# Lecture 5

## The classical linear regression model

# Review of key concepts from Lecture 3 and 4

▶ Simple linear regression model

    ▶ ARM = B0 + B1 (age – 6) + e, e~N(0,$\sigma^2$), independent

Systematic
component

↳ random
component

B0 = average ARM
among 6 month
old children

B1 = difference in
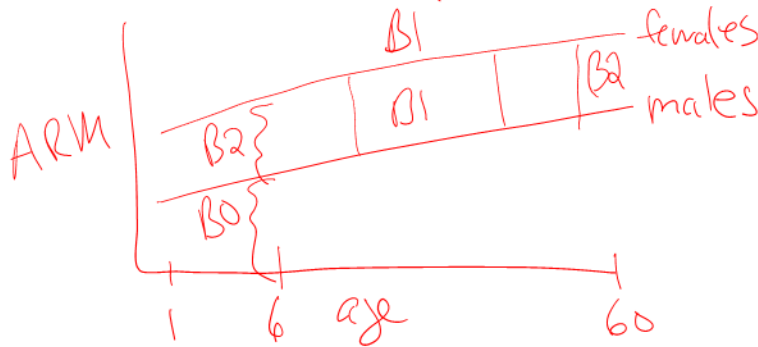average ARM
comparing children
who vary in a age
by 1-month

ARM

B1

B0

1    6    Age    60

# Review of key concepts from Lecture 3 and 4

▶ Sex adjusted relationship between ARM and age

  ▶ ARM = B0 + B1 (age – 6) + B2 Female + e, e~N(0,$\sigma^2$), independent

$\left\{ \right.$ = Controlling for

= adjusting for

male : ARM = B0 + B1 (age-6) + e

female : ARM = B0 + B1 (age-6) + B2 + e

= (B0+B2) + B1 (age+6) + e

B0 = Average ARM among male children who are 6 months of age

& B1 = difference in mean ARM comparing two children who have the same sex but differ in age by 1 month

ARM

B2 B1 B1 B2 females

B2 B1 males

B0

1   6   age   60

B2 : Difference in mean ARM comparing females to males who are the same age

3

# Review of key concepts from Lecture 3 and 4

▶ Height adjusted relationship between ARM and age
  ▶ ARM = B0 + B1 (age – 6) + B2 (HT – 62) + e, e~N(0,σ²), independent

# Review of key concepts from Lecture 3 and 4

▶ Effect modification: Is the ARM vs. age relationship the same or different by sex

  ▶ ARM = B0 + B1 (age – 6) + B2 Female + B3 (age – 6) Female + e, e~N(0,σ²), independent

*Handwritten annotations:*

male: $ARM = B0 + B1(age-6) + e$

female: $ARM = (B0 + B2) + (B1 + B3)(age-6) + e$

$E(\$) = B0 + B1(older) + B2(low income) + B3(older) \times (low income)$

ARM vs. Age plot with lines labeled "female" (slope B1+B3), "males", intercepts B0, B2, at age 6 and 60.

$E(\$)$ chart: B0, B1, B2 — younger higher income, older only, low income only, Both

# Multiple Linear Regression Model

▶ Y is a random variable representing the outcome of interest in the population

$y$ = observations

▶ The explanatory variables, $X_1$, $X_2$, ..., $X_p$ are fixed/known (not random or measured with error)

▶ Sample of size n is observed, data are:

$$(y_i, X_{1i}, X_{2i}, ..., X_{pi})$$

$$\longrightarrow Y_i = \mu_i(\beta, X_i) + \varepsilon_i$$

— random component

systematic

$E(Y_i | X_i)$

$$X = \begin{bmatrix} 1 & X_{11} & X_{21} & \cdots & X_{p1} \\ 1 & X_{12} & X_{22} & & X_{p2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{1n} & X_{2n} & & X_{pn} \end{bmatrix}$$

▶ X is the design matrix

▶ $X_i$ is the row of the design matrix corresponding to subject i

$X_n \leftarrow \begin{bmatrix} 1 & X_{1n} & X_{2n} & X_{pn} \end{bmatrix}$

# Multiple Linear Regression Model

$$Y_i = \mu_i(\beta, X_i) + \varepsilon_i \qquad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}_{P+1 \times 1}$$

▶ Systematic component:

  ▶ $\mu_i(\beta, X_i) = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \ldots + \beta_p X_{pi}$

▶ $\varepsilon_i$ is the random components: $\varepsilon_i \sim N(0, \sigma^2), \quad Cov(\varepsilon_i, \varepsilon_j) = 0$ for all $i \neq j$

▶ The least squares solution finds the values of $\beta$ that minimize:

$$\sum_{i=1}^{n} \left( y_i - \mu_i(\beta, X_i) \right)^2$$

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \beta_1 X_{1i} - \beta_2 X_{2i} - \ldots - \beta_p X_{pi} \right)^2$$

# Least squares solution: simple linear regression

$$SLR = p = 1 \qquad Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (y_i - \bar{y})(X_i - \bar{X})}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$

$$\frac{\text{covariance between } y \text{ and } X}{\text{variance in } X}$$

Standardize $Y$ and $X$ to have mean $0$ and $SD = 1$

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

$$\hat{\beta}_1 = r \quad \text{Pearson correlation coefficient}$$

# Maximum likelihood inference in MLR

▶ Start with the MLR:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_p X_{pi} + \varepsilon_i$$
$$\varepsilon_i \sim N(0, \sigma^2) \text{ and independent}$$
$$\text{Data}: \ (y_i, X_i) \text{ for } i = 1, \dots, n$$

