

Inference about slope β_1

Issue: Is the predictor X "significant", i.e., does it really help in predicting the response Y ?

Approach 1: Test $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$. This is equivalent to testing $H_0 : \rho = 0$ vs. $H_1 : \rho \neq 0$. (Why?)

Test statistic:

$$T = \frac{\hat{\beta}_1 - 0}{\hat{se}(\hat{\beta}_1)} \sim t_{n-2} \text{ when } H_0 \text{ is true.}$$

Null distribution:

- A two-sided t -test.
- Level & p -val. reject H_0 when $|T| > t_{n-2, \alpha/2}$
- p -value; p_{pivot} .

$100(1 - \alpha)\%$ Confidence Interval for β_1 :

$$\frac{\hat{\beta}_1 - \beta_1}{\hat{se}(\hat{\beta}_1)} \sim t_{n-2}$$

\uparrow
pivot.

Recall:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim N[0, \sigma^2]$$

β_1 really help $\hat{\beta}_1, n-1$

Recall:

$$\hat{\beta}_1 \sim N[\beta_1, \frac{\sigma^2}{n-2}]$$

Approach 2: Test for model significance. In simple linear regression, this is equivalent to testing $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$.

Test statistic:

$$F = \frac{MS_{\text{REG}}}{MS_{\text{ERR}}} = \frac{(n-2) \frac{R^2}{1-R^2}}{MS_{\text{ERR}}} \quad \begin{array}{l} \bullet \quad S_{\text{TOT}} = SS_{\text{REG}} + SS_{\text{ERR}} \\ \rightarrow \text{verify.} \end{array}$$

$$\boxed{\frac{Recall}{\bullet \quad R^2 = \frac{SS_{\text{REG}}}{SS_{\text{TOT}}},}}$$

Null distribution: This F statistic follows an F distribution with numerator d.f. 1 and denominator d.f. $n - 2$.

- An F -test. • $\text{reject } H_0 \text{ if } F \gg F_{1, n-2, \alpha}$ • $p\text{-value} - \text{know.}$
- Equivalent to the t -test seen before because $T^2 = F$ (verify).

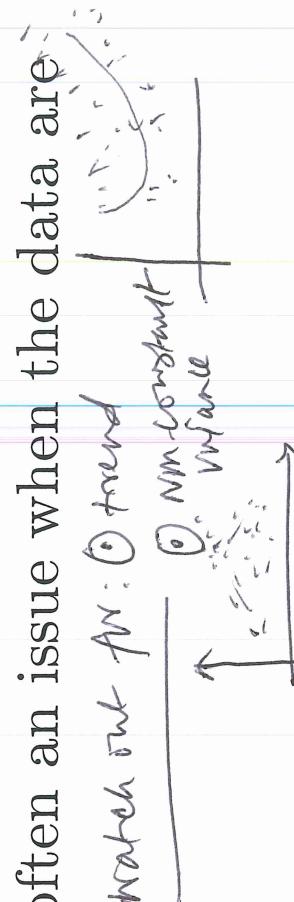
Model evaluation

e_i = random error
 e_i = residual — "residual" of e_i

Issue: Is the fitted model a good representation of the data?

Approach: Examine the residuals, $e_i = Y_i - \hat{Y}_i$, $i = 1, \dots, n$, and verify the key assumptions, namely,

- Errors have mean zero and constant variance
- Errors are normally distributed
- Errors are independent — often an issue when the data are collected over time.



Key Graphical Tools:

- **Residual plot:** Plot of residuals e_i against fitted values \hat{Y}_i . In the ideal plot, the points are scattered around zero and there is no pattern. This verifies the first assumption.
- **Normal QQ plot:** This verifies the normality assumption.
- **Time series plot:** Plot e_i against i . In the ideal plot, there should be no dependence, which verifies the independence assumption. More sophisticated tools exist.

