



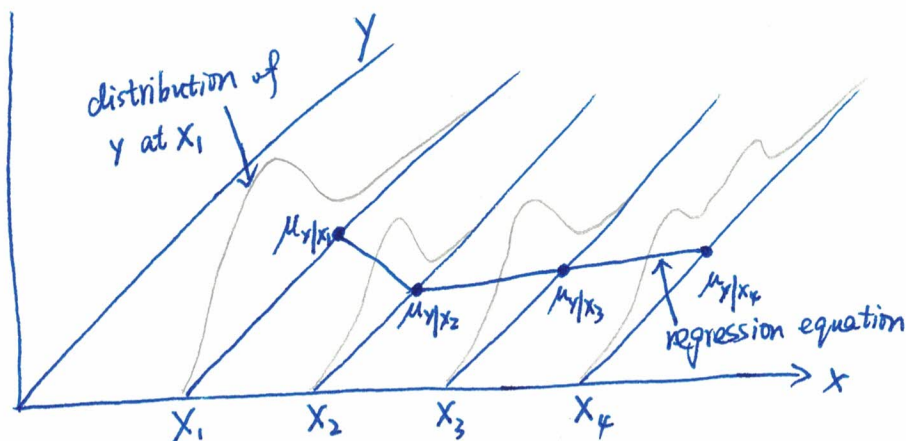
Simple Linear Regression

1. Statement of Assumptions

1.1. Existence: applies to any regression model

For any fixed value of the variable X , Y is a random variable with a certain probability distribution having finite mean and variance $\mu_{Y/X}$ $\sigma^2_{Y/X}$

Y/X : mean and variance of the r.v. Y depend on X .



1.2. Independence

The Y -values are statistically independent of one another.

violated examples: different observations are made on the same individual at different times. (longitudinal data).

mixed model \rightarrow involving repeated measurements

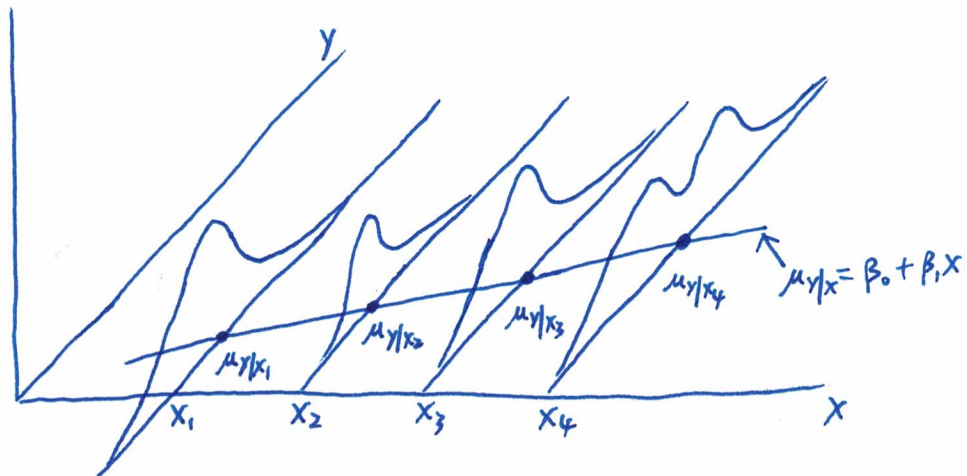
multivariate linear model

generalized estimating equations (GEE) \rightarrow analyzing correlated response data.



1.3 Linearity

The mean value of Y , $\mu_{Y|X}$, is a straight-line function of X .



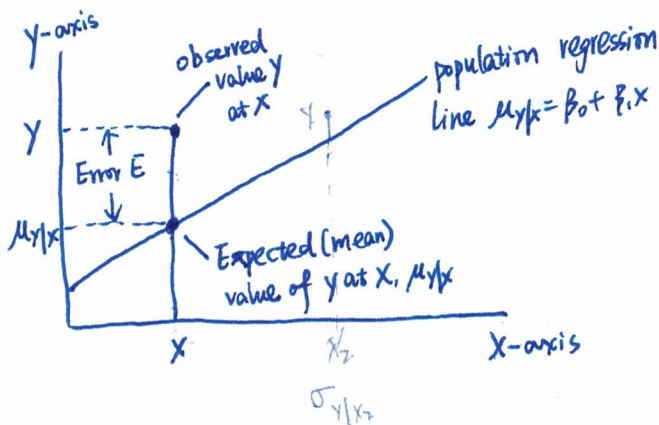
$$\mu_{Y|X} = \beta_0 + \beta_1 X$$

\uparrow intercept \uparrow slope

$$Y = \beta_0 + \beta_1 X + E$$

\uparrow r.v. \uparrow fixed not random. \uparrow r.v.

E is a r.v. with mean 0 at fixed X .
i.e. $\mu_{E|X} = 0$ for any X .
 \uparrow the error component.



The variable E describes how distant an individual's response can be from the population regression line.

What we observe at a given X (namely y) is in error from that expected on the average (namely $\mu_{Y|X}$) by an amount E , which is random and varies from individual to individual.



1.4 homoscedasticity

The variance of y is the same for any x .

example of violation: Page 2. y/x_1 vs y/x_2 , y/x_1 more spread $\sigma_{y/x_1}^2 > \sigma_{y/x_2}^2$

$$\sigma_{y/x}^2 \equiv \sigma^2$$

1.5. Normal distribution

For any fixed value of x , y has a normal distribution.

