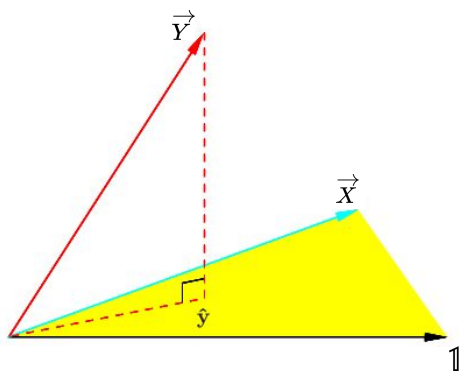


# Sampling distributions of $\hat{\beta}_0$ and $\hat{\beta}_1$

*14.2 of Rice*

08/04/2021

In the previous lecture,



- Simple linear regression (SLR):

$$Y_i = \underbrace{\beta_0} + \underbrace{\beta_1 X_i} + \underbrace{\epsilon_i}_{\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)}.$$

- Method of least squares;
- Maximum likelihood estimation;
- Projections onto the hyperplane  $\text{span}(1, \vec{X})$ .

$$\Rightarrow \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- Check model assumptions:

- Linearity - Plotting  $Y_i$  vs  $X_i$
- Normality - QQ plot
- Zero mean in error terms
- Homoscedasticity
- Independence

} Residual plot

## Mean and variance

$$E(\bar{Y}_n) = \frac{E(Y_1) + \dots + E(Y_n)}{n} = \frac{\sum_{i=1}^n \beta_0 + \beta_1 X_i}{n} = \beta_0 + \beta_1 \bar{X}_n$$

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \sim N \left( \begin{pmatrix} \beta_0 + \beta_1 X_1 \\ \vdots \\ \beta_0 + \beta_1 X_n \end{pmatrix}, \sigma^2 I \right)$$

**Proposition D.** Under the SLR assumptions,  $(\hat{\beta}_0, \hat{\beta}_1)$  is bivariate Normal with:

$$E(\hat{\beta}_0) = \beta_0, \text{var}(\hat{\beta}_0) = \frac{n^{-1} \sum_i X_i^2}{\sum_i (X_i - \bar{X}_n)^2} \sigma^2,$$

$$E(\hat{\beta}_1) = \beta_1, \text{var}(\hat{\beta}_1) = \frac{1}{\sum_i (X_i - \bar{X}_n)^2} \sigma^2,$$

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\bar{X}_n}{\sum_i (X_i - \bar{X}_n)^2} \sigma^2.$$

Proof.  $\hat{\beta}_0 = \bar{Y}_n - \hat{\beta}_1 \bar{X}_n = \frac{\sum_i X_i^2 \sum_i Y_i - \sum_i X_i \sum_i X_i Y_i}{n \sum_i X_i^2 - (\sum_i X_i)^2} = \frac{\sum_i [\sum_i X_i^2 - X_i \sum_i X_i] Y_i}{n \sum_i X_i^2 - (\sum_i X_i)^2} = \sum_i b_i Y_i$

$$\hat{\beta}_1 = \frac{\sum_i (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sum_i (X_i - \bar{X}_n)^2} = \frac{n \sum_i X_i Y_i - \sum_i X_i \sum_i Y_i}{n \sum_i X_i^2 - (\sum_i X_i)^2} = \frac{\sum_i [n X_i - \sum_i X_i] Y_i}{n \sum_i X_i^2 - (\sum_i X_i)^2} = \sum_i a_i Y_i$$

$$E(\hat{\beta}_1) = E\left(\sum_i a_i Y_i\right) = \sum_i a_i E(Y_i) = \sum_i \frac{n X_i - \sum_i X_i}{n \sum_i X_i^2 - (\sum_i X_i)^2} \times (\beta_0 + \beta_1 X_i)$$

$$= \beta_0 \underbrace{\sum_i \frac{n X_i - \sum_i X_i}{n \sum_i X_i^2 - (\sum_i X_i)^2}}_{=0} + \beta_1 \underbrace{\sum_i \frac{n X_i^2 - X_i \sum_i X_i}{n \sum_i X_i^2 - (\sum_i X_i)^2}}_{=1} = \beta_1$$

$$\begin{aligned} E(\hat{\beta}_0) &= E(\bar{Y}_n + \hat{\beta}_1 \bar{X}_n) \\ &= \beta_0 - \beta_1 \bar{X}_n + \beta_1 \bar{X}_n \\ &= \beta_0 \end{aligned}$$

## Mean and variance

$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x}_n) (y_i - \bar{y}_n)}{\sum_i (x_i - \bar{x}_n)^2} = \frac{\sum_i (x_i - \bar{x}_n) y_i - \bar{y}_n \underbrace{\sum_i (x_i - \bar{x}_n)}_{=0}}{\sum_i (x_i - \bar{x}_n)^2} = \frac{\sum_i (x_i - \bar{x}_n) y_i}{\sum_i (x_i - \bar{x}_n)^2}$$