Ex: House price data, continued.

```
x <- house$size  
y <- house$price
```

```
house.reg <- lm(y ~ x)
```

```
# ANOVA table
```

```
> (anova(house.reg))  
Analysis of Variance Table
```

MSE_{reg.}

Response:	y	Df	Sum Sq	Mean Sq	F value	Pr(>F)
x		1	71534	71534	184.62	< 2.2e-16 ***
Residuals	56	21698	387			

```
Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

>

```
# Testing for zero slope
```

```
> summary(house.reg)
```

Call:

```
lm(formula = y ~ x)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-38.489	-14.512	-1.422	14.919	54.389

$\frac{Est}{se}$

Coefficients:

(Intercept)	Estimate	Std. Error	t value	Pr(> t)
β_0	5.432	8.191	0.663	0.51
$x - \beta_1$	56.083	4.128	13.587	<2e-16 ***

$H_0: \beta_0 = 0$

$H_1: \beta_0 \neq 0$

$H_0: \beta_1 = 0$

$H_1: \beta_1 \neq 0$

```
Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

$$R^2 \rightarrow \hat{e}$$

Residual standard error: 19.68 on 56 degrees of freedom

Multiple R-squared: 0.7673, Adjusted R-squared: 0.7631

F-statistic: 184.6 on 1 and 56 DF, p-value: < 2.2e-16

>

```
# Confidence interval for slope
```

```
> confint(house.reg)
2.5 % 97.5 %
(Intercept) -10.97619 21.83933
x ~ beta_1
>
```

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

```
# Prediction at a new x
```

```
x.new <- data.frame(x=3) → e(3,5)
```

```
> (predict(house.reg, newdata=x.new))
```

```
1
```

```
173.6814
```

```
>
```

$$\hat{y}_i$$

```
# Use fitted(house.reg) to get the fitted values
```

```
# Use resid(house.reg) to get the residuals
```

$$e_i = y_i - \hat{y}_i$$

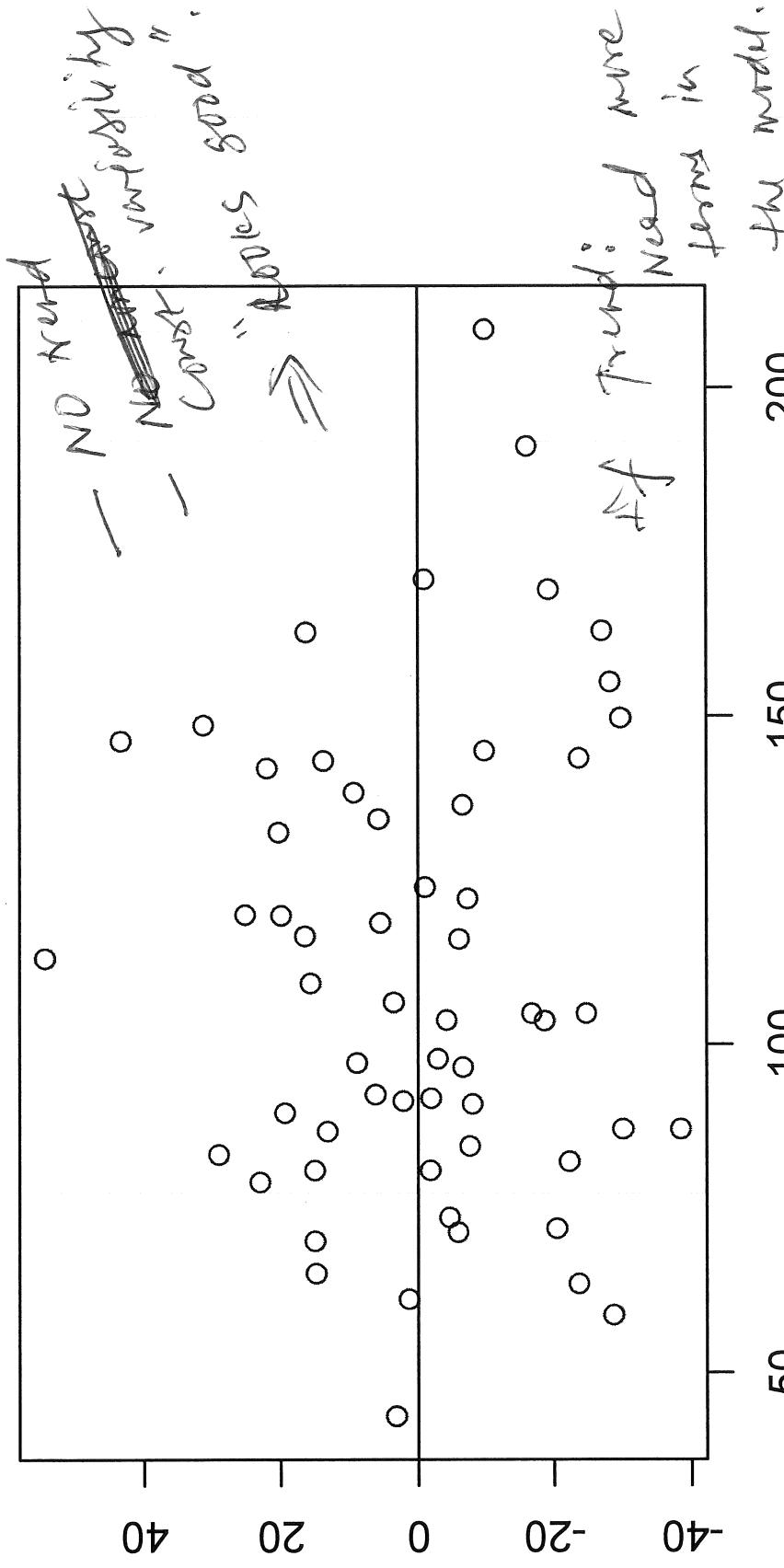
```
# Residual plot
```

