

# Miscellaneous Notes on Regression

Based on SJS and KNN

© Tingfeng Xia

Fall 2019, modified on Tuesday 12<sup>th</sup> November, 2019

## Preface

Notes for STA302H1F fall offering, 2019 with Prof. Shivon Sue-Chee. These notes are based on the KNN and SJS text, in an aim for better understanding of the course material.

This work is licensed under a Creative Commons “Attribution-NonCommercial-ShareAlike 4.0 International” license.



## Contents

<b>1</b>	<b>Preliminaries</b>	<b>3</b>
1.1	Distribution Theories . . . . .	3
1.2	Matrix Calculus . . . . .	3
1.2.1	Lemma I (Real Valued Fcn Matrix Differentiation) . . . . .	3
1.2.2	Lemma II . . . . .	4
1.2.3	Matrix Idempotency . . . . .	4
<b>2</b>	<b>Simple Linear Regression</b>	<b>4</b>
2.1	Ordinary Least Square . . . . .	4
2.1.1	Simple Linear Regression Models . . . . .	4
2.2	Inferences on Slope and Intercept . . . . .	5
2.2.1	Inference Assumptions . . . . .	5
2.2.2	Inference of Slope . . . . .	5
2.2.3	Inference of Intercept . . . . .	7
2.3	CI for <i>Unknown</i> Population Regression Line . . . . .	8
2.4	Prediction Intervals for Actual Value of Y . . . . .	9
2.5	Analysis of Variance (ANOVA) . . . . .	10
2.5.1	Sum of Squares Decomposition . . . . .	10
2.5.2	Test for Zero Slope . . . . .	11
2.5.3	Coefficient of Determination . . . . .	11
2.5.4	The ANOVA Table . . . . .	11
2.6	SLR in Matrix Form . . . . .	12
2.6.1	Set-Up . . . . .	12
2.6.2	The Design Matrix . . . . .	12

2.6.3	Normal Error Regression Model . . . . .	12
2.6.4	OLS in Matrix Form . . . . .	12
2.6.5	Properties of OLS Regressors . . . . .	13
2.6.6	The Hat Matrix . . . . .	13
2.6.7	Properties of the Residuals in Matrix Form . . . . .	14
2.6.8	ANOVA in Matrix Form . . . . .	14
2.6.9	ANOVA Table in Matrix Form . . . . .	15
<b>3</b>	<b>Diagnostics and Transformations for SLR</b>	<b>15</b>
3.1	Valid and Invalid Data . . . . .	15
3.1.1	Residuals . . . . .	15
3.1.2	Reading Residual Plots . . . . .	16
3.2	Regression Diagnostics . . . . .	16
3.2.1	Leverage Point . . . . .	17
3.2.2	Standardized Residuals . . . . .	18
3.2.3	Recommendations for Handling Outliers & Leverage . . . . .	20
3.2.4	Influence of Certain Cases . . . . .	21
3.2.5	Normality of the Errors . . . . .	21
3.2.6	Constant Variance (Homoscedasticity) . . . . .	23
3.3	Transformation . . . . .	24
3.3.1	Variance Stabilizing Transformations . . . . .	24
3.3.2	Logarithms to Estimate Percentage Effects . . . . .	24
<b>4</b>	<b>Weighted Least Square Regression</b>	<b>25</b>
4.1	Motivation and Set-Up . . . . .	25
4.2	Deriving LS Regressors . . . . .	25
<b>5</b>	<b>Multiple Linear Regression (Under Construction)</b>	<b>26</b>
<b>6</b>	<b>Selected Properties, Formulae, and Theorems</b>	<b>26</b>
6.1	Properties of Fitted Regression Line . . . . .	26
6.2	Rules of Expectation . . . . .	27
6.3	Variance and Covariance . . . . .	27
6.4	The Theorem of Gauss-Markov . . . . .	27
6.5	Matrix Form Rules . . . . .	28
6.5.1	Summations . . . . .	28
6.5.2	Transpositions . . . . .	28
6.5.3	Inversions . . . . .	28
6.5.4	Covariance Matrix . . . . .	28

# 1 Preliminaries

## 1.1 Distribution Theories

- Suppose  $Y_1, Y_2, \dots, Y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , and consider  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$ . Then,

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{(df=n-1)}$$

- Under the Normal Error SLR model, where  $e_i \stackrel{iid}{\sim} N(0, \sigma^2)$ , and  $S^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{e}_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$  (Different from above!). Then,

$$\frac{(n-2)S^2}{\sigma^2} = \frac{(n-2)S^2/SXX}{\sigma^2/SXX} = \sum_{i=1}^n \left( \frac{Y_i - \hat{Y}_i}{\sigma} \right)^2 \sim \chi^2_{(df=n-2)}$$

- Let  $Z \sim N(0, 1)$ , and  $V \sim \chi^2_{(df=v)}$ . Assume further that  $Z \perp\!\!\!\perp V$ , then

$$\frac{Z}{\sqrt{V/v}} \sim t_{(df=v)}$$

- Under the Normal Error SLR model,

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{S^2}{SXX}}} \sim t_{(df=n-2)}$$

- Suppose  $V \sim \chi^2_{(df=v)}$ ,  $W \sim \chi^2_{(df=w)}$  and  $V \perp\!\!\!\perp W$ . Then,

$$\frac{V/v}{W/w} \sim F_{(v,w)}$$

- Suppose  $Q \sim t_{(df=v)}$ , then

$$Q^2 \sim F_{(1,v)}$$

## 1.2 Matrix Calculus

Manuel to matrix calculus provided by professor: <http://biostat.mc.vanderbilt.edu/wiki/pub/Main/CourseBios312/chap2.pdf>

### 1.2.1 Lemma I (Real Valued Fcn Matrix Differentiation)

If  $\theta' = (\theta_1, \theta_2, \dots, \theta_k)$  and  $c' = (c_1, c_2, \dots, c_k)$  is a vector of constant, such that

$$f(\theta) = c'\theta = \theta'c = \sum_i c_i \theta_i$$

is a scalar, then

$$\frac{\partial f(\theta)}{\partial \theta} = c$$

### 1.2.2 Lemma II

Let  $\mathbf{A}$  be a  $k \times k$  symmetric matrix. Suppose  $f(\boldsymbol{\theta}) = \boldsymbol{\theta}'\mathbf{A}\boldsymbol{\theta}$ . Then,

$$\frac{\partial f(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 2\mathbf{A}\boldsymbol{\theta}$$

### 1.2.3 Matrix Idempotency

## 2 Simple Linear Regression

### 2.1 Ordinary Least Square

#### 2.1.1 Simple Linear Regression Models

The cost function to use in this case is the RSS, defined as

$$\text{RSS} = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

We will now derive the OLS estimators as follows.

