

Chapter 1 Intro

Monday, January 6, 2025 11:27 AM



slides1

ECON 3040 – Intro to Econometrics

Lecture 1 – Course outline, RStudio, “What is Econometrics?”

Course Description

The principal objective of this course is to provide a basic introduction to econometric theory and its application. Much of the emphasis of the course is on the linear multiple regression model, under standard assumptions. The course begins with a review of probability and statistics, and ordinary least squares (OLS).

Required Textbook

Godwin, R. T., *Introduction to Econometrics*

Recommended Textbook

Introduction to Econometrics, 3rd Edition Update, by Stock and Watson.

Course Website

Course resources (including lecture notes, past exams, assignments, and computer labs) are available on rtgodwin.com/3040

Evaluation

Assignments:	15%
Midterm 1 (Feb. 3):	20%
Midterm 2 (Mar. 10):	20%
Final Exam:	45%

Assignments

You will use RStudio and work with data in order to complete your assignments.

Midterm and final examination

These will be closed book/closed notes. The final examination will cover all of the material presented in the course.

Grading scale

A+	93 – 100
A	87 – 93
B+	80 – 87
B	72 – 80
C+	64 – 72
C	57 – 64
D	50 – 57
F	0 – 50

- A missed assessment will result in make-up work, or reweighting of your grade.
- Mar. 19 is the last day for Voluntary Withdrawal from courses.

Ignorance is not a defense. Familiarize yourself with section 2.5 of [Academic Misconduct Procedures](#).

I own the copyright to all course content. Sharing my content (e.g. on Course Hero) is illegal!

All course material is copyrighted by Ryan Godwin, 2020¹⁵. No audio or video recording of this material, lectures, or presentations is allowed in any format, openly or surreptitiously, in whole or in part without permission of Ryan Godwin. Course materials are for the participant's private study and research, and must not be shared. Violation of these and other Academic Integrity principles, will lead to serious disciplinary action.

- Mar. 19 is the last day for Voluntary Withdrawal from courses.

Violation of these and other Academic Integrity principles, will lead to serious disciplinary action.

Academic Integrity

- All assignments and exams must be completed independently.
- Do not engage in "contract" cheating.
- Do not provide your UM Learn login information to anyone else. This is "personation", a serious form of academic misconduct.

Tentative Course Topics

- Review of Probability
- Review of Statistics
- Linear Regression with One Regressor
- Hypothesis Tests
- Linear Regression with Multiple Regressors
- Hypothesis Tests in Multiple Regression
- Nonlinear Regression Functions
- Instrumental Variables
- Heteroskedasticity

DiD

Student Accessibility Services

Students with disabilities should contact Student Accessibility Services to facilitate the implementation of accommodations, and meet with me to discuss the accommodations recommended by Student Accessibility Services.

Academic Supports

Sample Lecture

What is Econometrics?

- Econometrics is a subset of statistics
- Science of testing economic theories
- Used to estimate causal effects *
- Used to forecast or predict (not covered in this course)
- Often characterized by "observational data"

AI
computer science
statistics

Causal Effects

Economic models often suggest that one variable causes another. This often has policy implications. The economic models, however, do not provide quantitative magnitudes of the causal effects.

for

For example:

- How would a change in the price of alcohol or cigarettes effect the quantity consumed?
- ~~If income increases, how much of the increase will be consumed?~~
- If an additional fireplace is added to a house, how much will the price of the house increase? houses w/ fireplace vs. without
- How does another year of education change earnings? throughout

Using data to estimate causal effects

An experiment would be best.

- How would you determine the effect of fertilizer on crop yield?
- How would you use an experiment to determine the above four causal effects (on the previous slide)?
- What is the advantage of experiments? controls, for

the previous slide)?

- What is the advantage of experiments?

controls for
confounding
factors

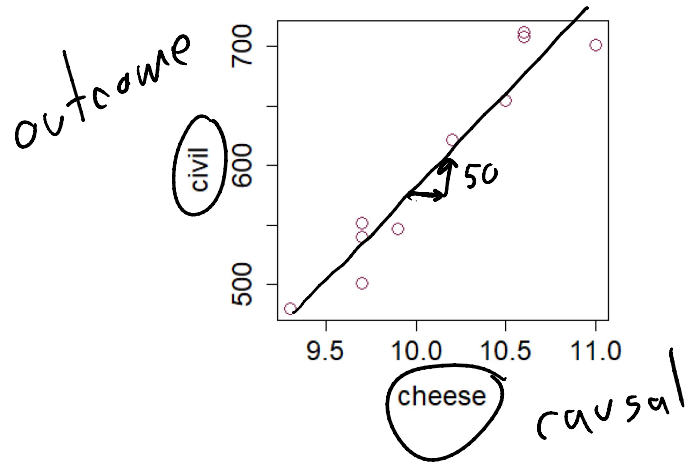
Economic experiments are usually unethical and/or too expensive.

We usually don't have experimental data in econometrics – we have observational data.

There are issues when dealing with observational data:

- Omitted variables
- Simultaneous causality
- Correlation vs. causation

Civil engineering PhDs awarded, and per-capita consumption of cheese, from 2000-2009 in the U.S. (Spurious correlations, Tyler Vigen)



What is wrong with the above picture?

Shouldn't exist

Objectives of this course

- Learn a method for estimating causal effects (least squares, “LS”)
- Understand some theoretical properties of LS
- Learn about hypothesis testing
- Practice LS using data sets


R and RStudio

The theory and concepts presented in this course will be illustrated by analysing several data sets. Data analysis will be accomplished through the R Statistical Environment and RStudio. Both are free, and R is fast becoming the best and most widely used statistical software.

First, install R

- Go to <https://muug.ca/mirror/cran/>
- Choose Windows or Mac

← → ↻ muug.ca/mirror/cran/ ☆ ⚙ R ⋮



Cran
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

About R
[R Homepage](#)
[The R Journal](#)

Software
[R Sources](#)
[R Binaries](#)
[Packages](#)
[Other](#)

Documentation
[Manuals](#)

The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms


Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2020-06-22, Taking Off Again) [R-4.0.2.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in

- Click “install R for the first time”

← → ↻ muug.ca/mirror/cran/ ☆ ⚙ R ⋮

R for Windows



Subdirectories:

base	Binaries for base distribution. This is what you want to install R for the first time .
contrib	Binaries of contributed CRAN packages (for R \geq 2.13.x; managed by Uwe Ligges). There is also information on third party software available for CRAN Windows services and corresponding environment and make variables.
old.contrib	Binaries of contributed CRAN packages for outdated versions of R (for R $<$ 2.13.x; managed by Uwe Ligges).
Rtools	Tools to build R and R packages. This is what you want to build your own packages on Windows, or to build R itself.