```
plot(fitted(house.reg), resid(house.reg))  
abline(h=0)
```

```
# QQ plot
```

```
qqnorm(resid(house.reg))  
qqline(resid(house.reg))  
  
# Time series plot of residuals  
  
plot(resid(house.reg), type="l")  
abline(h=0)
```

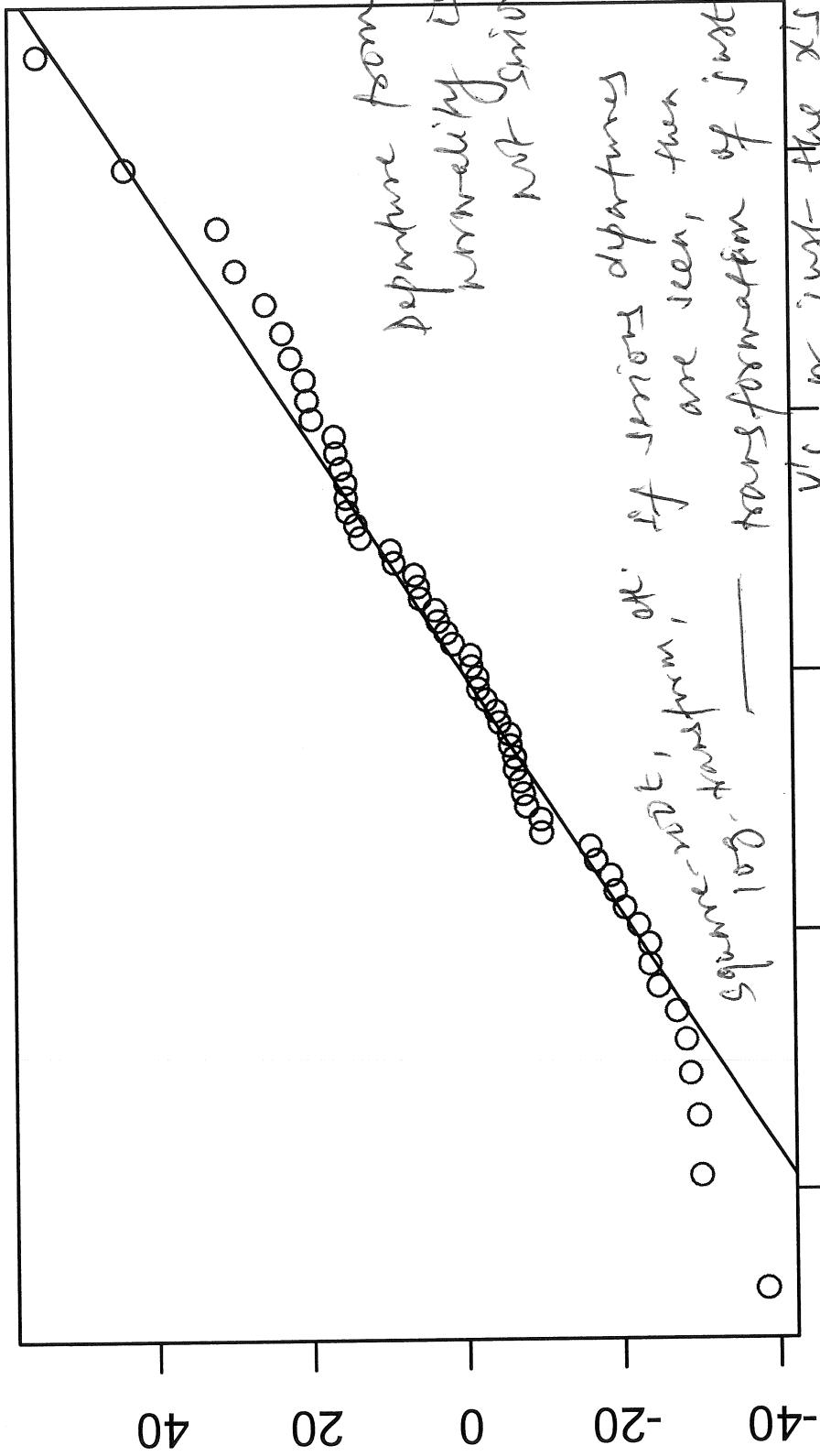
Residual plot



resid(house.reg)

fitted(house.reg)

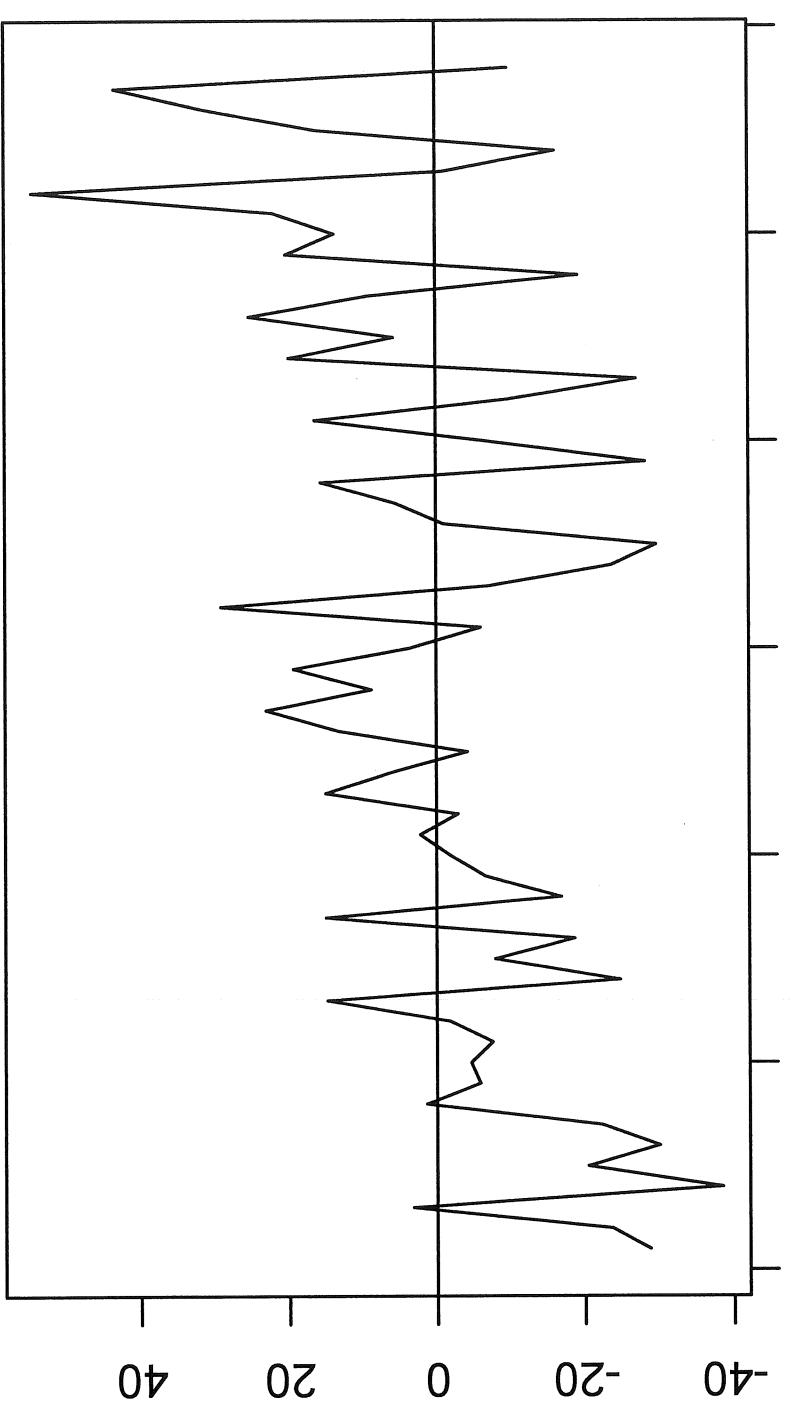
Normal Q-Q Plot



Sample Quantiles

Theoretical Quantiles

Time series plot



resid(house.reg)

Q. Is there any ordering in the data? If yes, we need to
some dependence /~~dependence~~ Index
depends in random over time.

Q. Is there any ordering in the data? If yes, we need to
some dependence /~~dependence~~ Index
depends in random over time.

Multiple Linear Regression

Simple linear regression: One predictor — X

Multiple linear regression: Several predictors — X_1, \dots, X_k

Linear (regression) model:

$E(Y|X_1 = x_1, \dots, X_k = x_k) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$ — models mean response as a function of predictors

Examples:

- $E(Y|x) = \beta_0 + \beta_1 x - \psi_{12}$
- $E(Y|\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x^2 - \psi_{12}$
- $E(Y|\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3(x_1 * x_2) - \psi_{12}$
- $E(Y|\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3(x_1 * x_2) - \psi_{12}$
- $E(\log(Y)|x) = \beta_0 + \beta_1 \log(x) - \psi_{12}$ *predictor*
