



Statistics Review

- A statistic is a **function** of a **sample of data**
- An **estimator** is a statistic
- **Population parameter** \rightarrow **unknown** (μ, σ^2)
- **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**

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

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

y is random!
 y is a sample of values
 \hookrightarrow like in Assign 1 die rolls
 anything I calculate using y is also random!

3.1 Random Sampling from the Population

\leftarrow holds an unknown truth

- 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**

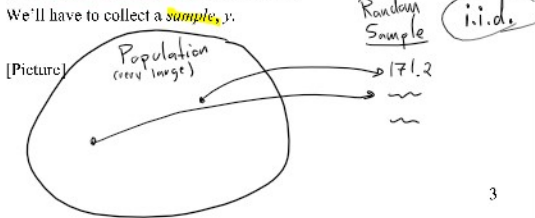
Example: suppose we want to know the mean height of a U of M student

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

- Population: all **male** students
- Population parameter of interest: μ_y \rightarrow **mean/expected height**

We can't afford to observe the whole pop.

We'll have to collect a **sample, y** .



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 all come from the **correct pop. (no mini-U)**
- **Independently** \rightarrow no connection/link btw people (no basketball teams)
- **Distributed**

anything that follows is also random!

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	170.2	187.2	178.3	170.3	179.4	181.2	180.0	178.9
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	180.9	173.4	168.7	180.0	175.1
175.7	180.8	176.2	170.8	177.3	168.4	180.3	177.1	191.2
171.0	180.3	169.5	167.2	178.0	172.9	176.0	176.5	174.0
175.1	184.2	165.3	180.2	178.3	183.4	179.9	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

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

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

basketball
mini-U

Example: mean income of ~~Canada~~ ^{people in Canada}. How could i.i.d. be violated?

sample mean to using these
average cannot everyone has one

How should we estimate the mean height?

We want μ_y . Use \bar{y} to estimate μ_y .

3.2 Estimators and Sampling Distributions

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

We collect our sample, $y = \{173.9, 171.7, 182.6, 181.3, 169.1, 174.0, 165.7, 182.2, 171.7, 168.1, 180.9, 175.7, 163.4, 186.3, 169.5, 171.9, 173.9, 172.0, 172.7, 172.0\}$. How should we use this sample to estimate the mean height?

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)

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

in reality don't know this (3.1)

- From heights example: $\bar{y} = 174.1$, $\mu_y = 176.8$ sample variance to estimate σ_y^2
- There are many ways to estimate μ_y . Examples? mode/median/geometric average/
- Why is (3.1) so popular? harmonic average/
- How good is \bar{y} at estimating μ_y in general? $\text{Var}(\bar{y}) = \text{Var}(y)$
- To answer these questions: idea of a sampling distribution

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 $\bar{y} = \{175.9, 175.3, 182.2, 178.6, 175.2, 180.3, 178.3, 183.7, 176.0, 167.4, 178.7, 178.7, 186.0, 175.6, 180.0, 168.7, 178.6, 173.1, 173.2, 182.1\}$, where the # in \bar{y} denotes that we are in the parallel universe. In this parallel universe, we got $\bar{y} = 177.6$. But in every universe, the population (table 3.1), is the same ($\mu_y = 176.8$).

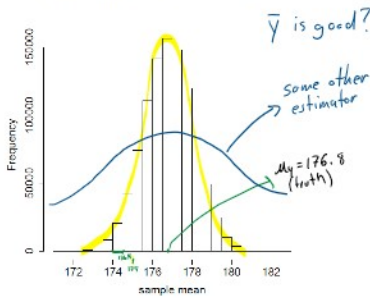
- Randomly sample from the population \rightarrow get y
 - y is random
- Use y to calculate \bar{y}
 - \bar{y} is random
 - could have gotten a different sample \rightarrow could have gotten a different \bar{y}
 - population is always the same (μ_y)

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 \bar{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



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

CLT

10

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

\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

\bar{y} , mode, median, max, etc.

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

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

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

Biased if it gives the wrong answer on average.

