# Lecture 5

## The classical linear regression model

► Simple linear regression model

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



systematic component → random component / error

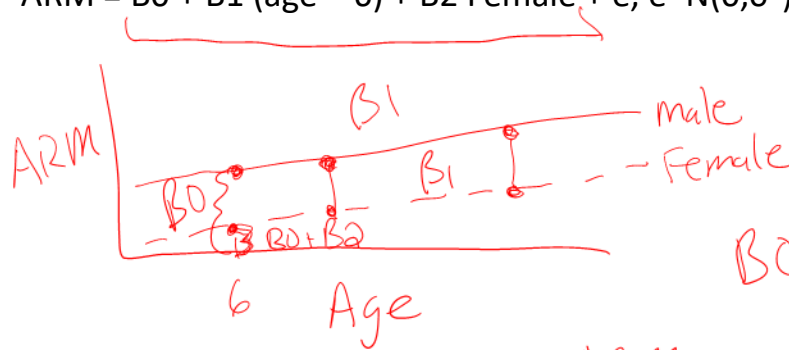E(ARM | age)

BO = average ARM at age = 6 months

B1 = expected difference in ARM comparing two children who differ in age by 1 month

ARM

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,σ²), independent



$$B0 + B1 (age-6) + B2$$
$$(B0 + B2) + B1(age-6)$$

B0 = Average ARM among male children 6 months of age

B1 = difference in mean ARM comparing children of the same gender but whom differ by 1 month of age

B2 = Difference in average ARM comparing female to males at any Age

# 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~$N(0,\sigma^2)$, independent



B1 = difference in the mean ARM comparing children who differ in age by 1 month but whom have the same HT.

► 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,$\sigma^2$), independent

$\beta_1 + \beta_3$ — Female

$B_1$ — male

$B_2$

$B_0$

6   Age

$B1 =$ expected change in mean ARM per month among males

$B1 + B3 =$ " among females

$B3 =$ difference in monthly growth of ARM comparing females to males

# Multiple Linear Regression Model

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

▶ 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: $\left(Y_i, X_{1i}, X_{2i}, ..., X_{pi}\right)$

$$\Rightarrow \quad Y_i = \mu_i(\beta, X_i) + \varepsilon_i \quad \rightarrow \text{random component}$$

outcome random variable

systematic component

▶ X is the design matrix — or table combining all the explanatory variables $[1_s, X_1, X_2, ..., X_p]$

▶ $X_i$ is the row of the design matrix corresponding to subject i — column vector of length n

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

# Multiple Linear Regression Model

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

▶ Systematic component:

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

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

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

$$\sum_{i=1}^{n} \left( Y_i - \underbrace{\mu_i(\beta, X_i)}_{\beta_0 + \beta_1 X_{1i} + \dots + \beta_p X_{pi}} \right)^2$$

Find $\beta$ s.t. the residual sums of squares are minimized

find $\beta_0, \beta_1$ that minimize SSR

# Least squares solution: simple linear regression

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Least squares solution is:

$$\hat{\beta}_1 = \frac{\sum (Y_i - \bar{Y})(X_i - \bar{X})}{\sum (X_i - \bar{X})^2} \quad \Bigg\} \quad \frac{\text{Covariance between } Y \text{ and } X}{\text{Variance of } X}$$

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

If we were to standardize $Y$ and $X$,

$$\text{mean}(Y) = \text{mean}(X) = 0$$

$$SD(Y) = SD(X) = 1$$

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

# Maximum likelihood inference in MLR

▶ Start with the MLR:

$$\Rightarrow Y_i = \underbrace{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{pi}} + \varepsilon_i \quad, \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$\text{independent}$$

$$\text{Data:} \quad (y_i, X_i)$$

▶ Other notation:

$$Y_i = \text{Random variable} \quad, \quad y_i = \text{observation}$$

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \quad Y, \quad y = \text{vectors (lists) of independent random variables or observations} \quad \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

$$X_i = \text{row vector containing explanatory variables for subject } i$$

$$\hat{Y}_i = RV = \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_2 X_{pi}$$

$$R_i = RV = Y_i - \hat{Y}_i \quad \text{residual}$$

# Likelihood function definition

▶ **Model:** under the assumption $e_i \sim N(0, \sigma^2)$, $Y_i = RV$, $X_i = fixed$

