

Regression Analysis

Q: What are the goals of regression analysis?

A: Estimation, Testing, and Prediction

- Estimation of the **effect** of one variable (exposure), called the predictor of interest, after **adjusting**, or controlling for other measured variables
 - ★ Remove confounding variables
 - ★ Remove bias
- Testing whether variables are associated with the response
- Prediction of a response variable given a collection of covariates

1

2

Regression Analysis

- **Classification of Variables**
- Response variable
 - ★ Dependent variable
 - ★ Outcome variable
- Predictor of interest (POI)
 - ★ Exposure variable
 - ★ Treatment assignment
- Confounding variables
 - ★ Associated with response and POI
 - ★ Not intermediate
- Precision variables
 - ★ Associated with response and not with POI
 - ★ Reduces response uncertainty

3

Example

Impact of maternal smoking on low birth

- Measured variables
 - birth weight (g)*
 - maternal smoking (yes/no)*
 - maternal age (yrs)*
 - maternal weight at last menses (kg)*
 - race*
 - history of premature labor (yes/no)*
 - history of hypertension*

Response:

Predictor of interest:

Confounder:

Precision:

4

Linear regression model

- y_{ij} , $j = 1, \dots, n$, $i = 1, \dots, m$
- t_j corresponding times at which the measurements are taken on each unit

- x_{ijk} , $k = 1, \dots, p$ explanatory variables

$$y_{ij} = \beta_1 x_{ij1} + \dots + \beta_p x_{ijp} + \epsilon_{ij}$$

- in the classical linear regression model

$$\epsilon_{ij} \sim N(0, \sigma^2), \text{ cor}(\epsilon_{ij}, \epsilon_{ik}) = 0$$

- ordinary least square estimation

5

Review of linear model

$$y_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$$

$$= \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

$$\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$$

$$\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$$

$$\epsilon_i \sim N(0, \sigma^2), i = 1, \dots, m, \text{ iid}$$

$$x_{i1} = 1 \text{ then } \beta_1 \text{ is the intercept}$$

$$E(\epsilon_i) = 0, V(\epsilon_i) = \sigma^2$$

$$Y = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \sim MVN(0, \sigma^2 I)$$

$$\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_m)$$

6

$$y_1 = \beta_1 x_{11} + \dots + \beta_p x_{1p} + \epsilon_1$$

$$y_2 = \beta_1 x_{21} + \dots + \beta_p x_{2p} + \epsilon_2$$

$$\vdots = \vdots$$

$$y_m = \beta_1 x_{m1} + \dots + \beta_p x_{mp} + \epsilon_m$$

$$Y = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- ϵ_i iid $N(0, \sigma^2)$ if and only if $\boldsymbol{\epsilon} \sim MVN(0, \sigma^2 I)$

- $[cov(\boldsymbol{\epsilon})] = cov(\epsilon_i, \epsilon_j) = \sigma^2 I$

$$\bullet \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

7

Review of linear model

- Linear model includes as special cases:

1. the analysis of variance, x_{ij} are dummy variables
2. multiple regression, x_{ij} are continuous
3. the analysis of covariance, x_{ij} are dummy and continuous variables

- β is the expected value of the response variable Y , per unit change of its corresponding explanatory variable x , all other variables held fixed.

8

Estimation of β

- Principle of maximum likelihood
RA Fisher 1925
- *Estimate β by the values that make the observations maximally likely*
- Likelihood $P(\text{data} \mid \beta)$ - as a function of β