This assumption makes it possible to evaluate the statistical significance (eg. C.I., testing)

Violation: y transformation, eg: $\log(y)$, \sqrt{y} , ... box-cox procedure, ...

Summary

Maintain distinctions among random variables, parameters, and point estimates.

Y : r.v.

X : assumed to be measured without error. fixed value.

β_0, β_1 : parameters (unknown) for a population.

E/ε random, ^{unobservable} ~~unobserved~~ variable.

$\hat{\beta}_0$: point estimate for β_0 .

$\hat{E} = y - \hat{y} = y - (\hat{\beta}_0 + \hat{\beta}_1 x)$: point estimate of E at value x . \rightarrow residual

Normal distribution = Gaussian distribution.



2. Estimate the coefficients

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ n observation pairs

Goal: find β_0, β_1 such that the resulting line is as close as possible to the data points
different measures of closeness

2.1 The Least-squares Method

$$(\hat{\beta}_0, \hat{\beta}_1)_{LS} = \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$= \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n e_i^2$$

$e_i = y_i - \hat{y}_i$: i^{th} residual

$\sum_{i=1}^n e_i^2$: residual sum of squares (RSS), or sum of squares about the regression line
sum of squares due to error (SSE), Error sum of squares (SSE)

2.2. The minimum-variance method.

$\hat{\beta}_0, \hat{\beta}_1$: unbiased for β_0, β_1 , and have minimum variance among all unbiased estimators.

2.3. Maximum Likelihood method.

$\hat{\beta}_0, \hat{\beta}_1$: maximize the likelihood function.

$$L(Y; \beta_0, \beta_1, \sigma^2) = \prod_{i=1}^n f(y_i; \beta_0, \beta_1, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2 \right\}$$

Under the assumption of Gaussian and mutual independence of y_i 's,
maximum-likelihood estimates, LS estimates, and minimum-variance estimates are all the same.



2.4 LS estimates

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

The value of β_0, β_1 that minimize $Q(\beta_0, \beta_1)$ satisfy:

$$\frac{\partial Q(\beta_0, \beta_1)}{\partial \beta_0} = 0, \quad \frac{\partial Q(\beta_0, \beta_1)}{\partial \beta_1} = 0$$

This leads to the normal equations:

$$\begin{cases} n\beta_0 + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \\ \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \end{cases}$$

useful identities:

$$\sum (x_i - \bar{x})^2 = \sum x_i^2 - n\bar{x}^2$$

$$\begin{aligned} \sum (x_i - \bar{x})(y_i - \bar{y}) &= \sum (x_i - \bar{x})y_i \\ &= \sum x_i y_i - n\bar{x}\bar{y} \end{aligned}$$

$$\Rightarrow \begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = r_{xy} \cdot \frac{s_y}{s_x} \end{cases}$$

The fitted regression line (LS line):

$$y = \hat{\beta}_0 + \hat{\beta}_1 x = \bar{y} + \hat{\beta}_1 (x - \bar{x}) \quad , \quad \text{pass through the point } (\bar{x}, \bar{y}) \text{ — center of data}$$

The fitted value for the i^{th} case:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x}), \quad i=1, 2, \dots, n$$

Residual $e_i = y_i - \hat{y}_i = y_i - \bar{y} - \hat{\beta}_1 (x_i - \bar{x})$, an "estimator" of the error term ε_i

Properties of residuals:

$$1) \sum_{i=1}^n e_i = 0 \quad \sum (y_i - \bar{y}) - \hat{\beta}_1 \sum (x_i - \bar{x}) = 0$$

$$2) \sum_{i=1}^n x_i e_i = 0 \quad \sum x_i (y_i - \bar{y}) - \hat{\beta}_1 \sum x_i (x_i - \bar{x}) = \sum x_i y_i - n\bar{x}\bar{y} - \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} (\sum x_i^2 - n\bar{x}^2)$$

$$3) \sum_{i=1}^n \hat{y}_i e_i = 0 \quad \sum (\hat{\beta}_1 (x_i - \bar{x}) + \bar{y}) e_i = \sum \bar{y} e_i + \hat{\beta}_1 (\sum x_i e_i - \bar{x} \sum e_i)$$



2.5. Estimate of Error Variance.

$$\text{Var}(\varepsilon_i) = \sigma^2$$

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

$$E(SSE) = (n-2)\sigma^2$$

Mean squared error (MSE), unbiased estimator of σ^2 .

$$MSE = \frac{SSE}{n-2}, \quad E(MSE) = \sigma^2$$

Analogy

$$\text{usually } \text{var} = \frac{1}{n-1} \sum (y_i - \bar{y})^2$$

y_i 's are indep. with the same mean μ and variance σ^2 .

$n-1$: to estimate σ^2 , we need to estimate μ first.

our case: the population mean $\mu_{y|x}$ changes with x .

\hat{y} is the estimate of $\mu_{y|x}$, so we have $(y_i - \hat{y}_i)^2$

$n-2$: we need to estimate β_0, β_1 first

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

3. Properties of LS Estimators

3.1. LS estimators are linear functions of the responses y_i 's.

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \frac{(x_i - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} y_i = \sum_{i=1}^n k_i y_i$$

$$\hat{\beta}_0 = \sum_{i=1}^n \left(\frac{1}{n} - \bar{x} k_i \right) y_i$$

\hat{y}_i and e_i are also linear functions of y_i 's.