#### Derivatives

$$\begin{aligned} \frac{\partial \text{RSS}}{\partial b_0} &= -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0 \\ \frac{\partial \text{RSS}}{\partial b_1} &= -2 \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) = 0 \end{aligned}$$

**Normal Equations** are obtained by rearranging

$$\sum_{i=1}^n y_i = b_0 n + b_1 \sum_{i=1}^n x_i \tag{1}$$

$$\sum_{i=1}^n x_i y_i = b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 \tag{2}$$

#### OLS Regressor

$$\begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{XY}}{S_{XX}} \end{aligned}$$

#### Estimating Variance of Error Term (Using Residuals)

$$\text{Unbiased Estimator } \hat{\sigma}^2 = S^2 = \frac{\text{RSS}}{n-2} = \frac{1}{n-2} \sum_{i=1}^n \hat{e}_i^2$$

## Notes

- $\bar{\hat{e}} = 0$ , since  $\sum \hat{e}_i = 0$  as the least square estimates minimizes RSS. (This is like a minimization goal where derivatives are taken w.r.t  $\hat{e}_i$ .)
- $S^2$  has  $n - 2$  degrees of freedom since we have estimated two parameters, namely  $\beta_0$  and  $\beta_1$ .

## 2.2 Inferences on Slope and Intercept

### 2.2.1 Inference Assumptions

The following assumptions need to be made in order to perform inference

- $Y$  is explained by  $x$  through a simple linear regression model

$$Y_i = \beta_0 + \beta_1 x_i + e_i (i = 1, \dots, n), \text{ i.e., } E(Y|X = x_i) = \beta_0 + \beta_1 x_i$$

- Independent Errors,  $e_i \perp e_j, \forall i \neq j$
- Homoscedasticity,  $\text{Var}(e_i) = \sigma^2, \forall i$
- Normal Error:  $e|X \sim N(0, \sigma^2)$

### 2.2.2 Inference of Slope

#### Distribution

$$\hat{\beta}_1|X \sim N\left(\beta_1, \frac{\sigma^2}{SXX}\right)$$

#### Standardized Test Statistic (Var Known)

$$Z = \frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{SXX}} \sim N(0, 1)$$

**Test Statistic (Var Unknown)** Recall that degrees of freedom = sample size - number of mean parameters estimated. Then,

$$T = \frac{\hat{\beta}_1 - \beta_1}{S/\sqrt{SXX}} = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim t_{(df=n-2)} \quad \text{where } S^2 = \frac{\sum_i \hat{e}_i^2}{n-2}$$

**Confidence Interval (Var Unknown)** The  $100(1 - \alpha)\%$  CI is

$$CI \leftarrow \hat{\beta}_1 \pm t_{(\alpha/2, df=n-2)}^* \times SE(\hat{\beta}_1) \equiv \hat{\beta}_1 \pm t_{(\alpha/2, df=n-2)}^* \frac{S}{\sqrt{SXX}}$$

**Distribution Proof** Recall OLS regressor for  $\beta_1$  is

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i \quad \text{where } c_i = \frac{x_i \bar{x}}{SXX}$$

The expectation could be derived as<sup>1</sup>

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

and the variance

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

Then, since  $e_i | X$  are normally distributed, then  $y_i = \beta_0 + \beta_1 x_i + e_i$ ,  $Y_i | X$  is normally distributed. Since  $\hat{\beta}_1 | X$  is a linear combination of  $y_i$ 's,  $\hat{\beta}_1 | X$  is normally distributed.  $\mathcal{Q.E.D.}^\dagger$

**t-test using PMCC on Bi-variate Normal** Recall that

$$\hat{\rho}_{MLE} = r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

---

<sup>1</sup>using the fact that  $\sum_{i=1}^n (x_i - \bar{x}) = 0$  and  $\sum_{i=1}^n (x_i - \bar{x}) x_i = \sum_{i=1}^n x_i^2 - n\bar{x}^2 = SXX$

Use the null hypothesis that  $H_0 : \rho_{XY} = 0$  and alternative  $H_1 : \rho_{XY} \neq 0$ , then

$$T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{b_1}{\sqrt{\frac{S^2}{SXX}}} \sim t_{(df=n-2)}$$

### 2.2.3 Inference of Intercept

#### Distribution

$$\hat{\beta}_0|X \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{SXX}\right)\right)$$

#### Standardized Test Statistic (Var Known)

$$Z = \frac{\hat{\beta}_0 - \beta_0}{\sigma \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SXX}}} \sim N(0, 1)$$

#### Test Statistic (Var Unknown)

$$Z = \frac{\hat{\beta}_0 - \beta_0}{S \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SXX}}} = \frac{\hat{\beta}_0 - \beta_0}{SE(\hat{\beta}_0)} \sim t_{(df=n-2)} \quad \text{where } S^2 = \frac{\sum_i \hat{e}_i^2}{n-2}$$

#### Confidence Interval (Var Unknown)

$$CI \leftarrow \hat{\beta}_0 \pm t_{(\alpha/2, df=n-2)}^* \times SE(\hat{\beta}_0) \equiv \hat{\beta}_0 \pm t_{(\alpha/2, df=n-2)}^* S \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SXX}}$$

**Distribution Proof** Recall that the OLS regressor of  $\beta_0$  is given by

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

The expectation,

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

and the variance,

$$\begin{aligned} \text{Var}(\hat{\beta}_0|X) &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}|X) \\ &= \text{Var}(\bar{y}|X) + \bar{x}^2 \text{Var}(\hat{\beta}_1|X) - 2\bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1|X) \end{aligned}$$

where

$$\begin{aligned}\text{Var}(\bar{y}|X) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n y_i | X = x_i\right) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \\ \text{Var}(\hat{\beta}_1|X) &= \frac{\sigma^2}{SXX} \\ \text{Cov}(\bar{y}, \hat{\beta}_1|X) &= \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n y_i, \sum_{i=1}^n c_i y_i\right) = \frac{1}{n} \sum_{i=1}^n c_i \text{Cov}(y_i, y_i) = \frac{\sigma^2}{n} \sum_{i=1}^n c_i = 0\end{aligned}$$

Thus,

$$\text{Var}(\hat{\beta}_0|X) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)$$

*Q.E.D.*†

### 2.3 CI for *Unknown* Population Regression Line

Goal: Find a confidence interval for a unknown population regression line at  $X = x^*$ . The population regression line is given by

$$E(Y|X = x^*) = \beta_0 + \beta_1 x^*$$

**Distribution** To get an estimate of the y value at  $X = x^*$ , we can use the regression output (evaluate the estimated regression line at  $X = x^*$ )

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

where we claim it follows the distribution

$$\hat{y}^* = \hat{y}|X = x^* \sim N\left(\beta_0 + \beta_1 x^*, \sigma^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)\right)$$

**Proof of Distribution** The expectation follows directly from definition, and we will now show that the variance has the claimed value. Notice that  $\text{Var}(\hat{\beta}_0|X) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)$ ,  $\text{Var}(\hat{\beta}_1|X) = \frac{\sigma^2}{SXX}$  and  $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1|X) = \frac{-\bar{x}\sigma^2}{SXX}$ , then

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

*Q.E.D.*†



**Standardized Test Statistic (Var Known)**

$$Z = \frac{\hat{y}^* - (\beta_0 + \beta_1 x^*)}{\sigma \sqrt{\left(\frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX}\right)}} \sim N(0, 1)$$

**Test Statistic (Var Known)**

$$T = \frac{\hat{y}^* - (\beta_0 + \beta_1 x^*)}{S \sqrt{\left(\frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX}\right)}} \sim t_{(df=n-2)} \quad \text{where } S^2 = \frac{\sum_i \hat{e}_i^2}{n-2}$$

**Confidence Interval** A  $100(1 - \alpha)\%$  CI for  $E(Y|X = x^*) = \beta_0 + \beta_1 x^*$  is given by

$$\begin{aligned} CI &\leftarrow \hat{y}^* \pm t_{(\alpha/2, df=n-2)}^* S \sqrt{\left(\frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX}\right)} \\ &\equiv \hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{(\alpha/2, df=n-2)}^* S \sqrt{\left(\frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX}\right)} \end{aligned}$$

Do note that this is only valid for  $x^*$  values in the range of the original data values of  $X$ . Avoid extrapolation.

