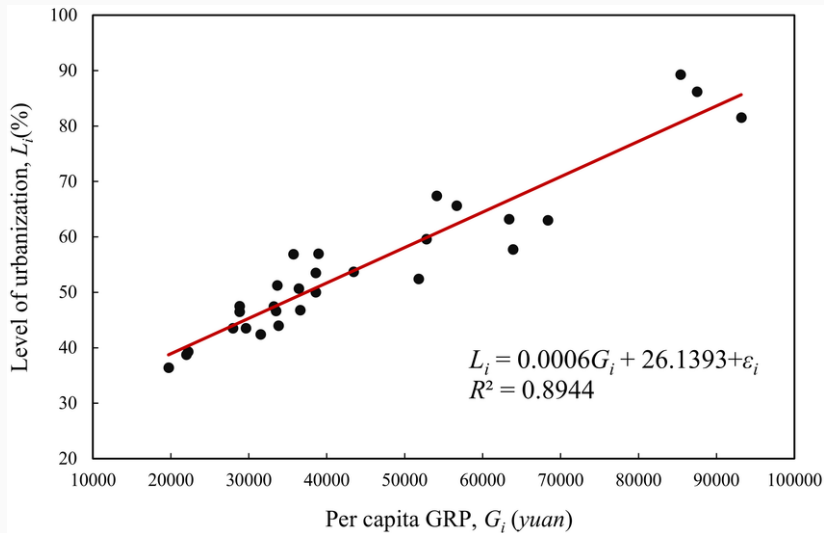
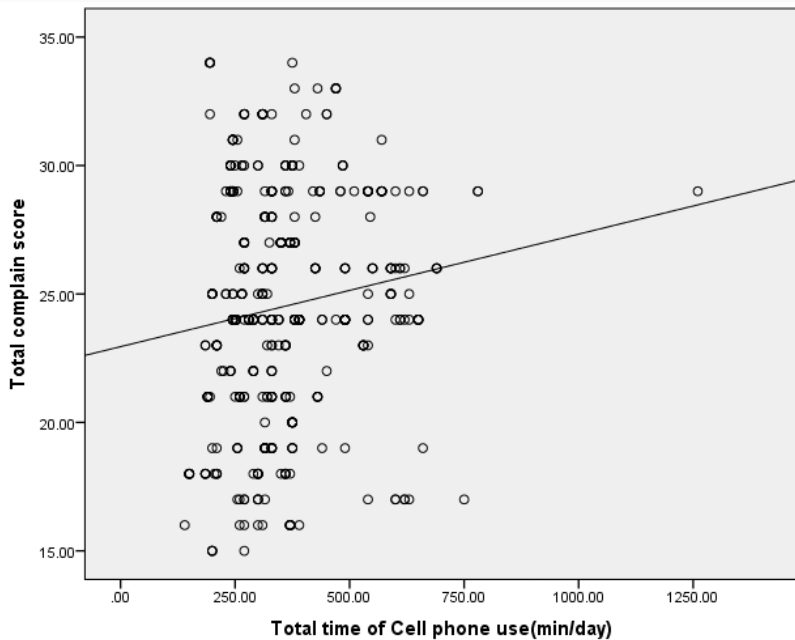


Introduction to Regression







THE REGRESSION ALGORITHM

Idea: Draw a line given a sample of data

What does this line *mean*?

Under what assumptions?

THE MECHANICS OF OLS REGRESSION

- The OLS algorithm
- Statistical assumptions
- Statistical guarantees

Unit Plan

GOALS OF THIS WEEK

At the end of this week, you will:

1. Understand that regression is the plug-in estimator of the best linear predictor (BLP).
2. Understand overall model fit and use an F-test to assess whether a candidate model is performing better than a baseline model.
3. Understand how to appropriately interpret regression coefficients.
4. Lay a foundation for careful interpretation.

Reading: The Golem of Prague

READING: THE GOLEM OF PRAGUE

This is a placeholder for a reading call. We're just placing it here for organization.