12

$$E[\bar{y}] = E\left[\frac{1}{n} \sum_{i=1}^n y_i\right]$$

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

Rules of the mean
 (1) $E[cy] = c \cdot E[y]$
 (2) $E[x+y] = E[x] + E[y]$

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

13

$\bar{y} = \frac{1}{n} \sum y_i$
 $E[\bar{y}] = \mu_y$ if unbiased
 $E\left[\frac{1}{n} \sum y_i\right] = \frac{1}{n} E[\sum y_i] = \frac{1}{n} E[y_1 + y_2 + \dots + y_n]$
 $= \frac{1}{n} \{E[y_1] + E[y_2] + \dots + E[y_n]\}$ *sample is i.i.d.*
 $= \frac{1}{n} \{\mu_y + \mu_y + \dots + \mu_y\} = \frac{1}{n} n\mu_y = \mu_y$ *identical*

Rules of the mean
 (i) $E[cY] = cE[Y]$
 (ii) $E[X+Y] = E[X] + E[Y]$

3.2.4 Efficiency

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} .

$Var(\bar{y}) = Var\left(\frac{1}{n} \sum y_i\right)$
 $= \frac{1}{n^2} Var(\sum y_i) = \frac{1}{n^2} Var(y_1 + y_2 + \dots + y_n)$
 $= \frac{1}{n^2} \{Var(y_1) + Var(y_2) + \dots + Var(y_n)\}$ *i.i.d.*
 $= \frac{1}{n^2} \{\sigma_y^2 + \sigma_y^2 + \dots + \sigma_y^2\} = \frac{n\sigma_y^2}{n^2}$
 $= \frac{\sigma_y^2}{n}$

Rules of var
 $Var(cY) = c^2 Var(Y)$
 $Var(X+Y) = Var(X) + Var(Y) + 2Cov(X,Y)$

if 2 variables are independent: then: \emptyset cov & corr.

\hookrightarrow independent
 \hookrightarrow if we random sample

14

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

- 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

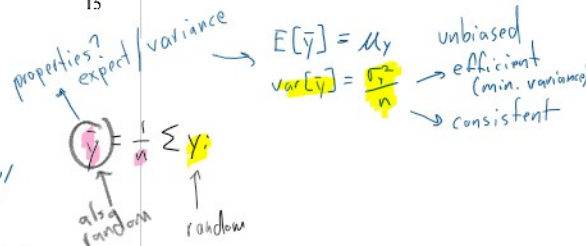
Suppose we had a lot of information. ($n \rightarrow \infty$)

What value should we get for our estimator? \rightarrow truth w/ prob. 1

How would state this mathematically?
 $\lim_{n \rightarrow \infty} Var(\bar{y}) \rightarrow 0$ and $\lim_{n \rightarrow \infty} E[\bar{y}] \rightarrow \mu_y$

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

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



16

"Null" \rightarrow reject / fail to reject
 3.3 Hypothesis tests (known σ_y^2)
 $H_0: \mu_y = \mu_{y,0}$ # we pick
 very unrealistic assumption
 almost all tests in Econ

3.3 Hypothesis tests (known σ_y^2)

$H_0: \mu_y = \mu_{y,0}$ # we pick

$H_A: \mu_y \neq \mu_{y,0}$ (2-sided alternative) (3.4)

almost all tests in Econ

- 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

In assign #1: $H_0: \mu_y = 3.5$

Example:

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

$H_0: \mu_y = 173$ (3.5)

$H_A: \mu_y \neq 173$

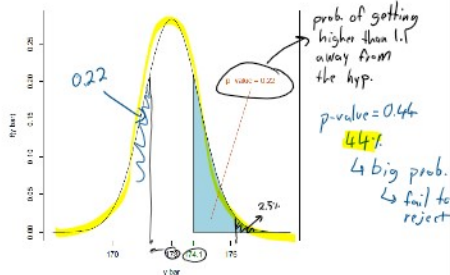
- Collect a sample: $y = \{173.9, 171.7, \dots, 172.0\}$
- Calculate $\bar{y} = 174.1$
- Suppose (very unrealistically that we know that) $\sigma_y^2 = 39.7$
- What now?