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

▶ **Probability density function:**

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

pdf is a function of $\underset{\sim}{y}$ with $\mu_i(\beta, X_i)$ and $\sigma^2$ fixed

▶ **Likelihood function:**

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

likelihood function is viewed as a
function of $\mu_i(\beta, X_i)$ and $\sigma^2$ for fixed $y_i$
Identify the values of $\beta$ and $\sigma^2$ that maximize
$L$ given fixed $\underset{\sim}{y}$

# Maximum likelihood estimation under gaussian residuals

▶ Likelihood function , rely on the normality assumption and independence

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

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

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}\left(y_i - \beta_0 - \beta_1 X_{1i} - \beta_2 X_{2i} - \ldots - \beta_p X_{p_i}\right)^2\right)$$

# Maximum likelihood estimation under gaussian residuals

▶ Log Likelihood Function

$$\ell\left(\hat{\beta}, \sigma^2 \mid \underset{\sim}{y}\right) = Log \, L\left(\hat{\beta}, \sigma^2 \mid \underset{\sim}{y}\right)$$

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

Find $\hat{\beta}$ and $\sigma^2$ that maximize $\ell\left(\hat{\beta}, \sigma^2 \mid \underset{\sim}{y}\right)$ by differentiating with respect to $\hat{\beta}$ and $\sigma^2$, setting these derivatives $= 0$ and solving for $\hat{\beta}$ and $\sigma^2$

# Maximum likelihood estimation under gaussian residuals

▶ Solution for $\beta_j$

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

$j = 0, \cdots, p$

wrt $\beta$

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

Score equation for $\beta_j$

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

$$= \frac{d}{d\beta_j} \sum_{i=1}^{n} \frac{-1}{2\sigma^2}\left(y_i - (\beta_0 + \beta_1 X_{1i} + \cdots + \beta_p X_{pi})\right)^2$$

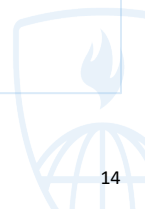$$= \frac{-2}{2\sigma^2} \sum_{i=1}^{n} \left(y_i - \mu_i(\beta, X_i)\right)\left(-X_{ij}\right)$$

↳ a score equation for each $\beta$ ⟹ p+1 score equations

⟹ 0 solve for $\beta_j$   p+1 unknowns

13

# Maximum likelihood estimation under gaussian residuals

▶ Solution for $\beta_j$

$$U_\beta = \sum_{i=1}^{n} \left( y_i - \mu_i(\beta, X_i) \right) \begin{bmatrix} 1 \\ X_{1i} \\ X_{2i} \\ \vdots \\ X_{pi} \end{bmatrix} = 0$$

$$p+1 \times 1$$

$$p+1 \quad \times 1$$

# Maximum likelihood estimation under gaussian residuals

▶ Solution for $\sigma^2$

Given the MLEs $\hat{\beta}$, derive score equation for $\sigma^2$

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

$\Rightarrow 0$ and solve for $\sigma^2$

$$\Rightarrow \quad \hat{\sigma}^2_{MLE} = \frac{1}{n}\sum_{i=1}^{n}\left(y_i - \mu_i\left(\hat{\beta}, X_i\right)\right)^2$$

In practice, we use

$$E\left(\hat{\sigma}^2_{MLE}\right) = \frac{n-(p+1)}{n}\sigma^2 \qquad \tilde{\sigma}^2 = \frac{1}{n-(p+1)}\sum\left(y_i - \mu_i\left(\hat{\beta}, X_i\right)\right)^2$$
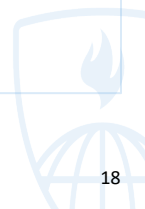
# MLEs for simple linear regression

# 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