# Statistics Review

- A statistic is a *function* of a *sample of data*
- An *estimator* is a statistic
- Population parameter → unknown
- Estimator → 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

contains 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.
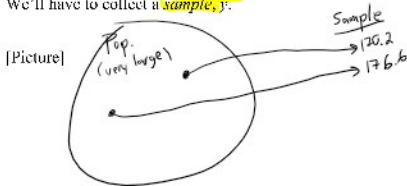- We'll use the sample, $y$, to estimate that something

2

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

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

- Population: all ~~male~~ students
- Population parameter of interest: $\mu_Y$

We can't afford to observe the whole pop.

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

[Picture]

Pop. (very large)    Sample
→ 120.2
→ 176.6

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 → come from pop. of U of M students (no mini-U students)
- Independently → 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 | 170.2 | 187.2 | 178.3 | 170.3 | 179.4 | 181.2 | 180.0 | 173.9 |
| 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 | 173.4 | 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 | 173.9 | 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, 162.1, 174.9, 165.7, 182.2, 171.7, 168.1, 189.9, 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?

sample mean / sample average / average

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

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*)

$$\left( \tilde{y} + \frac{1}{n} \left( \sum \right) \right) \sim N \qquad (3.1)$$

$$\prod_{i=1}^{n} y_i^{1/n} \quad \text{geometric average / thumb avg.}$$

- From heights example: $\tilde{y} = 174.1$, $\mu_y = 176.8$
- There are many ways to estimate $\mu_y$. Examples? median / mode
- Why is (3.1) so popular? it's the best
- How good is $\tilde{y}$ at estimating $\mu_y$ in general?
- To answer these questions: idea of a *sampling distribution*

$$\tilde{y}_{0.5} \quad \frac{min(y) + max(y)}{2}$$

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, 186.0, 175.6, 180.0, 168.7, 178.6, 173.1, 173.2, 187.1}, where the $*$ in $y^*$ denotes that we are in the parallel universe. In this parallel universe, we get $\tilde{y} = 177.6$. But in every universe, the population (table 3.1), is the same. $\tilde{y} = 174.1$

- Randomly sample from the population $\rightarrow$ get $y$
  - $y$ is random
- Use $y$ to calculate $\tilde{y}$
  - $\tilde{y}$ is random
  - could have gotten a different sample $\rightarrow$ could have gotten a different $\tilde{y}$
  - population is always the same ($\mu_y$)
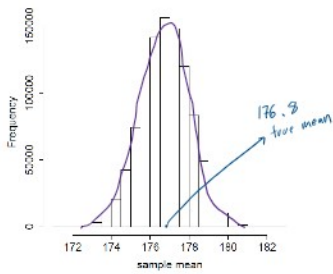
8

## 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 $\bar{y}$ 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



176.8
true mean

10

---

*Normal dist$^n$*

### 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 mean)

$\bar{y}$ 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) $E[X+Y] = E[X] + E[Y]$

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

12

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

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

$$= \frac{1}{n}E[y_1 + y_2 + \cdots + y_n] \qquad (3.2)$$

$$= \frac{1}{n}(E[y_1] + E[y_2] + \cdots + E[y_n])$$

$$= \frac{1}{n}(\mu_y + \mu_y + \cdots + \mu_y)$$

$$= \frac{n\mu_y}{n} = \mu_y$$

13

### 3.2.4 Efficiency → 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**
(i) $\text{var}(cY) = c^2\text{var}(Y)$
(ii) $\text{var}(X+Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X,Y)$

If 2 variables are independent $\Rightarrow \emptyset$ cov, $\emptyset$ corr.

$$\text{Var}(\bar{y}) = \text{var}\left(\frac{1}{n}\sum y_i\right)$$
$$= \frac{1}{n^2}\text{var}(\sum y_i) = \frac{1}{n^2}\text{var}(y_1 + y_2 + \cdots + y_n)$$
$$= \frac{1}{n^2}\{\text{var}(y_1) + \text{var}(y_2) + \cdots + \text{var}(y_n)\} \quad \longrightarrow \text{ i.i.d.} \quad \hookrightarrow \text{"independent"}$$
$$= \frac{1}{n^2}\{\sigma_Y^2 + \sigma_Y^2 + \cdots + \sigma_Y^2\}$$
$$= \frac{1}{n^2}n\sigma_Y^2 = \boxed{\frac{\sigma_Y^2}{n}}$$