- $E(Y|x) = \beta_0 + (\beta_1 x)^{-1} - \psi_{12}$ \rightarrow *Not* $E(Y|x)$ is a nonlinear fn. of β_1 .

Note: “Linear” refers to linear in regression coefficients

Linear model: $E(Y|\underline{\mathbf{x}}) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$

Interpretation of $k+1$ regression coefficients:

- $\beta_0 = E(Y|\mathbf{x} = \mathbf{0})$ — intercept
- $\beta_j = E(Y|x_1, \dots, x_j + 1, \dots, x_k) - E(Y|x_1, \dots, x_j, \dots, x_k)$
 - slope of x_j , i.e., change in mean response when j th predictor increases by 1, while keeping other predictors fixed, $j = 1, \dots, k$.

Data: n independent subjects, i th subject gives

$$(Y_i, \mathbf{X}_{1i}, \mathbf{X}_{2i}, \dots, \mathbf{X}_{ki}), i = 1, \dots, n.$$

Linear model for data: For $i = 1, \dots, n$,

$$E(Y_i|x_{i1}, \dots, x_{ik}) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

$$\text{Alternative form: } Y_i = \underbrace{\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}}_{\rightarrow E[Y_i | x_i]} + \epsilon_i$$

Assumptions:

- $E(\epsilon_i) = 0$, $\text{var}(\epsilon_i) = \sigma^2$, and ϵ_i are independent.
- $k+1 < n$ — i.e., have more observations than the number of regression coefficients
- The predictors are considered fixed and are measured without error

Same as before.

These imply:

- $E(Y_i | x_{i1}, \dots, x_{ik}) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$