Please do not submit binaries to CRAN. Package developers might want to contact Uwe Ligges directly in case of questions / suggestions related to Windows binaries.

You may also want to read the [R FAQ](#) and [R for Windows FAQ](#).

Note: CRAN does some checks on these binaries for viruses, but cannot give guarantees. Use the normal precautions with downloaded executables.

CRAN
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

About R
[R Homepage](#)
[The R Journal](#)

Software
[R Sources](#)
[R Binaries](#)
[Packages](#)
[Other](#)

Documentation
[Manuals](#)

- Click “Download R 4.4.1 for Windows” (or Mac)
- Run the “.exe” file
- Click “Next” a bunch of times
- Don’t download RTools!

Second, install RStudio

- Go to <https://rstudio.com/products/rstudio/download/>
- Scroll down until you see the download button “Download RStudio Desktop for Windows (Mac)”. Click it.

Step 2: Install RStudio Desktop

DOWNLOAD RSTUDIO DESKTOP FOR WINDOWS

Size: 202.76MB | [SHA-256: FD8EA4B4](#) | Version: 2022.12.0+353 |
Released: 2022-12-15

- Run the “.exe”
- Keep clicking “next” / “install”
- Find RStudio on your computer and open it. It should look something like this:

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Project: (None)

Environment History Connections Tutorial

Global Environment

Environment is empty

Files Plots Packages Help Viewer

New Folder Delete Rename More

Home

	Name	Size	Modified
<input type="checkbox"/>	Custom Office Templates		
<input type="checkbox"/>	desktop.ini	402 B	Sep 7, 2020, 12:47 PM
<input type="checkbox"/>	My Music		
<input type="checkbox"/>	My Pictures		
<input type="checkbox"/>	My Videos		
<input type="checkbox"/>	R		
<input type="checkbox"/>	Zoom		

Console

```
R version 4.0.2 (2020-06-22) -- "Taking Off Again"
Copyright (c) 2020 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |
```



Probability Review – 2.1 Fundamental Stuff

2.1.1 Randomness

- Unpredictability
- Outcomes we can't predict are random
- Represents an inability to predict
- Example: rolling two dice

Sample Space

- Set of all outcomes of interest
- Dice example

1 die: $S = \{1, 2, \dots, 6\}$

1

Event

- Subset of outcomes
- Example: rolling higher than a 10

2.1.2 Probability

- Between 0 and 1 (or a percentage)
- “The probability of an event is the proportion of times it occurs in the long run”
- Probability of rolling 7, 12, or higher than 10?

2

2.2 Random Variables

- Translates random outcomes into numerical values
- Die roll has numerical meaning
- RVs are human-made
- Example: temperature in Celsius, Fahrenheit, Kelvin
- RVs can be discrete or continuous
- A continuous RV always has an infinite number of possibilities
- Probability of temp. being -20 tomorrow? = 0
- Random variable vs. the realization of a random variable

countable

uncountable

3

2.3 Probability function

Probability function = probability distribution = probability distribution function (PDF) = probability mass function (PMF) = probability function

- Usually an equation
- Probability function: (i) lists all possible numerical values the RV can take; (ii) assigns a probability to each value.
- Prob. function contains all possible knowledge we can have about an RV
- 2.3.1 Example: die roll

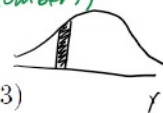
$$Pr(Y = y) = \frac{1}{6} \quad y = 1, \dots, 6 \quad (2.2)$$

4

- 2.3.2 Example: a normal RV

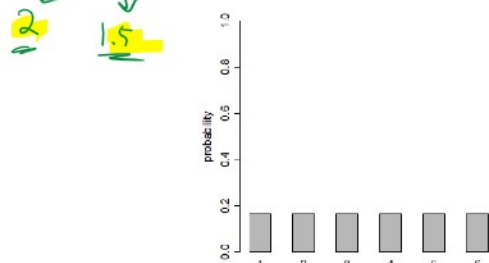
$$f(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$

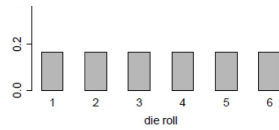
(2.3)



- Probability function for die roll in a picture:

Figure 2.1: Probability function for the result of a die roll





5

2.3.3 Probabilities of events

Probability function can be used to calculate the probability of events occurring.

Example. Let Y be the result of a die roll. What is the probability of rolling higher than 3?

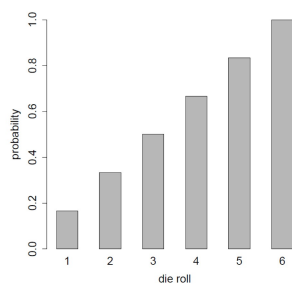
$$Pr(Y > 3) = Pr(Y = 4) + Pr(Y = 5) + Pr(Y = 6) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}$$

6

2.3.4 Cumulative distribution function (CDF)

- CDF is related to the probability function
- It's the prob. that the RV is *less than or equal to* a particular value
- In a picture:

Figure 2.2: Cumulative density function for the result of a die roll



7

2.4 Moments of a random variable

- “Moment” refers to a concept in physics
- 1st moment is the mean
- 2nd (central) moment is the variance
- 3rd is skewness
- 4th is kurtosis
- Covariance and correlation is a mixed moment

Moments summarize information about the RV. Moments are obtained from the probability function

8

2.4.1 Mean (expected value)

- Value that is expected
 - Average through repeated realizations of the RV
 - Determined from the probability function (do some math to it)
 - Mean is summarized info that is already contained in the prob. function
- Let Y be the RV
 - Mean of Y = expected value of $Y = \mu_Y = E[Y]$
 - If Y is discrete:

The mean is the weighted average of all possible outcomes, where the weights are the probabilities of each outcome.

9

The equation for the mean of Y (Y is discrete):

$$E[Y] = \sum_{i=1}^K p_i Y_i \quad (2.5)$$

where p_i is the probability of the i^{th} event, Y_i is the value of the i^{th} outcome, and K is the total number of outcomes (K can be infinite). Study this equation. It is a good way of understanding what the mean is.

Exercise: calculate the mean die roll. $E[Y] = 3.5$

What are the *properties* of the mean?

4 properties

Let Y be the result of a die roll.

$$E[Y] = \frac{1}{6}(1) + \frac{1}{6}(2) + \frac{1}{6}(3) + \dots + \frac{1}{6}(6) = 3.5$$

Properties of Expected Value

$$E[cY] = cE[Y] \rightarrow \text{Let } Z = 2Y \quad E[Z] = E[2Y] = 2 \times 3.5 = 7$$

$$E[c+Y] = c + E[Y] \rightarrow \text{Let } W = 1+Y \quad E[W] = 4.5$$

$$E[c] = c$$

$$E[X+Y] = E[X] + E[Y] = 7$$

↗ another die

10

The equation for the mean of y (y is continuous):

Let y be a random variable. The mean of y is

$$E[y] = \int y f(y) dy$$

If y is normally distributed, then $f(y)$ is equation (2.3), and the mean of y turns out to be μ . You do not need to integrate for this course, but you should have some idea about how the mean of a continuous random variable is determined from its probability function.