9

Likelihood Inference

- Likelihood Inference is based on the specification of the probability density for the observed data

$$L(\theta \mid \mathbf{y}) = f(\mathbf{y}; \theta)$$

- Maximum likelihood estimate $\hat{\theta}$ is defined as

$$L(\theta \mid \mathbf{y}) \leq L(\hat{\theta} \mid \mathbf{y})$$

- $\hat{\theta}$ is then regarded as the value of θ which is most strongly supported by the observed data
- θ is obtained by direct maximization of $L(\theta)$ or $\log L(\theta)$ by solving

$$S(\theta) = \frac{\partial \log L(\theta)}{\partial \theta} = 0$$

- $S(\theta)$ is called score equation for θ

10

Linear Model – IID Normal Errors

$$\begin{aligned} L(\beta; Y) &= \prod_{i=1}^m g(y_i; \beta, \sigma^2) \\ &= \prod_{i=1}^m \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(y_i - \mathbf{x}_i' \beta)^2}{2\sigma^2}\right) \end{aligned}$$

Maximizing likelihood = maximizing log likelihood

$$\begin{aligned} \log L &= l(\beta, \sigma^2, Y) = \\ &= \\ &= \sum_{i=1}^m \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y_i - \mathbf{x}_i' \beta)^2}{2\sigma^2} \right\} \\ &= \\ &= c(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \mathbf{x}_i' \beta)^2 \\ &= \\ &= c(\sigma^2) - SS(\beta) \end{aligned}$$

Maximizing normal likelihood = minimizing sum of squares

11

Geometry of least squares

$$\begin{aligned} \mathbf{x}_1 &= \begin{pmatrix} x_{11} \\ \vdots \\ x_{m1} \end{pmatrix}, \dots, \mathbf{x}_p = \begin{pmatrix} x_{1p} \\ \vdots \\ x_{mp} \end{pmatrix} \\ X &= (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p), \dim(X) = m \times p \\ X\beta &= \mathbf{x}_1\beta_1 + \mathbf{x}_2\beta_2 + \dots + \mathbf{x}_p\beta_p \\ SS(\beta) &= \sum_{i=1}^m (y_i - \mathbf{x}_i' \beta)^2 \\ &= (Y - X\beta)'(Y - X\beta) \\ &= \|Y - X\beta\|^2 \end{aligned}$$

$SS(\beta)$ is the LENGTH² of the residual vector $(Y - X\beta)$.

Choose $\hat{\beta}$ so that the distance from Y to $X\beta$ is as small as possible!

12

Methods of Least Squares

$\hat{\beta}$ that minimizes the length of residual vector
 $Y - X\hat{\beta}$ makes $Y - X\hat{\beta}$ orthogonal to x_1, \dots, x_p

$$\rightarrow X'(Y - X\hat{\beta}) = 0$$

$$\rightarrow X'Y - X'X\hat{\beta} = 0$$

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

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

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

$$= HY$$

$$H \cdot H = H \text{ (you check)}$$

H project Y onto space spanned by columns of X

13

Distribution Theory for Linear model

$$Y = X\beta + \epsilon, \epsilon \sim MVN(0, \sigma^2 I)$$

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

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

β is unbiased

$$\begin{aligned} Var\hat{\beta} &= Var\left((X'X)^{-1}X'Y\right) \\ &= (X'X)^{-1}X'VarYX(X'X)^{-1} \\ &= (X'X)^{-1}\sigma^2 \end{aligned}$$

$$\bullet E\hat{Y} = EHY = HEY = HX\beta = X\beta$$

$$\bullet Var\hat{Y} = HVarYH' = \sigma^2 H^2 = \sigma^2 H$$

$$\bullet E\hat{\epsilon} = E(Y - \hat{Y}) = X\beta - X\beta = 0$$

$$\bullet Var\hat{\epsilon} = (I - H)\sigma^2 I(I - H)' = \sigma^2(I - H)$$

14

Statistical Properties of $\hat{\beta}$

- It is an unbiased estimator: $E(\hat{\beta}) = \beta$
- $Var(\hat{\beta}) = \sigma^2(X'X)^{-1}$
- For any vector a of known coefficients, $\phi = a'\beta$ then $\hat{\phi} = a'\hat{\beta}$ has the smallest possible variances amongst all unbiased estimators for ϕ which are linear combination of the Y_i . (Gauss-Markov theorem)