

Final concepts of SLR

Author: Nicholas G Reich, Jeff Goldsmith

*This material is part of the **statsTeachR** project*

Made available under the Creative Commons Attribution-ShareAlike 3.0 Unported License: http://creativecommons.org/licenses/by-sa/3.0/deed.en_US

Today's lecture

- Simple Linear Regression Continued
- Multiple Regression Intro

Simple linear regression model

- Observe data (y_i, x_i) for subjects $1, \dots, I$. Want to estimate β_0, β_1 in the model

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i; \epsilon_i \stackrel{iid}{\sim} (0, \sigma^2)$$

- Note the assumptions on the variance:

- $E(\epsilon | x) = E(\epsilon) = 0$
- Constant variance
- Independence
- [Normally distributed is not needed for least squares, but is needed for inference]

$E[y|x]$

Normality

Some definitions / SLR products

$:=$

- Fitted values: $\hat{y}_i := \hat{\beta}_0 + \hat{\beta}_1 x_i$ *is fitted*
- Residuals / estimated errors: $\hat{\epsilon}_i := y_i - \hat{y}_i$
- Residual sum of squares: $RSS := \sum_{i=1}^n \hat{\epsilon}_i^2$
- Residual variance: $\hat{\sigma}^2 := \frac{RSS}{n-2}$
- Degrees of freedom: $n - 2$

Notes: residual sample mean is zero; residuals are uncorrelated with fitted values.

R^2

Looking for a measure of goodness of fit.

- RSS by itself doesn't work so well:

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Coefficient of determination (R^2) works better:

prop of var. explained by model

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

prop of all variance described by errors

R^2



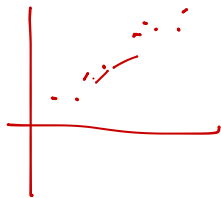
Some notes about R^2

- Interpreted as proportion of outcome variance explained by the model.
- Alternative form

$$R^2 = \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2}$$

- R^2 is bounded: $0 \leq R^2 \leq 1$
- For simple linear regression only, $R^2 = \rho^2$

R^2_2



$$R^2_1 = R^2_2$$

$$\rho_1 \neq \rho_2$$

ANOVA

Lots of sums of squares around.

- Regression sum of squares $SS_{reg} = \sum(\hat{y}_i - \bar{y})^2$
- Residual sum of squares $SS_{res} = \sum(y_i - \hat{y}_i)^2$
- Total sum of squares $SS_{tot} = \sum(y_i - \bar{y})^2$
- All are related to sample variances

Analysis of variance (ANOVA) seeks to address goodness-of-fit by looking at these sample variances.

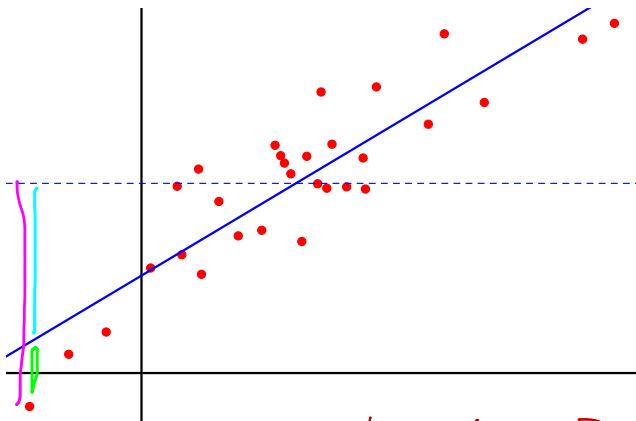
ANOVA

ANOVA is based on the fact that $SS_{tot} = SS_{reg} + SS_{res}$

$$\begin{aligned} SS_{tot} &= \sum (y_i - \bar{y})^2 \\ &= \sum_i ((y_i - \hat{y}) + (\hat{y} - \bar{y}))^2 \end{aligned}$$

ANOVA

ANOVA is based on the fact that $SS_{tot} = SS_{reg} + SS_{res}$



see textbook?

