

STAT 420: Methods of Applied Statistics

Midterm I review

Ruoqing Zhu, Ph.D. <rqzhu@illinois.edu>

Course website: <https://sites.google.com/site/teazrq/teaching/STAT420>

University of Illinois at Urbana-Champaign
February 28, 2017

- Midterm 1 is testing your mathematical understanding of the simple and multiple linear regressions, including:
 - Estimating the coefficients (by direct formula or manipulating matrices)
 - Model fitting assessments
 - Distributions of parameter estimates, confidence intervals and hypothesis testing
 - Prediction intervals
 - Testing multiple parameters
 - R implementation
- A calculator (cannot directly fit linear regression), and a one-sided cheat-sheet (A4 size) are allowed.

- For simple linear regression, one predictor x with intercept:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = r_{xy} \frac{s_y}{s_x}$$

an alternative formula

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{(n-1) s_x^2}$$

- For multiple linear regression, any type of design matrix \mathbf{X} , as long as its full rank,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

This is derived from the normal equation.

- Matrix operations are very important in learning linear regressions, especially for MLR, where no direct formula is given for each individual parameter. They have to be solved jointly using the $(\mathbf{X}^T \mathbf{X})^{-1}$ matrix.
- Another important result is the linear transformation of MVN random variable. If $X \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and $Z = \mathbf{A}_{q \times p} X + \mathbf{b}_{q \times 1}$,

$$Z \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T).$$

- This is the foundation for deriving the distribution of $\hat{\boldsymbol{\beta}}$

- Variance breakdown: $SST = SSR + SSE$.
- The goodness-of-fit statistic: coefficient of determination R^2

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

- Estimating the error variance:

$$\hat{\sigma}^2 = \frac{SSE}{\text{degrees of freedom}}$$

What is the degrees of freedom (if intercept is included)?

- The key result:

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$$

- So the marginal variances of each $\hat{\beta}$ are just the diagonal elements of $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$.
- This can be calculated directly in SLR:

$$\sigma^2(\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 \begin{pmatrix} \frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} & -\frac{\bar{x}}{(n-1)s_x^2} \\ -\frac{\bar{x}}{(n-1)s_x^2} & \frac{1}{(n-1)s_x^2} \end{pmatrix}$$

- To test a single parameter in SLR, replace σ^2 with $\hat{\sigma}^2$:

$$\frac{\hat{\beta}_0 - \beta_0}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2}}} \sim t(n-2)$$

$$\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma} / (\sqrt{(n-1)s_x})} \sim t(n-2)$$

- These results can be used to build CI:

$$\hat{\beta}_0 \pm t_{\alpha/2, n-2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2}}$$

$$\hat{\beta}_1 \pm t_{\alpha/2, n-2} \frac{\hat{\sigma}}{\sqrt{(n-1)s_x}}$$

- For testing a single parameter MLR, find the corresponding variance estimate for the parameter in the $\sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$, and do similar constructions of the t distribution.
- CI's are constructed in a similar fashion.

- We will only test prediction intervals of SLR. There are two types of predictions. Predicting the mean response

$$\hat{\mu}_{\text{new}} \sim \mathcal{N} \left(\beta_0 + \beta_1 x_{\text{new}}, \sigma^2 \left(\frac{1}{n} + \frac{(x_{\text{new}} - \bar{x})^2}{(n-1)s_x^2} \right) \right)$$

with CI

$$\hat{\mu}_{\text{new}} \pm t_{\alpha/2, n-2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x_{\text{new}} - \bar{x})^2}{(n-1)s_x^2}}$$

- CI for a future observed value Y_{new}

$$\hat{\mu}_{\text{new}} \pm t_{\alpha/2, n-2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_{\text{new}} - \bar{x})^2}{(n-1)s_x^2}}$$

Testing multiple parameters

- The F test statistic for testing multiple parameters is given by

$$F = \frac{(\text{SSE}_R - \text{SSE}_F) / q}{\text{SSE}_F / (n - p - 1)}$$

where q is the number of restrictions in the hypothesis test.

- The `R lm()` summary output gives the overall F test for all predictors.
- Then this overall F -statistic for testing all predictors is essentially

$$\frac{\text{SSR}/p}{\text{SSE}/(n - p - 1)} = \frac{R^2}{1 - R^2} \frac{n - p - 1}{p}$$

Testing linear constraints

- A linear constraint is specified by

$$H_0 : \mathbf{A}\beta = c$$

- Apply the linear transformation theorem

$$\mathbf{A}\hat{\beta} \sim \mathcal{N}(\mathbf{A}\beta, \sigma^2 \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A})$$

- Under the Null, we have $\mathbf{A}\beta = c$, so a t -statistic for testing $\mathbf{A}\beta = c$ is

$$\frac{\mathbf{A}\hat{\beta} - c}{\sqrt{\hat{\sigma}^2 \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}}}$$

```
1 > summary( fit )
2 Call:
3 lm(formula = Species ~ Area + Elevation + Nearest + Scrub + Adjacent,
4     data = gala)
5 Coefficients:
6             Estimate Std. Error t value Pr(>|t|)
7 (Intercept)  7.068221   19.154198   0.369  0.715351
8 Area        -0.023938    0.022422  -1.068  0.296318
9 Elevation    0.319465    0.053663   5.953 3.82e-06 ***
10 Nearest     0.009144    1.054136   0.009  0.993151
11 Scrub       -0.240524    0.215402  -1.117  0.275208
12 Adjacent    -0.074805    0.017700  -4.226  0.000297 ***
13 ---
14 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
15                  1
16 Residual standard error: 60.98 on 24 degrees of freedom
17 Multiple R-squared:  0.7658, Adjusted R-squared:  0.7171
18 F-statistic: 15.7 on 5 and 24 DF, p-value: 6.838e-07
```

- Questions?