- $\text{var}(Y_i) = \text{var}[\epsilon_i] = \sigma^2$.
- Y_1, \dots, Y_n are independent.

Linear Model in Matrix Notation

Define:

$$\text{Define: } \mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \mathbf{X}_{n \times (k+1)} = \begin{bmatrix} 1 & \cdots & x_{11} & \cdots & x_{1k} \\ 1 & \cdots & x_{21} & \cdots & x_{2k} \\ \vdots & & \vdots & & \vdots \\ 1 & \cdots & x_{n1} & \cdots & x_{nk} \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

values of X

values of y

regression matrix

response

- $Y_i = (\mathbf{X} + \epsilon) \beta + \epsilon_i$

- $\hat{\mathbf{Y}} = \mathbf{X} \hat{\beta} + \epsilon$

- $E(\mathbf{Y} | \mathbf{X}) = \mathbf{X} \beta$

- rank of X is full, i.e., $(\mathbf{X}' \mathbf{X})^{-1}$ exists.

- $\hat{\beta}$ = estimator of β

- $\hat{Y} = \mathbf{X} \hat{\beta}$ = fitted (or predicted) response

Predicted response when $\mathbf{x} = \mathbf{x}_0: \hat{Y}_0 = \mathbf{x}'_0 \hat{\beta}$

values of (x_1, x_2, \dots, x_k)

Least Squares Estimation of β

X' = transpose
of X .

- As before: Minimize $\sum_{i=1}^n \epsilon_i^2$ with respect to $\beta_0, \beta_1, \dots, \beta_k$ to get $\hat{\beta}$
- Least squares estimator: $\hat{\beta} = (\underline{X}' \underline{X})^{-1} \underline{X}' \underline{Y}$
- Minimum value of $\sum_{i=1}^n \epsilon_i^2$ is $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = (\underline{Y} - \hat{\underline{Y}})'(\underline{Y} - \hat{\underline{Y}}) = S_{\text{ERR}} - \text{error (or residual) sum of squares}$

