Linear Regression Models P8111

Lecture 09

10:00 am - 5:00 pm

Milterm!!

3/8

3/10

d

3R (([

Jeff Goldsmith February 18, 2016

Columbia University

THE DEPARTMENT OF **BIOSTATISTICS**

Today's Lecture

$$MLR : LSE \rightarrow \hat{\beta} = (x^{T}x)^{T}x^{T}y$$

$$L\eta = x\beta + \epsilon$$

Sampling distribution of β
 Hypothesis tests for individual coefficients
 Global tests

Circle of Life



Statistical inference

- We have LSEs $\hat{\beta}_0, \hat{\beta}_1, \ldots$; we want to know what this tells us about β_0, β_1, \ldots .
- Two basic tools are confidence intervals and hypothesis tests
 - Confidence intervals provide a plausible range of values for
 - the parameter of interest based on the observed data
 - Hypothesis tests ask how probable are the data we gathered under a null hypothesis about the data generating distribution

A quick word about p-values



Motivation

Recall the MLB data:

```
> setwd("'/Desktop")
> download.file("http://www.openintro.org/stat/data/mlb11.RData", destfile = "mlb11.RData")
> load("mlb11.RData")
> mlb11 %>% tbl_df
Source: local data frame [30 x 12]
```

	team	runs	at_bats	hits	homeruns	bat_avg	strikeouts	stolen_bases	wins	new
	(fctr)	(int)	(int)	(int)	(int)	(dbl)	(int)	(int)	(int)	
1	Texas Rangers	855	5659	1599	210	0.283	930	143	96	
2	Boston Red Sox	875	5710	1600	203	0.280	1108	102	90	
3	Detroit Tigers	787	5563	1540	169	0.277	1143	49	95	
4	Kansas City Royals	730	5672	1560	129	0.275	1006	153	71	
5	St. Louis Cardinals	762	5532	1513	162	0.273	978	57	90	
6	New York Mets	718	5600	1477	108	0.264	1085	130	77	
7	New York Yankees	867	5518	1452	222	0.263	1138	147	97	
8	Milwaukee Brewers	721	5447	1422	185	0.261	1083	94	96	
9	Colorado Rockies	735	5544	1429	163	0.258	1201	118	73	
10	Houston Astros	615	5598	1442	95	0.258	1164	118	56	
Variables not shown: new_slug			g (dbl),	new_ob	s (dbl)					

Motivation



Residual standard error: 26.85 on 25 degrees of freedom Multiple R-squared: 0.9087,Adjusted R-squared: 0.894 F-statistic: 62.17 on 4 and 25 DF, p-value: 1.26e-12

Motivation

• Can we say anything about whether the effect of >> stolen_bases is "significant" after adjusting for other variables? $\iint_{a} : \hat{\beta}_{a} = 0$ Can we compare this model to a model with only hits 7 and homeruns? β y = pot β, ot + p2 h. ts + β3 hR + β4 5B + € by yin pot pi hater & pahr + E Ho: Pre pi enter & pahr + E Ho: Pre pi = 0 Vg

Sampling distribution

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{x}^{\mathsf{T}}\boldsymbol{x})^{\mathsf{T}}\boldsymbol{x}^{\mathsf{T}}\boldsymbol{j}, \quad \hat{\boldsymbol{\beta}} \sim (\boldsymbol{\beta}, \ \boldsymbol{\varepsilon}^{\mathsf{C}}(\boldsymbol{x}^{\mathsf{T}}\boldsymbol{x})^{\mathsf{T}})$$
If our usual assumptions are satisfied and $\boldsymbol{\epsilon} \stackrel{iid}{\sim} \mathbf{N} [0, \sigma^2]$ then
$$\hat{\boldsymbol{\beta}} \sim \mathbf{N} [\boldsymbol{\beta}, \sigma^2 (\boldsymbol{X}^T \boldsymbol{X})^{-1}].$$

$$\hat{\boldsymbol{\beta}}_{j} \sim \mathbf{N} [\boldsymbol{\beta}, \sigma^2 (\boldsymbol{X}^T \boldsymbol{X})_{jj}^{-1}].$$
• This will be used for inference.

Asymptotic distribution

Assume that

- $E(\epsilon_i | \mathbf{x}_i) = \underline{0 \ \forall i};$
- $Var(\epsilon_i | \mathbf{x}_i) = \sigma^2 \forall i;$

• $n \xrightarrow{\lim} \infty \frac{X^T X}{n} \rightarrow Q$ where *Q* is a finite non-singular matrix.

Then

$$\underbrace{\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}_{\mathbf{V}} \rightarrow \underbrace{N\left[0, \sigma^2 Q^{-1}\right]}_{\mathbf{V}}$$

(This is essentially an extension of the central limit theorem)

Simulations exploring distributions

Look at SLR

$$y_i = 0 + 1x_i + \epsilon_i$$

under various conditions.