**2.4 Prediction Intervals for Actual Value of Y**

**Goal** Find a prediction interval for the actual value of  $Y$  at  $x^*$ , a given value of  $X$ .

**Important Notes**

- $E(Y|X = x^*)$ , the expected value or average value of  $Y$  for a given value  $x^*$  of  $X$ , is what one would expect  $Y$  to be in the long run when  $X = x^*$ .  $E(Y|X = x^*)$  is therefore a fixed but unknown quantity whereas  $Y$  can take a number of values when  $X = x^*$ .
- $E(Y|X = x^*)$ , the value of the regression line at  $X = x^*$ , is entirely different from  $Y^*$ , a single value of  $Y$  when  $X = x^*$ . In particular,  $Y^*$  need not lie on the population regression line.
- A confidence interval is always reported for a parameter (e.g.,  $E(Y|X = x^*) = b_0 + b_1 x^*$ ) and a prediction interval is reported for the value of a random variable (e.g.,  $Y^*$ ).

**Difference Between CI and PI (My Thoughts)** The intrinsic difference is that: The CI we found above for  $E(Y|X = x^*)$  is a CI for a fixed value. We are trying to find, in the long run where can we expect the regression line to lie given infinite samples. However, PI is trying to report for a specific value, possibly a not-already-observed new value, what is the range that it may appear in. PI captures more variability.

**Distribution**

$$Y^* - \hat{y}^* \sim N \left( 0, \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right] \right)$$

**Test Statistic (Var Unknown)**

$$T = \frac{Y^* - \hat{y}^*}{S \sqrt{\left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)}} \sim t_{(df=n-2)}$$

**Derivation of Distribution** We base our prediction of  $Y$  at  $X = x^*$  (which is  $Y^*$ ) on

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

The error (deviation, to be precise) of our prediction is

$$Y^* - \hat{y}^* = \beta_0 + \beta_1 x^* + e^* - \hat{y}^* = (E(Y|X = x^*) - \hat{y}^*) + e^*$$

that is, the deviation between  $E(Y|X = x^*)$  and  $y^*$  plus the random fluctuation  $e^*$  (which represents the deviation of  $Y^*$  from  $E(Y|X = x^*)$ ). Thus the variability in the error for predicting a single value of  $Y$  will exceed the variability for estimating the expected value of  $Y$  (because of the random error  $e^*$ ). We have

$$E(Y^* - \hat{y}^*) = E(Y - \hat{y}|X = x^*) = 0$$

and

$$\text{Var}(Y^* - \hat{y}^*) = \text{Var}(Y - \hat{y}|X = x^*) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right]$$

**Prediction Interval** A  $100(1 - \alpha)\%$  prediction interval for  $Y^*$  (the value of  $Y$  at  $X = x^*$ ) is given by

$$\begin{aligned} PI &\leftarrow \hat{y}^* \pm t_{(\alpha/2, df=n-2)} S \sqrt{\left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)} \\ &\equiv \hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{(\alpha/2, df=n-2)} S \sqrt{\left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{SXX} \right)} \end{aligned}$$

**2.5 Analysis of Variance (ANOVA)****2.5.1 Sum of Squares Decomposition**

Define **Total Sample Variability**

$$\text{SST} = \text{SYY} = \sum_i^n (y_i - \bar{y})^2$$

and recall the familiar residual squared sum (Unexplained (or error) variability)

$$\text{RSS} = \sum_i^n (y_i - \hat{y}_i)^2$$

and we define (Sum of Squares explained by the regression model)

$$\text{SSreg} = \sum_i^n (\hat{y}_i - \bar{y})^2$$

Then, the decomposition is

$$\text{SST} = \text{RSS} + \text{SSreg}$$

### 2.5.2 Test for Zero Slope

**t-test** Note that for SLR, this is equivalent to the F-test outlined below. Consider the null  $H_0 : \beta_1 = 0$  against alternative  $H_1 : \beta_1 \neq 0$ , then

$$T = \frac{\hat{\beta}_1 - 0}{\text{se}(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - 0}{S/\sqrt{SXX}} \stackrel{H_0}{\sim} t_{(df=n-2)}$$

This is just a specific application of the general  $t$ -test that we mentioned earlier, not very interesting.

**F-test** Assume that  $e_i \perp e_j, \forall i \neq j \wedge e_i \sim N(0, \sigma^2), \forall i$ . Consider the null  $H_0 : \beta_1 = 0$  against alternative  $H_1 : \beta_1 \neq 0$ , then

$$F = \frac{\text{SSreg}/1}{\text{RSS}/(n-2)} \stackrel{H_0}{\sim} F_{1,n-2}$$

### 2.5.3 Coefficient of Determination

The Coefficient of Determination ( $R^2$ ) of a regression line is defined as the proportion of the total sample variability in the  $Y$ 's explained by the regression model, that is

$$R^2 = \frac{\text{SSreg}}{\text{SST}} = 1 - \frac{\text{RSS}}{\text{SST}}$$

### 2.5.4 The ANOVA Table

The above  $F$ -test, as well as the sum of squares decomposition, could be summarized using the following handy table.

Source of Variation	$df$	SS	$MS = SS/df$	$F$
Regression	1	SSreg	SSreg/1	$F = \frac{\text{SSreg}/1}{\text{RSS}/(n-2)}$
Residual	$n - 2$	RSS	$\text{RSS}/(n - 2)$	
Total	$n - 1$	SST		

## 2.6 SLR in Matrix Form

### 2.6.1 Set-Up

The simple linear regression model is

$$Y = X\beta + e$$

where  $Y \in M_{n \times 1}(\mathbb{R})$ ,  $X \in M_{n \times 2}(\mathbb{R})$ ,  $\beta \in M_{2 \times 1}(\mathbb{R})$ ,  $e \in M_{n \times 1}(\mathbb{R})$ .