Proof cont'd.  $\text{var}(\hat{\beta}_1) = \frac{\sum_i (x_i - \bar{x}_n)^2 \text{var}(y_i)}{\left[ \sum_i (x_i - \bar{x}_n)^2 \right]^2} = \frac{\cancel{\sum_i (x_i - \bar{x}_n)^2} b^2}{\left[ \sum_i (x_i - \bar{x}_n)^2 \right]} = \frac{1}{\sum_i (x_i - \bar{x}_n)^2} b^2$

$$\text{var}(\hat{\beta}_0) = \text{var}(\bar{y}_n - \hat{\beta}_1 \bar{x}_n)$$

$$= \text{var}(\bar{y}_n) + \bar{x}_n^2 \text{var}(\hat{\beta}_1) - 2\bar{x}_n \text{cov}(\bar{y}_n, \hat{\beta}_1)$$

$$= \frac{b^2}{n} + \frac{\bar{x}_n^2}{\sum_i (x_i - \bar{x}_n)^2} b^2 - 2\bar{x}_n \text{cov}\left(\frac{\sum_i y_i}{n}, \frac{\sum_i (x_i - \bar{x}_n) y_i}{\sum_i (x_i - \bar{x}_n)^2}\right)$$

$$= \frac{b^2}{n} + \frac{\bar{x}_n^2}{\sum_i (x_i - \bar{x}_n)^2} b^2 - 2\bar{x}_n \sum_{i=1}^n \frac{1}{n} \cdot \frac{x_i - \bar{x}_n}{\sum_i (x_i - \bar{x}_n)^2} \underbrace{\text{var}(y_i)}$$

$$= \frac{b^2}{n} + \frac{\bar{x}_n^2}{\sum_i (x_i - \bar{x}_n)^2} b^2 - \frac{2\bar{x}_n b^2}{n} \sum_{i=1}^n \frac{x_i - \bar{x}_n}{\sum_i (x_i - \bar{x}_n)^2} = 0$$

$$= \frac{\frac{1}{n} \sum_i x_i^2 - \bar{x}_n^2 + \bar{x}_n^2}{\sum_i (x_i - \bar{x}_n)^2} b^2 = \frac{n^{-1} \sum_i x_i^2}{\sum_i (x_i - \bar{x}_n)^2} b^2$$

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \text{cov}\left(\sum_i b_i Y_i, \sum_i a_i Y_i\right) = \sum_i a_i b_i \text{var}(Y_i) = b^2 \sum_i a_i b_i$$

$$= \frac{-\bar{X}_n}{\sum (X_i - \bar{X}_n)^2} b^2$$

$$Y \sim N(\mu, \Sigma)$$

$$AY \sim N(A\mu, A\Sigma A^T)$$

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \sim N\left(\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \begin{pmatrix} \text{var}(\hat{\beta}_0) & \text{cov}(\hat{\beta}_0, \hat{\beta}_1) \\ \text{cov}(\hat{\beta}_0, \hat{\beta}_1) & \text{var}(\hat{\beta}_1) \end{pmatrix}\right)$$

recall that  $\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \underbrace{(X^T X)^{-1} X^T \vec{Y}}_{\text{in which } X = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{pmatrix}, \vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}}$

$$\text{cov}\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = (X^T X)^{-1} X^T (b^2 I) X (X^T X)^{-1}$$

$$= b^2 (X^T X)^{-1} \cancel{X^T X (X^T X)^{-1}}$$

$$= b^2 (X^T X)^{-1}$$

$H = X(X^T X)^{-1} X^T$  projection matrix  $\hat{\epsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i$   
 $= Y_i - (2, X_i) \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \Rightarrow \begin{pmatrix} \hat{\epsilon}_1 \\ \vdots \\ \hat{\epsilon}_n \end{pmatrix} = \vec{Y} - X \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}$

## Sampling distribution

**Lemma D.** Under the SLR assumptions,  $\vec{Y} \sim N(\mu, \Sigma)$ ,  $\text{cov}(A\vec{Y}, B\vec{Y}) = A \Sigma B^T$ .

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \perp\!\!\!\perp \vec{Y} - X \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix},$$

and  $\text{RSS} = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \sim \chi_{n-2}^2$ .

1. A matrix  $R$  is idempotent if  $R^2 = R$ .
2. If  $Z \sim N(0, I)$  and  $R$  is symmetric and idempotent of rank  $r$ , then  $Z^T R Z \sim \chi_r^2$ .

