the mean.
- The regression fallacy: when the regression effect is mistaken for a real effect.
- Ecological fallacy/Correlation: inference is made about an individual based on aggregate data for a group.

### 3.2 Validity of SLR Model: Model Linearity

#### 3.2.1 Residuals

Recall that the residuals is  $\hat{e}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$ .

**Property 3.1.**  $\mathbb{E}[\hat{e}_i] = 0$ .

*Proof.* We have

$$\mathbb{E}[\hat{e}_i] = \frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0.$$

□

**Property 3.2.**  $\text{Var}[\hat{e}_i] = (1 - h_{ii})\sigma^2$ .

*Proof.* We have

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1 = \frac{\sum_{j=1}^n (x_j - \bar{x})y_j}{SXX}, \bar{y} = \frac{1}{n} \sum_{j=1}^n y_j.$$

Thus

$$\begin{aligned} \hat{y}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i = \frac{1}{n} \sum_{j=1}^n y_j + \frac{1}{SXX} \sum_{j=1}^n (x_j - \bar{x})(x_i - \bar{x})y_j \\ &= \sum_{j=1}^n \left( \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \right) y_j := \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j. \end{aligned}$$

Hence,

$$\hat{e}_i = y_i - \hat{y}_i = (1 - h_{ii})y_i + \sum_{j \neq i} h_{ij} y_j,$$

and thus

$$\begin{aligned} \text{Var}[\hat{e}_i] &= \text{Var} \left[ (1 - h_{ii})y_i + \sum_{j \neq i} h_{ij} y_j \right] = (1 - h_{ii})^2 \sigma^2 + \sum_{j \neq i} h_{ij}^2 \sigma^2 + 0 \\ &= \left( 1 - 2h_{ii} + \sum_{j=1}^n h_{ij}^2 \right) \sigma^2. \end{aligned}$$



Since

$$\begin{aligned}\sum_{j=1}^n h_{ij}^2 &= \sum_{j=1}^n \left( \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \right)^2 = \frac{1}{n} + 0 + \sum_{j=1}^n \frac{(x_i - \bar{x})^2 (x_j - \bar{x})^2}{SXX^2} \\ &= \frac{1}{n} + \frac{(x_i - \bar{x})^2}{SXX} = h_{ii},\end{aligned}$$

then  $\text{Var}[\hat{e}_i] = (1 - h_{ii})\sigma^2$ , where  $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{SXX}$ . □

**Property 3.3.**  $\sum_{j=1}^n h_{ij} = 1$ .

*Proof.* We have

$$\sum_{j=1}^n h_{ij} = \sum_{j=1}^n \left( \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \right) = \frac{n}{n} + \frac{x_i - \bar{x}}{SXX} \sum_{j=1}^n (x_j - \bar{x}) = 1 + 0 = 1.$$

□

**Property 3.4.**  $\sum_{j=1}^n h_{ij}^2 = h_{ii}$ .

**Property 3.5.**  $\sum_{j=1}^n h_{ij}x_j = x_i$ .

*Proof.* We have

$$\sum_{j=1}^n h_{ij}x_j = \sum_{j=1}^n \left( \frac{x_j}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})x_j}{SXX} \right) = \bar{x} + \frac{(x_i - \bar{x})SXX}{SXX} = x_i.$$

□

**Property 3.6.**  $\text{Var}[\hat{y}_i] = h_{ii}\sigma^2$ .

*Proof.* We have

$$\text{Var}[\hat{y}_i] = \text{Var} \left[ \sum_{j=1}^n h_{ij}y_j \right] = \sum_{j=1}^n h_{ij}^2 \text{Var}[y_j] = h_{ii}\sigma^2.$$

□

**Property 3.7.**  $\text{Cov}(\hat{e}_i, \hat{e}_j) = -h_{ij}\sigma^2$ , for  $i \neq j$ .

**Property 3.8.**  $\hat{e}_i \sim \mathcal{N}(0, (1 - h_{ii})\sigma^2)$ .

We can plot the residuals in three way:

- $\hat{e}_i$  against observation order/time  $i$ ;
- $\hat{e}_i$  against predictor value  $x_i$ ;
- $\hat{e}_i$  against fitted value  $\hat{y}_i$ .

Residuals plot can be used to assess whether an appropriate model has been fit to the data: if no pattern is found, then the model provides an adequate summary of the data; or the shape of the pattern provides information on the function of  $x$  that is missing from the model.

### 3.2.2 Diagnostics

We can use **scatter plot**, **residual plot** and **added-variable plot** to check for SLR model.

- Residual plot: Using residuals against  $x_i$  or  $\hat{y}_i$  will yield the same information. If the model is linear, there should be no pattern.
- Added-Variable plot: Using residuals against other potential predictors. Any pattern indicates that the other predictor should be included in the model.

If the assumption is violated, we can add additional predictors or transform  $X$  and/or  $Y$ .