### 2.6.2 The Design Matrix

$$X_{n \times 2} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \implies X\beta (n \times 1) \begin{bmatrix} \beta_0 + X_1\beta_1 \\ \beta_0 + X_2\beta_1 \\ \vdots \\ \beta_0 + X_n\beta_1 \end{bmatrix}$$

### 2.6.3 Normal Error Regression Model

#### Gauss - Markov Conditions

- The errors have zero mean,  $E(\mathbf{e}) = \mathbf{0}$
- The errors have constant variance,  $\sigma^2$
- The errors are uncorrelated,  $V(\mathbf{e}) = \sigma^2 \mathbf{I}$

**Jointly Normal** The error terms follow a multivariate normal,

$$\mathbf{e} \sim N_n(\Sigma = \mathbf{0}, \sigma^2 I)$$

### 2.6.4 OLS in Matrix Form

Consider  $\beta = [\beta_0, \beta_1]'$ , and the cost function

$$\begin{aligned} RSS(\beta) &= \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2 \\ &= (Y - X\beta)(Y - X\beta)' \\ &= (Y' - \beta'X')(Y - X\beta) \\ &= Y'IY + \beta(X'X)\beta' - Y'X\beta - \beta'X'Y \end{aligned}$$

#### Derivative

$$\frac{\partial RSS(\beta)}{\partial \beta} = 0 - 2X'Y\beta + 2X'X\beta$$

**Normal Equation** obtained by setting to zero and re-arrange

$$2X'X\hat{\beta} = 2X'Y$$

**OLS Regressor**

$$\hat{\beta} = (X'X)^{-1}X'Y$$

**2.6.5 Properties of OLS Regressors****Expectation: Unbiased Estimator**

$$\begin{aligned} E(\hat{\beta}) &= E((X'X)^{-1}X'Y) \\ &= E(Y)(X'X)^{-1}X' \\ &= (X'X)^{-1}X'E(X\beta + e) \\ &= (X'X)^{-1}X'(X\beta + 0) \\ &= I\beta = \beta \end{aligned}$$

Q.E.D.†

**Variance - Covariance Matrix**

$$\begin{aligned} VAR(\hat{\beta}) &= VAR((X'X)^{-1}X'Y) \\ &= [(X'X)^{-1}X']\sigma^2 I[(X'X)^{-1}X']' \\ &= \sigma^2 I(X'X)^{-1} \underbrace{X'X(X'X)^{-1}}_{\equiv I} \\ &= \sigma^2 (X'X)^{-1} \\ &= \begin{bmatrix} VAR(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{nSXX} & COV(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 \bar{x}}{SXX} \\ COV(\hat{\beta}_1, \hat{\beta}_0) = \frac{-\sigma^2 \bar{x}}{SXX} & VAR(\hat{\beta}_1) = \frac{\sigma^2}{SXX} \end{bmatrix} \end{aligned}$$

**2.6.6 The Hat Matrix****Defining the Hat** We have

$$\hat{\mathbf{e}} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} \quad \text{and} \quad \hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

where

$$\text{HAT Matrix } \mathbf{H} := \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

Thus,  $\hat{\mathbf{e}} = \mathbf{Y} - \mathbf{H}\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$ **Facts About the Hat Matrix**

- $\mathbf{H}$  is symmetric

$$H' = (X(X'T)^{-1}X')' = X(X'X)^{-1}X' = H$$

- $\mathbf{H}$  is idempotent

$$\begin{aligned} H^2 &= HH = X(X'X)^{-1} \underbrace{X'X(X'X)^{-1}}_{\equiv I} X' \\ &= X(X'X)^{-1} X' = H \end{aligned}$$

### 2.6.7 Properties of the Residuals in Matrix Form

$$\begin{aligned} \hat{e} &= (I - H)Y = (I - H)(X\beta + e) \\ &= (I - H)X\beta + (I - H)e \\ &= IX\beta - HX\beta + (I - H)e \\ &= X\beta - \underbrace{X(X'X)^{-1}X}_{\equiv I} \beta + (I - H)e \\ &= X\beta - X\beta + (I - H)e \\ &= (I - H)e \end{aligned}$$

#### Expectation

$$E(\hat{e}|X) = E((I - H)Y|X) = E((I - H)e|X) = \underbrace{E(e|X)}_{\equiv 0} (I - H) = 0$$

#### Variance - Covariance

$$\begin{aligned} VAR(\hat{e}|X) &= VAR((I - H)e|X) \\ &= (I - H)VAR(e|X)(I - H)' \\ &= (I - H)\sigma^2 I(I - H)' \\ &= (I - H)\sigma^2 I(I - H) \quad \text{since } (I - H) \text{ is symmetric} \\ &= \sigma^2(I - H) \quad \text{since } (I - H) \text{ is idempotent} \end{aligned}$$

### 2.6.8 ANOVA in Matrix Form

#### Total Sum Squared (SST)

$$\begin{aligned} SST &= \sum (y_i - \bar{y})^2 = \sum y_i^2 - n\bar{y}^2 \\ &= Y'Y - \frac{1}{n}Y'JY \\ &= Y' \underbrace{\left(I - \frac{1}{n}J\right)}_{\substack{\text{symmetric square matrix} \\ \text{quadratic form}}} Y \end{aligned}$$

To find the corresponding df, we first check the idempotency of  $I - \frac{1}{n}J$ . Indeed,

$$\left(I - \frac{1}{n}J\right)\left(I - \frac{1}{n}J\right) = I^2 - \frac{2}{n}J + \frac{1}{n^2}J^2 = I - \frac{1}{n}J$$

Then, applying  $\text{Rank}(\text{idempotent mat}) = \text{Tr}(\text{idempotent mat})$ , we have

$$\text{Rank}(I - \frac{1}{n}J) = n - 1 = \text{degrees of freedom of } SST$$

**Residual Sum Squared (RSS)** Notice that  $(I - H)$  is symmetric and idempotent,

$$RSS = \sum \hat{e}_i^2 = \hat{e}'\hat{e} = Y'(I - H)(I - H)Y = Y'(I - H)Y$$

and the corresponding degrees of freedom is

$$\text{Rank}(I - H) = \text{Rank}(I) - \text{Rank}(H) = \text{Tr}(I) - \text{Tr}(H) = n - 2 = \text{df of MSE}$$

**Regression Sum Squared (SSReg)**

$$\begin{aligned} SSReg &= SST - RSS \\ &= Y'(I - 1/nJ)Y - Y'(I - H)Y \\ &= Y'IY - Y'1/nJY - Y'IY + Y'HY \\ &= Y'(H - \frac{1}{n}J)Y \end{aligned}$$

and the corresponding degrees of freedom is

$$\text{Rank}(H - \frac{1}{n}J) = \text{Rank}(H) - \text{Rank}(\frac{1}{n}J) = \sum h_{ii} - \sum \frac{1}{n} = 2 - 1 = 1 = \text{df of MSReg}$$

### 2.6.9 ANOVA Table in Matrix Form

Source	SS	df
Regression	$\mathbf{Y}'(\mathbf{H} - \frac{1}{n}\mathbf{J})\mathbf{Y}$	1
Error	$\mathbf{Y}'(\mathbf{I} - \mathbf{H})\mathbf{Y}$	$n - 2$
Total	$\mathbf{Y}'(\mathbf{I} - \frac{1}{n}\mathbf{J})\mathbf{Y}$	$n - 1$

## 3 Diagnostics and Transformations for SLR

### 3.1 Valid and Invalid Data

#### 3.1.1 Residuals

One tool that we can use to validate a regression model is one or more plots of residuals (or standardized residuals). These plots will enable us to assess visually whether an appropriate model has been fit to the data no matter how many predictor variables are used.

