

4.3 CENTRAL LIMIT THEORY

WK 7-2

MOST IMPORTANT THEORY IN STATS! IT SAYS! ALMOST EVERYTHING EVENTUALLY BECOMES NORMAL!

STANDARD DEFINITION:

SUPPOSE THAT THE RANDOM VARIABLES $X_1, X_2, X_3, \dots, X_n$ ARE DRAWN FROM A LARGER POPULATION INDEPENDENTLY (therefore $X_1, X_2, X_3, \dots, X_n$ are identically distributed) WITH A COMMON MEAN μ AND THE STANDARD DEVIATION σ . LET

WE GET 3 PIECES OF INFO \rightarrow

SAMPLE TOTAL

①

\rightarrow have identical distributions

$$T_0 = X_1 + X_2 + \dots + X_n$$

REQUIREMENT: LARGE SAMPLE SIZE!

② EXPECTED VALUE

③ SD

THEN $E(T_0) = n\mu$ AND $\sigma_{T_0} = \sigma\sqrt{n}$ (standard error)

FURTHERMORE, IF n IS LARGE (≥ 30 OR ≥ 40), THEN WE MAY CLOSELY APPROXIMATE

\uparrow
IF YOU MAKE A RANDOM VARIABLE FROM MANY, IT WOULD BE A NORMAL DISTRIBUTION

THE DISTRIBUTION OF T_0 BY A NORMAL DISTRIBUTION.

\rightarrow WE DON'T SELECT SAMPLE W/ REPLACEMENT, BUT FOR THE THEORY, THEY NEED TO BE INDEPENDENT. SO! WE NEED LARGE SAMPLE SIZE THAT THE PROBABILITY OF REPEATING IS LOW SO WE CAN SELECT SAMPLE W/ REPLACEMENT.

\hookrightarrow MAGICALLY, THE THEOREM PROVES $X_1 + X_2 + \dots + X_n$ WILL BE NORMAL

\circ BUT APPLICATION REALLY IS CASE BY CASE. IF DATA IS SKEWED, YOU WOULD NEED EVEN LARGER SAMPLE SIZES

\circ IRL IF YOU DON'T KNOW COMMON μ + σ YOU CAN EVEN USE THE AVERAGE μ + σ

WAIT.. WE HAVE 4 CLT'S \sim THERE ARE 4 CASES

2 SUPPOSE THAT THE RANDOM VARIABLES X_1, X_2, \dots, X_n ARE DRAWN WITH REPLACEMENT FROM A LARGE POPULATION INDEPENDENT + IDENTICALLY DISTRIBUTED WITH COMMON MEAN μ + STANDARD DEVIATION σ . LET:

PAY ATTENTION TO TOTAL / AVERAGE IN HW

SAMPLE AVERAGE (MEAN)

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

THEN

basically derived from SAMPLE TOTAL

$$E(\bar{X}) = \mu \text{ AND } \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \text{ (standard error)}$$

FURTHERMORE, IF n IS LARGE (≥ 30 OR ≥ 40) THEN WE MAY CLOSELY APPROXIMATE THE DISTRIBUTION OF \bar{X} BY A NORMAL DISTRIBUTION

3 A SPECIAL CASE OF 1

\rightarrow identical

ADDING BERNOLLI'S! \rightarrow BINOMIAL DIST.!

SUPPOSE THAT THE RANDOM VARIABLES X_1, X_2, \dots, X_n ARE DRAWN WITH REPLACEMENT. SUPPOSE ALSO THAT $P(X_i = 1) = p$ AND $P(X_i = 0) = 1 - p$

1 = SUCCESS 0 = FAILURE

SAMPLE TOTAL COUNT

THEN X HAS A BINOMIAL DISTRIBUTION WITH $E(X) = np$ AND $\sigma_x = \sqrt{np(1-p)}$ ALSO CALLED STANDARD ERROR

TOTAL SUCCESS

FURTHERMORE, IF n IS LARGE ($np \geq 10$ AND $n(1-p) \geq 10$), THEN WE MAY APPROXIMATE CLOSELY THE DISTRIBUTION OF X BY A NORMAL DISTRIBUTION

THIS APPROXIMATES A BINOMIAL DISTRIBUTION INTO A NORMAL ONE (BUT IT'S NOT EXACT ANSWER)

BUT IN HW DON'T USE NORMAL TO APPROXIMATE BINOMIAL

4 WE HAVE 3 TO GET THE PERCENTAGE!

SUPPOSE THAT THE RANDOM VARIABLES X_1, X_2, \dots, X_n ARE DRAWN WITH REPLACEMENT. SUPPOSE ALSO THAT $P(X_i = 1) = p$ AND $P(X_i = 0) = 1 - p$

