

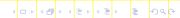
Recap: Sampling Distributions of Sums of Squares (SS)

Under the Normal error model:

•
$$SSE$$
 and SSR are independent.

• $SSE \sim \sigma^2 \chi^2_{(n-p)}$.

• If $\beta_1 = \cdots = \beta_{p-1} = 0$, then $SSR \sim \sigma^2 \chi^2_{(p-1)}$.



Mean squares (MS): MS = SS/d.f.(SS). MSE: $MSE = \frac{SSE}{n-p}, E(MSE) = \sigma^2.$ MSE is an estimator of the error variance σ^2 . MSR: $MSR = \frac{SSR}{p-1}$. if $\beta_1 = \cdots = \beta_{p-1} = 0$ if otherwise E(MSR) =4□ > 4₫ > 4 ≧ > 4 ≧ > ½ 900 €

F Test of Regression Relation

Under the Normal error model:

 Test whether there is a regression relation between the response variable Y and the set of X variables:

F ratio and its null distribution:

where
$$F_{p-1,n-p}$$
 denotes the F distribution with $(p-1,n-p)$ degrees of freedom.

Decision rule at level α : reject H_0 if $F^* >$

ANOVA Table

Source of Vari		S\$		d.f.	MS	F*
R <mark>egre</mark> ssion	SS	$R = \mathbf{Y}'(\mathbf{H} -$	$\frac{1}{n}J_n)Y$	p – 1	$MSR = \frac{SSR}{p-1}$	$F^* = \frac{MSR}{MSE}$
Error		$SE=\mathbf{Y}'(\mathbf{I}_n$ -	/	n – p	$MSE = \frac{SSE}{n-p}$	
Total	SST	$O = \mathbf{Y}' \left(\mathbf{I}_n - \mathbf{I}_n \right)$	$-\frac{1}{n}\mathbf{J}_{n}\mathbf{Y}$	n – 1		

Example Model 2: n = 30, p = 5.

Source of Variation	SS	d.f. M\$	F*
Regression	SSR = 366.4846	4 $MSR = 91.62116$	$F^* = 87.03703$
Error	SSE = 26.31672	25 $MSE = 1.052669$	
Total	SSTO = 392.8013	29	

Pvalue = $P(F_{4,25} > 87.037) \approx 0$, so there is a significant regression relation between Y and X_1, X_2, X_3, X_1X_2 .

Coefficient of Multiple Determination

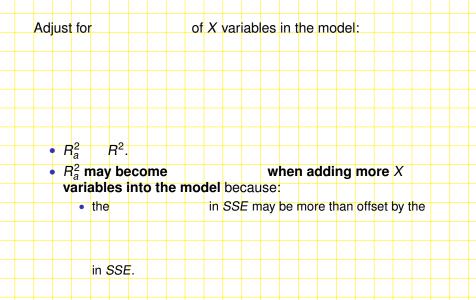
$$R^2 := \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO}$$

- R² is the of the total variation in Y by using the X variables to explain Y.
- $0 \le R^2 \le 1$.
- When $R^2 = 0$? When $R^2 = 1$?
- Adding more X variables to the model will always R² because:
 - SSTO
 - SSE

Since adding more X variables can only R^2 , does this mean we should use as many X variables as possible?

- With more X variables, the model does fit the observed data , indicated by SSE.
- However, a model with many X variables that are unrelated to the response variable and/or are highly correlated with each other tends to
 - the observed data and often do a job for prediction (i.e., generalize poorly for new cases) due to sampling variability.
 - make interpretation
 - make model maintenance more

Adjusted Coefficient of Multiple Determination



Example

$$R^2 = 0.8883, \quad R_a^2 = 0.8754$$

• Model 2 :
$$Y \sim X_1, X_2, X_3, X_1 X_2$$

$$R^2 = 0.933, R_a^2 = 0.9223.$$

$$R^2 = 0.937, \quad R_a^2 = 0.9205.$$

(i) For each model, $R^2 > R_a^2$; (ii) Adding more X variable(s) increases R^2 . The increase of R^2 is much more from Model 1 to Model 2 than from Model 2 to Model 3; (iii) Model 3 has a smaller R_a^2 than Model 2.

Inferences about Regression Coefficients

LS estimators:
$$\hat{\beta} = \begin{vmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_{p-1} \end{vmatrix}$$

$$\mathbf{E}\{\hat{\beta}\} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_{p-1} \end{vmatrix}$$

$$\mathbf{F}\{\hat{\beta}\} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_{p-1} \end{pmatrix}$$
The standard error of $\hat{\beta}_k$, $s(\hat{\beta}_k)$, is the

- Studentized pivotal quantity:
 - $\frac{\hat{\beta}_k \beta_k}{\hat{\beta}_k} \sim$
- (1α) -Confidence interval for β_k :

T statistic:

Two-sided T-Test: $H_0: \beta_k = \beta_k^0$ vs. $H_a: \beta_k \neq \beta_k^0$. At level α , the decision rule is to reject H_0 if and only if

What are decision rules for one-sided tests?