**Expected Behavior** We expect that the residual graph to have no discern-able pattern and centered at some value (0 in the case of standardized residual). Patterns such as curves, skewness et cetra indicates non-normal residuals. More on this in the below section.

### 3.1.2 Reading Residual Plots

**Criterion** One way of checking whether a valid simple linear regression model has been fit is to plot residuals versus  $x$  and look for patterns. If no pattern is found then this indicates that the model provides an adequate summary of the data, i.e., is a valid model. If a pattern is found then the shape of the pattern provides information on the function of  $x$  that is missing from the model.

**Rationale** Suppose that the true model is a straight line (which we never know) defined as

$$Y_i = E(Y_i|X_i = x_i) + e_i = \beta_0 + \beta_1 x_i + e_i \quad (3)$$

where

$$e_i = \text{Random error on } Y_i \quad \text{and} \quad E(e_i) = 0$$

and we fit a regression line

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Under the assumption that our regression line is very close to the true model, i.e.  $\beta_0 \approx b_0$  and  $\beta_1 \approx b_1$ , we see

$$\begin{aligned} \hat{e}_i &= y_i - \hat{y}_i \\ &= \beta_0 + \beta_1 x_i + e_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \\ &= (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) x_i + e_i \\ &\approx e_i \end{aligned}$$

which means that our residuals resembles the random error!

## 3.2 Regression Diagnostics

### Categorization

- **X-Direction Outlier, i.e. Leverage Point:** Away from the bulk of data in the  $x$ -direction.
  - **Good:** Not much change after removing the data point, i.e. the data point originally was quite close to the regression line although away from the bulk of data in the  $x$  direction. “A good leverage point is a leverage point which is NOT also an outlier.”
  - **Bad, Influential Point:** If its  $Y$ -value does not follow the pattern set by the other data points, i.e. a bad leverage point is a leverage point which is also an outlier.
- **Y-Direction Outlier** Trait: large residuals



### 3.2.1 Leverage Point

**Defining The Hat** The hat came from yet another representation of the  $\hat{y}_i$ . Recall that  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ , where  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ ,  $\hat{\beta}_1 = \sum_{j=1}^n c_j y_j$  and  $c_j = \frac{x_j - \bar{x}}{SXX}$ . Then we have

$$\begin{aligned}\hat{y}_i &= \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x}) \\ &= \frac{1}{n} \sum_{j=1}^n y_j + \sum_{j=1}^n \frac{(x_j - \bar{x})}{SXX} y_j (x_i - \bar{x}) \\ &= \sum_{j=1}^n \left[ \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \right] y_j = \sum_{j=1}^n h_{ij} y_j\end{aligned}$$

where we define

$$h_{ij} = \left[ \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \right]$$

**Property of The Hat** Recall that  $\sum_{j=1}^n [x_j - \bar{x}] = n\bar{x} - n\bar{x} = 0$ , then

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

Thus,

$$\hat{y}_i = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j \quad \text{where } h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

**Defining Leverage** The term  $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$  above is commonly known as the leverage of the  $i$ th data point. Notice the following in the definition of the leverage  $h_{ii}$

- The second term measures the proportion, in terms of squared deviation in  $x$ -direction over sum of square of total deviation in  $x$ -direction, of the  $i$ -th data point's deviation. When the second term tends to 1, meaning that  $i$ -th data point is some extreme outlier in the  $x$ -direction, then  $h_{ii}$  would close to one, signifying the 'leverage'-ness.
- Recall that  $\sum_{j=1}^n h_{ij} = 1$ , then when  $h_{ii} \cong 1$ ,  $h_{ij} \rightarrow 0$  and

$$\hat{y}_i = 1 \times y_i + \text{other terms} \cong y_i$$

which means  $\hat{y}_i$  will be very close to  $y_i$ , regardless of the rest dataset.

- A point of high leverage (or a leverage point) can be found by looking at just the values of the  $x$ 's and not at the values of the  $y$ 's

**Average of Leverage** For simple linear regression,

$$\text{average}(h_{ii}) = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{2}{n}$$

**Identifying Leverage** Rule:  $x_i$  is a high leverage (i.e., a leverage point) in a SLR model if

$$h_{ii} > 2 \times \text{average}(h_{ii}) = 2 \times 2/n = 4/n$$

### Dealing with ‘Bad’ Leverage

- **Remove invalid data points;** Question the validity of the data points corresponding to bad leverage points, that is: Are these data points unusual or different in some way from the rest of the data? If so, consider removing these points and refitting the model without them.
- **Fit a different regression model;** Question the validity of the regression model that has been fitted, that is: Has an incorrect model been fitted to the data? If so, consider trying a different model by including extra predictor variables (e.g., polynomial terms) or by transforming  $Y$  and/or  $x$  (which is considered later in this chapter).

### 3.2.2 Standardized Residuals

**Problem of Non-constant Variance** Recall that

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX}$$

and (we will show this later)

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

which is indeed non-constant for different data points. When  $h_{ii} \cong 1$  ( $h_{ii}$  is very close to 1), the  $i$ -th data point is a leverage point and

$$\text{Var}(\hat{e}_i) = \sigma^2 [1 - h_{ii}] \approx 0 \quad \text{and} \quad \hat{y}_i \cong y_i$$

The above results intuitively makes sense: When  $i$ -th data point is a leverage,  $\hat{e}_i$  will be small and it does not vary much (data point close to the estimated regression line).

**Derivation of Residual Variance (Not Important)** Recall that

$$\hat{y}_i = h_{ii}y_i + \sum_{j \neq i} h_{ij}y_j \quad \text{where} \quad h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX}$$

Then,

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

Hence,

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

Notice that

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

So,

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

and

$$\text{Var}(\hat{y}_i) = \text{Var}\left(\sum_{j=1}^n h_{ij}y_j\right) = \sum_{j \neq i} h_{ij}^2 \text{Var}(y_j) = \sigma^2 \sum_j h_{ij}^2 = \sigma^2 h_{ii}$$

**Overcome with Standardization** The above problem of each  $\hat{e}_i$  having different variances could be overcome by standardizing the residuals. The  $i$ -th standardized residual is defined as (notice that the  $s = \hat{\sigma}$  is the estimated variance in the SLR settings)

$$r_i = \frac{\hat{e}_i}{s\sqrt{1 - h_{ii}}} \quad \text{where } s = \sqrt{\frac{1}{n-2} \sum_{j=1}^n \hat{e}_j^2}$$

### Advantages of Standardization

- When points of high leverage exist, instead of looking at residual plots, it is generally more informative to look at plots of standardized residuals since plots of the residuals will have non-constant variance even if the errors have constant variance.

- When points of high leverage do not exist, there is generally little difference in the patterns seen in plots of residuals when compared with those in plots of standardized residuals.
- The other advantage of standardized residuals is that they immediately tell us how many estimated standard deviations any point is away from the fitted regression model.