Proof\*.  $\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = (X^T X)^{-1} X^T \vec{Y}$

$$\vec{Y} - X \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \vec{Y} - X (X^T X)^{-1} X^T \vec{Y} = \underbrace{(I - X(X^T X)^{-1} X^T)}_{H} \vec{Y}.$$

$$\begin{aligned} \text{Cov} \left[ (X^T X)^{-1} X^T \vec{Y}, (I - H) \vec{Y} \right] &= b^2 (X^T X)^{-1} X^T (I - H)^T \\ &= b^2 \underbrace{(X^T X)^{-1} X^T} \underbrace{(I - X(X^T X)^{-1} X^T)} \\ &= b^2 \left[ \underbrace{(X^T X)^{-1} X^T} - \underbrace{(X^T X)^{-1} X^T X}_{I} \underbrace{(X^T X)^{-1} X^T} \right] \\ &= 0. \end{aligned}$$

By definition,

$$RSS = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = \sum_{i=1}^n \hat{\epsilon}_i^2 = (\hat{\epsilon}_1, \dots, \hat{\epsilon}_n) \begin{pmatrix} \hat{\epsilon}_1 \\ \vdots \\ \hat{\epsilon}_n \end{pmatrix}$$

$$= [(\mathbf{I} - \mathbf{H})\vec{Y}]^T (\mathbf{I} - \mathbf{H})\vec{Y} = \vec{Y}^T \underbrace{(\mathbf{I} - \mathbf{H})^T (\mathbf{I} - \mathbf{H})}_{\mathbf{I} - \mathbf{H}} \vec{Y} = \vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}.$$

$$(\mathbf{I} - \mathbf{H})^T (\mathbf{I} - \mathbf{H}) = (\mathbf{I} - \mathbf{H})^2 = (\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) (\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)$$

$$= \mathbf{I} - \underbrace{2\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T} + \underbrace{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \cancel{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}}_{\substack{\Downarrow \\ H^2 = H}}$$

$$= \mathbf{I} - \mathbf{H}$$

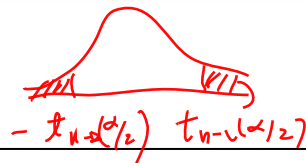
$$\text{rank}(\mathbf{I} - \mathbf{H}) = \text{trace}(\mathbf{I} - \mathbf{H}) = n - \text{trace}(\mathbf{H}) = n - \text{rank}(\mathbf{H}) \xrightarrow{\substack{\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \parallel \\ \text{rank}(\mathbf{X})}} n - \text{rank}(\mathbf{X}) = n - 2.$$

$$RSS = \vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y} \sim \underline{\underline{b^2}} \chi_{n-2}^2 \quad \text{and} \quad \underbrace{\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}}_{\perp RSS}.$$

$$\downarrow$$

$$\vec{Y} \sim N(\mathbf{X}\vec{\beta}, b^2 \mathbf{I})$$

# Sampling distribution



**Theorem D.** Under the SLR assumptions,

$$\frac{\hat{\beta}_k - \beta_k}{\text{se}(\hat{\beta}_k)} \sim t_{n-2}, k = 0, 1.$$

95% exact confidence interval for  $\beta_k$ :



$$\hat{\beta}_k \pm t_{n-2}(\alpha/2) \cdot \text{se}(\hat{\beta}_k)$$

$$\hat{\beta}_0 \sim N\left(\beta_0, \frac{n^{-1} \sum_i x_i^2}{\sum_i (x_i - \bar{x}_n)^2} b^2\right)$$

$$RSS \sim b^2 \chi_{n-2}^2$$

$$\hat{\sigma}^2 = \frac{RSS}{n-2} \Rightarrow E(\hat{\sigma}^2) = b^2.$$

$$\Rightarrow \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\frac{n^{-1} \sum_i x_i^2}{\sum_i (x_i - \bar{x}_n)^2} \hat{\sigma}^2}} \sim t_{n-2}$$

$$\Rightarrow \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\frac{n^{-1} \sum_i x_i^2}{\sum_i (x_i - \bar{x}_n)^2} \hat{\sigma}^2}} \sim t_{n-2}$$

$$RSS = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$\hat{\sigma}^2 = \frac{RSS}{n-2}$$

$$\text{se}(\hat{\beta}_0) = \sqrt{\frac{n^{-1} \sum_i X_i^2}{\sum_i (X_i - \bar{X}_n)^2} \hat{\sigma}^2}$$

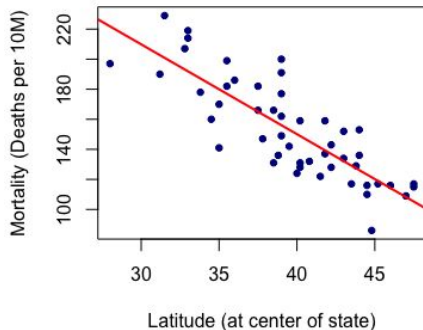
$$\text{se}(\hat{\beta}_1) = \sqrt{\frac{1}{\sum_i (X_i - \bar{X}_n)^2} \hat{\sigma}^2}$$

$\sim t_{n-2}$  Similarly for  $\hat{\beta}_1$ .