## 3.3 Validity of SLR Model: Uncorrelated Errors

It is based on the design of the study and we can use randomization, where possible, to satisfy the assumption and widen the scope of inferences. Since  $\text{Cov}(\hat{e}_i, \hat{e}_j) = -h_{ij}\sigma^2, i \neq j$ , residuals are not uncorrelated even if errors are independent. However the covariance is usually so small and can be ignored in practice.

We can use **residual plot** to check, using residuals against observation order. If the errors are uncorrelated, there should be no pattern.

If the predictor is time-dependent, auto-correlation may exist. If the assumption is violated, we can fit a time series model or do longitudinal data analysis.

## 3.4 Validity of SLR Model: Homoscedasticity

### 3.4.1 Standardized Residuals

**Definition 3.1** (Standardized Residuals). The  $i$ th standardized residual is

$$r_i = \frac{\hat{e}_i}{\sqrt{S^2(1 - h_{ii})}},$$

where  $\hat{\sigma}^2 = S^2 = \frac{1}{n-2} \sum_{j=1}^n \hat{e}_j^2$ .

### 3.4.2 Diagnostics

We can use **residual plot** to check, using  $\sqrt{|\text{Residuals}|}$  or  $\sqrt{|\text{Standardized Residuals}|}$  against  $x_i$ . If the variance is constant, there should be no pattern.

If the assumption is violated, we can transform  $X$  and/or  $Y$ , do weighted least squares or fit a generalized linear model (models variance as a function of the mean).

## 3.5 Normality

### 3.5.1 Q-Q Plots

Probability plots are useful tools to graphically assess goodness-of-fit: we plot the observed order statistics against the expected theoretical quantiles, and if the data follows the particular

distribution, the plot should look roughly linear. The most common probability plot is the normal Q-Q plot.

### 3.5.2 The Shapiro-Wilk Test

We test  $H_0$  : the data follow a normal distribution, and the test statistic is

$$W = \frac{\left( \sum_{i=1}^n a_i x_{(i)} \right)^2}{\sum_{i=1}^2 (x_i - \bar{x})^2},$$

where  $a_i = \frac{\mathbf{m}^T V^{-1}}{(\mathbf{m}^T V^{-1} V^{-1} \mathbf{m})^{1/2}}$ ,  $\mathbf{m}^T = (m_1, \dots, m_n)$  are expected values of standard normal order statistics and  $V$  is the covariance matrix of those normal order statistics.

If the  $p$ -value is less than the significance level, there is evidence that the data is not normal.

In R, we can use

---

```
shapiro.test(data)
```

---

### 3.5.3 Diagnostics

We can use **normal Q-Q plot** of residuals or standardized residuals. If the errors are normally distributed, then there should be a linear relationship.

Recall that  $\sum_{j=1}^n h_{ij} = 1$  and  $\sum_{j=1}^n h_{ij} x_j = x_i$  and thus

$$\begin{aligned} \hat{e}_i &= y_i - \sum_{j=1}^n h_{ij} y_j = \beta_0 + \beta_1 x_i + e_i - \sum_{j=1}^n (\beta_0 + \beta_1 x_j + e_j) h_{ij} \\ &= \beta_0 + \beta_1 x_i + e_i - \beta_0 - \beta_1 x_i - \sum_{j=1}^n h_{ij} e_j = e_i - \sum_{j=1}^n h_{ij} e_j. \end{aligned}$$

With large samples, by the CLT, linear combinations of random variables are approximately normally distributed, no matter what their original distribution is. and thus for larger sample, the residuals can be used to assess the normality of errors.

Therefore, CIs and tests for  $\beta_0, \beta_1$  and  $\mathbb{E}[Y|X]$  are robust against non-normality (i.e., have approximately the correct coverage or approximately the correct  $p$ -value) but PIs are not robust against departure from normality because they are for one point.

Note that do not bother to check normality in the presence of other issues.

## 3.6 Outliers, Leverage Points, and Influential Points

Outliers are the points that do not follow the pattern of the data. Outliers with respect to the explanatory variable (in the  $x$  direction) are called leverage points. If we remove the leverage points from the data, the fitted model is substantially different, then we call it influential points.

A good leverage point has no adverse effect on the estimated regression coefficients, decreases the standard errors, and increases  $R^2$ .

### 3.6.1 Quantifying Leverage

Note that  $h_{ii}$  is the leverage of the  $i$ th data point and it varies only by the squared distance of  $x_i$  from its mean but not the values of the  $y$ 's. Recall that

$$\hat{y}_i = h_{ii}y_i + \sum_{j \neq i} h_{ij}y_j,$$

and  $h_{ii}$  shows how  $y_i$  affects  $\hat{y}_i$ .