$\bar{y} - \mu_{y,0} = 174.1 - 173 = 1.1$

18

$\text{var}(\bar{y}) = \frac{\sigma^2}{n}$

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



19

Significance level
Pre-determined p-value that decided if you reject/fail to reject
10% / 5% / 1%

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

\hookrightarrow prob. of getting a \bar{y} that is more adverse to H_0 , compared to what we just observed.

3.3.1 Significance of a test

$\alpha = 10\% / 5\% / 1\%$ $P(H_0 \text{ is true}) = 0 \text{ or } 1$

3.3.2 Type I error

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

3.3.3 Type II error (and power)

$\text{pr}(\text{fail to reject } H_0 \mid H_0 \text{ is false}) = \beta = ?$

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

How big is $\mu_{y,0} - \mu_y$ truth?

20

3.3.4 Test statistics

\hookrightarrow z-test statistic
t-test stat.

- Just a more convenient way of getting the p-value for the test
- Each hypothesis test would present us with a new normal curve

$H_0: \mu_y = 1000$
 $\bar{y} = 1022.3$

21

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.

$$z = \frac{\text{estimate} - \text{hyp.}}{\sqrt{\text{var}(\text{estimator})}} = \frac{\bar{y} - \mu_{y,0}}{\sqrt{\frac{\sigma_y^2}{n}}}$$

$$H_0: \mu_y = 1000$$

$$\bar{y} = 1022.3$$

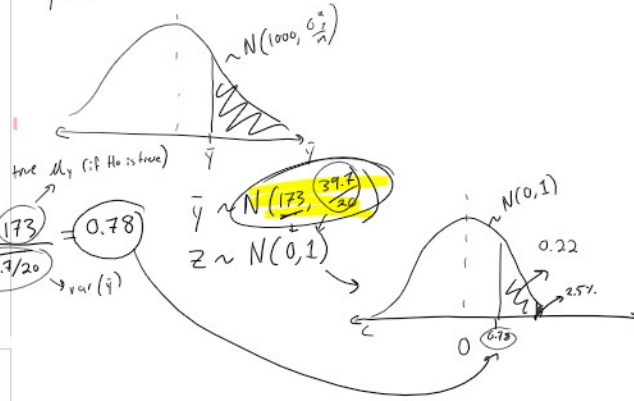


Table 3.2 Areas under the standard normal curve to the left of z

z	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5518	.5558	.5598	.5638	.5677	.5717	.5757
0.2	.5797	.5837	.5877	.5917	.5957	.5997	.6037	.6077	.6117	.6157
0.3	.6197	.6237	.6277	.6317	.6357	.6397	.6437	.6477	.6517	.6557
0.4	.6597	.6637	.6677	.6717	.6757	.6797	.6837	.6877	.6917	.6957
0.5	.6997	.7037	.7077	.7117	.7157	.7197	.7237	.7277	.7317	.7357
0.6	.7397	.7437	.7477	.7517	.7557	.7597	.7637	.7677	.7717	.7757
0.7	.7797	.7837	.7877	.7917	.7957	.7997	.8037	.8077	.8117	.8157
0.8	.8197	.8237	.8277	.8317	.8357	.8397	.8437	.8477	.8517	.8557
0.9	.8597	.8637	.8677	.8717	.8757	.8797	.8837	.8877	.8917	.8957
1.0	.8997	.9037	.9077	.9117	.9157	.9197	.9237	.9277	.9317	.9357
1.1	.9397	.9437	.9477	.9517	.9557	.9597	.9637	.9677	.9717	.9757
1.2	.9797	.9837	.9877	.9917	.9957					

3.3.5 Critical values

max z-stat before you reject
 1.96 for 5% significance
 if $|z| > 1.96$ we reject the H_0 at 5% level

3.3.6 Confidence intervals

What is the probability that our statistic will be within a certain interval if the null hypothesis is true? For example, what is the following probability?