The mean is different from the median and the mode, although all are measures of central tendency.

The mean is different from the sample mean or sample average.

The mean comes from the probability function. The sample mean/average comes from a sample of data.

11

2.4.3 Variance

- Measure of the spread or dispersion of a RV
- Denoted by σ^2 . The variance of y would be σ_y^2 and the variance of X would be σ_x^2
- Variance is the expected squared difference of a variable from its mean
- Equation:

$$E(Y - \mu_Y)^2 = \text{var}(Y)$$

$E[Y]$

2.4.3 Variance

- Measure of the *spread* or *dispersion* of a RV
- Denoted by σ^2 . The variance of y would be σ_y^2 and the variance of X would be σ_X^2
- Variance is the expected squared difference of a variable from its mean
- Equation:

$$\text{Var}(Y) = E[(Y - E[Y])^2] \quad (2.6)$$

When Y is a discrete random variable, then equation (2.6) becomes

$$\begin{aligned} E(Y) &= \sum_{i=1}^K p_i Y_i \\ \text{Var}(Y) &= \sum_{i=1}^K p_i \times (Y_i - E[Y])^2 \end{aligned} \quad (2.7)$$

- For variance (the ~~2nd moment~~), we are taking the expectation of a squared term
- For skewness (the ~~3rd moment~~), we would take the expectation of a cubed term, etc.

Exercise: calculate the variance of a die roll

$$\text{var}(Y) = \frac{1}{6} (1 - 3.5)^2 + \frac{1}{6} (2 - 3.5)^2 + \dots + \frac{1}{6} (6 - 3.5)^2 \approx 2.92$$

What are the *properties* of the variance?

4 (on board)

Exercise: I change the sides of the die to equal 2,4,6,8,10,12. What is the mean and variance of the die roll?

$$\text{var} = 2^2 \text{var}(Y)$$

Exercise: What is the mean and variance of the sum of two dice?

2.4.5 Covariance

- Measures the relationship between two random variables
- Random variables Y and X have a joint probability function
- Joint prob. func.: (i) lists all possible combos of Y and X ; (ii) assign a probability to each combination
- A useful summary of a joint probability function is the covariance
- The covariance between Y and X is the expected difference of Y from its mean, multiplied by the expected difference of X from its mean
- Covariance tells us something about how two variables are related, or how they move together
- Tells us about the direction and strength of the relationship between two variables

15

$$\text{Cov}(Y, X) = E[(Y - \mu_Y)(X - \mu_X)] = E[(Y - \mu_Y)^2]$$

$$\text{Cov}(Y, X) = E[(Y - \mu_Y)(X - \mu_X)] \quad (2.8)$$

The covariance between Y and X is often denoted as σ_{YX} . Note the following properties of σ_{YX} :

- σ_{YX} is a measure of the linear relationship between Y and X . Non-linear relationships will be discussed later.
- $\sigma_{YX} = 0$ means that Y and X are linearly independent.
- If Y and X are independent (neither variable causes the other), then $\sigma_{YX} = 0$. The converse is not necessarily true (because of non-linear relationships).
independence \Rightarrow 0 cov/corr
- The $\text{Cov}(Y, Y)$ is the $\text{Var}(Y)$.
- A positive covariance means that the two variables tend to differ from their mean in the same direction.
- A negative covariance means that the two variables tend to differ from their mean in the opposite direction.

16

2.4.6 Correlation

- Correlation usually denoted by ρ ^{$\sigma_r h_o''$}
- Similar to covariance, but is easier to interpret

$$\rho_{YX} = \frac{\text{Cov}(Y, X)}{\sqrt{\text{Var}(Y)\text{Var}(X)}} = \frac{\sigma_{YX}}{\sigma_Y \sigma_X} \quad (2.9)$$

The difficulty in interpreting the value of covariance is because $-\infty < \sigma_{YX} < \infty$. Correlation transforms covariance so that it is bound between -1 and 1. That is, $-1 \leq \rho_{YX} \leq 1$.

- $\rho_{YX} = 1$ means perfect positive linear association between Y and X .
- $\rho_{YX} = -1$ means perfect negative linear association between Y and X .
- $\rho_{YX} = 0$ means no linear association between Y and X (linear independence).

17

2.4.7 Conditional distribution

- Joint distribution – 2 RVs
- Conditional distribution – fix (condition on) one of those RVs
- Condition expectation – the mean of one RV after the other RV has been “fixed”

Let Y be a discrete random variable. Then, the conditional mean of Y given some value for X is

$$\mathbb{E}(Y|X=x) = \sum_{i=1}^K (p_i|X=x) Y_i \quad (2.10)$$

- If the two RVs are independent, the conditional distribution is the same as the ~~marginal~~ distribution

18

Example: Blizzard and cancelled midterm

Suppose that you have a midterm tomorrow, but there is a possibility of a blizzard. You are wondering if the midterm might be cancelled.

Table 2.1: Joint distribution for snow and a canceled midterm

	Midterm ($Y = 1$)	No Midterm ($Y = 0$)	
Blizzard ($X = 1$)	0.05	0.20	$0.05 + 0.2 = 0.25$
No Blizzard ($X = 0$)	0.72	0.03	$0.75 = 0.75$
	.77	.23	$\frac{0.2}{0.25} = 0.8$

- What are the *marginal* probability distributions?
- What is $E[Y]$? What is $E[Y|X=1]$? $= .2(1) + .8(0) = .2$
- What is the covariance and correlation between X and Y ?
- More exercises in the "Review Questions"

19

2.5 Some special probability functions

2.5.1 The normal distribution

- Common because of the "central limit theorem" (in a few slides)

$$f(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y-\mu)^2}{2\sigma^2}\right] \quad (2.3)$$

- Mean of y is μ
- Variance of y is σ^2

20

2.5.2 The standard normal distribution

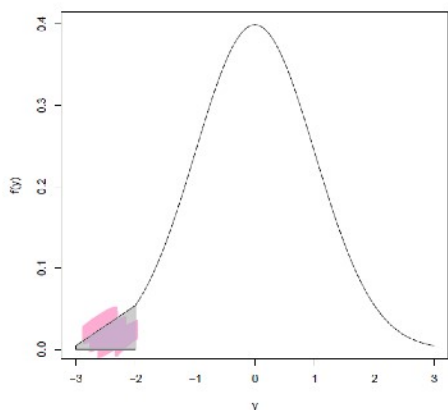
- Special case of a normal distribution, where $\mu = 0$ and $\sigma^2 = 1$
- Equation 2.3 becomes:

$$f(y) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{y^2}{2} \right) \quad (2.11)$$