For SLR,

$$\bar{h}_{ii} = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{n} + \frac{(x_i - \bar{x})^2}{SXX} \right) = \frac{1}{n} \left( \frac{n}{n} + \frac{SXX}{SXX} \right) = \frac{2}{n}$$

and we define a point is a leverage point if

$$h_{ii} > 2\bar{h}_{ii} = \frac{4}{n}.$$

If  $h_{ii} \approx 1$ , then for  $i \neq j$ ,  $h_{ij} \approx 0$  (since  $\sum_{j=1}^n h_{ij} = 1$ ) and thus  $\hat{y}_i \approx y_i$ , so the  $i$ th point is a leverage point, i.e., the fitted line is attracted by the point.  $\hat{e}_i$  has small variance and

$$\text{Var}[\hat{y}_i] \approx \text{Var}[y_i].$$

When leverage points do not exist, there is little or no difference in the plots of residuals when compared to plots of standardized residuals. In small to moderate size data sets, an influential point is one if  $|r_i| > 2$ . In very large data sets, an influential point is one if  $|r_i| > 4$ .

In R, we have:

---

```
# Calculate  $h_{ii}$ .
lm.influence(model)$hat
hatvalues(model)

# Calculate  $r_i$ .
rstandard(model)

# Calculate  $\hat{e}_i$ .
model$residuals
```

---

### 3.6.2 Cook's Distance

**Definition 3.2** (Cook's Distance). Cook's distance for the  $i$ th point is given by

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_{j(i)} - \hat{y}_j)^2}{2S^2},$$

where  $\hat{y}_{j(i)}$  is the  $j$  fitted value based on the fit obtained when the  $i$ th case has been removed from the fit.

Cook's distance measures influence of the  $i$ th observation.

**Property 3.9.**  $D_i = \frac{r_i^2}{2} \left( \frac{h_{ii}}{1-h_{ii}} \right)$ , where  $r_i$  is the  $i$ th standardized residual and  $h_{ii}$  is the leverage of the  $i$ th point.

Hence a large  $D_i$  may be due to large  $r_i$ , large  $h_{ii}$  or both.

If the largest  $D_i$  is much less than 1, deletion of a case will not change the estimate of  $\hat{\beta}$  by much.

**Definition 3.3.** A point is noteworthy if

$$D_i > \frac{4}{n-2}.$$

In practice, we look for gaps in the values of Cook's distance and not just whether values exceed the suggested cut-off.

In R, we have:

---

```
cooks.distance(model)
```

---

### 3.6.3 Recommendations

- Base estimates and confidence intervals only on valid model.
- Unusual points should be thoroughly investigated and should not be routinely deleted from an analysis.
- Outliers often point to important features of the problem not considered before.

## 3.7 Diagnostic Plots from R

In R, we have:

---

```
plot(model)
```

---

We will have 4 plots:

- $\hat{e}_i$  against  $\hat{y}_i$ .

- Normal Q-Q plot of  $r_i$ .
- $\sqrt{r_i}$  against  $\hat{y}_i$ .
- $r_i$  against  $h_{ii}$  with Cook's distance.

### 3.8 Transformations

With transformations, we can overcome problems due to non-constant variance, estimate percentage effects, overcome problems due to nonlinearity, and remedy non-normality.

#### 3.8.1 Common Transformations

Common monotonic transformations are  $X^2, \ln X, \sqrt{X}$ .

- If relationship is non-linear but variance of  $Y$  is nearly constant - transform  $X$ .
- If relationship is non-linear and variance is non-constant - transform  $Y$ .
- If relationship is linear but variance is non-constant - transform both  $X$  and  $Y$ .

Note that transforming changes the relative spacing of the observations.

#### 3.8.2 Variance Stabilizing Transformations

**Theorem 3.1** (Delta Method). Suppose  $Y$  has a distribution with mean  $\mu$  and variance  $\sigma_Y^2$ , and  $Z = f(Y)$ . We have

$$\mathbb{E}[Z] \approx f(\mu), \text{Var}[Z] \approx \sigma_Y^2 [f'(\mu)]^2.$$

*Proof.* Recall the first order Taylor series expansion of  $Z$ , we have

$$Z = f(Y) = f(\mu) + (Y - \mu)f'(\mu) + o(Y - \mu)$$

and thus

$$\mathbb{E}[Z] = \mathbb{E}[f(\mu)] + \mathbb{E}[(Y - \mu)f'(\mu)] + \mathbb{E}[o(Y - \mu)] \approx f(\mu).$$

Besides,  $Z - f(\mu) = (Y - \mu)f'(\mu) + o(Y - \mu)$  and thus

$$\text{Var}[Z] = \mathbb{E}[(Z - f(\mu))^2] \approx \mathbb{E}[(Y - \mu)f'(\mu)]^2 = \sigma_Y^2 [f'(\mu)]^2.$$

□

**Example 3.1.** Suppose  $Y \sim \text{Poisson}(\mu)$ , then  $\mathbb{E}[Y] = \mu = \text{Var}[Y]$ . Hence variance is linearly related to expectation and we want a  $f(Y)$  s.t.  $\text{Var}[f(Y)]$  will be constant (independent of  $\mu$ ), i.e., we want

$$\text{Var}[f(Y)] \approx \mu [f'(\mu)]^2 = C \Rightarrow [f'(\mu)]^2 \propto \frac{1}{\mu} \Rightarrow f'(\mu) \propto \frac{1}{\sqrt{\mu}}.$$

Thus,

$$f(\mu) \propto \sqrt{\mu}.$$