$$\Pr(-1.96 \leq \frac{\bar{y} - \mu_{y,0}}{\sqrt{\frac{\sigma_y^2}{n}}} \leq 1.96) = 95\% \quad (3.12)$$

$$\Pr(-1.96 \leq \frac{\bar{y} - \mu_{y,0}}{\sqrt{\frac{\sigma_y^2}{n}}} \leq 1.96) = 0.95 \quad (3.13)$$

Finally, we solve equation 3.13 so that the null hypothesis $\mu_{y,0}$ is in the middle of the probability statement:

$$\Pr(\bar{y} - 1.96 \sqrt{\frac{\sigma_y^2}{n}} \leq \mu_{y,0} \leq \bar{y} + 1.96 \sqrt{\frac{\sigma_y^2}{n}}) = 0.95 \quad (3.14)$$

$$\bar{y} \pm 1.96 \times \text{s.e.}(\bar{y})$$

95%

In R , to do a t-test
 unbiased
 efficient
 consistent
 properties that are desirable
 \bar{y} is an example

3 ways to decide on H_0

- compare p-value to α
- compare z-stat to crit. value (1.96)
- see if H_0 is inside a confidence interval

Interpretation

μ_y
 $[-\bar{y} - \dots]$
 (i) 95% of such intervals contain the truth
 (ii) contains all null hypotheses you fail to reject

Ch 3

Assign #1 $\mu_{y,0}$

$$H_0: \mu_y = 173$$

$$H_A: \mu_y \neq 173$$

$$H_0: \mu_y = 3.5$$

$$H_A: \mu_y \neq 3.5$$

$$z = \frac{\bar{y} - \mu_{y,0}}{\sqrt{\text{var}(\bar{y})}}$$

$\text{var}(\bar{y}) = \frac{\sigma_y^2}{n}$ sample size
 $\sqrt{\text{var}(\bar{y})}$ s.e. (\bar{y}) "standard error"

3.4 Hypothesis Tests (unknown σ_y^2)

- Much more realistically, σ_y^2 (variance of y) will be unknown.
- Recall that: $\text{Var}[\bar{y}] = \frac{\sigma_y^2}{n}$
- $z = \frac{\bar{y} - \mu_{y,0}}{\text{s.e.}(\bar{y})} = \frac{\bar{y} - \mu_{y,0}}{\sqrt{\frac{\sigma_y^2}{n}}}$
- So, we need to estimate σ_y^2 in order to perform hypothesis tests.

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

$$\text{mean}(y) = E(\bar{y})$$

$$\text{var}(y) = E[(\bar{y} - \mu_y)^2]$$

3.4.1 Estimating σ_y^2

- A "natural" estimator:

$$E\left[\frac{1}{n} \sum (y_i - \bar{y})^2\right] = \frac{n-1}{n} \sigma_y^2$$

BIASED

$$E\left[\frac{1}{n-1} \sum (y_i - \bar{y})^2\right] = \sigma_y^2$$

3.4.1 Estimating σ^2

- A "natural" estimator:

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

$\frac{1}{n} \sum (y_i - \bar{y})^2$
 $E[s^2] = \frac{n-1}{n} \sigma^2$ (3.15) **BIASED**
 $E\left[\frac{1}{n-1} \sum (y_i - \bar{y})^2\right] = \sigma^2$

- Is this a good estimator? Why or why not? → it's biased

- A better estimator:

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

- Degrees-of-freedom correction → after \bar{y} , there are $n-1$ d.o.f. → lose 1 piece of information

$y = \{1, 3, \dots\}$
 $\bar{y} = 3$
 \downarrow
 5

So:

Estimated variance of $\bar{y} = \frac{\sigma^2}{n}$

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

$$z = \frac{\bar{y} - \mu_0}{\frac{s}{\sqrt{n}}}$$

replace σ^2 w/ s^2
 → z becomes t
 $z \sim N(0,1)$
 $t \sim t_{n-1}$

Note: for large n , the t test is equivalent to the z test

↓ determines "slope" of curve

