

STAT 331: Applied Linear Regression

Professor Kun Liang
L^AT_EXer Iris Jiang

Spring 2020

Contents

| | | |
|----------|---|----------|
| 1 | Introduction to Regression | 2 |
| 1.1 | What is regression | 2 |
| 1.2 | Why linear model? | 3 |
| 1.3 | Sample vs. population | 3 |
| 2 | Simple Linear Regression (SLR) | 4 |
| 2.1 | Population model | 4 |
| 2.2 | Assumptions | 5 |
| 2.3 | Least Square Estimation (LSE) | 5 |
| 2.3.1 | Task | 5 |
| 2.3.2 | Goal and Derivation | 5 |

Chapter 1

Introduction to Regression

1.1 What is regression

Definition 1.1 (Regression analysis). Regression analysis is a statistical methodology that models the **functional relationship** between a response variable y and one or more explanatory variables x_1, x_2, \dots, x_p .

A typical regression model is:

$$y = f(x_1, x_2, \dots, x_p) + \epsilon$$

- y : dependent variable or **response** variable
- x_1, x_2, \dots, x_p : **covariates**, explanatory variables, independent variables, or **predictors**
- ϵ : random error term

Regression models can be used to:

- Identify important predictors
- Estimate regression coefficients
- Estimate the response for given values of predictors
- Predict of future values of response

In STAT 331, we focus on the simplest form of regression: **linear models**

$$\begin{aligned} y &= f(x_1, x_2, \dots, x_p) + \epsilon \\ &= \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon \end{aligned}$$

where the β 's are the regression parameters(coefficients).

Linear in the parameter (not predictor). Linear model is the basic building block of more complicated models

Remark. *We refer to the model as linear in the parameters β 's ($\frac{\partial f}{\partial \beta_i}$ do not depend on the parameters)*

Example 1.1. Are the following models linear?

1. $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$

2. $f(x) = \beta_0 + \beta_1 e^{\beta_2 x}$
3. $f(x_1, x_2) = \beta_0 + \beta_1 x_1 x_2$

Solution.

1. This is a linear model. The predictor is x , this is not a linear model on the predictor but we define the linear model as to parameter, $\beta_0, \beta_1, \beta_2$ in this case.
2. This is not a linear model. If taking derivative to β_1 , the result involves β_2 .
3. This is a linear model.

□

1.2 Why linear model?

- Linear model is easy to implement and interpret
- All functions can be approximated locally by a linear function
- The simplest starting model to fit

1.3 Sample vs. population

Definition 1.2 (Sample). A sample is the collection of units (people, animals, cities, whatever you study) that is actually measure or surveyed.

Definition 1.3 (Population). The population is the large group of unites we are interested in, from which the sample was selected.

Remark. *We assume the data we have a representative sample (random sample) from a larger population*

Chapter 2

Simple Linear Regression (SLR)

2.1 Population model

$$y = \beta_0 + \beta_1 x + \epsilon$$

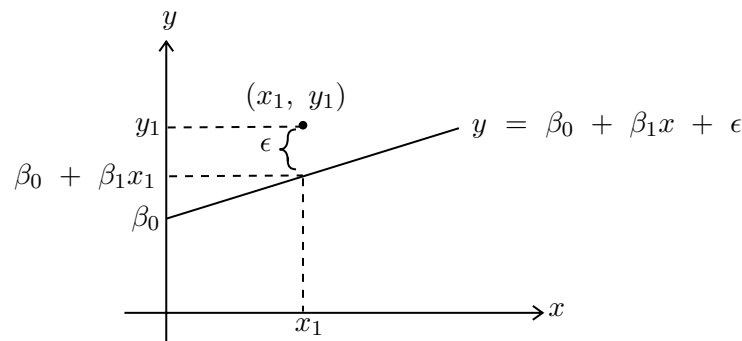
- y : response
- β_0, β_1 : regression Coefficients
- x : predictor
- ϵ : random error
- $\beta_0 + \beta_1 x$: systematic (deterministic) part

Observed sample: suppose we have n pairs of observations (x_i, y_i) , $i = 1, \dots, n$. Then

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

- x_i : fixed and known (for this course)
- β : fixed and unknown
- ϵ_i : random and unknown
- y_i : random and known

“known” means we can observe them



- β_0 : intercept
- β_1 : slope

2.2 Assumptions

1. $E(\epsilon_i) = 0$
2. $\epsilon_1, \dots, \epsilon_n$ are statistically independent
3. Constant variance: $\text{Var}(\epsilon_i) = \sigma^2 \implies \text{Var}(y_i) = \sigma^2$
The randomness of y_i comes from ϵ_i
4. ϵ_i is normally distributed. $\epsilon_i \sim N(0, \sigma^2)$ and $y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

Note. Assumption 1 to 3 are called Gauss-Markov assumptions.

Assumption 4 is stronger than all 3 assumptions combined.

1 to 3 are useful if you want to point estimate β .

4 is useful for further results.

There is no guarantee that the assumptions are correct. We will talk about model diagnostic and checking these assumptions.

2.3 Least Square Estimation (LSE)

2.3.1 Task

Given the sample observation (x_i, y_i) , $i = 1, \dots, n$, estimate (β_0, β_1) as $(\hat{\beta}_0, \hat{\beta}_1)$ such that the values of

$$r_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i = y_i - \hat{y}_i$$

are “small”.

- r_i : residual
- \hat{y}_i : fitted value

We define discrepancy function

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = \sum_{i=1}^n r_i^2$$

The reason we use square: the least square method provides an elegant solution; ϵ 's follow normal distribution the Least Squared Estimation has the equivalence with Maximum Likelihood estimation.

2.3.2 Goal and Derivation

Minimize $S(\beta_0, \beta_1)$, i.e. solve
$$\begin{cases} \frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = 0 \\ \frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = 0 \end{cases}.$$

$$\begin{cases} \frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-1) = 0 \\ \frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = 0 \end{cases}$$

This is called Normal Equation

$$\Rightarrow \begin{cases} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0 & \dots (1) \\ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)(x_i) = 0 & \dots (2) \end{cases}$$

$$(1) \Rightarrow n\hat{\beta}_0 = \sum_{i=1}^n y_i - \hat{\beta}_1 \sum_{i=1}^n x_i \text{ or } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \text{ where } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \text{ and } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$(2) \Rightarrow \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0$$

we can replace $\hat{\beta}_0$ from previous results.

$$\Rightarrow \sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i - \hat{\beta}_1 \bar{x} \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \text{ where } S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) \text{ and } S_{xx} = \sum (x_i - \bar{x})^2$$

$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$ is the fractions of the sample covariance between x and y over the sample variance of x
Thus, the solution is

$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \end{cases}$$