Properties of $\hat{\beta}$:

- Linear in \mathbf{Y}
- Unbiased, i.e., $E[\hat{\beta}] = (\underline{X}' \underline{X})^{-1} \underline{X}' \underline{\varepsilon} = \beta$
- $\text{var}(\hat{\beta}) = \sigma^2 (\underline{X}' \underline{X})^{-1} \rightarrow \text{matrix } (k+1) \times (k+1)$
- $\text{var}(\hat{\beta}_0) = \sigma^2 \times \text{first diagonal element of } (\underline{X}' \underline{X})^{-1}$
- $\text{var}(\hat{\beta}_j) = \sigma^2 \times (j+1)\text{th diagonal element of } (\underline{X}' \underline{X})^{-1}$
- $\hat{\sigma}^2 = S_{\text{ERR}} / (n - k - 1) = M S_{\text{ERR}}$ is unbiased for σ^2 .

ANOVA table

As before:

- $SS_{TOT} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (\mathbf{Y} - \bar{\mathbf{Y}})'(\mathbf{Y} - \bar{\mathbf{Y}})$, where

$$\bar{\mathbf{Y}} = \bar{Y} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

- $SS_{REG} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = (\hat{\mathbf{Y}} - \bar{\mathbf{Y}})'(\hat{\mathbf{Y}} - \bar{\mathbf{Y}})$
- $SS_{ERR} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = (\mathbf{Y} - \hat{\mathbf{Y}})'(\mathbf{Y} - \hat{\mathbf{Y}})$

Source	SS	d.f.	MS	F
Model	SS_{REG}	k	$MS_{REG} = \frac{SS_{REG}}{k}$	$\frac{MS_{REG}}{MS_{ERR}}$
Error	SS_{ERR}	$n - k - 1$	$MS_{ERR} = \frac{SS_{ERR}}{n-k-1}$	
Total	SS_{TOT}	$n - 1$		