

Review of Normal Linear Regression¹

STA312 Fall 2023

¹See last slide for copyright information.

Multiple Linear Regression with normal errors

$$y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i,$$

where

$\beta_0, \dots, \beta_{p-1}$ are unknown constants.

$x_{i,j}$ are known constants.

$\epsilon_1, \dots, \epsilon_n$ are independent $N(0, \sigma^2)$ random variables.

σ^2 is an unknown constant.

y_1, \dots, y_n are observable random variables.

This implies y_i are independent $N(\mu_i, \sigma^2)$, with

$$\mu_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_{p-1} x_{i,p-1}.$$

Another way to think about it

- y_1, \dots, y_n are independent $N(\mu_i, \sigma^2)$
- $\mu_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_{p-1} x_{i,p-1}$
- We have substituted a regression function for the location parameter.
- Anything that makes the regression function larger or smaller shifts the distribution to the right or left.
- In normal regression we are always talking about μ_i as the expected value (which it is), but more generally we mean the location.

Log Likelihood

$$\begin{aligned}\ell(\boldsymbol{\beta}, \sigma^2) &= \log \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}} \\ &= -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2\end{aligned}$$

To make this quantity as *large* as possible over all $\beta_0, \dots, \beta_{p-1}$,

- Make $\sum_{i=1}^n (y_i - \mu_i)^2$ as *small* as possible.
- That is, minimize $\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i,1} - \dots - \beta_{p-1} x_{i,p-1})^2$.
- This is a familiar problem – least squares.
- So the least-squares estimates for multiple regression are the same as the MLEs.

Vocabulary

$$y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i$$

- Explanatory variables are x
- Response variable is y .

“Control” means hold constant

- Regression model with four explanatory variables.
- Hold x_1 , x_2 and x_4 constant at some fixed values.

$$\begin{aligned}E(Y|\mathbf{X} = \mathbf{x}) &= \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \beta_4x_4 \\ &= (\beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_4x_4) + \beta_3x_3\end{aligned}$$

- The equation of a straight line with slope β_3 .
- Values of x_1 , x_2 and x_4 affect only the intercept.
- So β_3 is the rate at which $E(Y|\mathbf{x})$ changes as a function of x_3 with all other variables held constant at fixed levels.
- *According to the model.*

More vocabulary

$$E(Y|\mathbf{X} = \mathbf{x}) = (\beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_4x_4) + \beta_3x_3$$

- If $\beta_3 > 0$, describe the relationship between x_3 and (expected) y as “positive,” controlling for the other variables. If $\beta_3 < 0$, negative.
- Useful ways of saying “controlling for” or “holding constant” include
 - Allowing for
 - Correcting for
 - Taking into account

Categorical Explanatory Variables (Unordered categories)

Example: $Y = \beta_0 + \beta_1 x + \epsilon$

- $x = 1$ means Drug, $x = 0$ means Placebo.
- Population mean is $E(Y|x) = \beta_0 + \beta_1 x$.
- For patients getting the drug, mean response is $E(Y|x = 1) = \beta_0 + \beta_1$
- For patients getting the placebo, mean response is $E(Y|x = 0) = \beta_0$
- And β_1 is the difference between means, the average treatment effect.

More than Two Categories

Suppose a study has 3 treatment conditions. For example Group 1 gets Drug 1, Group 2 gets Drug 2, and Group 3 gets a placebo, so that the Explanatory Variable is Group (taking values 1,2,3) and there is some Response Variable Y (maybe response to drug again).

Why is $E[Y|X = x] = \beta_0 + \beta_1 x$ (with $x = \text{Group}$) a silly model?

Indicator Dummy Variables With Intercept

Example: $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$

- $x_1 = 1$ if Drug A, zero otherwise
- $x_2 = 1$ if Drug B, zero otherwise
- $E[Y|\mathbf{X} = \mathbf{x}] = \beta_0 + \beta_1 x_1 + \beta_2 x_2$.
- Fill in the table.

Drug	x_1	x_2	$E(Y \mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
A			$\mu_1 =$
B			$\mu_2 =$
Placebo			$\mu_3 =$

Answer

- $x_1 = 1$ if Drug A, zero otherwise
- $x_2 = 1$ if Drug B, zero otherwise
- $E[Y|\mathbf{X} = \mathbf{x}] = \beta_0 + \beta_1x_1 + \beta_2x_2$.

Drug	x_1	x_2	$E(Y \mathbf{x}) = \beta_0 + \beta_1x_1 + \beta_2x_2$
A	1	0	$\mu_1 = \beta_0 + \beta_1$
B	0	1	$\mu_2 = \beta_0 + \beta_2$
Placebo	0	0	$\mu_3 = \beta_0$

Regression coefficients are contrasts with the category that has no indicator – the *reference category*.

Indicator dummy variable coding with intercept

- With an intercept in the model, need $p - 1$ indicators to represent a categorical explanatory variable with p categories.
- If you use p dummy variables and an intercept, trouble.
- Regression coefficients are contrasts with the category that has no indicator.
- Call this the *reference category*.