- Any normal random variable can be “standardized”
- How to standardize? *subtract mean, divide by standard deviation*
- Standardizing has long been used in hypothesis testing (as we shall see)

21

Figure 2.3: Probability function for a standard normal variable, $p_{y < -2}$ in gray



22

2.5.3 The central limit theorem

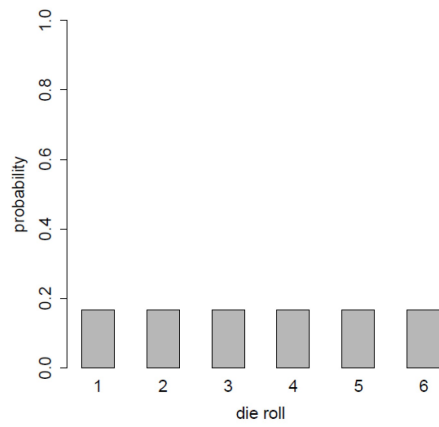
- There are hundreds of different probability functions
- Examples: Poisson, Binomial, Generalized Pareto, Nakagami, Uniform
- So why is the normal distribution so important? Why are so many RVs normal?
- Answer: CLT
- CLT (loosely speaking) if we add up enough RVs, the resulting sum tends to be normal

Exercise: draw the probability function for one die roll, then for the sum of two dice.

Exercise: draw the probability function for one die roll, then for the sum of two dice.

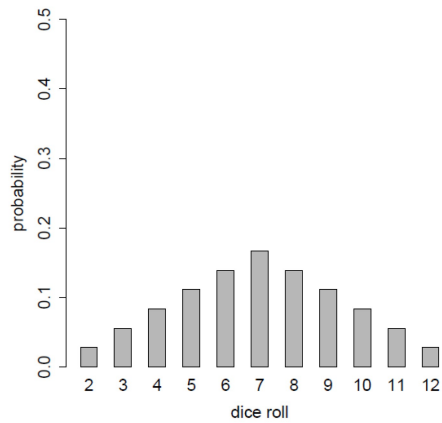
23

Figure 2.1: Probability function for the result of a die roll



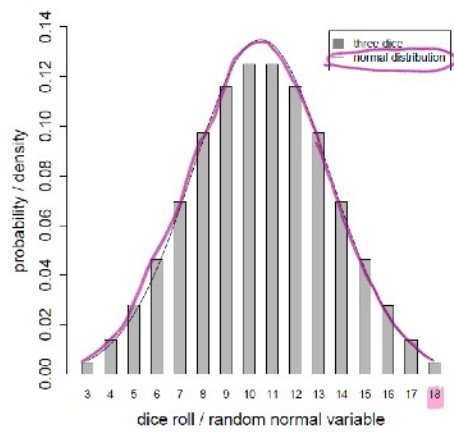
24

Figure 2.4: Probability function for the sum of two dice



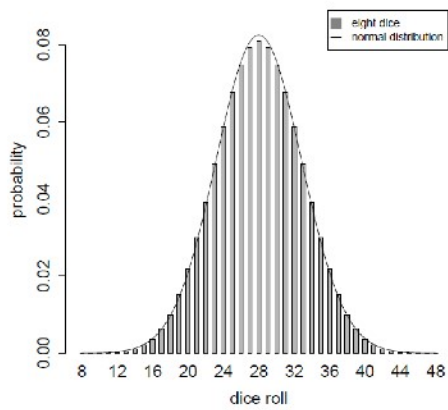
25

Figure 2.5: Probability function for three dice, and normal distribution



26

Figure 2.6: Probability function for eight dice, and normal distribution



27

2.5.4 The chi-square distribution

- Add to a normal RV – still normal
- Multiply a normal RV – still normal
- Square a normal RV – now it is chi-square distributed
- We will use the chi-square distribution for the F-test in a later chapter



Statistics Review

- A statistic is a *function* of a *sample* of data
- An *estimator* is a statistic
- Population parameter \rightarrow unknown
- Estimator \rightarrow used to estimate an unknown population parameter
- The sample, y , will be considered random
- Since y is random, estimators using y will be random

sample (like the die rolls in assign 1)

Since estimators are random, they have a probability function, given a special name: sampling distribution.

We will obtain properties of the sampling distribution to see if the estimator is "good" or not.

1

3.1 Random Sampling from the Population

- Typically, we want to know something about a population
- The population is considered to be very large (infinite), and contains some unknown "truth"
- We likely won't observe the whole population, but a sample from the pop.
- We'll use the sample, y , to estimate that something

contains truth

2

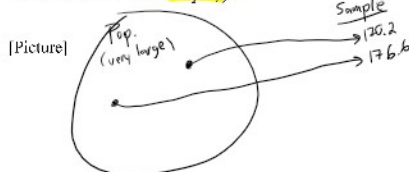
Example: suppose we want to know the mean height of a
U of M student

Let y = height of a ~~single~~ student

- Population: all ~~single~~ students
- Population parameter of interest: μ_y

We can't afford to observe the whole pop.

We'll have to collect a sample, y .



3

We want the sample to reflect the population.

Question: How should the sample be selected from the population?
randomly

In particular we want the sample to be i.i.d.

- Identically: \rightarrow come from pop. of U of M students (no min: = U students)
- Independently: \rightarrow no link/connection (entire basketball team)
- Distributed

4

So, the sample y is random!!

- Could have gotten a different y
- Parallel universe

Table 3.1: Entire population of heights (in cm). The true (unobservable) population mean and variance are $\mu_y = 176.8$ and $\sigma_y^2 = 39.7$.

177.3	176.2	187.2	178.3	176.3	179.4	181.2	180.0	175.9
178.7	171.7	160.5	183.9	175.7	175.9	182.6	181.7	180.2
181.5	176.5	162.1	180.3	175.6	174.9	165.7	172.7	178.9
175.3	178.7	175.6	166.4	173.1	173.2	175.6	183.7	181.3
174.2	180.9	179.9	171.2	171.0	178.6	181.4	175.2	182.2
171.7	178.4	168.1	186.0	189.9	173.1	168.7	180.0	175.1
175.7	180.8	176.2	170.8	177.3	163.4	186.3	177.1	191.2
171.0	180.3	169.5	167.2	178.0	172.9	176.0	176.5	171.9
175.1	184.2	165.3	180.2	178.3	183.4	178.9	178.6	177.9
184.5	184.1	180.9	187.1	179.9	167.1	172.0	167.4	172.7
171.6	186.6	182.4	185.5	174.8	178.8	192.8	179.3	172.0

5

How could i.i.d. be violated in the heights example?