### Recognizing Outliers Using Standardized Residuals

- An **outlier** is a point whose standardized residual falls outside the interval from -2 to 2, i.e.  $|r_i| > 2$
- A **Bad Leverage Point** is a leverage point whose standardized residual falls outside the interval from -2 to 2, i.e.  $|r_i| > 2 \wedge h_{ii} > \frac{4}{n}$
- A **Good Leverage Point** is a leverage point whose standardized residual falls inside the interval from -2 to 2, i.e.  $|r_i| \leq 2 \wedge h_{ii} > \frac{4}{n}$
- **Dealing with large datasets:** In this case, we should change the above criterion to  $|r_i| > 4$  and  $|r_i| \leq 4$  respectively. This is to give allowance for more occurrence of rare events in a large data set.

**Correlation Between Residuals** Even if the errors are independent (homogeneous), i.e.  $e_i \perp e_j$  ( $i \neq j$ ), the residuals are still correlated. It can be shown that the covariance and the correlation is given by

$$\begin{aligned}\text{Cov}(\hat{e}_i, \hat{e}_j) &= -h_{ij}\sigma^2 (i \neq j) \\ \text{Corr}(\hat{e}_i, \hat{e}_j) &= \frac{-h_{ij}}{\sqrt{(1-h_{ii})(1-h_{jj})}} (i \neq j)\end{aligned}$$

Such correlation could be safely ignored in practice. They are usually given rise by inherent correlation such as data collected over time.

**Variance of Residuals** Above we discussed ‘inter-correlation’ of residuals. **The variance of a single residual is**

$$\text{Var}(\hat{e}_i) = (1 - h_{ii})\sigma^2$$

### 3.2.3 Recommendations for Handling Outliers & Leverage

We have discussed multiple ways of assessing outliers and talked about the way to deal with them by removing them. However, it is not always a good idea to delete them for the following reasons:

- Points should not be routinely deleted from an analysis just because they do not fit the model. Outliers and bad leverage points are signals, flagging potential problems with the model.
- Outliers often point out an important feature of the problem not considered before. They may point to an alternative model in which the points are not an outlier. In this case it is then worth considering fitting an alternative model.

### 3.2.4 Influence of Certain Cases

It can sometimes be the case that certain data points in a data set are drastically controlling the entire regression model (the model has paid too much attention to them). We now develop methods where we measure the “importance” of a data point.

**Cook’s Distance** First, define (recall if already defined) the following notation

- $\hat{y}_{j(i)}$  means the fitted value of the  $j$ -th data point on the regression line obtained by removing the  $i$ -th case.
- $S^2 = \frac{1}{n-2} \sum_{j=1}^n \hat{e}_j^2$  is the variance (Original MSE) of the **total regression model**.
- $r_i = \frac{\hat{e}_i}{s\sqrt{1-h_{ii}}}$  where  $s = \sqrt{\frac{1}{n-2} \sum_{j=1}^n \hat{e}_j^2}$

Then, the Cook’s Distance of the  $i$ -th data point is given by

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

where we should note that  $D_i$  may be large due to large  $r_i$ , or large  $h_{ii}$  or both.

#### Rule: Cook’s Distance

- A point is noteworthy if

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

- In practice, look for gaps in the values of Cook’s Distance and not just whether one value exceeds the suggested cut off.

### 3.2.5 Normality of the Errors

The assumption of normal errors is (especially) needed in small samples for the validity of  $t$ -dist based tests and inferences. This assumption is generally checked by looking at the distribution of the residuals or standardized residuals. Recall that the  $i$ -th least squares residuals is given by  $\hat{e}_i = y_i - \hat{y}_i$ . We will now show  $\hat{e}_i = e_i - \sum_{j=1}^n h_{ij}e_j$ . First, in the derivation we will need these two facts

$$\sum_{i=1}^n h_{ij} = 1$$

and

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

We then proceed as follows

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

*Q.E.D.*<sup>†</sup>

The above result showed that the  $i$ -th least squares residual is equal to  $e_i$  minus a weighted sum of all the  $e$ 's. There are two cases to consider,

- In small to moderate samples, the second term could dominate the first and first and the residuals can look like they come from a normal distribution even if the errors do not.
- When  $n$  is large, the second term in the derived result (thistle colored) has a much smaller variance than that of the first term and as such the first term dominates the last equation.

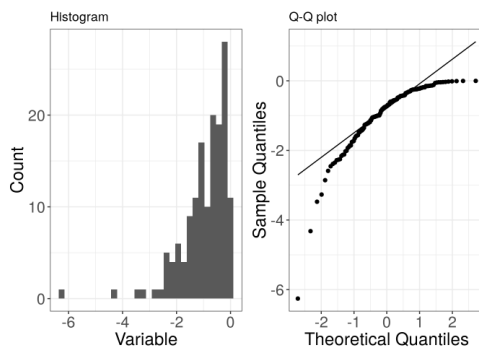
**Conclusion:** For large samples, the residuals can be used to assess normality of the errors.

**Assessment Using Normal Q-Q** A normal probability plot of the standardized residuals is obtained by plotting the ordered standardized residuals on the vertical axis against the expected order statistics from a standard normal distribution on the horizontal axes. If the resulting plot produces points “close” to a straight line then the data are said to be consistent with that from a normal distribution. On the other hand, departures from linearity provide evidence of non-normality.

#### Left-skewed data

Below is an example of data (150 observations) that are drawn from a distribution that is **left-skewed** (in this case it is a negative exponential distribution). Left-skew is also known as **negative skew**.

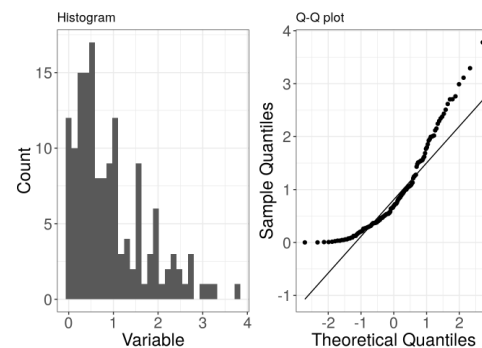
On a Q-Q plot left-skewed data appears curved (the opposite of right-skewed data).



#### Right-skewed data

Below is an example of data (150 observations) that are drawn from a distribution that is **right-skewed** (in this case it is the exponential distribution). Right-skew is also known as **positive skew**.

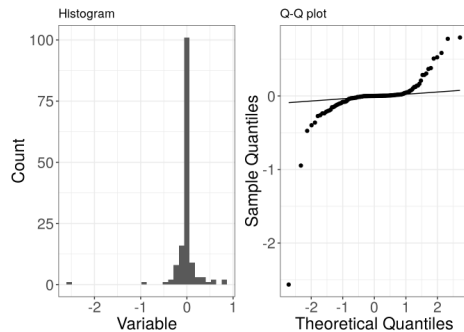
On a Q-Q plot right-skewed data appears curved.



### Over-dispersed data

Below is an example of data (150 observations) that are drawn from a distribution that is **over-dispersed** relative to a normal distribution (in this case it is a Laplace distribution). Over-dispersed data has an increased number of outliers (i.e. the distribution has fatter tails than a normal distribution). Over-dispersed data is also known as having a **leptokurtic distribution** and as having **positive excess kurtosis**.