- Read Sections 1.0 and 1.1 of *Statistical Rethinking*, which we have provided a copy of in PDF form from the publisher.
- *Statistical Rethinking* is a great book and reference that you should consider later in your data science and statistics path.

Regression, a Statistical Golem

REGRESSION, A STATISTICAL GOLEM

- Regression—like all models—is a tool.
- We put tools to use toward a data scientific purpose; however, tools are only tools.
- Use of a straight-edge, scale, and T-square doesn't make one an architect any more than use of {insert language} or {insert technique} makes one a data scientist.

THE MACHINERY OF OLS REGRESSION

- OLS regression is a plug-in estimator for the best linear predictor (BLP).
- The BLP is the lowest mean squared error (MSE) estimator, out of all linear functions.

APPLICATIONS OF OLS REGRESSION

Regression is fantastically versatile

- Under some circumstances, regression has explainable internal weights (coefficients) that are of interest.
- Under other circumstances, regression identifies causal effects.
- Under many circumstances, regression is the *de facto* baseline estimator.

Elements of a Linear Model

ELEMENTS IN A LINEAR MODEL

Linear model A.K.A. linear predictor

A **linear model** is a representation of a random variable (Y) as a linear function of other random variables

$$\mathbf{X} = (X_1, X_2, \dots, X_k).$$

Y	X_1, X_2, \dots, X_k
Target	Features
Outcome	Predictors
Dependent variable	Independent variables
Output	Inputs
Response	Controls
Left-hand side (LHS)	Right-hand side (RHS)
\vdots	\vdots

THE LINEAR MODEL FORMULA

The linear model formula

$$\begin{aligned}\hat{Y} &= g(X_1, X_2, \dots, X_k) \\ &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k\end{aligned}$$

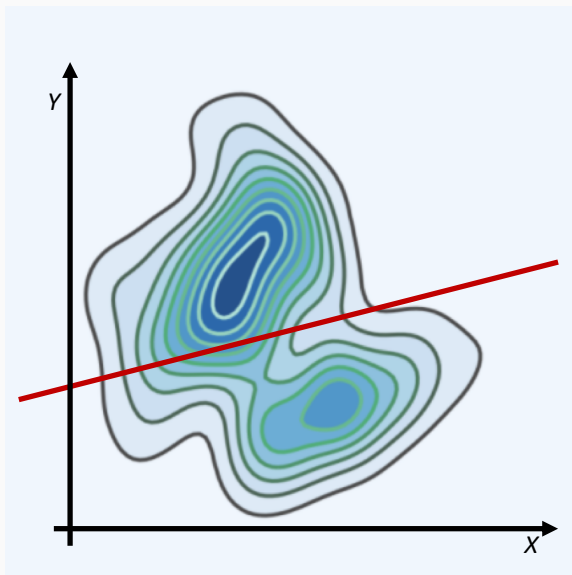
A LINEAR MODEL EXAMPLE

Brunch in Berkeley

$$\widehat{\text{Avocados}} = 2 + 1 \cdot \text{Lemons} + 2 \cdot \text{Loaves_Bread}$$

Review: Outcome, Prediction, and Error

REVIEW: OUTCOME, PREDICTION, AND ERROR



Concept Check: Making Predictions with a Linear Model

CONCEPT CHECK: MAKING PREDICTIONS WITH A LINEAR MODEL

- Students will be given a model and (x,y) and compute prediction and error.

Metric Inputs

INTERPRETING MODEL WEIGHTS (COEFFICIENTS)

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

Interpretation of coefficients

- If X_i changes by ΔX_i units, the predicted value of the target, \hat{Y} changes by $\beta_i \cdot \Delta X_i$ units.
- If X_i and X_j change by ΔX_i and ΔX_j respectively, then the predicted value of the target, \hat{Y} changes by $(\beta_i \Delta X_i) + (\beta_j \Delta X_j)$.