Example: mean income of Canadians. How could i.i.d. be violated?

How should we estimate the mean height?

3.2 Estimators and Sampling Distributions

An estimator uses the sample y to "guess" something about the pop.

We collect our sample, $y = \{173.9, 171.7, 182.6, 181.5, 169.1, 174.9, 165.7, 182.2, 171.7, 168.1, 180.9, 175.7, 163.4, 186.3, 169.3, 171.9, 173.9, 173.9, 172.9, 172.0\}$. How should we use this sample to estimate the mean height?

sample mean / sample average / average

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

6

3.2.1 Sample mean

A popular choice for estimating a population mean is by using a sample mean (or sample average or just average)

- From heights example: $y = 174.1$, $\mu_y = 176.8$

- There are many ways to estimate μ_y . Examples?

- Why is (3.1) so popular? It's the best

- How good is \bar{y} at estimating μ_y in general?

- To answer these questions: idea of a sampling distribution

(3.1)

$\frac{1}{n} \sum_{i=1}^n y_i$

median / mode / geometric average / harmonic avg.

$\frac{\min(y) + \max(y)}{2}$

7

Recall that the sample, y , is random. Each element of y was selected randomly from the population. We could have selected a different sample of size $n = 20$. For example, in a parallel universe, we could have gotten $y' = \{175.9, 175.3, 182.2, 178.6, 175.2, 180.3, 178.3, 183.7, 176.0, 167.4, 178.7, 178.7, 185.0, 175.6, 180.0, 168.7, 178.6, 173.1, 173.2, 187.4\}$, where the y' in y' denotes that we are in the parallel universe. In this parallel universe, we get $\bar{y}' = 177.6$. But in every universe, the population (table 3.1), is the same, $\bar{y} = 176.1$.

- Randomly sample from the population \rightarrow get y

y is random

- Use y to calculate \bar{y}

\bar{y} is random

\rightarrow could have gotten a different sample \rightarrow could have gotten a different \bar{y}

\rightarrow population is always the same (μ_y)

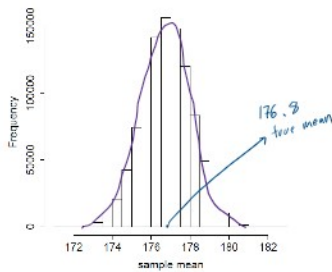
8

3.2.2 Sampling distribution of the sample mean

- \bar{y} is random variable (it's an estimator, all estimators are random)
- random variables usually have probability functions
- \bar{y} has a sampling distribution (probability function for an estimator)
- sampling distribution imagine all possible values for y that you could get – plot a histogram
- Using a computer, I drew 1 mil. different random samples of $n=20$ from table 3.1. Calculate \bar{y} each time. Plot histogram:

9

Figure 3.1: Histogram for 1 million \bar{y} s



10

Normal distⁿ

Which probability function is right for \bar{y} ? Why?

- Look at figure 3.1
- Notice the summation operator in equation 3.1
- Answer: Normal Reason: CLT (adding in sample means)

\bar{y} is random. We'll derive its:

- mean
- variance

Use these to determine if it's a "good" estimator via three statistical properties:

- Bias
- Efficiency
- Consistency

11

3.2.3 Bias

An estimator is unbiased if its expected value is equal to the population parameter it's estimating.

That is, \bar{y} is unbiased if $E[\bar{y}] = \mu_y$.

Unbiased if it gives "the right answer on average".

Biased if it gives the wrong answer on average.

Rules of the mean

$$(i) E[cY] = c E[Y]$$

$$(ii) E[X+Y] = E[X] + E[Y]$$

$$\begin{aligned} E[\bar{y}] &= E\left[\frac{1}{n} \sum y_i\right] \\ &= \frac{1}{n} E\left[\sum y_i\right] = \frac{1}{n} E[y_1 + y_2 + \dots + y_n] \\ &= \frac{1}{n} \{E[y_1] + E[y_2] + \dots + E[y_n]\} \\ &= \frac{1}{n} \{\mu_y + \mu_y + \dots + \mu_y\} \quad \text{i.i.d.} \\ &= \frac{1}{n} n \mu_y = \mu_y \quad \hookrightarrow \text{"identical"} \\ &\quad \bar{y} \text{ is unbiased} \end{aligned}$$

12

$$\begin{aligned}
 E[\bar{y}] &= E\left[\frac{1}{n} \sum_{i=1}^n y_i\right] \\
 &= \frac{1}{n} E\left[\sum_{i=1}^n y_i\right] \\
 &= \frac{1}{n} E[y_1 + y_2 + \dots + y_n] \\
 &= \frac{1}{n} (E[y_1] + E[y_2] + \dots + E[y_n]) \\
 &= \frac{1}{n} (\mu_y + \mu_y + \dots + \mu_y) \\
 &= \frac{n\mu_y}{n} = \mu_y
 \end{aligned}
 \tag{3.2}$$

13

3.2.4 Efficiency \rightarrow accuracy / spread of estimator

An estimator is **efficient** if it has the **smallest variance** among all **other potential estimators** (for us, potential = linear, unbiased)

Need to get the variance of \bar{y} .

Rules of variance

$$\begin{aligned}
 \text{Var}(\bar{y}) &= \text{Var}\left(\frac{1}{n} \sum y_i\right) \\
 &= \frac{1}{n^2} \text{Var}\left(\sum y_i\right) = \frac{1}{n^2} \text{Var}(y_1 + y_2 + \dots + y_n) \\
 &= \frac{1}{n^2} \{ \text{Var}(y_1) + \text{Var}(y_2) + \dots + \text{Var}(y_n) \} \\
 &= \frac{1}{n^2} \{ \sigma_y^2 + \sigma_y^2 + \dots + \sigma_y^2 \} \\
 &= \frac{1}{n^2} n \sigma_y^2 = \frac{\sigma_y^2}{n}
 \end{aligned}$$

If 2 variables are independent \Rightarrow no cov. coll.