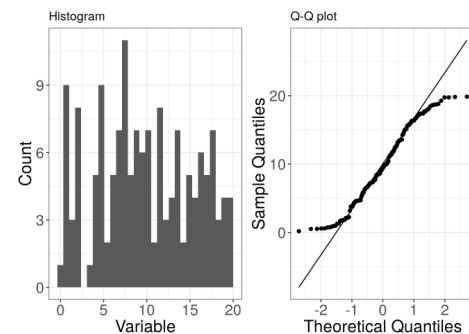
On a Q-Q plot over-dispersed data appears as a flipped S shape (the opposite of under-dispersed data).



### Under-dispersed data

Below is an example of data (150 observations) that are drawn from a distribution that is **under-dispersed** relative to a normal distribution (in this case it is the uniform distribution). Under-dispersed data has a reduced number of outliers (i.e. the distribution has thinner tails than a normal distribution). Under-dispersed data is also known as having a **platykurtic distribution** and as having **negative excess kurtosis**.

On a Q-Q plot under-dispersed data appears S shaped.



The above four are figures that I borrowed from [http://www.ucd.ie/ecomodel/Resources/QQplots\\_WebVersion.html](http://www.ucd.ie/ecomodel/Resources/QQplots_WebVersion.html) which illustrates how to interpret QQ plots with non-normal behavior.

### 3.2.6 Constant Variance (Homoscedasticity)

A crucial assumption in any regression analysis is that errors have constant variance. **Notice the difference between error and residual, we have demonstrated in (3.2.2 Standardized Residuals) that residuals are not of constant variance.** There are two general methods that we can adopt to overcome this issue, namely (both of which will be discussed later)

- Transformations
- Weighted Least Squares

**Important:** Ignoring non-constant variance when it exists invalidates all inferential tools, including *p-values*, *CI*, *PI*, et cetera!

**Behavior of Non-Homoscedasticity** For example, on the plot explanatory var against standardized residuals, we might see that as  $x$  increases, the residuals are more spread out, indicating an increasing trend in the variance.

**Checking for Constant Variance** To check this, check the plot of

$$| \text{Residuals} |^{0.5} \text{ against } x \quad \text{or} \quad | \text{Standardized Residuals} |^{0.5} \text{ against } x$$

The power of 0.5 here is used to reduce skewness in the absolute values. In the above mentioned example where the residuals become more spread out as  $x$  increases, the plot  $| \text{Standardized Residuals} |^{0.5}$  against  $x$  will have an overall increasing trend! **This is essentially mirroring all the points to the positive side (and de-skew) to observe a general trend.**

### 3.3 Transformation

#### 3.3.1 Variance Stabilizing Transformations

**Goal** When non-constant variance exists, it is often possible to transform one or both of the regression variables to produce a model in which the error variance is constant.

**Delta Method, Poisson** Suppose that  $Y \sim \text{Poi}(\mu = \lambda)$  and we want to find the appropriate transformation of  $Y$  for stabilizing variance. In this case, square root is the appropriate transformation to apply. We will now justify this choice. Consider the McLaurin Series expansion

$$f(Y) = f(E(Y)) + f'(E(Y))(Y - E(Y)) + \dots$$

According to the delta rule, the first order variance term is obtained by taking variance on both sides of the above equation, which yields

$$\text{Var}(f(Y)) \simeq [f'(E(Y))]^2 \text{Var}(Y)$$

Using the proposed transformation  $f(Y) = Y^{0.5}$  and recall from properties of Poisson Random Variable that  $\text{Var}(Y) = \lambda = E(Y)$ , then

$$\text{Var}(Y^{0.5}) \simeq [0.5(E(Y))^{-0.5}]^2 \text{Var}(Y) = [0.5\lambda^{-0.5}]^2 \lambda = \text{constant}$$

**Rule of Thumb:** When both  $Y$  and  $X$  are measured in the same units then it is often natural to consider the same transformation for both  $X$  and  $Y$

Hence in this case our regression model would be

$$Y = \beta_0 + \beta_1 x + e$$

where

$$Y \leftarrow \sqrt{Y} \quad \text{and} \quad x \leftarrow \sqrt{x}$$

#### 3.3.2 Logarithms to Estimate Percentage Effects

Consider the regression model

$$\log(Y) = \beta_0 + \beta_1 \log(x) + e$$

The slope,<sup>2</sup>

$$\begin{aligned} \beta_1 &= \frac{\Delta \log(Y)}{\Delta \log(x)} = \frac{\log(Y_2) - \log(Y_1)}{\log(x_2) - \log(x_1)} = \frac{\log(Y_2/Y_1)}{\log(x_2/x_1)} \\ &\cong \frac{Y_2/Y_1 - 1}{x_2/x_1 - 1} \quad (\text{using } \log(1+z) \cong z \text{ and assuming } \beta_1 \text{ is small}) \\ &= \frac{100(Y_2/Y_1 - 1)}{100(x_2/x_1 - 1)} = \frac{\% \Delta Y}{\% \Delta x} \end{aligned}$$

---

<sup>2</sup>Notice that the first step is possible since here we are considering the regression straight line



**Interpretation** We showed above that  $\% \Delta Y \simeq \beta_1 \times \% \Delta x$ . Thus for every 1% increase in  $x$ , the model predicts a  $\beta_1\%$  increase in  $Y$  (provided  $\beta_1$  is small).

## 4 Weighted Least Square Regression

### 4.1 Motivation and Set-Up

Consider the straight line (simple) linear regression model

$$Y_i = \beta_0 + \beta_1 x_i + e_i \quad \text{where } e_i \sim N\left(0, \frac{\sigma^2}{w_i}\right)$$

For the weight  $w_i$ , we should note the following

- $w_i \rightarrow \infty \implies \text{Var}(e_i) \rightarrow 0$ . In this case, the estimates of the regression parameters  $\beta_0, \beta_1$  should be such that the fitted line at  $x_i$  should be very close to  $y_i$ . (Small variance means more strict in terms of deviation from the regression line, corresponding to a larger emphasis on the  $i$ -th data point.)
- If  $w_i$  is some small value, then the variance of the  $i$ -th data point would be quite large. In this case, we have a loose restriction of the deviation of the  $i$ -th data point from the regression line meaning that little emphasis is taken for this data point.
- $w_i \rightarrow 0 \implies \text{Var}(e_i) \rightarrow \infty$ . In this case, we have the variance tending to infinity. Meaning that there is absolutely no restriction/emphasis on the  $i$ -th data point and it could be simply removed from the set.

We define the cost function, WRSS as

$$\text{WRSS} = \sum_{i=1}^n w_i (y_i - \hat{y}_{W_i})^2 = \sum_{i=1}^n w_i (y_i - b_0 - b_1 x_i)^2$$

and the estimators  $\mathbf{b} = [b_0, b_1]^T$  are derived using MLE.