Ceteris paribus: all else equal

INTERPRETING MODEL COEFFICIENTS: EXAMPLE

Does this model say peacocks with longer tails fly slower?

$$\text{air_speed} = 4.3 - 1.2 \cdot \text{tail_length} + 0.8 \cdot \text{muscle_mass}$$



Photo by Thimindu Goonatillake CC BY-SA 2.0

Categorical Inputs

CATEGORIAL INPUTS, PART I

What if, rather than *numeric* inputs, we had *categorical* inputs?

- Information that says it belongs to one category or another, but doesn't provide a value to that category?
- **Reminder:** There are four “levels” of information – (1) Categorical; (2) Ordinal; (3) Interval; (4) Ratio

CATEGORICAL INPUTS, PART II

How is this represented in a model?

- The model aims only to distinguish one category from another, and so switches from stacked *labels* to *one-hot encoded*, or *dummy* variables.
- Practically, in this model, one of the labels is identified as the *baseline*, *default*, or *omitted* category
- Indicators for *alternate* levels mark changes from the baseline to the alternate level.

CATEGORICAL INPUTS, PART II

Learnosity: Interpreting Model Weights

This is a placeholder for a Learnosity activity.

- have to change this, so it's about interpretation, not prediction
- In this activity, students will be given a fitted model that conforms with the data that they have for the peacock
- They will make predictions *first* from the data, and then
- Second, from newly created data to see how the predictions change

Part 2: Selecting a Linear Model with OLS

OLS is Regression for Estimating the BLP

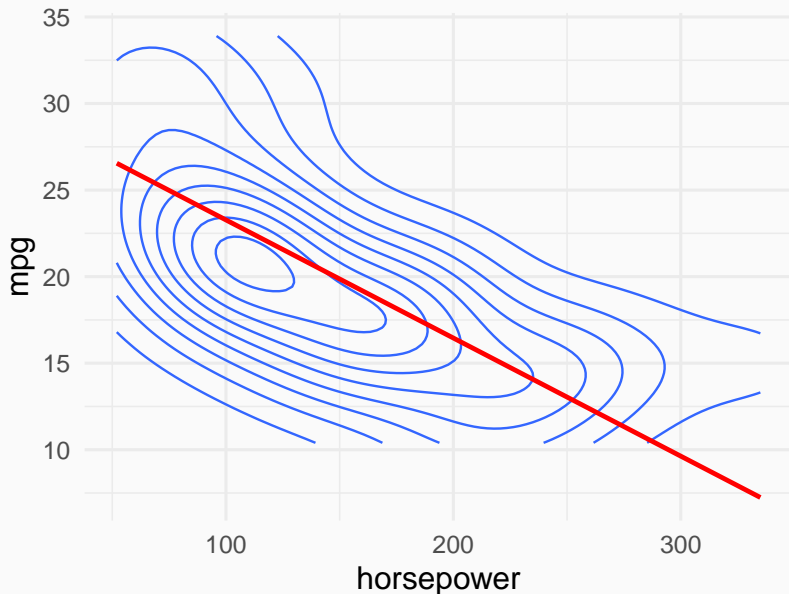
Linear regression: an algorithm for fitting a linear model given a sample of data

- Ordinary least squares (OLS) regression
- Quantile regression
- Regularized regression
 - Lasso
 - Ridge regression

Ordinary least squares (OLS) regression

- The most well-known type of linear regression
- A foundation for many other types of regression
- Key goal: estimating the best linear predictor (BLP)

THE BLP MINIMIZES EXPECTED SQUARED ERROR



THE BEST LINEAR PREDICTOR

The BLP (population regression function)

The best linear predictor is defined by the function

$$g(\mathbf{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k$$

where $\beta = (\beta_0, \beta_1, \dots, \beta_k)$ are chosen to minimize the expected squared error.