ANOVA and R^2

- Both take advantage of sums of squares
- Both are defined for more complex models
- ANOVA can be used to derive a “global hypothesis test” based on an F test (more on this later)

R^2 never used for hypothesis test

R example

```
require(alr3)
data(heights)
linmod <- lm(Dheight ~ Mheight, data = heights)
linmod

##
## Call:
## lm(formula = Dheight ~ Mheight, data = heights)
##
## Coefficients:
## (Intercept)      Mheight
##      29.917         0.542
```

R example

```
summary(linmod)

##
## Call:
## lm(formula = Dheight ~ Mheight, data = heights)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.397 -1.529  0.036  1.492  9.053
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   29.917     1.623    18.4   <2e-16 ***
## Mheight        0.542     0.026    20.9   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.27 on 1373 degrees of freedom
## Multiple R-squared:  0.241, Adjusted R-squared:  0.24
## F-statistic: 435 on 1 and 1373 DF, p-value: <2e-16
```

R example

*class(linmod)
"lm"*

```
names(linmod)
```

```
## [1] "coefficients" "residuals"      "effects"         "rank"  
## [5] "fitted.values" "assign"          "qr"              "df.residual"  
## [9] "xlevels"       "call"           "terms"           "model"
```

*linmod\$coef
/lm*

R example

```
head(linmod$residuals)
```

```
##      1      2      3      4      5      6  
## -7.160 -4.947 -6.747 -6.001 -7.397 -2.084
```

```
head(resid(linmod))
```

```
##      1      2      3      4      5      6  
## -7.160 -4.947 -6.747 -6.001 -7.397 -2.084
```

```
head(linmod$fitted.values)
```

```
##      1      2      3      4      5      6  
## 62.26 61.45 62.75 62.80 63.40 59.98
```

```
head(fitted(linmod))
```

```
##      1      2      3      4      5      6  
## 62.26 61.45 62.75 62.80 63.40 59.98
```

R example

```
names(summary(linmod))
```

```
## [1] "call"      "terms"      "residuals"  "coefficients"  
## [5] "aliased"    "sigma"      "df"         "r.squared"  
## [9] "adj.r.squared" "fstatistic" "cov.unscaled"
```

```
summary(linmod)$coef
```

[1, 4]

```
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept)  29.9174    1.62247   18.44 5.212e-68  
## Mheight      0.5417     0.02596   20.87 3.217e-84
```

```
summary(linmod)$r.squared
```

```
## [1] 0.2408
```

R example

```
anova(linmod)

## Analysis of Variance Table
##
## Response: Dheight
##              Df Sum Sq Mean Sq F value Pr(>F)
## Mheight         1    2237      2237    435 <2e-16 ***
## Residuals    1373    7052         5
## --- SSreg --- SSres
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{SS_{res}}{SS_{res} + SS_{reg}}$$

R example

```
anova(linmod)

## Analysis of Variance Table
##
## Response: Dheight
##           Df Sum Sq Mean Sq F value Pr(>F)
## Mheight      1   2237    2237    435 <2e-16 ***
## Residuals 1373   7052         5
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(r2 <- 1 - 7052/(7052 + 2237))

## [1] 0.2408
```

Note on interpretation of β_0

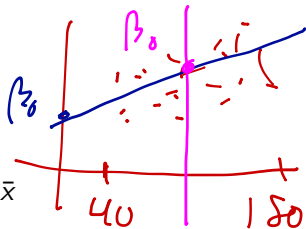
Recall $\beta_0 = E(y|x=0)$

- This often makes no sense in context
- "Centering" x can be useful: $x^* = x - \bar{x}$
- Center by mean, median, minimum, etc
- Effect of centering on slope:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_1^* = \frac{\sum (y_i - \bar{y})}{\sum (x_i - c - \bar{x}^*)}$$

$$= \hat{\beta}_1$$



$$\begin{aligned}\bar{x}^* &= \frac{\sum x_i^*}{n} \\ &= \frac{\sum (x_i - c)}{n} \\ &= \dots\end{aligned}$$

Note on interpretation of β_0, β_1

- The interpretations are sensitive to the scale of the outcome and predictors (in reasonable ways)
- You can't get a better model fit by rescaling variables

$$X_i^* = C \cdot X_i$$

$$\hat{\beta}_1 \neq \hat{\beta}_1^*$$

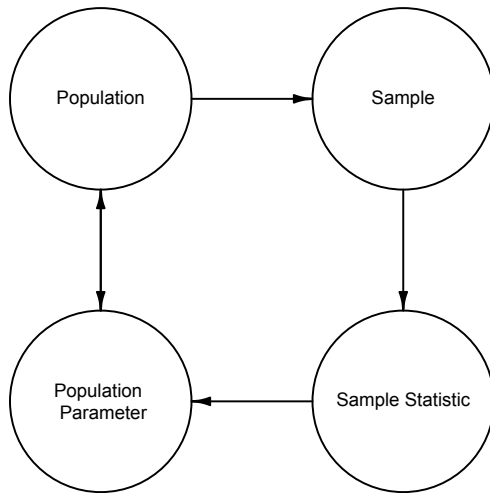
but inference won't
change

R example

