

A bit of regression: Quick and very applied¹
STA2101 Fall 2019

¹See last slide for copyright information.

Fixed Effects Linear Regression

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- \mathbf{X} is an $n \times p$ matrix of known constants.
- $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown constants.
- $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, where $\sigma^2 > 0$ is an unknown constant.

- $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$
- $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$
- $\mathbf{e} = (\mathbf{y} - \hat{\mathbf{y}})$

Comparing scalar and matrix form

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

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}$$
$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 14.2 & \cdots & 1 \\ 1 & 11.9 & \cdots & 0 \\ 1 & 3.7 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 6.2 & \cdots & 1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Vocabulary

- 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

Partitioning Sums of Squares

$$\begin{aligned} SST &= SSR + SSE \\ \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \end{aligned}$$

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

Categorical Explanatory Variables

Unordered categories

- $X = 1$ means Drug, $X = 0$ means Placebo.
- Population mean is $E(Y|X = x) = \beta_0 + \beta_1 x$.
- 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

- $x_1 = 1$ 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 differences from 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 if there is random assignment to drug?
- Good covariates reduce MSE, make testing of categorical variables more sensitive.

Hypothesis Testing

Standard tests when errors are normal

- Overall F -test for all the explanatory variables at once
 $H_0 : \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$
- 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*, or the *Unrestricted Model*.
- Fit a model with just A : Call it the *Restricted Model*.
 $R_F^2 \geq R_R^2$.
- The F -test is a likelihood ratio test (exact).

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

By how much? Basis of the F test.

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

Strength of Relationship: Change in R^2 is not enough

$$\begin{aligned} F &= \frac{(R_F^2 - R_R^2)/r}{(1 - R_F^2)/(n - p)} \\ &= \left(\frac{n - p}{r} \right) \left(\frac{a}{1 - a} \right) \end{aligned}$$

where

$$a = \frac{R_F^2 - R_R^2}{1 - R_R^2} = \frac{rF}{n - p + rF}$$

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.

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/2101f19>