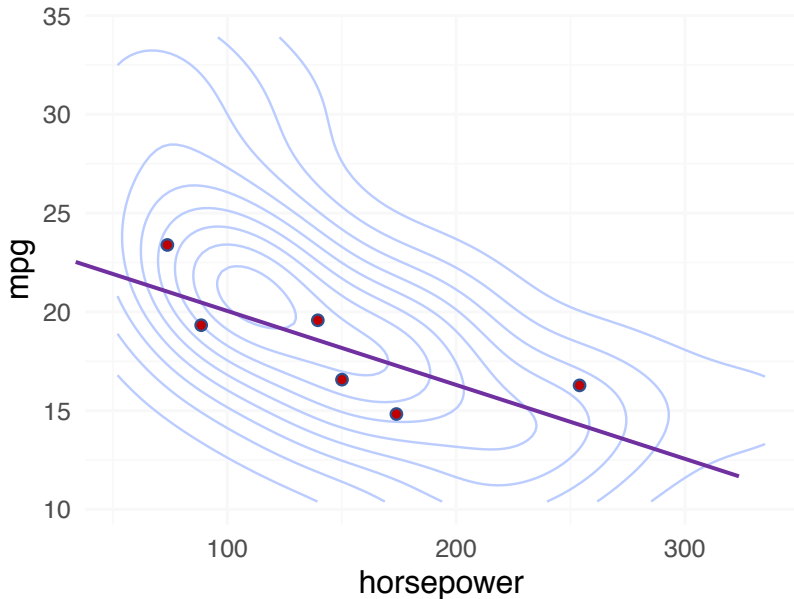
$$\min_{(b_0, \dots, b_k)} E[(Y - (b_0 + b_1 X_1 + \dots + b_k X_k))^2]$$

GREAT THINGS ABOUT THE BLP

The best linear predictor...

- minimizes MSE out of all linear models.
- captures an infinitely complex distribution in a few parameters.
- can be estimated with much less data compared to a probability density.
- is easy to reason about.
- is easy to communicate to others, helping knowledge advance.
- has a closed form solution that is relatively easy to work with.

APPLYING THE PLUG-IN PRINCIPLE



OLS REGRESSION IS THE BLP PLUG-IN ESTIMATOR

Plug-in strategy

The OLS regression line is given by

$$\hat{g}(\mathbf{X}) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \cdots + \hat{\beta}_k X_k$$

where $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k)$ are given by

$$\min_{(b_0, \dots, b_k)} \frac{1}{n} \sum_{i=1}^n (Y_i - (b_0 + b_1 X_{[1]i} + \dots + b_k X_{[k]i}))^2.$$

AN ANALOGY WITH THE MEAN

Given only Y	Given Y and X
$E[Y]$ minimizes MSE out of all numbers.	the BLP minimizes MSE out of all linear models.
We can't compute $E[Y]$ without knowing the distribution.	We can't compute the BLP without knowing the distribution.
\bar{X} is the plug-in estimator for $E[Y]$.	OLS is the plug-in estimator for the BLP.

COMING UP SOON...

We still have to:

- solve the minimization problem
- show that OLS is *consistent* for the BLP

Learnosity: You Minimize It!

LEARNOSITY: YOU MINIMIZE IT!

Note: This is a Learnosity Activity. We're just placing it here for organization.

This is the activity that is currently coded

`regression_fit_2d_exercise/`. **Note that we would like to expand this to ask students to work through several examples in a row. Presently we have a single example made; expanding this to a broader set is relatively easy. In the expanded set, we should ensure that we have some complexity in the data—e.g., a sine curve.**

Reading: OLS Regression Estimates the BLP

READING: LINEAR REGRESSION IS A PLUG-IN ESTIMATOR FOR THE BLP

Note: this is a reading call, we're just placing it here for organization.

Read pages 143–147 of *Foundations of Agnostic Statistics*.

Choosing Assumptions for OLS Regression

WHEN DOES OLS "WORK"?