## CIs of $\beta_0$ and $\beta_1$

**Example 2 cont'd.** During the 50s, data were collected to examine the relationship between the mortality rate due to skin cancer (number of deaths per 10 million people) and the latitude at the center of each of 48 states in the United States (Alaska and Hawaii were not yet states. And, Washington, D.C. was included in the data set even though it is not technically a state.)



```
> beta1 <- cov(Mort, Lat)/var(Lat)
> beta0 <- mean(Mort)-beta1*mean(Lat)
> c('beta0' = beta0, 'beta1' = beta1)
      beta0      beta1
389.189351 -5.977636
```

```
> n <- length(Mort)
> RSS <- sum(Res^2)
> s2_hat <- RSS/(n-2)
> se_beta0 <- sqrt(mean(Lat^2)*s2_hat/((n-1)*var(Lat)))
> se_beta1 <- sqrt(s2_hat/((n-1)*var(Lat)))
>
> alpha <- 0.05
> t_alpha <- qt(alpha/2, n-2, lower.tail = FALSE)
> c(beta0 - t_alpha*se_beta0, beta0 + t_alpha*se_beta0)
[1] 341.2852 437.0936
> c(beta1 - t_alpha*se_beta1, beta1 + t_alpha*se_beta1)
[1] -7.181404 -4.773867
```

$$RSS = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$\hat{\sigma}^2 = \frac{RSS}{n-2}$$

$$se(\hat{\beta}_0) = \sqrt{\frac{n^{-1} \sum_i X_i^2}{\sum_i (X_i - \bar{X}_n)^2} \hat{\sigma}^2}$$

$$se(\hat{\beta}_1) = \sqrt{\frac{1}{\sum_i (X_i - \bar{X}_n)^2} \hat{\sigma}^2}$$

```
> confint(fit, 1, level = 0.95)
              2.5 %    97.5 %
(Intercept) 341.2852 437.0936
> confint(fit, 'Lat', level = 0.95)
              2.5 %    97.5 %
Lat -7.181404 -4.773867
```

## SE of the regression line

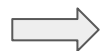
$$\beta_0 + \beta_1 x$$

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \perp \text{RSS}$$

**Theorem E.** Under the SLR assumptions, we can use  $(\hat{\beta}_0, \hat{\beta}_1)$  to predict at any  $X = x$  :

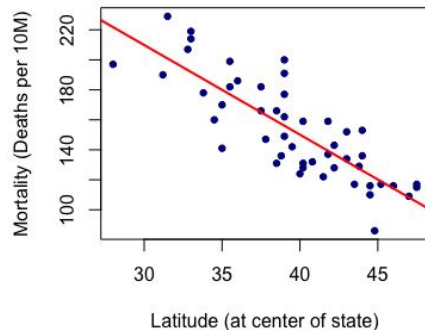
$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x,$$

Then  $E(\hat{y}) = \beta_0 + \beta_1 x$ ,  $\text{var}(\hat{y}) = \sigma^2 \left[ \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \right]$  and  $\hat{y}$  is independent of RSS.



95% exact confidence interval for  $\beta_0 + \beta_1 x$  :

$$\hat{y} \pm t_{n-2}(\alpha/2) \cdot \text{se}(\hat{y})$$

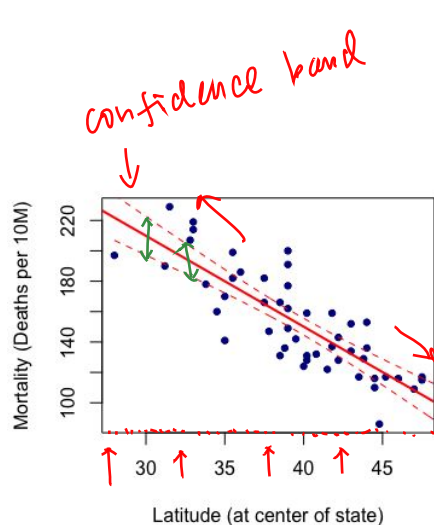


$$\begin{aligned} \text{var}(\hat{y}) &= \text{var}(\hat{\beta}_0 + \hat{\beta}_1 x) \\ &= \text{var}(\hat{\beta}_0) + \text{var}(\hat{\beta}_1) x^2 + 2x \text{cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= b^2 \left[ \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \right] \end{aligned}$$

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \\ \hat{\sigma}^2 &= \frac{\text{RSS}}{n-2} \\ \text{se}(\hat{y}) &= \sqrt{\left[ \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_i (X_i - \bar{X}_n)^2} \right] \hat{\sigma}^2} \end{aligned}$$

