

STAT 420: Methods of Applied Statistics

Inference of Linear Regressions

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

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

University of Illinois at Urbana-Champaign
February 16, 2017

Linear Regressions

- We have already learned how to estimate the parameters β of a linear regression, and to assess the goodness-of-fit R^2 .
- Now we are going to learn how make inference on the parameters.
- Recall that we used simulation study to estimate the p -values, now we will derive a rigorous formula for testing:
 - Distribution of the parameter estimates
 - Testing for individual coefficient
 - Confidence interval for individual coefficient
 - Confidence interval for predicting new subject
 - Testing for multiple coefficients and the entire model
- We will derive the results for both SLR and MLR.
- Simulation study can still help us understand these results.

- For now, we will make several assumption:
- The error terms ϵ 's follow i.i.d. Normal distribution, with mean 0 and variance σ^2 :

$$\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2).$$

- We also assume that the covariates X_i 's are fixed (not a random variable), with no measurement errors.
- The random errors ϵ are independent of the covariates \mathbf{X} .
- The design matrix \mathbf{X} has full rank.

- Some consequences:
 - The outcome variable $Y_i = X_i^T \beta + \epsilon_i$ also follows a normal distribution:

$$Y_i \sim \mathcal{N}(X_i^T \beta, \sigma^2),$$

- The error vector ϵ follows a multivariate normal distribution:

$$\epsilon \sim \mathcal{N}(\mathbf{0}_{n \times 1}, \sigma^2 \mathbf{I}_{n \times n})$$

where $\mathbf{I}_{n \times n}$ is a diagonal matrix with size n , and all diagonal elements are 1.

- **Question:** What is the distribution of the outcome vector \mathbf{Y} ?

Distribution of $\hat{\beta}$

- Consider the vector \mathbf{Y} , we want to find its distribution.
- Since $\mathbf{Y} = \mathbf{X}\beta + \epsilon$, we have

$$\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2\mathbf{I}),$$

independent normal variables, with means $X_i^T\beta$, $i = 1, \dots, n$ and variance all equal to σ^2 .

- Then the next step is to derive the distribution of $\hat{\beta}$, which is a **linear transformation** of \mathbf{Y} .
- Recall the multivariate Normal distribution results, since

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

is a linear transformation of the vector \mathbf{Y} , where we let the transformation matrix $\mathbf{A} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$. Then, what is the distribution of $\mathbf{A}\mathbf{Y}$?

- When the mean of \mathbf{Y} is $\boldsymbol{\mu}$, and variance-covariance matrix is $\boldsymbol{\Sigma}$, the distribution of $\mathbf{A}\mathbf{Y}$ is

$$\mathcal{N}(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

- In our case, $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta}$, $\boldsymbol{\Sigma} = \sigma^2\mathbf{I}$, and $\mathbf{A} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top$, so

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{Y} \\ &\sim \mathcal{N}\left((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}\boldsymbol{\beta}, (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\sigma^2\mathbf{I}((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top)^\top\right) \\ \implies \hat{\boldsymbol{\beta}} &\sim \mathcal{N}\left(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}\right)\end{aligned}$$

- **Properties** of the distribution $\mathcal{N}(\beta, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1})$?
 1. $E(\hat{\beta}) = \beta$, so $\hat{\beta}$ is an unbiased estimator of the parameter β .
 2. $\mathbf{X}^T\mathbf{X}$ grows in the same rate as the sample size n (essentially each element of $\mathbf{X}^T\mathbf{X}$ is a sum of n terms. So $(\mathbf{X}^T\mathbf{X})^{-1}$ is in the order of $\mathcal{O}(1/n)$. Hence the variance of goes to 0, i.e., $\hat{\beta}$ becomes more accurate as we collect more samples.
- However, we still don't know what σ^2 is. What to do?
 - We have $\hat{\sigma}^2$
 - After replacing σ^2 with $\hat{\sigma}^2$, is it still normal? z -test vs. t -test?

Distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$ in SLR

- In a SLR, we can explicitly write out the inverse matrix:

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

- What is the variance of $\hat{\beta}_1$?

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{(n-1)s_x^2}$$

$$\Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\sigma/(\sqrt{(n-1)}s_x)} \sim \mathcal{N}(0, 1)$$

- Try our simulation study to confirm this.

Distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$ in SLR

- Replace σ^2 with $\hat{\sigma}^2 = \frac{SSE}{n-2}$. Why $n-2$? What's its distribution?

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2$$

- Then we have

$$\begin{aligned}\frac{\mathcal{N}(0,1)}{\sqrt{\chi_r^2/r}} &\sim t(r) \\ \frac{\hat{\beta}_1 - \beta_1}{\sigma/(\sqrt{(n-1)s_x})} / \sqrt{\frac{\hat{\sigma}^2}{\sigma^2}} &\sim t(n-2) \\ \frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}/(\sqrt{(n-1)s_x})} &\sim t(n-2)\end{aligned}$$

- Similarly, we have

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

CI for $\hat{\beta}_0$ and $\hat{\beta}_1$ in SLR

- Now we know the distributions of both $\hat{\beta}_0$ and $\hat{\beta}_1$, we can construct the distributions of them accordingly.
- The $(1 - \alpha)100\%$ two-sided confidence intervals for $\hat{\beta}_0$ and $\hat{\beta}_1$:

$$\hat{\beta}_0 \pm t_{\alpha/2, n-2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2}}$$
$$\hat{\beta}_1 \pm t_{\alpha/2, n-2} \frac{\hat{\sigma}}{\sqrt{(n-1)s_x}}$$

- Similarly, p -values can be computed.