14

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

$$= \frac{1}{n^2}\text{Var}\left[\sum_{i=1}^{n} y_i\right]$$

$$\frac{1}{n^2}\text{Var}[y_1 + y_2 + \cdots + y_n]$$

$$= \frac{1}{n^2}(\text{Var}[y_1] + \text{Var}[y_2] + \cdots + \text{Var}[y_n]) \qquad (3.3)$$

any other estimator for $\mu$

$$= \frac{1}{n}(\sigma_y^2 + \sigma_y^2 + \cdots + \sigma_y^2)$$

$$= \frac{n\sigma_y^2}{n^2} = \frac{\sigma_y^2}{n}$$

$\text{var}(\bar{y}) < \text{var}(\hat{\mu}_y)$
$\longrightarrow \bar{y}$ 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. $(n \to \infty)$
What value should we get for our estimator? right answer, every time
How would state this mathematically?

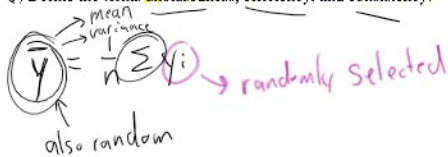$\lim\limits_{n\to\infty} \text{var}(\bar{y}) \to 0$ ✓     $\lim\limits_{n\to\infty} \text{bias}(\bar{y}) \to 0$ ✓

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

$\text{var}(\bar{y}) = \frac{\sigma^2}{n}$     $\lim\limits_{n\to\infty} \frac{\sigma^2}{n} \to 0$

"variance goes away"

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

mean
variance
$\bar{y} = \frac{1}{n}\sum y_i \longrightarrow$ randomly selected
also random

$E[\bar{y}] = \mu_y$     unbiased
$\text{var}[\bar{y}] = \frac{\sigma_y^2}{n} \to$ efficient
$\hookrightarrow$ consistent

16

very unrealistic

null   ### 3.3 Hypothesis tests (known $\sigma_y^2$)

a number
true
alternative   $H_0: \mu_y = \mu_{y,0} \qquad$ almost always 2-sided in Econ
$H_A: \mu_y \neq \mu_{y,0}$ (2-sided alternative) $\qquad (3.4)$

## 3.3 Hypothesis tests (known $\sigma_y^2$)

null

alternative

$H_0: \mu_y = \mu_{y,0} \rightarrow$ a number $\rightarrow$ almost always 2-sided in Econ

true

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

- Estimate $\mu_y$ (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 → fail to reject
- If it's far → reject

17

---

Example:

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

z-test
t-test

$H_0: \mu_y = 173$

$H_A: \mu_y \neq 173$   (3.5)

$(174.1 - 173) = 1.1 \text{ cm}$

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

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

$\sim N\left(173, \frac{39.7}{20}\right)$



173   174.1
       $\bar{y}$

18

---

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.



0.22

p-value = 0.22

prob. of getting a $\bar{y}$ (than further away (than what we just calc.)

p-val = 44%

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 $H_0$. → fail to reject

### 3.3.1 Significance of a test

↳ pre-determined p-value that decides reject/fail to reject

$\alpha = 10\%, \boxed{5\%}, 1\%. \rightarrow$ if p-val < 5% ⇒ 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)

$Pr(\text{fail to reject } H_0 \mid H_0 \text{ is false})$

depends on "how" false $H_0$

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

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

20

- 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.

$\sim N(173, 2)$     $\sim N(3.5, \frac{2.12}{n})$

174.1     4.1

$N(0,1)$

$\bar{y} \sim N(173, \frac{34.7}{20})$

$Z \sim N(0,1)$

$\sim N(0,1)$

$Z = \dfrac{\text{estimate} - \text{hypothesis}}{\sqrt{Var(\text{estimator})}} = \dfrac{\bar{y} - \mu_{y,0}}{\sqrt{\sigma^2}/n} = \dfrac{174.1 - 173}{\sqrt{\frac{34.7}{20}}} = 0.78$

$\bar{y} \sim N$

0.22

$\sim N(0,1)$

0.22

$\bar{4} - 0.025$

$p\text{-val} = 0.22 \times 2$
$= 0.44 > .05$
$\hookrightarrow$ fail to reject

0   0.78

Table 3.2: Area under the standard normal curve, to the right of $z$.