What null hypotheses would you test?

Drug	x_1	x_2	$E(Y \mathbf{x}) = \beta_0 + \beta_1x_1 + \beta_2x_2$
A	1	0	$\mu_1 = \beta_0 + \beta_1$
B	0	1	$\mu_2 = \beta_0 + \beta_2$
Placebo	0	0	$\mu_3 = \beta_0$

- Is the effect of Drug A different from the placebo? $H_0 : \beta_1 = 0$
- Is Drug A better than the placebo? $H_0 : \beta_1 = 0$
- Did Drug B work? $H_0 : \beta_2 = 0$
- Did experimental treatment have an effect? $H_0 : \beta_1 = \beta_2 = 0$
- Is there a difference between the effects of Drug A and Drug B?
 $H_0 : \beta_1 = \beta_2$

Now add a quantitative explanatory variable (covariate)

Covariates often come first in the regression equation

- $x_1 = 1$ if Drug A, zero otherwise
- $x_2 = 1$ if Drug B, zero otherwise
- $x_3 = \text{Age}$
- $E[Y|\mathbf{X} = \mathbf{x}] = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3$.

Drug	x_1	x_2	$E(Y \mathbf{x}) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3$
A	1	0	$\mu_1 = (\beta_0 + \beta_1) + \beta_3x_3$
B	0	1	$\mu_2 = (\beta_0 + \beta_2) + \beta_3x_3$
Placebo	0	0	$\mu_3 = \beta_0 + \beta_3x_3$

Parallel regression lines.

More comments

Drug	x_1	x_2	$E(Y \mathbf{x}) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3$
A	1	0	$\mu_1 = (\beta_0 + \beta_1) + \beta_3x_3$
B	0	1	$\mu_2 = (\beta_0 + \beta_2) + \beta_3x_3$
Placebo	0	0	$\mu_3 = \beta_0 + \beta_3x_3$

- If more than one covariate, parallel regression planes.
- Non-parallel (interaction) is testable.
- “Controlling” interpretation holds.
- In an experimental study, quantitative covariates are usually just observed.
- Could age be related to drug?
- Good covariates reduce MSE, make testing of categorical variables more sensitive.

Partitioning sums of squares

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_{p-1} x_{i,p-1}$$

- Variation to explain: **Total Sum of Squares**

$$\text{SSTO} = \sum_{i=1}^n (y_i - \bar{y})^2$$

- Variation that is still unexplained: **Error (or Residual) Sum of Squares**

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

- Variation that is explained: **Regression (or Model) Sum of Squares**

$$\text{SSR} = \text{SSTO} - \text{SSE} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

R^2 : Proportion of variation in y that is explained

- $SSTO = SSR + SSE$
- Proportion of variation in the response variable that is explained by the explanatory variable

$$R^2 = \frac{SSR}{SSTO}$$

- For a simple regression, same as the square of the correlation coefficient: $r^2 = R^2$

Hypothesis Testing (Standard tests with normal errors)

$$y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i$$

- Overall F -test for all the explanatory variables at once:
 $H_0 : \beta_1 = \beta_2 = \cdots = \beta_{p-1} = 0$.
- t -tests for each regression coefficient: Controlling for all the others, does that explanatory variable matter? $H_0 : \beta_j = 0$.
- Test a collection of explanatory variables controlling for another collection: $H_0 : \beta_2 = \beta_3 = \beta_5 = 0$.
- Example: Controlling for mother's education and father's education, are (any of) total family income, assessed value of home and total market value of all vehicles owned by the family related to High School GPA?
- Most general: Testing whether sets of linear combinations of regression coefficients differ from specified constants. $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{h}$.

Full versus Restricted Model

Restricted by H_0

- You have 2 sets of variables, A and B . Want to test B controlling for A .
- Fit a model with both A and B : Call it the *Full Model*.
- Fit a model with just A : Call it the *Restricted Model*. It's restricted by the null hypothesis.
 $R_F^2 \geq R_R^2$.
- The F -test for full versus restricted is a likelihood ratio test (exact).

When you add the r explanatory variables in set B to the model, R^2 can only go up

By how much? Basis of the F test.

$$\begin{aligned} F &= \frac{(R_F^2 - R_R^2)/r}{(1 - R_F^2)/(n - p)} \\ &= \frac{(SSR_F - SSR_R)/r}{MSE_F} \stackrel{H_0}{\sim} F(r, n - p) \end{aligned}$$

General Linear Test of $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{h}$

\mathbf{L} is $r \times p$, rows linearly independent

$$F = \frac{(\mathbf{L}\hat{\boldsymbol{\beta}} - \mathbf{h})^\top (\mathbf{L}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{L}^\top)^{-1} (\mathbf{L}\hat{\boldsymbol{\beta}} - \mathbf{h})}{r \text{MSE}_F}$$

$$\stackrel{H_0}{\sim} F(r, n - p)$$

- Equal to full-restricted formula.
- Numerator looks like a Wald statistic, and it's no accident.

Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The \LaTeX source code is available from the course website:

<http://www.utstat.toronto.edu/brunner/oldclass/312f23>