Example

- Use the `cheddar` data to fit SLR with Lactic, and calculate:
 - The variance of both parameter estimates
 - The 95% confidence intervals
 - p -value for testing $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$
- See `R` code.
- Fit the model using Acetic and construct a 90% CI for the parameter estimate. Does the interval include 8?

A: Yes; B: No;

Example

- The `R` function for calculating the CIs:

```
1 > fit = lm(taste ~ Lactic, data= cheddar)
2 > confint(fit, level = 0.95)
3           2.5 %      97.5 %
4 (Intercept) -51.53573 -8.181935
5 Lactic       22.99928 52.440613
```

- The p -values can be found in

```
1 > summary(fit)
2 Coefficients:
3           Estimate Std. Error t value Pr(>|t|)
4 (Intercept)  -29.859     10.582  -2.822  0.00869 **
5 Lactic        37.720      7.186   5.249 1.41e-05 ***
6 ---
7 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- You still need to know how to calculate them by hand if given sufficient information

Example

- Calculate the t -statistic and p -value based on the following `lm()` fitting results:

```
1 Coefficients:
2             Estimate Std. Error t value Pr(>|t|)
3 (Intercept)  4.84691    0.37422  12.952  <2e-16 ***
4 light       -0.10712    0.07419       
5 -----
6 Signif. codes:
7 0   ***  0.001  **  0.01  *  0.05  .  0.1      1
8 Residual standard error: 0.2875 on 45 degrees of freedom
9 Multiple R-squared:  0.04427, Adjusted R-squared:  0.02304
10 F-statistic: 2.085 on 1 and 45 DF,  p-value: 0.1557
```

- Calculate the sum of squares total, SST.

Inference about predictions: μ_{new}

- Suppose we have a new subject with Lactic = x_{new} , can we predict its **mean** taste score? How accurate is that prediction?
- Let $\mu_{\text{new}} = \beta_0 + \beta_1 x_{\text{new}}$ be the mean taste score of cheddar if the Lactic level is x
- The prediction is $\hat{\mu}_{\text{new}} = \hat{\beta}_0 + \hat{\beta}_1 x_{\text{new}}$, i.e.

$$\hat{\mu}_{\text{new}} = (1, x_{\text{new}}) \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}$$

and we already know the distribution of $\hat{\beta}$. This is yet another linear transformation of MVN variables.

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

- The CI for μ_{new} is

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

Example

- Construct a 90% CI for mean taste with Lactic = 1

```
1 > xnew <- data.frame(Lactic = 1)
2 > predict.lm(fit, xnew, interval = c("confidence"), level =
    0.9)
3      fit      lwr      upr
4 1 7.861116 1.341625 14.38061
```

- Construct a 95% CI for mean taste with Lactic = 2

```
1 > xnew <- data.frame(Lactic = 2)
2 > predict.lm(fit, xnew, interval = c("confidence"), level =
    0.95)
3      fit      lwr      upr
4 1 45.58106 36.26624 54.89588
```

- The data `star` in the `faraway` package is trying to model the surface temperature of a star with its light intensity. The stars are in the star cluster CYG OB1, which is in the direction of Cygnus.
- Fit linear regression and test for $x_{\text{new}} = 6$ at 95% confidence level:

$$H_0 : \mu_{\text{new}} = 4.06 \quad \text{vs.} \quad H_1 : \mu_{\text{new}} > 4.06$$

- Do we reject the hypothesis?

A: Yes; B: No;

- What if we want to predict the observed value Y_{new} for a new subject with x_{new} , there will be additional variations.
- Since $\hat{Y}_{\text{new}} = \hat{\mu}_{\text{new}} + \epsilon_{\text{new}}$, this is only adding additional σ^2 to the variance component. Also ϵ_{new} is independent of everything else.
- Hence the CI for a future observed value Y_{new}

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

Example

- Construct a 90% CI for future taste score with Lactic = 2

```
1 > xnew <- data.frame(Lactic = 2)
2 > predict.lm(fit, xnew, interval = c("prediction"), level =
    0.95)
3      fit      lwr      upr
4 1 45.58106 19.78216 71.37997
```

Example

- Suppose the following data are observed.

```
1 > x = c(-1, -1, -1, -1, 1, 1, 1, 1)
2 > y = c(1.3, 0.8, 1.2, 0.6, 2.3, 2.5, 1.8, 1.6)
```

- Construct a 90% confidence interval for the $\hat{\beta}_1$ estimation.
- Predict a future outcome value at $x = 0$.

The Gauss-Markov Theorem

- We know that $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ and derived its properties. However, is this the best estimator? Is it possible that we can have some other estimators that are also unbiased and more accurate (smaller variance)?
- Turns out that this is the best we can do if $\hat{\beta}$ is unbiased and also a “linear estimator”, i.e., $\hat{\beta} = \mathbf{C} \mathbf{y}$ for some matrix \mathbf{C} .
- This is guaranteed by the Gauss-Markov Theorem.
- $\hat{\beta}$ is the BLUE (best linear unbiased estimate).