## CIs of the population mean

**Example 2** *cont'd*. During the 50s, data were collected to examine the relationship between the mortality rate due to skin cancer (number of deaths per 10 million people) and the latitude at the center of each of 48 states in the United States (Alaska and Hawaii were not yet states. And, Washington, D.C. was included in the data set even though it is not technically a state.)



```
x_vec <- seq(28, 49, by=0.1) ←  
fit_vec <- beta0 + beta1*x_vec ←  
denom <- {(n-1)*var(Lat)} ←  
se_vec <- sqrt({1/n + (x_vec-mean(Lat))^2/denom}*s2_hat)  
  
lower_vec <- fit_vec - t_alpha*se_vec ←  
upper_vec <- fit_vec + t_alpha*se_vec ←  
lines(x_vec, lower_vec, lty=2, col='red') ←  
lines(x_vec, upper_vec, lty=2, col='red')
```

```
new_dat <- data.frame(Lat=x_vec) ←  
CIs <- predict(fit, newdata = new_dat,  
               interval = 'confidence', level = 0.95) ←  
lines(x_vec, CIs[, 'lwr'], lty=2, col='red') ←  
lines(x_vec, CIs[, 'upr'], lty=2, col='red')
```

$$\text{RSS} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$
$$\hat{\sigma}^2 = \frac{\text{RSS}}{n-2}$$
$$\text{se}(\hat{y}) = \sqrt{\left[ \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_i (X_i - \bar{X}_n)^2} \right] \hat{\sigma}^2}$$

# Prediction interval of the response

$\mu, b^2, \dots$

$$\frac{y - \hat{y}}{\sqrt{b^2 + \text{var}(\hat{y})}} \sim N(0, 1)$$

**Theorem F.** Under the SLR assumptions, the new observation at  $X = x$  is a random variable:

$$y = \beta_0 + \beta_1 x + \epsilon, \epsilon \sim N(0, \sigma^2).$$

$$RSS/b^2 \sim \chi_{n-2}^2$$

Then  $E(y - \hat{y}) = 0$ ,  $\text{var}(y - \hat{y}) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \right]$  and  $y$  is independent of  $\hat{y}$  and RSS.

95% ~~exact confidence interval for  $\beta_0 + \beta_1 x$ :~~ prediction interval for the response  $y$ :

$$\hat{y} \pm t_{n-2}(\alpha/2) \cdot \text{se}(y - \hat{y})$$

$$y \sim N(\beta_0 + \beta_1 x, b^2) \Leftarrow$$

$$\hat{y} \sim N(\beta_0 + \beta_1 x, \text{var}(\hat{y})) \Leftarrow$$

indep  $\Rightarrow y - \hat{y} \sim N(0, b^2 + \text{var}(\hat{y}))$

$$y \perp \hat{y}$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$\uparrow \quad \uparrow$   
 $\begin{pmatrix} Y_1 & X_1 \\ \vdots & \vdots \\ Y_n & X_n \end{pmatrix}$

$$RSS = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \Leftarrow$$

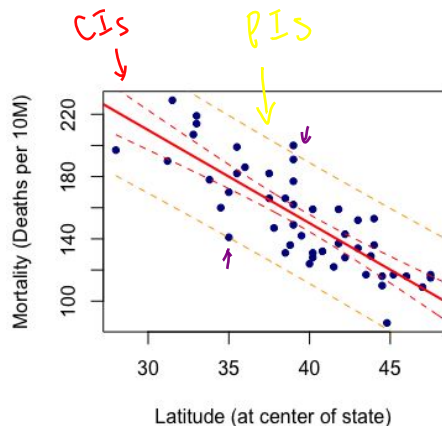
$$\hat{\sigma}^2 = \frac{RSS}{n-2}$$

$$\text{se}(y - \hat{y}) = \sqrt{\left[ 1 + \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_i (X_i - \bar{X}_n)^2} \right] \hat{\sigma}^2}$$

More uncertainty

## Prediction interval of the response

**Example 2 cont'd.** During the 50s, data were collected to examine the relationship between the mortality rate due to skin cancer (number of deaths per 10 million people) and the latitude at the center of each of 48 states in the United States (Alaska and Hawaii were not yet states. And, Washington, D.C. was included in the data set even though it is not technically a state.)



```
x_vec <- seq(28, 49, by=0.1)
fit_vec <- beta0 + beta1*x_vec
denom <- {(n-1)*var(Lat)}
se_pred_vec <- sqrt({1+1/n + (x_vec-mean(Lat))^2/denom}*s2_hat)

lower_pred_vec <- fit_vec - t_alpha*se_pred_vec
upper_pred_vec <- fit_vec + t_alpha*se_pred_vec
lines(x_vec, lower_pred_vec, lty=2, col='orange')
lines(x_vec, upper_pred_vec, lty=2, col='orange')
```

```
new_dat <- data.frame(Lat=x_vec)
CIs_pred <- predict(fit, newdata = new_dat,
  interval = 'prediction', level = 0.95)
lines(x_vec, CIs_pred[, 'lwr'], lty=2, col='orange')
lines(x_vec, CIs_pred[, 'upr'], lty=2, col='orange')
```

$$\text{RSS} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$
$$\hat{\sigma}^2 = \frac{\text{RSS}}{n-2}$$
$$\text{se}(y - \hat{y}) = \sqrt{\left[1 + \frac{1}{n} + \frac{(x - \bar{X}_n)^2}{\sum_i (X_i - \bar{X}_n)^2}\right] \hat{\sigma}^2}$$