\rightarrow i.i.d. \hookrightarrow "independent"

14

$$\begin{aligned}
 \text{Var}[\bar{y}] &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n y_i\right] \\
 &= \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n y_i\right] \\
 &= \frac{1}{n^2} \text{Var}[y_1 + y_2 + \dots + y_n] \\
 &= \frac{1}{n^2} (\text{Var}[y_1] + \text{Var}[y_2] + \dots + \text{Var}[y_n]) \\
 &= \frac{1}{n^2} (\sigma_y^2 + \sigma_y^2 + \dots + \sigma_y^2) \\
 &= \frac{n\sigma_y^2}{n^2} = \frac{\sigma_y^2}{n}
 \end{aligned}
 \tag{3.3}$$

any other estimator for μ

$\text{Var}(\bar{y}) < \text{Var}(\hat{\mu}_1)$

\bar{y} is BLUE (best linear unbiased estimator)

- Gauss-Markov theorem proves this is minimum variance
- We'll also need this to **prove consistency**, and for **hyp. testing**

15

3.2.5 Consistency

Suppose we had a lot of information ($n \rightarrow \infty$)

What value should we get for our estimator? **right answer, every time**

How would state this mathematically?

$$\lim_{n \rightarrow \infty} \text{Var}(\bar{y}) \rightarrow 0 \quad \checkmark \quad \lim_{n \rightarrow \infty} \text{bias}(\bar{y}) \rightarrow 0 \quad \checkmark$$

Q) Prove that the **sample mean** is a **consistent estimator** for the population mean.

$$\text{Var}(\bar{y}) = \frac{\sigma^2}{n} \quad \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} \rightarrow 0$$

"variance goes away"

Q) Define the terms **unbiasedness**, **efficiency**, and **consistency**.

$$\bar{y} = \frac{1}{n} \sum y_i$$

mean

variance

randomly selected

also random

16

$$\begin{aligned}
 E[\bar{y}] &= \mu_y \quad \text{unbiased} \\
 \text{Var}[\bar{y}] &= \frac{\sigma_y^2}{n} \rightarrow \text{efficient} \\
 &\rightarrow \text{consistent}
 \end{aligned}$$

very unrealistic

3.3 Hypothesis tests (known σ_y^2)

null $H_0: \mu_y = \mu_{y,0}$

alternative $H_A: \mu_y \neq \mu_{y,0}$ (2-sided alternative)

\rightarrow almost always 2-sided in Econ

$$\tag{3.4}$$

3.3 Hypothesis tests (known σ_y^2)

Null: $H_0: \mu_y = \mu_{y,0}$ (almost always 2-sided in Econ)

Alternative: $H_A: \mu_y \neq \mu_{y,0}$ (2-sided alternative)

(3.4)

- Estimate μ_y (using \bar{y} for example)
- See if \bar{y} appears "close" to $\mu_{y,0}$
 - Remember, \bar{y} is random! (and Normal)
- If it's close \rightarrow fail to reject
- If it's far \rightarrow reject

17

Example:

- Hypothesize that mean height of a U of M student is 173cm

$$H_0: \mu_y = 173 \quad (3.5)$$

$$H_A: \mu_y \neq 173 \quad (174.1 - 173) = 1.1 \text{ cm}$$

- Collect a sample: $y = \{173.9, 171.7, \dots, 172.0\}$

- Calculate $\bar{y} = 174.1$

- Suppose (very unrealistically that we know that) $\sigma_y^2 = 39.7$

- What now?

$$\bar{y} \sim N(173, \frac{39.7}{20})$$

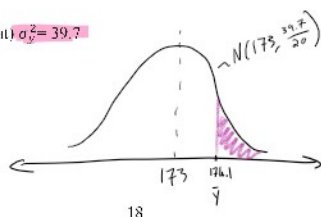
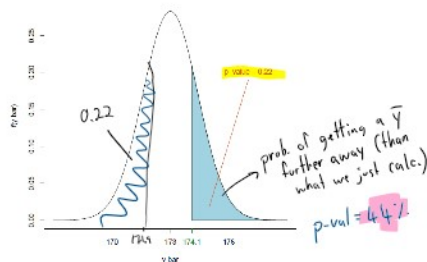


Figure 3.2: Normal distribution with $\mu = 173$ and $\sigma^2 = 39.7/20$. Shaded area is the probability that the normal variable is greater than 174.1.



19

The p-value for the above test is 0.44. How to interpret this?

44% chance of getting a \bar{y} that is more adverse to H_0 . \rightarrow fail to reject

3.3.1 Significance of a test

\rightarrow pre-determined p-value that decides reject/fail to reject

$\alpha = 10\%, 5\%, 1\% \rightarrow$ if $p\text{-val} < 5\% \Rightarrow$ reject

3.3.2 Type I error

$$Pr(\text{reject } H_0 \mid H_0 \text{ is true}) = \alpha$$

3.3.3 Type II error (and power)

(fail to reject H_0) H_0 is false

$$\text{power} = 1 - \text{type II} = Pr(\text{reject } H_0 \mid H_0 \text{ is false}) = ?$$

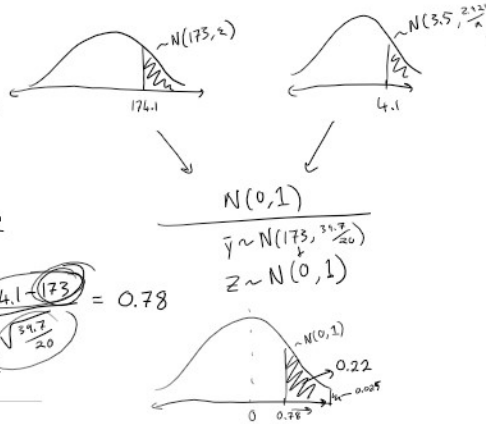
depends on how false H_0

$H_0: \mu_y = 20$
in reality 20.01 vs. 1,000,000

20

3.3.4 Test statistics

- Just a more convenient way of getting the p-value for the test
- Each hypothesis test would present us with a new normal curve that we would have to draw, and calculate a new area (see fig. 3.2)
- Instead: **standardize**
- This gives us **one curve for all testing problems** (the standard normal curve)
- Calculate a bunch of areas under the curve, and tabulate them
- Not an issue with modern computers, but this is still the way we do things
- How to get a z test statistic?
- Do a z test for our heights example.