1% sig $\to$ 2.56
10% sig $\to$ 1.65

**3.3.5 Critical values** $\to$ pre-determined max $z$-stat before ($t$-stat) you reject

1.96 for a 5% sig level   if $|z| > 1.96 \to$ reject

**3.3.6 Confidence intervals**

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

5% crit-value

$$Pr(-1.96 \leq z \leq 1.96)? \qquad (3.12)$$

$$Pr\left(-1.96 < \frac{\bar{y} - \mu_{y,0}}{\sqrt{\sigma^2/n}} < 1.96\right) = 0.95 \qquad (3.13)$$

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

$$Pr\left(\bar{y} - 1.96 \times \sqrt{\frac{\sigma^2}{n}} \leq \mu_{y,0} \leq \bar{y} + 1.96 \times \sqrt{\frac{\sigma^2}{n}}\right) = 0.95 \qquad (3.14)$$

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

23

- **unbiased** ⎫
- **efficient** ⎬ desirable properties for an estimator
- **consistent** ⎭

$\bar{y}$

**3 ways to decide about $H_0$**

(i) compare p-value to $\alpha$
(ii) compare $z$-stat to crit. value (1.96)
(iii) check if $H_0$ is inside C.I.
    $\hookrightarrow$ fail to reject

$[- \bar{y} -]$

$\uparrow$ $H_0$ inside $\to$ fail to reject

**Ch.3**

$H_0: \mu_y = 173$
$H_a: \mu_y \neq 173$

**Assign #1**
$H_0: \mu_y = 3.5 \leftarrow \mu_{y,0}$
$H_a: \mu_y \neq 3.5$

$Z = \dfrac{\bar{y} - \mu_{y,0}}{\sqrt{Var(\bar{y})}} \to$ "standard error" of $\bar{y}$
      s.e.$(\bar{y})$

**Interpretation of a C.I.**

$\mu_y$

$[- \bar{y} -]$
$[- \bar{y}^* -]$
$[- \bar{y}^{**} -]$

(i) 95% prob. that C.I contains truth

(ii) contains all null hypotheses that you fail to reject

$\mu_{y,0}$      $\mu_{y,0}$

## 3.4 Hypothesis Tests (unknown $\sigma_y^2$)

- Much more realistically, $\sigma_y^2$ (variance of $y$) will be unknown.

- Recall that: $Var[\bar{y}] = \dfrac{\sigma_y^2}{n}$

- $z = \dfrac{\bar{y} - \mu_{y,0}}{s.e.(\bar{y})} = \dfrac{\bar{y} - \mu_{y,0}}{\sqrt{\sigma_y^2/n}} \to s.e.(\bar{y})$

- So, we need to estimate $\sigma_y^2$ in order to perform hypothesis tests.

$\bar{y}$    24

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

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

## 3.4.1 Estimating $\sigma_y^2$

- A "natural" estimator:

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

$\ldots \frac{n-1}{?} \ldots$   **UNBIASED**

## 3.4.1 Estimating $\sigma_y^2$

$\text{var}(y) = E\left[(y - E(y))^2\right]$ estimate $\frac{1}{n}\sum(y-...)$

- A "natural" estimator:

$$\hat{\sigma}^2 = \frac{1}{n}\sum_{i}^{n}(y_i - \bar{y})^2 \qquad (3.15)$$

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

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

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

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

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

- Degrees-of-freedom correction

$$y = \{1, 3, ?\}$$
$$5$$
$$\bar{y} = 3$$
$$\rightarrow \text{calc. } \bar{y} \rightarrow \text{lose 1 piece info.}$$
$$25$$

---

So:

Estimated variance of $y \rightarrow \frac{\hat{\sigma}^2}{n} \rightarrow \frac{1}{n-1}\sum(y_i-\bar{y})^2$

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

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

$$\rightarrow \text{random } (\chi^2)$$

using $s_y^2$ instead of $\sigma^2$ makes $z$ become $t$

$$t \sim t_{n-1}$$

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

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

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



$\sim N(0,1)$
$t_n$
$t_{\text{large}}$

26