1 = SUCCESS 0 = FAILURE

SAMPLE PERCENTAGE (PROPORTION)

$$\hat{p} = \frac{X}{n} = \frac{X_1 + X_2 + X_3 + \dots + X_n}{n}$$

\rightarrow THIS IS NOT BINOMIAL, YOU HAVE TO USE CLT TO GET PERCENTAGE

$$\text{THEN } E(\hat{p}) = p \text{ AND } \sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} \text{ + STANDARD ERROR}$$

FURTHERMORE, IF n IS LARGE ($np \geq 10$ AND $n(1-p) \geq 10$), THEN WE MAY APPROXIMATE CLOSELY THE DISTRIBUTION OF \hat{p} BY A NORMAL DISTRIBUTION

ALL THE CLT'S

SAMPLE TOTAL

$$E(T_0) = n\mu$$

$$\sigma_{T_0} = \sigma\sqrt{n}$$

SAMPLE AVERAGE

$$E(\bar{X}) = \mu$$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

TOTAL SUCCESS

$$E(X) = np$$

$$\sigma_x = \sqrt{np(1-p)}$$

SAMPLE PERCENTAGE

$$E(p) = p$$

$$\sigma_p = \sqrt{\frac{p(1-p)}{n}}$$

FOLLOWS NORMAL DISTRIBUTION IF

n IS LARGE

ACTUALLY BINOMIAL

DON'T USE NORMAL APPROX. WHEN
COMPUTING PROBABILITY

EXAMPLE 1: FINDING TOTAL APPROXIMATION

THIRTY-SIX JETS WAIT TO TAKE OFF FROM AN AIRPORT. THE AVERAGE TAXI & TAKE-OFF TIME FOR EACH JET IS 8.5 MINUTES, WITH AN SD OF 2.5 MINUTES. WHAT IS THE PROBABILITY THAT THE TOTAL TAXI & TAKE OFF TIME FOR THE 36 JETS IS LESS THAN 320 MIN?

WHAT IS THE VARIABLE? TIME

WHAT DISTRIBUTION DOES IT FOLLOW? WE DON'T KNOW, BUT WE CAN
FIND TOTAL
IT IS CONTINUOUS

↓

TOTAL = SAMPLE SUM

$V_{tot} = \text{time taxi} + \text{take off}$

$$T_0 = X_1 + X_2 + \dots + X_{36}$$

↑
time per jet

$$\mu = 8.5 \quad \sigma = 2.5 \quad n = 36$$

$$E(T_0) = \mu_{T_0} = 8.5 \cdot 36$$

$$= 306$$

$$SD(T_0) = \sigma_{T_0} = \sqrt{n} \sigma = \sqrt{36} \cdot 2.5$$

$$= 15$$

WE KNOW TOTAL HAS NORMAL DISTRIBUTION

$$P(T_0 < 320) = \text{normalcdf}(-1000, 320, 306, 15) \approx 0.8247$$

EXAMPLE 2:

THE COOKIE MACHINE AT CHIPS AHOY ADDS A RANDOM # OF CHIPS TO EACH COOKIE. THE NUMBER OF CHIPS IS A RANDOM # WITH AVERAGE 28.5 AND SD 5.3. FIND THE PROBABILITY THAT, IN A BAG OF 50 COOKIES, THE AVERAGE # OF CHIPS PER COOKIE IS AT LEAST 30.

WHAT IS THE VARIABLE? # OF CHIPS

WHAT DISTRIBUTION DOES IT FOLLOW? WE DON'T KNOW

IT IS DISCRETE THO.

$$\mu = 28.5$$

$$\sigma = 5.3$$

$$n = 50$$

↑
each cookie is a sample

$$E(\bar{X}) = \mu = 28.5$$

$$SD = \frac{\sigma}{\sqrt{n}} = \frac{5.3}{\sqrt{50}}$$

$$\bar{X} \sim \text{NORMAL}(28.5, \frac{5.3}{\sqrt{50}})$$

$$P(\bar{X} \geq 30) = \text{NORMALCDF}(30, 1000, 28.5, \frac{5.3}{\sqrt{50}})$$

$$= 0.8227$$



↑
PUT IN EXACT,
NO ROUNDING ERROR

EXAMPLE 3: BINOMIAL SUCCESS

IT IS KNOWN THAT 40% OF PEOPLE IN THE CITY ARE INFECTED BY COVID-19 (SOME NO, SOME MILD, SOME SEVERE SYMPTOMS).

A RANDOM SAMPLE OF 50 PEOPLE ARE SELECTED.

a) WHAT IS THE PROBABILITY THAT BETWEEN 200 + 210 ARE INFECTED? (USE BINOMIAL THEN NORMAL APPROX. TO SEE HOW ACCURATE)

b) WHAT IS THE PROBABILITY THAT THE SAMPLE PORTION IS BETWEEN 39% AND 43%?