The OLS regression line is

$$\hat{g}(\mathbf{X}) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \cdots + \hat{\beta}_k X_k$$

where $\hat{\beta}$ are chosen to minimize squared residuals.

WHEN DOES OLS "WORK"?

The OLS regression line is

$$\hat{g}(\mathbf{X}) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \cdots + \hat{\beta}_k X_k$$

where $\hat{\beta}$ are chosen to minimize squared residuals.

Different assumptions



Different statistical guarantees

More data \implies less restrictive assumptions

More data \implies easier to assess assumptions

OLS IN A LARGE SAMPLE

The large-sample model (not an official name)

Just two assumptions:

- I.I.D.
- Unique BLP exists

Asymptotic behavior as $n \rightarrow \infty$ provides considerable guarantees.

The classical linear model

- A parametric model—fully specifies $f_{Y|X}$
- Traditional starting point for regression
- Even with extensive transformations, may be hard to justify assumptions

Guarantees come from strict assumptions.

OLS IN VERY SMALL SAMPLES

Special difficulties when $n < \sim 15$

- No help from asymptotics
- Not enough data to assess CLM

Randomization inference

- A framework for testing (restrictive) null hypotheses

RULES OF THUMB FOR OLS ASSUMPTIONS

Rules of thumb

In general, you might reason about data and regression models in the following way.

Sample size	Required assumptions
$100 \leq n$	Large-sample linear model
$15 \leq n < 100$	Classical linear model
$5 \leq n < 15$	Randomization inference

The Bivariate OLS Solution

Existing Content (in distinct format)

9.5 Deriving the Bivariate OLS Estimators

Consistency of Bivariate OLS Under the Large-Sample Model

Continuous mapping theorem: Let $(S^{(1)}, S^{(2)}, S^{(3)}, \dots)$ and $(T^{(1)}, T^{(2)}, T^{(3)}, \dots)$ be sequences of random variables. Let $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a function that is continuous at $(a, b) \in \mathbb{R}^2$. If $S^{(n)} \xrightarrow{p} a$ and $T^{(n)} \xrightarrow{p} b$, then $g(S^{(n)}, T^{(n)}) \xrightarrow{p} g(a, b)$

Assumptions: 1) I.I.D. 2) Unique BLP exists ($V(X) > 0$)

Continuous mapping theorem: Let $(S^{(1)}, S^{(2)}, S^{(3)}, \dots)$ and $(T^{(1)}, T^{(2)}, T^{(3)}, \dots)$ be sequences of random variables. Let $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a function that is continuous at $(a, b) \in \mathbb{R}^2$. If $S^{(n)} \xrightarrow{p} a$ and $T^{(n)} \xrightarrow{p} b$, then $g(S^{(n)}, T^{(n)}) \xrightarrow{p} g(a, b)$

Assumptions: 1) I.I.D. 2) Unique BLP exists ($V(X) > 0$)

$$x^{(n)} = (x_1, x_2, \dots, x_n) \quad y^{(n)} = (y_1, y_2, \dots, y_n)$$

$$S^{(n)} = \widehat{\text{cov}}(x^{(n)}, y^{(n)}) \quad T^{(n)} = \widehat{V}(x^{(n)})$$

$$\hat{\beta}_1^{(n)} = S^{(n)} / T^{(n)}$$

$$S^{(n)} \xrightarrow{p} \text{cov}[X, Y], \quad T^{(n)} \xrightarrow{p} V[X]$$

$g(c, d) = c/d$ is continuous where $d \neq 0$

$$\hat{\beta}_1^{(n)} = g(S^{(n)}, T^{(n)}) \xrightarrow{p} g(\text{cov}[X, Y], V[X, Y]) = \beta_1$$

The Matrix Formulation of a Linear Model

THE MATRIX FORMULATION OF A LINEAR MODEL

- Insert content from previous version of course: 10.7 Matrix Form of the Linear Model
- This content leads into the next lightboard of the derivation of the OLS normal equations

Reading: The Matrix Solution For OLS Regression

READING: THE MATRIX SOLUTION FOR OLS REGRESSION

Read section 4.1.3, which is on pages 147 - 151.