```
heights$centeredMheight <- heights$Mheight - mean(heights$Mheight)
centeredLinmod <- lm(Dheight ~ centeredMheight, data = heights)
summary(centeredLinmod)
```

```
##
## Call:
## lm(formula = Dheight ~ centeredMheight, data = heights)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.397 -1.529  0.036   1.492  9.053
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    63.7511    0.0611  1043.1  <2e-16 ***
## centeredMheight  0.5417    0.0260    20.9  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.27 on 1373 degrees of freedom
## Multiple R-squared:  0.241, Adjusted R-squared: 0.24
## F-statistic: 435 on 1 and 1373 DF, p-value: <2e-16
```

Properties of $\hat{\beta}_0, \hat{\beta}_1$



Properties of $\hat{\beta}_0, \hat{\beta}_1$

Estimates are unbiased:

$$E(\hat{\beta}_0) = \beta_0$$

$$E(\hat{\beta}_1) = \beta_1$$

Properties of $\hat{\beta}_0, \hat{\beta}_1$

Variances of estimates

$$\text{Var}(\hat{\beta}_0) = \frac{\bar{x}\sigma^2}{\sum x^2}$$

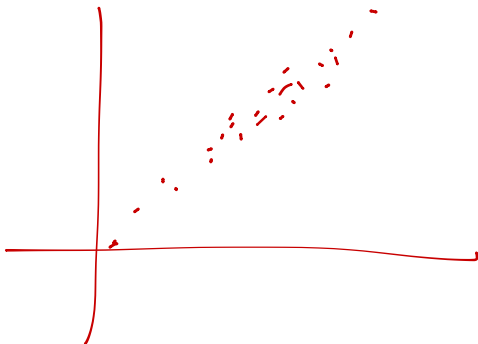
$$\text{Var}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{S_{xx}}$$

where $S_{xx} = \sum (x - \bar{x})^2$

Properties of $\hat{\beta}_0, \hat{\beta}_1$

Note about the variance of β_1 :

- Denominator contains $SS_x = \sum (x_i - \bar{x})^2$
- To decrease variance of $\hat{\beta}_1$, increase variance of x



One slide on multiple linear regression

- Observe data $(y_i, \underline{x_{i1}}, \dots, x_{ip})$ for subjects $1, \dots, n$. Want to estimate $\beta_0, \beta_1, \dots, \beta_p$ in the model

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i; \quad \epsilon_i \stackrel{iid}{\sim} (0, \sigma^2)$$

- Assumptions (residuals have mean zero, constant variance, are independent) are as in SLR
- Notation is cumbersome. To fix this, let

- $\mathbf{x}_i = [1, x_{i1}, \dots, x_{ip}]$ $1 \times p$

- $\boldsymbol{\beta}^T = [\beta_0, \beta_1, \dots, \beta_p]$

- Then $y_i = \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i$

$\boldsymbol{\beta} = p \times 1$ vector

$$= [1 \times p] [p \times 1]$$

Summary

Today's big ideas

- ▶ Simple linear regression definitions
- ▶ Properties of least squares estimates

Coming up soon

- ▶ More on MLR