

Properties of Least Square Estimates

DS 301

Iowa State University

Today's Agenda

- Implementation of multiple linear regression in R
- Properties of least square estimates

Multiple Linear Regression (Recap)

$$Y = f(X) + \varepsilon$$

• If we assume a linear relationship btwn X & Y

$$f(X) = B_0 + B_1 X_1 + B_2 X_2 + \dots + B_p X_p.$$

\downarrow
 X_1, X_2, \dots, X_p

$$\hookrightarrow Y = B_0 + B_1 X_1 + B_2 X_2 + \dots + B_p X_p + \varepsilon$$

(true population regression line)

$B_0, B_1, B_2, \dots, B_p$ are unknown parameters.

\Rightarrow estimate them from (training) data

$$\hat{B}_0, \hat{B}_1, \hat{B}_2, \dots, \hat{B}_p$$

(least squares criterion) $\left\{ \sum_{i=1}^n (y_i - \hat{y}_i)^2 : \text{RSS} \right.$

$$\hat{y}_i = \hat{B}_0 + \hat{B}_1 X_{i1} + \hat{B}_2 X_{i2} + \dots + \hat{B}_p X_{ip}$$

\uparrow predicted value / fitted value

let's find \hat{y}_i that minimizes RSS

Interpretation of least squares coefficients

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_p X_{ip}$$

is an estimate for $E(Y)$, not Y .

$\hat{\beta}_j$ can be interpreted as the average change in Y associated with a 1 unit change in X_j , *holding all other predictors constant*.

Assumptions

(1) There is a linear relationship
btwn X_1, X_2, \dots, X_p and Y

(2) $E(\varepsilon_i) = 0$

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

(4) ε_i 's are uncorrelated

See R script IntroMLR.R

Properties of least square estimators

- Remember in real applications, the true parameters $\beta_0, \beta_1, \dots, \beta_p$ are unknown to us.
- Ideally, we hope our least square estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ are close to the true values of $\beta_0, \beta_1, \dots, \beta_p$.
- We can quantify how 'close' our estimates are to the truth using the following concepts:
 - Bias
 - Standard error

Properties of least square estimators

How good are our estimates?

The least sq. estimates $\hat{B}_0, \hat{B}_1, \dots, \hat{B}_p$
are unbiased estimates of B_0, B_1, \dots, B_p
(respectively).

conceptually, the property of unbiasedness
says that if we took the average of $\hat{B}_0, \hat{B}_1, \dots, \hat{B}_p$
obtained over a huge number of
datasets, then these averages would
exactly equal the true parameters B_0, B_1, \dots, B_p

formal definition:
 $E(\hat{B}_0) = B_0, E(\hat{B}_1) = B_1,$
 $\dots E(\hat{B}_p) = B_p$

Properties of least square estimators

$$\hat{B}_0, \hat{B}_1, \dots, \hat{B}_p \Rightarrow \hat{Y}$$

\hat{Y} is an unbiased estimate for $E(Y)$.

$$\begin{aligned} E(\hat{Y}) &= E(\hat{B}_0 + \hat{B}_1 X_1 + \dots + \hat{B}_p X_p) \\ &= B_0 + B_1 X_1 + \dots + B_p X_p = E(Y). \end{aligned}$$

$$E(\hat{Y}) = E(Y)$$

↳ on average, our predictions (\hat{Y}) will equal the true mean of Y : $E(Y)$

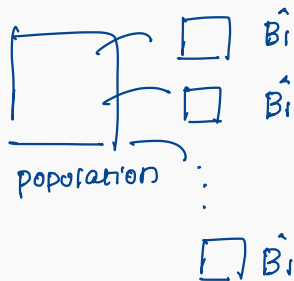
$$\hookrightarrow f(x)$$

⇒ This gives us some sense that we can trust our model.

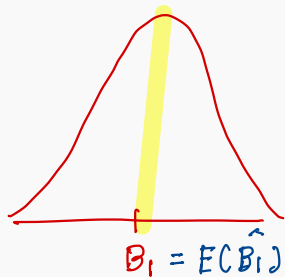
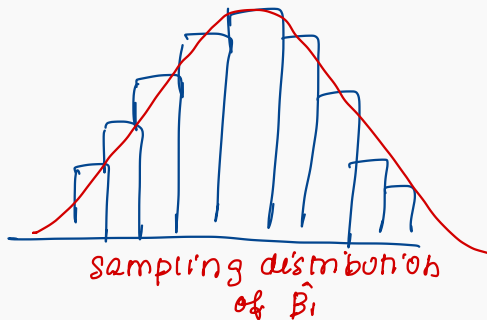
Accuracy of $\hat{\beta}$

- The unbiasedness property tell us that the average of our estimates from many many datasets will be very close to the true population parameter β .
- But we don't have access to many many datasets. For a particular data set, the single estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ may be a substantial underestimate or overestimate of the true $\beta_0, \beta_1, \dots, \beta_p$.
- How far off will our single estimate $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ be?

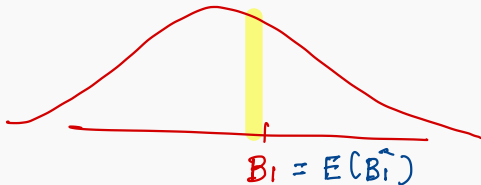
Accuracy of $\hat{\beta}$ (standard error)



$n=10,000$



$n=10,000$



Accuracy of $\hat{\beta}$: How far off is a single estimate?

\Rightarrow formalize this with the standard error of $\hat{B}_0, \hat{B}_1, \dots, \hat{B}_p$.

$$se(\hat{B}_0) = \sqrt{\sigma^2 / \sum (x_i - \bar{x})^2}$$

$$se(\hat{B}_1) = \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)}$$

$$\begin{matrix} \vdots & \hat{B}_2 \\ & \vdots \\ \vdots & \hat{B}_p \end{matrix}$$

σ^2 is a parameter. It represents the variance of γ .

$$Y = B_0 + B_1 X_1 + \epsilon$$

$$\text{var}(Y) = \text{var}(B_0 + B_1 X_1 + \epsilon)$$

$$= \underbrace{\text{var}(B_0)}_0 + \underbrace{\text{var}(B_1 X_1)}_0 + \underbrace{\text{var}(\epsilon)}_{\sigma^2}$$

$\hookrightarrow \sigma^2$ needs to be estimated from data.