The Multiple OLS Solution

THE MULTIPLE OLS SOLUTION

- Pull in the lightboard called *Matrix Derivation of the OLS Estimator*.
- In the next iteration of the course, pull in a geometric derivation of the OLS coefficients.

Sample Moment Conditions

REVIEW: POPULATION MOMENT CONDITIONS

Population: Let ϵ represent error from the BLP.

Version 1: $E[\epsilon] = 0$, $E[X_j\epsilon] = 0$ for all j .

Version 2: $E[\epsilon] = 0$, $\text{cov}[X_j, \epsilon] = 0$ for all j .

SAMPLE MOMENT CONDITIONS

Sample: Let $\mathbf{e} = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$ represent OLS residuals.

SAMPLE MOMENT CONDITIONS

Sample: Let $\mathbf{e} = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$ represent OLS residuals.

$$\mathbb{X}^T \mathbf{Y} = \mathbb{X}^T \mathbb{X} \boldsymbol{\beta}, \quad \mathbf{0} = \mathbb{X}^T (\mathbf{Y} - \mathbb{X} \boldsymbol{\beta}) = \mathbb{X}^T \mathbf{e}$$

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{[1]1} & X_{[1]2} & \dots & X_{[1]n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{[k]1} & X_{[k]2} & \dots & X_{[k]n} \end{bmatrix} \cdot \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

$$\sum e_i = 0, \sum X_{[j]i} e_i = 0. \text{ or } \widehat{\text{cov}}(\mathbf{X}_{[j]}, \mathbf{e}) = \mathbf{0}$$

Consistency of Multiple OLS

CONSISTENCY OF MULTIPLE OLS

Assumptions: 1) I.I.D. 2) Unique BLP exists

In population: $\beta = E[X^T X]^{-1} E[X^T Y]$

CONSISTENCY OF MULTIPLE OLS

Assumptions: 1) I.I.D. 2) Unique BLP exists

In population: $\beta = E[X^T X]^{-1} E[X^T Y]$

$$\hat{\beta}^{(n)} = \left(\frac{1}{n} \mathbb{X}^{(n)T} \mathbb{X}^{(n)} \right)^{-1} \frac{1}{n} \mathbb{X}^{(n)T} \mathbf{Y}^{(n)}$$

$$\frac{1}{n} \mathbb{X}^{(n)T} \mathbf{Y}^{(n)} = \frac{1}{n} \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \frac{1}{n} \sum_{i=1}^n x_i^T y_i$$

$$\frac{1}{n} \mathbb{X}^T \mathbb{X} = \frac{1}{n} \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} \text{---} & x_1 & \text{---} \\ \text{---} & x_2 & \text{---} \\ & \vdots & \\ \text{---} & x_n & \text{---} \end{bmatrix} = \frac{1}{n} \sum_{i=1}^n x_i^T x_i$$

CONSISTENCY OF OLS

$$\text{WLLN} \implies \frac{1}{n} \sum_{i=1}^n x_i^T x_i \xrightarrow{p} E[X^T X] \quad \frac{1}{n} \sum_{i=1}^n x_i^T y_i \xrightarrow{p} E[X^T Y]$$

$$\text{CMT} \implies \hat{\beta} = \beta + \left(\frac{1}{n} \sum_{i=1}^n x_i^T x_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i^T \epsilon_i \right) \xrightarrow{p} \beta + \mathbf{0} = \beta$$

Unique Variation and Regression Anatomy

HOW CAN WE UNDERSTAND A SPECIFIC $\hat{\beta}_i$?