p-val = 0.22 * 2
= 0.44 > 0.05
→ fail to reject

Table 3.3 Area under the standard normal curve to the right of z

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.4960	0.4961	0.4963	0.4964	0.4965	0.4966	0.4967	0.4968	0.4969	0.4970
0.1	0.4970	0.4971	0.4972	0.4973	0.4974	0.4975	0.4976	0.4977	0.4978	0.4979
0.2	0.4979	0.4980	0.4981	0.4982	0.4983	0.4984	0.4985	0.4986	0.4987	0.4988
0.3	0.4988	0.4989	0.4990	0.4991	0.4992	0.4993	0.4994	0.4995	0.4996	0.4997
0.4	0.4997	0.4998	0.4999	0.5000	0.5001	0.5002	0.5003	0.5004	0.5005	0.5006
0.5	0.5006	0.5007	0.5008	0.5009	0.5010	0.5011	0.5012	0.5013	0.5014	0.5015
0.6	0.5015	0.5016	0.5017	0.5018	0.5019	0.5020	0.5021	0.5022	0.5023	0.5024
0.7	0.5024	0.5025	0.5026	0.5027	0.5028	0.5029	0.5030	0.5031	0.5032	0.5033
0.8	0.5033	0.5034	0.5035	0.5036	0.5037	0.5038	0.5039	0.5040	0.5041	0.5042
0.9	0.5042	0.5043	0.5044	0.5045	0.5046	0.5047	0.5048	0.5049	0.5050	0.5051
1.0	0.5051	0.5052	0.5053	0.5054	0.5055	0.5056	0.5057	0.5058	0.5059	0.5060
1.1	0.5060	0.5061	0.5062	0.5063	0.5064	0.5065	0.5066	0.5067	0.5068	0.5069
1.2	0.5069	0.5070	0.5071	0.5072	0.5073	0.5074	0.5075	0.5076	0.5077	0.5078
1.3	0.5078	0.5079	0.5080	0.5081	0.5082	0.5083	0.5084	0.5085	0.5086	0.5087
1.4	0.5087	0.5088	0.5089	0.5090	0.5091	0.5092	0.5093	0.5094	0.5095	0.5096
1.5	0.5096	0.5097	0.5098	0.5099	0.5100	0.5101	0.5102	0.5103	0.5104	0.5105
1.6	0.5105	0.5106	0.5107	0.5108	0.5109	0.5110	0.5111	0.5112	0.5113	0.5114
1.7	0.5114	0.5115	0.5116	0.5117	0.5118	0.5119	0.5120	0.5121	0.5122	0.5123
1.8	0.5123	0.5124	0.5125	0.5126	0.5127	0.5128	0.5129	0.5130	0.5131	0.5132
1.9	0.5132	0.5133	0.5134	0.5135	0.5136	0.5137	0.5138	0.5139	0.5140	0.5141
2.0	0.5141	0.5142	0.5143	0.5144	0.5145	0.5146	0.5147	0.5148	0.5149	0.5150
2.1	0.5150	0.5151	0.5152	0.5153	0.5154	0.5155	0.5156	0.5157	0.5158	0.5159
2.2	0.5159	0.5160	0.5161	0.5162	0.5163	0.5164	0.5165	0.5166	0.5167	0.5168
2.3	0.5168	0.5169	0.5170	0.5171	0.5172	0.5173	0.5174	0.5175	0.5176	0.5177
2.4	0.5177	0.5178	0.5179	0.5180	0.5181	0.5182	0.5183	0.5184	0.5185	0.5186
2.5	0.5186	0.5187	0.5188	0.5189	0.5190	0.5191	0.5192	0.5193	0.5194	0.5195
2.6	0.5195	0.5196	0.5197	0.5198	0.5199	0.5200	0.5201	0.5202	0.5203	0.5204
2.7	0.5204	0.5205	0.5206	0.5207	0.5208	0.5209	0.5210	0.5211	0.5212	0.5213
2.8	0.5213	0.5214	0.5215	0.5216	0.5217	0.5218	0.5219	0.5220	0.5221	0.5222
2.9	0.5222	0.5223	0.5224	0.5225	0.5226	0.5227	0.5228	0.5229	0.5230	0.5231
3.0	0.5231	0.5232	0.5233	0.5234	0.5235	0.5236	0.5237	0.5238	0.5239	0.5240
3.1	0.5240	0.5241	0.5242	0.5243	0.5244	0.5245	0.5246	0.5247	0.5248	0.5249
3.2	0.5249	0.5250	0.5251	0.5252	0.5253	0.5254	0.5255	0.5256	0.5257	0.5258
3.3	0.5258	0.5259	0.5260	0.5261	0.5262	0.5263	0.5264	0.5265	0.5266	0.5267
3.4	0.5267	0.5268	0.5269	0.5270	0.5271	0.5272	0.5273	0.5274	0.5275	0.5276
3.5	0.5276	0.5277	0.5278	0.5279	0.5280	0.5281	0.5282	0.5283	0.5284	0.5285
3.6	0.5285	0.5286	0.5287	0.5288	0.5289	0.5290	0.5291	0.5292	0.5293	0.5294
3.7	0.5294	0.5295	0.5296	0.5297	0.5298	0.5299	0.5300	0.5301	0.5302	0.5303
3.8	0.5303	0.5304	0.5305	0.5306	0.5307	0.5308	0.5309	0.5310	0.5311	0.5312
3.9	0.5312	0.5313	0.5314	0.5315	0.5316	0.5317	0.5318	0.5319	0.5320	0.5321
4.0	0.5321	0.5322	0.5323	0.5324	0.5325	0.5326	0.5327	0.5328	0.5329	0.5330
4.1	0.5330	0.5331	0.5332	0.5333	0.5334	0.5335	0.5336	0.5337	0.5338	0.5339
4.2	0.5339	0.5340	0.5341	0.5342	0.5343	0.5344	0.5345	0.5346	0.5347	0.5348
4.3	0.5348	0.5349	0.5350	0.5351	0.5352	0.5353	0.5354	0.5355	0.5356	0.5357
4.4	0.5357	0.5358	0.5359	0.5360	0.5361	0.5362	0.5363	0.5364	0.5365	0.5366
4.5	0.5366	0.5367	0.5368	0.5369	0.5370	0.5371	0.5372	0.5373	0.5374	0.5375
4.6	0.5375	0.5376	0.5377	0.5378	0.5379	0.5380	0.5381	0.5382	0.5383	0.5384
4.7	0.5384	0.5385	0.5386	0.5387	0.5388	0.5389	0.5390	0.5391	0.5392	0.5393
4.8	0.5393	0.5394	0.5395	0.5396	0.5397	0.5398	0.5399	0.5400	0.5401	0.5402
4.9	0.5402	0.5403	0.5404	0.5405	0.5406	0.5407	0.5408	0.5409	0.5410	0.5411
5.0	0.5411	0.5412	0.5413	0.5414	0.5415	0.5416	0.5417	0.5418	0.5419	0.5420
5.1	0.5420	0.5421	0.5422	0.5423	0.5424	0.5425	0.5426	0.5427	0.5428	0.5429
5.2	0.5429	0.5430	0.5431	0.5432	0.5433	0.5434	0.5435	0.5436	0.5437	0.5438
5.3	0.5438	0.5439	0.5440	0.5441	0.5442	0.5443	0.5444	0.5445	0.5446	0.5447
5.4	0.5447	0.5448	0.5449	0.5450	0.5451	0.5452	0.5453	0.5454	0.5455	0.5456
5.5	0.5456	0.5457	0.5458	0.5459	0.5460	0.5461	0.5462	0.5463	0.5464	0.5465
5.6	0.5465	0.5466	0.5467	0.5468	0.5469	0.5470	0.5471	0.5472	0.5473	0.5474
5.7	0.5474	0.5475	0.5476	0.5477	0.5478	0.5479	0.5480	0.5481	0.5482	0.5483
5.8	0.5483	0.5484	0.5485	0.5486	0.5487	0.5488	0.5489	0.5490	0.5491	0.5492
5.9	0.5492	0.5493	0.5494	0.5495	0.5496	0.5497	0.5498	0.5499	0.5500	0.5501
6.0	0.5501	0.5502	0.5503	0.5504	0.5505	0.5506	0.5507	0.5508	0.5509	0.5510
6.1	0.5510	0.5511	0.5512	0.5513	0.5514	0.5515	0.5516	0.5517	0.5518	0.5519
6.2	0.5519	0.5520	0.5521	0.5522	0.5523	0.5524	0.5525	0.5526	0.5527	0.5528
6.3	0.5528	0.5529	0.5530	0.5531	0.5532	0.5533	0.5534	0.5535	0.5536	0.5537
6.4	0.5537	0.5538	0.5539	0.5540	0.5541	0.5542	0.5543	0.5544	0.5545	0.5546
6.5	0.5546	0.5547	0.5548	0.5549	0.5550	0.5551	0.5552	0.5553	0.5554	0.5555
6.6	0.5555	0.5556	0.5557	0.5558	0.5559	0.5560	0.5561	0.5562	0.5563	0.5564
6.7	0.5564	0.5565	0.5566	0.5567	0.5568	0.5569	0.5570	0.5571	0.5572	0.5573
6.8	0.5573	0.5574	0.5575	0.5576	0.5577	0.5578	0.5579	0.5580	0.5581	0.5582
6.9	0.5582	0.5583	0.5584	0.5585	0.5586	0.5587	0.5588	0.5589	0.5590	0.5591
7.0	0.5591	0.5592	0.5593	0.5594	0.5595	0.5596	0.5597	0.5598	0.5599	0.5600
7.1	0.5600	0.5601	0.5602	0.5603	0.5604	0.5605	0.5606	0.5607	0.5608	0.5609
7.2	0.5609	0.5610	0.5611	0.5612	0.5613	0.5614	0.5615	0.5616	0.5617	0.5618
7.3	0.5618	0.5619	0.5620	0.5621	0.5622	0.5623	0.5624	0.5625	0.5626	0.5627
7.4	0.5627	0.5628	0.5629	0.5630	0.5631	0.5632	0.5633	0.5634	0.5635	0.5636
7.5	0.5636	0.5637	0.5638	0.5639	0.5640	0.5641	0.5642	0.5643	0.5644	0.5645
7.6	0.5645	0.5646	0.5647	0.5648	0.5649	0.5650	0.5651	0.5652	0.5653	0.5654
7.7	0.5654	0.5655	0.5656	0.5657	0.5658	0.5659	0.5660	0.5661	0.5662	0.5663
7.8	0.5663	0.5664	0.5665	0.5666	0.5667	0.5668	0.5669	0.5670	0.5671	0.5672
7.9	0.5672	0.5673	0.5674	0.5675	0.5676	0.5677	0.5678	0.5679	0.5680	0.5681
8.0	0.5681	0.5682	0.5683	0.5684	0.5685	0.5686	0.5687	0.5688	0.5689	0.5690
8.1	0.5690	0.5691	0.5692	0.5693	0.5694	0.5695	0.5696	0.5697	0.5698	0.5699
8.2	0.5699	0.5700	0.5701	0.5702	0.5703	0.5704	0.5705	0.5706	0.5707	0.5708
8.3	0.5708	0.5709	0.5710	0.5711	0.5712	0.5713	0.5714	0.5715	0.5716	0.5717
8.4	0.5717	0.5718	0.5719	0.5720	0.5721	0.5722	0.5723	0.5724	0.5725	0.5726
8.5	0.5726	0.5727	0.5728	0.5729	0.5730	0.5731	0.5732	0.5733	0.5734	0.5735
8.6	0.5735	0.5736	0.5737	0.5738	0.5739	0.5740	0.5741	0.5742	0.5743	0.5744
8.7	0.5744	0.5745	0.5746	0.5747	0.5748	0.5749	0.5750	0.5751	0.5752	0.5753
8.8	0.5753	0.5754	0.5755	0.5756	0.5757	0.5758	0.5759	0.5760	0.5761	0.5762
8.9	0.5762	0.5763	0.5764	0.5765	0.5766	0.5767	0.5768	0.5769	0.5770	0.5771
9.0	0.5771	0.5772	0.5773	0.5774	0.5775	0.5776	0.5777	0.5778	0.5779	0.5780
9.1	0.5780	0.5781	0.5782	0.5783	0.5784	0.5785	0.5786	0.5787	0.5788	0.5789
9.2	0.5789	0.5790	0.5791	0.5792	0.5793	0.5794	0.5795	0.5796	0.5797	0.5798
9.3	0.5798	0.5799	0.5800	0.5801	0.5802	0.5803	0.5804	0		