# Multiple linear regression

*14.3 of Rice*

08/04/2021

# Multiple linear regression (MLR)

$$\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \rightarrow \vec{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

Model assumption ( $i = 1, \dots, n$ ):

Response variable	Predictor variable 1	Predictor variable 2	Predictor variable $p$	
$Y_1$	$X_{11}$	$X_{12}$	$\cdots$	$X_{1p}$
$Y_2$	$X_{21}$	$X_{22}$	$\cdots$	$X_{2p}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$Y_n$	$X_{n1}$	$X_{n2}$	$\cdots$	$X_{np}$



$$Y_i = \underbrace{\beta_0}_{\text{common mean level}} + \underbrace{\beta_1 X_{i1} + \cdots + \beta_p X_{ip}} + \underbrace{\epsilon_i}_{\substack{\text{iid} \\ \sim N(0, \sigma^2)}} .$$

1. Linearity - Plotting  $\vec{Y}$  vs  $\vec{X}_j$
  2. Normality - QQ plot
  3. Zero mean in error terms
  4. Homoscedasticity
  5. Independence
- } Residual plot

# Method of least squares

$$f(\vec{\beta}) \quad \frac{\partial f}{\partial \vec{\beta}} = \begin{pmatrix} \frac{\partial f}{\partial \beta_1} \\ \vdots \\ \frac{\partial f}{\partial \beta_p} \end{pmatrix}$$

$$f(\vec{\beta}) = \vec{a}^T \vec{\beta} \quad , \quad \frac{\partial f}{\partial \vec{\beta}} = \vec{a}$$

$$f(\vec{\beta}) = \vec{\beta}^T A \vec{\beta} \quad , \quad \frac{\partial f}{\partial \vec{\beta}} = (A^T + A) \vec{\beta}$$

**Proposition A'.** Under the SLR assumptions, find estimators for  $\vec{\beta}$  such that they minimize the sum of squared vertical deviations:

$$S(\vec{\beta}) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2.$$

$X^T X$  is invertible as long as  $X$  is full rank,

i.e.  $\text{rank}(X) = p+1$ .

$$X = \begin{pmatrix} 1 & X_{11} & \dots & X_{1p} \\ 1 & X_{21} & \dots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \dots & X_{np} \end{pmatrix} \leftarrow i\text{th row vector } \vec{X}^i$$

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i \\ &= (1, X_{i1}, \dots, X_{ip}) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \epsilon_i \end{aligned}$$

$$\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = X \vec{\beta} + \vec{\epsilon}$$

$$S(\vec{\beta}) = \sum_{i=1}^n [Y_i - \vec{X}^i \vec{\beta}]^2$$

$$= (Y_1 - \vec{X}^1 \vec{\beta}, \dots, Y_n - \vec{X}^n \vec{\beta})$$

$$\begin{pmatrix} Y_1 - \vec{X}^1 \vec{\beta} \\ \vdots \\ Y_n - \vec{X}^n \vec{\beta} \end{pmatrix}$$

$$= (\vec{Y} - X \vec{\beta})^T (\vec{Y} - X \vec{\beta}) \quad \vec{Y} - X \vec{\beta}$$

$$= \vec{Y}^T \vec{Y} - 2 \vec{Y}^T X \vec{\beta} + \vec{\beta}^T X^T X \vec{\beta}$$

$$\begin{aligned} \frac{\partial S}{\partial \vec{\beta}} &= -2 X^T \vec{Y} + (X^T X + X^T X) \vec{\beta} \Leftrightarrow X^T X \vec{\beta} = X^T \vec{Y} \\ &= -2 X^T \vec{Y} + 2 X^T X \vec{\beta} = 0 \Leftrightarrow \vec{\beta} = (X^T X)^{-1} X^T \vec{Y} \end{aligned}$$



# Maximum likelihood estimation

**Proposition B'.** Under the SLR assumptions, calculate  $\sup_{\theta} L(\vec{\beta}, \sigma^2)$  and find MLEs of  $\vec{\beta}$  and  $\sigma^2$ .

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{pmatrix}$$

$$L(\vec{\beta}, b^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi b^2}} e^{-\frac{(Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_p X_{ip})^2}{2b^2}}$$

$$= \left( \frac{1}{\sqrt{2\pi b^2}} \right)^n e^{-\frac{\sum_{i=1}^n (Y_i - \beta_0 - \cdots - \beta_p X_{ip})^2}{2b^2}}$$

$$\ell(\vec{\beta}, b^2) = -\frac{n}{2} \log(2\pi b^2) - \frac{1}{2b^2} \sum_{i=1}^n (Y_i - \beta_0 - \cdots - \beta_p X_{ip})^2$$

$$\hat{b}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \cdots - \hat{\beta}_p X_{ip})^2 \quad \frac{S(\vec{\beta})}{n}$$