**Intuition behind WRSS** This cost function may seem wierd at first glance, but it intuitively makes sense. Notice that when  $w_i$  is large, the  $i$ -th lost term  $w_i (y_i - \hat{y}_{W_i})^2$  is payed more emphasis on. On the contrary, when  $w_0 \rightarrow 0$ , the term  $\rightarrow 0$ . (Indeed, when Variance of the term  $\rightarrow \infty$  we just neglect it.)

### 4.2 Deriving LS Regressors

**Derivatives**

$$\frac{\partial \text{WRSS}}{\partial b_0} = -2 \sum_{i=1}^n w_i (y_i - b_0 - b_1 x_i) = 0 \quad (4)$$

$$\frac{\partial \text{WRSS}}{\partial b_1} = -2 \sum_{i=1}^n w_i x_i (y_i - b_0 - b_1 x_i) = 0 \quad (5)$$

**Normal Equations** Obtained from rearranging the above equations, we will call them Normal Eq1 and Normal Eq2 respectively for later reference.

$$\sum_{i=1}^n w_i y_i = b_0 \sum_{i=1}^n w_i + b_1 \sum_{i=1}^n w_i x_i \quad (6)$$

$$\sum_{i=1}^n w_i x_i y_i = b_0 \sum_{i=1}^n w_i x_i + b_1 \sum_{i=1}^n w_i x_i^2 \quad (7)$$

**Rearranging** Use Normal Eq1  $\times \sum_{i=1}^n w_i x_i$  and Normal Eq2  $\times \sum_{i=1}^n w_i$

$$\sum_{i=1}^n w_i x_i \sum_{i=1}^n w_i y_i = b_0 \sum_{i=1}^n w_i \sum_{i=1}^n w_i x_i + b_1 \left( \sum_{i=1}^n w_i x_i \right)^2 \quad (8)$$

$$\sum_{i=1}^n w_i \sum_{i=1}^n w_i x_i y_i = b_0 \sum_{i=1}^n w_i \sum_{i=1}^n w_i x_i + b_1 \sum_{i=1}^n w_i \sum_{i=1}^n w_i x_i^2 \quad (9)$$

**WLS Slope Regressor** <sup>3</sup>

$$\hat{\beta}_{1W} = \frac{\sum_{i=1}^n w_i \sum_{i=1}^n w_i x_i y_i - \sum_{i=1}^n w_i x_i \sum_{i=1}^n w_i y_i}{\sum_{i=1}^n \sum_{i=1}^n w_i x_i^2 - \left( \sum_{i=1}^n w_i x_i \right)^2} \quad (10)$$

$$= \frac{\sum_{i=1}^n x_i (x_i - \bar{x}_W) (y_i - \bar{y}_W)}{\sum_{i=1}^n w_i (x_i - \bar{x}_W)^2} \quad (11)$$

**WLS Intercept Regressor**

$$\hat{\beta}_{0W} = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} - \hat{\beta}_{1W} \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i} = \bar{y}_w - \hat{\beta}_{1W} \bar{x}_W \quad (12)$$

## 5 Multiple Linear Regression (Under Construction)

## 6 Selected Properties, Formulae, and Theorems

This section contains various properties mentioned in the slides/book. They may or may not have appeared in previous sections.

### 6.1 Properties of Fitted Regression Line

- $\sum_{i=1}^n \hat{e}_i = 0$
- $RSS = \sum_{i=1}^n \hat{e}_i^2 \neq 0$ , generally. Except for when we have perfect fit.
- $\sum_{i=1}^n \hat{e}_i x_i = 0$
- $\sum_{i=1}^n \hat{e}_i \hat{y}_i = 0$
- $\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n y_i$

---

<sup>3</sup>Note that  $\bar{x}_W = \sum_{i=1}^n w_i x_i / \sum_{i=1}^n w_i$  and  $\bar{y}_W = \sum_{i=1}^n w_i y_i / \sum_{i=1}^n w_i$

## 6.2 Rules of Expectation

- $E(a) = a, \forall a \in \mathbb{R}$
- $E(aY) = aE(Y)$
- $E(X \pm Y) = E(X) \pm E(Y)$
- $X \perp\!\!\!\perp Y \implies E(EY) = E(X)E(Y)$
- Tower Rule:  $E(Y) = E(E(Y|X))$

## 6.3 Variance and Covariance

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

## 6.4 The Theorem of Gauss-Markov

Under the conditions of the simple linear regression model, the OLS regressors are BLUE (“Best Linear Unbiased Estimators”).

- Best - obtains the minimum variance among all unbiased linear estimators.
- Linear - Linear in the parameter space. That is, feature maps are linear, although the actual curve of regression is ‘non-linear’.
- Unbiased - The estimators are unbiased, namely  $\hat{\beta}_0, \hat{\beta}_1$ .
- Estimator - Estimators  $\hat{\beta}_0, \hat{\beta}_1$  for  $\beta_0$  and  $\beta_1$  respectively.

## 6.5 Matrix Form Rules

### 6.5.1 Summations

Consider  $\mathbf{A}, \mathbf{B}$  as compatible matrices where appropriate and  $k \in \mathbb{R}$  then

- $\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$
- $(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$
- $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$
- $\mathbf{C}(\mathbf{A} + \mathbf{B}) = \mathbf{CA} + \mathbf{CB}$
- $k(\mathbf{A} + \mathbf{B}) = k\mathbf{A} + k\mathbf{B}$

### 6.5.2 Transpositions

- $(\mathbf{A}')' = \mathbf{A}$
- $(\mathbf{A} + \mathbf{B})' = \mathbf{A}' + \mathbf{B}'$
- $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$
- $(\mathbf{ABC})' = \mathbf{C}'\mathbf{B}'\mathbf{A}'$ , this properties is known as ‘cyclic’

### 6.5.3 Inversions

- $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$
- $(\mathbf{ABC})^{-1} = \mathbf{C}^{-1}\mathbf{B}^{-1}\mathbf{A}^{-1}$
- $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$
- $(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T$

### 6.5.4 Covariance Matrix

- The variance-covariance matrix of a random vector  $\mathbf{Y}$  is a symmetric, positive semi-definite matrix, defined as

$$\begin{aligned}\text{Var}(\mathbf{Y}) &= \mathbf{E}[(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))'] \\ &= \mathbf{E} \begin{pmatrix} (Y_1 - E(Y_1))^2 & (Y_1 - E(Y_1))(Y_2 - E(Y_2)) & \dots \\ (Y_2 - E(Y_2))(Y_1 - E(Y_1)) & (Y_2 - E(Y_2))^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}\end{aligned}$$

- Let  $\mathbf{A}$  be a square matrix of constants

$$\begin{aligned}\text{Var}(\mathbf{AY}) &= \mathbf{E}[(\mathbf{AY} - \mathbf{E}(\mathbf{AY}))(\mathbf{AY} - \mathbf{E}(\mathbf{AY}))'] \\ &= \mathbf{E}[\mathbf{A}(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))'\mathbf{A}'] \\ &= \mathbf{AE}[(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))(\mathbf{Y} - \mathbf{E}(\mathbf{Y}))']\mathbf{A}' \\ &= \mathbf{A}\text{Var}(\mathbf{Y})\mathbf{A}'\end{aligned}$$