3.4.1 Estimating σ_y^2

$$\text{var}[\bar{y}] = E[(\bar{y} - E(\bar{y}))^2] \xrightarrow{\text{estimate}} \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

- A "natural" estimator:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

(3.15)

$$E\left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2\right] = \frac{n-1}{n} \sigma^2 \quad \text{BIASED}$$

$$\hookrightarrow E\left[\frac{n}{n-1} \hat{\sigma}^2\right] = E\left[\frac{n}{n-1} \cdot \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2\right] = \frac{n}{n-1} E[\hat{\sigma}^2] = \frac{n}{n-1} \cdot \frac{n-1}{n} \sigma^2 = \sigma^2 \quad \text{UNBIASED}$$

- Is this a good estimator? Why or why not?
- A better estimator:

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \quad (3.17)$$

- Degrees-of-freedom correction

$$y = \{1, 3, ?\} \quad \bar{y} = 3$$

5

→ calc. $\bar{y} \rightarrow$ lose 1 piece info.

So:

Estimated variance of \bar{y} $\rightarrow \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$

We can implement hypothesis testing by replacing the unknown σ_y^2 with its estimator s_y^2 . The z test statistic now becomes:

$$\frac{\bar{y} - \mu_{y,0}}{\sqrt{s_y^2/n}} = t \quad \text{using } s_y^2 \text{ instead of } \sigma^2 \text{ makes } z \text{ become } t$$

random χ^2

only random thing (Normal) $\rightarrow z$ also Normal

$$Z = \frac{\bar{y} - \mu_{y,0}}{\sqrt{\sigma_y^2/n}}$$

Note: for large n , the t test is equivalent to the z test $t \sim t_{n-1}$