- First simulations: errors follow N [0, 1], let *n* vary
- Second simulations: errors follow $\frac{10}{3} * \text{Bern}(.1) \frac{1}{3}$, let *n* vary
- In both cases, $\epsilon \sim (0, 1)$









Non-normal errors



Testing procedure

Calculate the probability of the observed data (or more extreme data) under a null hypothesis.

- Often $H_0: \beta_1 = 0$ and $H_a: \beta_1 \neq 0$
- Set $\alpha = P($ falsely rejecting a true null hypothesis) (type I error rate) o, o s
- Calculate a test statistic assuming the null hypothesis is true
- Compute a p-value =

 $P(As \text{ or more extreme test statistic}|H_0)$

Reject or fail to reject H₀

Testing

$$\hat{\beta}_{j} \sim \mathcal{N}(\beta, \sigma^{2}(x^{T}x)^{T})$$

For real data we have to estimate σ^2 as well as β .

Recall our estimate of the error variance is

$$\hat{\sigma^2} = \frac{RSS}{n-p-1} = \frac{\sum_i (y_i - \hat{y}_i)^2}{n-p-1}$$

• With Normally distributed errors, it can be shown that

$$(n-p-1)\hat{\frac{\sigma^2}{\sigma^2}} \sim \chi^2_{n-p-1}$$

Implication is that test statistics follow a <u>t</u> distribution rather than Normal with df = n - p - 1

Individual coefficients



■ We can use the test statistic

T =

• For a two-sided test of size α , we reject if

$$|T| > t_{1-\alpha/2, n-p-1}$$
 clamps give P_{α}^{μ}

0.0Z

• The p-value gives $P(t_{n-p-1} > T_{obs}|H_0)$

Note that *t* is a symmetric distribution that converges to a Normal as n - p - 1 increases.

1025

Example revisited



Residual standard error: 26.85 on 25 degrees of freedom Multiple R-squared: 0.9087,Adjusted R-squared: 0.894 F-statistic: 62.17 on 4 and 25 DF, p-value: 1.26e-12 Inference for linear combinations $\begin{pmatrix} \rho_{1} & \rho_{2} \\ -(\rho_{1} & \rho_{2} & \rho_{3} \end{pmatrix}$ $\downarrow \downarrow_{i}: \rho_{i} - \rho_{2} = 0$?? $\begin{pmatrix} c_{1} & \rho_{2} \\ -(\rho_{1} & \rho_{2} & \rho_{3} \end{pmatrix}$ $\downarrow \downarrow_{i}: \rho_{0} + \rho_{1} \gamma_{2} \times \rho_{2} \gamma_{2} + \rho_{2} \gamma_{2} \times \rho_{3} + \rho_{3} \gamma_{3} + \rho_{3} \gamma_{3}$

• Define $H_0: \underline{c^T \beta} = c^T \beta_0$ or $H_0: \underline{c^T \beta} = 0$

• We can use the test statistic

ValaB

- This test statistic is asymptotically Normally distributed
- For a two-sided test of size *α*, we reject if

$$|T| > z_{1-\alpha/2}$$

Inference about multiple coefficients

Our model contains multiple parameters; often we want to perform multiple tests:



where each test has a size of α

• For any individual test, $P(\text{reject } H_{0i}|H_{0i}) = \alpha$

Inference about multiple coefficients

What about

 $P(\text{reject at least one } H_{0i}|\text{all } H_{0i} \text{ are true}) = \alpha$

Family-wise error rate

To calculate the FWER

- First note $P(\text{no rejections}|\text{all }H_{0i} \text{ are true}) = (1 \alpha)^k$
- It follows that

 $P(\text{at least one rejection}|\text{all } H_{0i} \text{ are true}) = 1 - (1 - \alpha)^k$

Further,

$$FWER = \underbrace{1 - (1 - \alpha)^k}_{\approx} = \underbrace{1 - \left(1 - \frac{k\alpha}{k}\right)^k}_{\approx}$$
$$\approx \underbrace{1 - exp(1 - k\alpha)}_{\approx}$$
$$\approx \underbrace{1 - (1 - k\alpha)}_{\approx}$$
$$= \underbrace{k\alpha}$$

Family-wise error rate



Addressing multiple comparisons

Three general approaches

Do nothing in a reasonable way

- • Define comparisons and expectations ahead of time
- Don't trust scientifically implausible results
 - Don't over-emphasize isolated findings
- Correct for multiple comparisons
 - C > Often, use the Bonferroni correction and use $\alpha_i = \alpha/k$ for each test
 - Thanks to the Bonferroni inequality, this gives an overall FWER < α
 - Control false discovery rate

Use a global test

Global tests

Compare a smaller "null" model to a larger "alternative" model

- Smaller model must be nested in the larger model
 - That is, the smaller model must be a special case of the larger model
 - For both models, the *RSS* gives a general idea about how well the model is fitting
 - In particular, something like

$$\frac{RSS_S - RSS_L}{RSS_L}$$
Sound = Reject
$$B_{i2} = fail + a_{R}$$

compares the relative RSS of the models

Nested models

• These models are nested:

Smaller = Regression of Y on
$$X_1$$

Larger = Regression of Y on X_1, X_2, X_3, X_4

• These models are not:

Smaller = Regression of Y on
$$X_2$$

Larger = Regression of Y on X_1, X_3

Global F tests

- Compute the test statistic $F_{obs} = \frac{(RSS_S - RSS_L)/(df_S - df_L)}{RSS_L/df_L}$
- If H_0 (the null model) is true, then $F_{obs} \sim F_{df_s df_L, df_L}$
- Note $df_s = n p_S 1$ and $df_L = n p_L 1$
- We reject the null hypothesis if the p-value is above α, where

$$p-value = P(F_{df_S} - df_L, df_L > F_{obs})$$

Global F tests

Cat $y_i = \beta_0 + \beta_1 g_{z_1} + \beta_2 g_{z_3} + \epsilon_i$ $H_0: \beta_1 = \beta_2 = 0$ $-y_i = \beta_0 + \epsilon_i$

There are a couple of important special cases for the *F* test The null model contains the intercept only

- When people say ANOVA, this is often what they mean (although all *F* tests are based on an analysis of variance)
- The null model and the alternative model differ only by one term
 - Gives a way of testing for a single coefficient
 - ► Turns out to be equivalent to a two-sided *t*-test: $t_{df_L}^2 \sim F_{1,df_L}$

MLB data

You can test multiple coefficient simultaneously using the F test

```
> linmod.null1 = lm(runs ~ hits + homeruns, data = mlb11)
> anova(linmod.null1, linmod)
Analysis of Variance Table
Model 1: runs ~ hits + homeruns
Model 2: runs ~ at_bats + hits + homeruns + stolen_bases
Res.Df RSS Df Sum of Sq F Pr(>F)
1 27 27128
2 25 18020 2 9107.8 6.3178 0.006015 **
---
Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

MLB data

The *F* test is equivalent to the *t* test when there's only one parameter of interest

```
> linmod.null2 = lm(runs ~ at_bats + hits + homeruns, data = mlb11)
> anova(linmod.null2, linmod)
Analysis of Variance Table
Model 1: runs ~ at_bats + hits + homeruns
Model 2: runs ~ at_bats + hits + homeruns + stolen_bases
Res.Df RSS Df Sum of Sq F Pr(>F)
1        26 24953
2        25 18020 1        6932.7         .618 0.004728 **
----
Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

MLB data

By default, R's summary function compares to an intercept-only null model

Test for "linearity"

- To test more flexible vs less flexible approaches to non-linearity, we can often use global tests
 - Polynomials and piecewise linear models have "linear" associations as nested model; B-splines don't
- Global *F* tests can be pretty useful here

Testing for linearity



Testing for linearity



Testing linearity

1 98 73.444

```
> piecewise.underfit = lm(y ~ x, data = data.nonlin)
> piecewise.fit = lm(y ~ x + spline_15 + spline_5 + spline_9, data = data.nonlin)
> anova(piecewise.underfit, piecewise.fit)
Analysis of Variance Table
Model 1: y ~ x
Model 2: y ~ x + spline_15 + spline_5 + spline_9
```

Res.Df RSS Df Sum of Sq F Pr(>F)

2 95 8.240 3 65.205 250.6 < 2.2e-16 *** ---Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1

Testing comparing twenty polynomials to four



Testing comparing twenty polynomials to four

```
> piecewise.overfit = lm(y ~ x + spline_1 + spline_15 + spline_2 + spline_25 + spline_3 + spl
+ spline_4 + spline_75 + spline_55 + spline_6 + spline_65
+ data = data.nonlin)
> anova(piecewise.fit, piecewise.overfit)
Analysis of Variance Table
Model 1: y ~ x + spline_15 + spline_5 + spline_9
Model 2: y ~ x + spline_15 + spline_22 + spline_25 + spline_3 +
spline_35 + spline_4 + spline_15 + spline_5 + spline_55 +
spline_6 + spline_65 + spline_75 + spline_75 + spline_8 +
spline_85 + spline_9
Res.Df RSS Df Sum of Sq F Pr(>F)
1 95 8.2395
2 81 6.7862 14 1.4533 1.239 0.2645
```

Testing comparing twenty polynomials to four

```
> anova(piecewise.underfit, piecewise.fit, piecewise.overfit)
Analysis of Variance Table
Model 1: y ~ x
Model 2: y ~ x + spline_15 + spline_5 + spline_9
Model 3: y ~ x + spline_1 + spline_15 + spline_2 + spline_25 + spline_3 +
    spline_35 + spline_4 + spline_45 + spline_5 + spline_55 +
    spline_6 + spline_6 + spline_7 + spline_75 + spline_8 +
    spline_85 + spline_9
Res.Df RSS Df Sum of Sq F Pr(>F)
1 98 73.444
2 95 8.240 3 65.205 259.427 <2e-16 ***
3 81 6.786 14 1.453 1.239 0.2645
---
Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1</pre>
```

Today's big ideas

Inference for multiple linear regression models

Suggested reading: Faraway Ch 3.1 - 3.3