## Geometric approach

**Proposition C'.** We generalize to any  $n \geq p+1$ :

$$\vec{Y} = \mathbf{X} \vec{\beta} + \vec{\epsilon}.$$

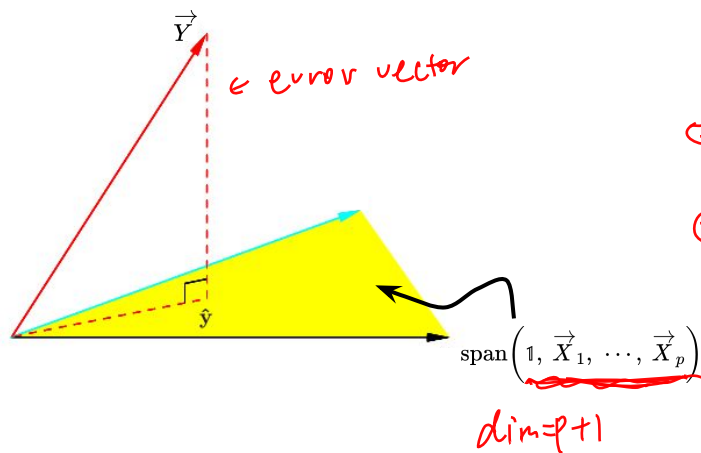
Find the best values of  $\vec{\beta}$ .

To ensure orthogonality,

$$\mathbf{X}^T (\vec{Y} - \mathbf{X} \vec{\beta}) = 0$$

$$\Leftrightarrow \mathbf{X}^T \vec{Y} - \mathbf{X}^T \mathbf{X} \vec{\beta} = 0$$

$$\Leftrightarrow \vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{Y}.$$



# Calculate in R

**Example 3.** Interested in answering whether person's brain size and body size predictive of his or her intelligence, some researchers (Willerman, et al, 1991) collected the following data on a sample of  $n = 38$  college students:

- **Response ( $y$ ):** Performance IQ scores (PIQ) from the revised Wechsler Adult Intelligence Scale. Potential
- $x_1$ : Brain size based on the count obtained from MRI scans (given as count/10,000).
- $x_2$ : Height in inches.
- $x_3$ : Weight in pounds

```
> dat <- read.table('~Downloads/iqsize.txt', header = TRUE)
> head(dat)
```

	PIQ	Brain	Height	Weight
1	124	81.69	64.5	118
2	150	103.84	73.3	143
3	128	96.54	68.8	172
4	134	95.15	65.0	147
5	110	92.88	69.0	146
6	131	99.13	64.5	138

$$\vec{\beta}_{\text{hat}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{Y}$$

```
> X <- cbind(1, dat$Brain, dat$Height, dat$Weight)
> y <- dat$PIQ
> solve(t(X)%*%X)%*%t(X)%*%y
```

[1,]	1.113536e+02
[2,]	2.060367e+00
[3,]	-2.731929e+00
[4,]	5.599371e-04

```
> fit <- lm(PIQ ~ Brain + Height + Weight, data=dat)
> fit$coefficients
```

	Brain	Height	Weight	
(Intercept)	1.113536e+02	2.060367e+00	-2.731929e+00	5.599371e-04

# Diagnostics

Model assumption ( $i = 1, \dots, n$ ):

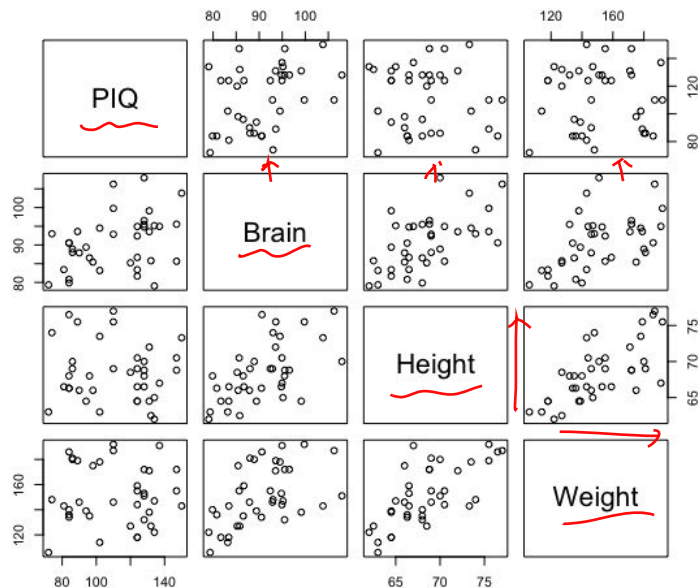
$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \underbrace{\epsilon_i}_{\substack{\text{iid} \\ \sim N(0, \sigma^2)}} . \quad \Rightarrow$$

$$\text{Residuals } \hat{\epsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_p X_{ip}$$

1. Linearity - Plotting  $Y_i$  vs  $X_i$
2. Normality - QQ plot
3. Zero mean in error terms
4. Homoscedasticity
5. Independence