▶ Other notation:

$$Y_i = RV \quad y_i = \text{observation}$$
$$\underset{\sim}{Y} = \text{vector of RVs} \quad \underset{\sim}{y} = \text{vector of observations} \quad \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}$$
$$X_i = i\text{th row of design matrix } X$$
$$\hat{Y}_i = \hat{\beta}_0 + \tilde{\beta}_1 X_{1i} + \dots + \hat{\beta} X_{pi}$$
$$\hat{R}_i = Y_i - \hat{Y}_i$$

# Likelihood function definition

▶ **Model:**

systematic component $\mu_i(\beta, X_i)$ plus random component

$Y_i \sim N\left(\mu_i(\beta, X_i), \sigma^2\right)$

$\varepsilon_i \sim N(0, \sigma^2)$

▶ **Probability density function:**

$$f\left(\underset{\sim}{y} \mid \underset{\sim}{\mu}(\beta, X), \sigma^2\right) = \prod_{i=1}^{n} f\left(y_i \mid \mu_i(\beta, X_i), \sigma^2\right)$$

↳ a function of $\underset{\sim}{y}$ given fixed mean $\mu_i(\beta, X_i)$ and variance $\sigma^2$

▶ **Likelihood function:**

$$L\left(\underset{\sim}{\mu}(\beta, X), \sigma^2 \mid \underset{\sim}{y}\right) = \prod_{i=1}^{n} L\left(\mu_i(\beta, X_i), \sigma^2 \mid y_i\right)$$

↳ a function of $\mu_i(\beta, X_i)$ and $\sigma^2$ given the observed data $y_i$

# Maximum likelihood estimation under gaussian residuals

▶ Likelihood function

$$Y_i \sim \mathcal{N}\left(\mu_i(\beta, X_i), \sigma^2\right)$$

$$L\left(\beta, \sigma^2 \mid y\right) = \prod_{i=1}^{n} L\left(\mu_i(\beta, X_i), \sigma^2 \mid y_i\right)$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}\left(y_i - \mu_i(\beta, X_i)\right)^2\right)$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}\left(y_i - \beta_0 - \beta_1 X_{1i} - \ldots - \beta_p X_{pi}\right)^2\right)$$

# Maximum likelihood estimation under gaussian residuals

▶ Log Likelihood Function

$$\ell\left(\beta, \sigma^2 \mid y\right) = \log L\left(\beta, \sigma^2 \mid y\right)$$

$$= \sum_{i=1}^{n} \left(-\frac{1}{2}\log(2\pi) - \log\sigma - \frac{1}{2\sigma^2}\left(y_i - \mu_i(\beta, X_i)\right)^2\right)$$

$=$ To find $\hat{\beta}$ and $\hat{\sigma}^2$ that maximize $\ell\left(\beta, \sigma^2 \mid y\right)$ w take the derivate w/t $\beta$ and $\sigma^2$, set the derivatives $= 0$ and solve for $\beta$ and $\sigma^2$

# Maximum likelihood estimation under gaussian residuals

▶ Solution for $\beta_j$

$j = 0, \ldots, p$

$$\ell\left(\beta, \sigma^2 \mid y\right) = \sum_{i=1}^{n} \left(-\frac{1}{2}\log(2\pi) - \log\sigma - \frac{1}{2\sigma^2}\left(y_i - \mu_i(\beta, X_i)\right)^2\right)$$

$$\ell\left(\beta, \sigma^2 \mid y\right) \quad \propto \quad \sum_{i=1}^{n} \frac{-1}{2\sigma^2}\left(y_i - \mu_i(\beta, X_i)\right)^2$$

as a function of $\beta$

Define a score equation for $\beta_j$

$$U_{\beta_j}\left(\beta \mid \sigma^2\right) = \frac{\partial}{\partial \beta_j} \ell\left(\beta, \sigma^2 \mid y\right)$$

$$= \frac{\partial}{\partial \beta_j} \sum_{i=1}^{n} \frac{-1}{2\sigma^2}\left(y_i - \beta_0 - \beta_1 X_{1i} - \ldots - \beta_p X_{p_i}\right)^2$$

= 0 and solve for $\beta_j$

$$= \frac{-1}{2\sigma^2} \sum_{i=1}^{n} 2 \times \left(y_i - \mu_i(\beta, X_i)\right)\left(-X_{ij}\right)$$

13

# Maximum likelihood estimation under gaussian residuals
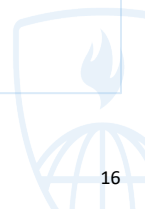
▶ Solution for $\beta_j$

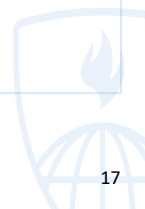# Maximum likelihood estimation under gaussian residuals
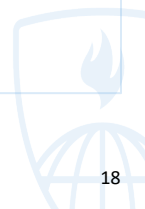
- Solution for $\sigma^2$

# MLEs for simple linear regression

# MLEs for simple linear regression

# MLEs for simple linear regression

# Take away messages

# Take away messages

# Next time....

- ▶ Vector / Matrix representation of MLR

- ▶ Geometry of least squares

- ▶ Distribution of MLEs for regression parameters