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

PARTIALLING OUT

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_k X_k$$

Step 1: Regress X_1 on other X s

$$\hat{X}_1 = \hat{\delta}_0 + \hat{\delta}_2 X_2 + \dots + \hat{\delta}_k X_k + r_1$$

Step 2: Regress Y on the residuals from Step 1

$$\hat{Y} = \hat{\gamma}_0 + \hat{\beta}_1 r_1$$

Regression anatomy: $\hat{\beta}_1 = \frac{\widehat{\text{cov}}(Y, r_1)}{\widehat{V}(r_1)}$

Deriving the Regression Anatomy Formula

DERIVING THE REGRESSION ANATOMY FORMULA

Use content from old course:
10.5 Regression Anatomy

Segment for Consideration: Applying the Regression Anatomy Formula

APPLYING THE REGRESSION ANATOMY FORMULA

Consider using the old concept check: 10.6 Applying the Regression Anatomy Formula

INTERPRETING MODEL COEFFICIENTS: WARNINGS

$$\widehat{Wage} = \beta_0 + \beta_1 Age + \beta_2 Birth_Year$$

What does it mean to hold *Age* constant while increasing *Birth_Year*?

Evaluating the Large-Sample Linear Model

THE LARGE-SAMPLE LINEAR MODEL

- I.I.D. data
- A unique BLP exists

WHAT DOES I.I.D. MEAN?



Imagine selecting each new datapoint...

- from the same distribution
- with no memory of any past datapoints

COMMON VIOLATIONS OF INDEPENDENCE

- Clustering
 - Geographic areas
 - School cohorts
 - Families
- Strategic Interaction
 - Competition among sellers
 - Imitation of species
- Autocorrelation
 - One time period may affect the next

How can observing one unit provide information about some other unit?

A UNIQUE BLP EXISTS

A BLP exists:

- $\text{cov}[X_i, X_j]$ and $\text{cov}[X_i, Y]$ are finite (no heavy tails)

The BLP is unique:

- No perfect collinearity
- $E[X^T X]$ is invertible

\implies No X_i can be written as a linear combination of the other X 's.

PERFECT COLLINEARITY EXAMPLE 1

$$\widehat{Price} = .5 \text{ Donuts} + 0.0 \text{ Dozens}$$

or

$$\widehat{Price} = 0.0 \text{ Donuts} + 6.0 \text{ Dozens}$$

PERFECT COLLINEARITY EXAMPLE 2

$$\widehat{Voters} = 200 \text{ Positive_Ads} + 100 \text{ Negative_Ads} + 0 \text{ Total_Ads}$$

or

$$\widehat{Voters} = 100 \text{ Positive_Ads} + 0 \text{ Negative_Ads} + 100 \text{ Total_Ads}$$

Goodness of Fit

Goodness of fit: How well does a model fit the data?

- R^2
- Adjusted R^2
- Akaike information criterion (AIC)
- Bayesian information criterion (BIC)

BREAKING DOWN VARIANCE

Total variance = explained variance + residual variance

DEFINING R^2

$$R^2 = 1 - \frac{\hat{V}(\hat{\epsilon})}{\hat{V}(\mathbf{Y})} = 1 - \frac{\text{residual variance}}{\text{total variance}}$$

$$\text{For OLS: } R^2 = \frac{\text{explained variance}}{\text{total variance}}$$

How much of the variation in the outcome does the model explain?

R IS CORRELATION

$$R^2 = \frac{\hat{V}(\hat{\mathbf{Y}})}{\hat{V}(\mathbf{Y})}$$

UNDERSTANDING SUMS OF SQUARES

Total sum of squares: $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$

Explained sum of squares: $ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

Residual sum of squares: $RSS = \sum_{i=1}^n (\hat{\epsilon}_i)^2$

For OLS: $TSS = ESS + RSS$

$$R^2 = 1 - \frac{RSS}{TSS}$$

THINGS TO REMEMBER ABOUT R^2

- Adding variables always makes R^2 go up.
- With many variables, consider alternatives.
 - Adjusted R^2
- R^2 is not a measure of practical significance.
 - For example, regress hospital admissions on being shot
- A low R^2 is a negative, but assess it in context.