Residual plot

```
> pairs(dat)
```



←

Height  
≈ 2 \* weight

# Diagnostics

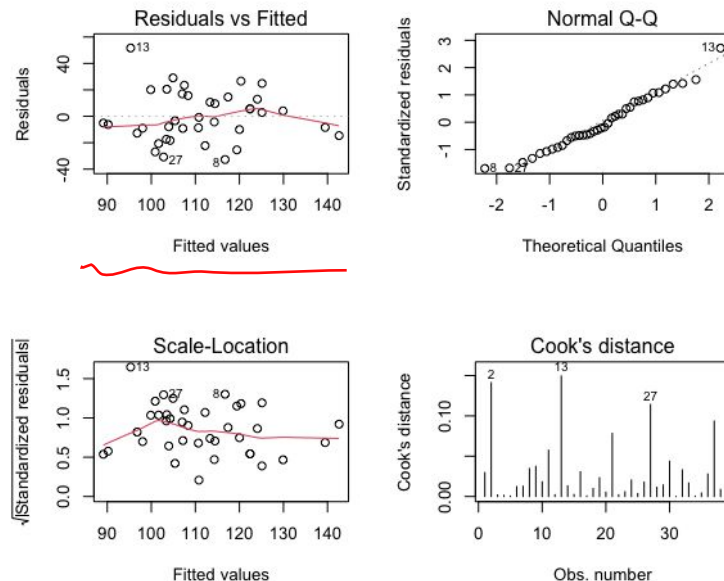
Model assumption ( $i = 1, \dots, n$ ):

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \underbrace{\epsilon_i}_{\substack{\text{iid} \\ \sim N(0, \sigma^2)}} . \quad \Rightarrow \quad \text{Residuals } \hat{\epsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_p X_{ip}$$

1. Linearity - Plotting  $Y_i$  vs  $X_i$
2. Normality - QQ plot
3. Zero mean in error terms
4. Homoscedasticity
5. Independence

} Residual plot

```
> par(mfrow=c(2,2))
> plot(fit, which = 1)
> plot(fit, which = 2)
> plot(fit, which = 3)
> plot(fit, which = 4)
> par(mfrow=c(1,1))
```



## Collinearity: redundant predictors

$$\vec{X}_1 = a \vec{X}_2$$

What if  $\vec{X}_1 \approx a \vec{X}_2$ ?

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{pmatrix}$$

$$R^2 = 1 - \frac{\text{RSS}}{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2} = \frac{SS_{\text{pred}}}{SS_{\text{tot}}}$$

$R^2$  = Percent of variation in  $\vec{Y}$  explained by  $\vec{X}_1, \dots, \vec{X}_p$

```
> n <- nrow(dat)
> RSS <- sum(fit$residuals^2)
> R_sq <- 1-RSS/[(n-1)*var(y)]
> R_sq
[1] 0.2949392
```

```
> summary(fit)$r.squared
[1] 0.2949392
```

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 \underbrace{X_{i1}}_{\downarrow} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip} + \epsilon_i \\ &= \beta_0 + \beta_1 a X_{i2} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip} + \epsilon_i \\ &= \beta_0 + (\beta_1 a + \beta_2) X_{i2} + \cdots + \beta_p X_{ip} + \epsilon_i \end{aligned}$$

$$\sum_i (Y_i - \bar{Y}_n)^2 = \sum_i (Y_i - \hat{\beta}_0 - \cdots - \hat{\beta}_p X_{ip})^2$$

↓

$SS_{\text{pred}}$     $SS_{\text{tot}}$

$SS_{\text{R}}$     $+ SS_{\text{pred}}$

$SS_{\text{R}}$     $SS_{\text{tot}}$

# Collinearity: redundant predictors

What if  $\vec{X}_1 \approx a\vec{X}_2$ ?

$$R^2 = 1 - \frac{\text{RSS}}{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}$$



$R^2$  = Percent of variation in  $\vec{Y}$  explained by  $\vec{X}_1, \dots, \vec{X}_p$

$R_j^2$  = Percent of variation in  $\vec{X}_j$  explained by all other predictors

$R_j^2 \uparrow$   $\vec{X}_j$  can be explained by  $\vec{X}_1, \vec{X}_2, \dots, \vec{X}_{j-1}, \vec{X}_{j+1}, \dots, \vec{X}_p$

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{pmatrix}$$

fit = lm(Res ~ Pred1 + Pred2

+ ... + Predp

```
> library(car)
Loading required package: carData
> vif(fit)
Brain Height Weight
1.578524 2.276641 2.021541
```



↑ ↑ ↑

- If  $R_j^2$  is large (>0.9), then estimation of  $\vec{\beta}$  will be difficult;
- The variance inflation factors (VIFs) are the most common diagnostics:

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

Rule of thumb:  $\text{VIF}_j > 10$  indicates strong collinearity in  $\vec{X}_j$ .

## Tomorrow ...

- Sampling distribution of  $\vec{\beta}_{\text{hat}}$
- Bayesian statistics