Multiple Regression: Example

n -	= 3	n c:	200	C 1	'ACI	าดท	22	var	iah	۱ ما	V a	nd ·	thre	o r)rei	dict	or۱	/ari	ahl	മഠ			
X_1	X_2	X_3	}.	ا ,	CS	JO1 1	30	vai	iau		a	·u		,), C.	JICL	0.	/an	αυι	C3			
	e –			X1		Х2		х3															
1 2								1.2 -0.															
3								0.4															
 30					12			0.6															
30			1.4	۷.	12	-0.	. 6 -	u . 0.															
														4	□ }	∢ 🗗	⊢ ∢	≣ ⊦	∢ ∄	Þ	=	99	(~

Example: Model 2

Nonadditive model with interaction between
$$X_1$$
 and X_2 :

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i1} X_{i2} + \epsilon_i, \quad i = 1, \cdots, 30.$$

($\rho = 5$)
Call:
Im(formula = Y \(^{-1} X_1 + X_2 + X_3 + X_1 : X_2, \text{ data} = \text{ data})

Coefficients:
Estimate Std. Error t value \(^{-1} Pr(>|t|) \)
(Intercept) 0.8832 0.2153 4.103 0.00038 ***
X1 1.5946 0.2421 6.587 6.69e-07 ***
X2 1.7091 0.2605 6.560 7.16e-07 ***
X3 2.1266 0.2687 7.916 2.85e-08 ***
X1:X2 1.0076 0.2467 4.084 0.00040 ***

X1:X2 1.0076 0.2467 4.084 0.00040 ***

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

Residual standard error: 1.026 on 25 degrees of freedom Multiple R-squared: 0.933, Adjusted R-squared: 0.9223 F-statistic: 87.04 on 4 and 25 DF, p-value: 2.681e-14



Test whether there is an interaction between X_1 and X_2 . Use

the null hypothesis and

interaction

, so

, vs., H_a:

$$\alpha = 0.01.$$
• H_0 :

•
$$n = 30, p = 5,$$

- conclude that there is
- effect between X_1 and X_2 .
- Alternatively, pvalue=
- H_0 .

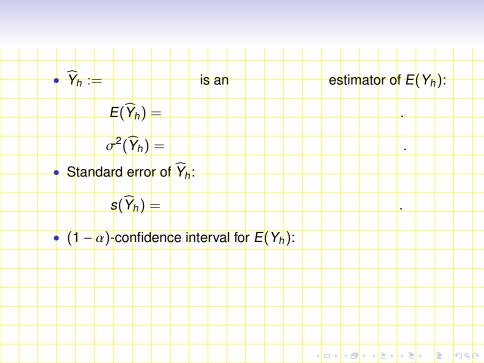
Notes: pvalue for the two-sided alternative is in the R output. What is a 99% confidence interval for β_4 ? How to test the

right-sided alternative?

Estimation of the Mean Response

• For a given set of values of the
$$X$$
 variables:
$$\mathbf{X}_h = \begin{bmatrix} 1 \\ X_{h1} \\ \vdots \\ X_{h,p-1} \end{bmatrix}$$
• Corresponding mean response:
$$E(Y_h) = \begin{bmatrix} 1 \\ X_{h1} \\ \vdots \\ X_{h,p-1} \end{bmatrix}$$

4 D > 4 B > 4 E > 4 E > E 9 Q C



Prediction of a New Observation

$$\bullet \ \ Y_{h(new)} = \mathbf{X}_h' \boldsymbol{\beta} + \epsilon_h : \qquad \text{with the observations } Y_i \mathbf{s}.$$

$$\bullet \ \ \text{Predicted value: } \widehat{Y}_h := \qquad .$$

$$\bullet \ \ \sigma^2(pred_h) := \qquad .$$

$$\bullet \ \ \text{Standard error for prediction:}$$

$$s(pred_h) = \qquad .$$

$$\bullet \ \ (1 - \alpha) \text{-prediction interval for } Y_{h(new)} :$$

4 D > 4 B > 4 E > 4 E > E 900

Example

Estimate the mean response when
$$X_1 = 0.8, X_2 = 0.5, X_3 = -1$$
 under Model 2.

•
$$n = 30, p = 5$$
:

$$\mathbf{X}'_{h}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}_{h} = 0.170, MSE = 1.053,$$

 $\widehat{Y}_h := \mathbf{X}_h' \hat{\boldsymbol{\beta}} = 1.290,$

$$s(\widehat{Y}_h) =$$

• A 99%-confidence interval for
$$E(Y_h)$$
: $t(0.995; 25) = 2.787$
1.290 $\pm 2.787 \times 0.423 = [0.111, 2.469]$.

- Predict a new observation when $X_1 = 0.8, X_2 = 0.5, X_3 = +1$
- under Model 2. Standard error for prediction:
 - s(pred) =
 - A 99%-prediction interval for Y_{hnew} :
 - $1.290 \pm 2.787 \times 1.1098 = [-1.803, 4.383].$ R codes.

 - > newX=data.frame(X1=0.8, X2=0.5, X3=-1) > predict.lm(fit2, newX, interval="confidence",
 - + level=0.99, se.fit=TRUE)
- > predict.lm(fit2, newX, interval="prediction", + level=0.99, se.fit=TRUE)

Hidden Extrapolations

- Recall that extrapolation occurs when predicting the response variable for values of the X variable(s) of the original data.
- The fitted model may when extended outside the range of the observations.
- With more than one X variables, the levels of define the region of the observations. One can not merely look at the ranges of each X variable.
- With two X variables, we can look at their scatter plot.
- Procedure to identify hidden extrapolation for more than two X variables will be discussed later.

