

- Based on our sample we hope to get a good "estimate" for a summary of the population
- A summary about the population is a parameter, and a summary for our sample is a statistic

- How good is our statistic in estimating our parameter?

Parameter

$\mu$  = mean

$p$  = proportion

Statistic

$\bar{x}$  = sample mean

$\hat{p}$  = sample proportion

In the case where we use a single estimate to represent the entire population, is referred to as a point estimate

Alternatively, we can use an Interval estimate to specify plausible values for our parameter of interest

- There will always be an error from our estimate.  
our goal is to reduce possible errors that can occur.  
such as

Sampling Error: Describes how much an estimate will vary from sample to sample. One goal is to quantify such error using the sample size " $n$ "

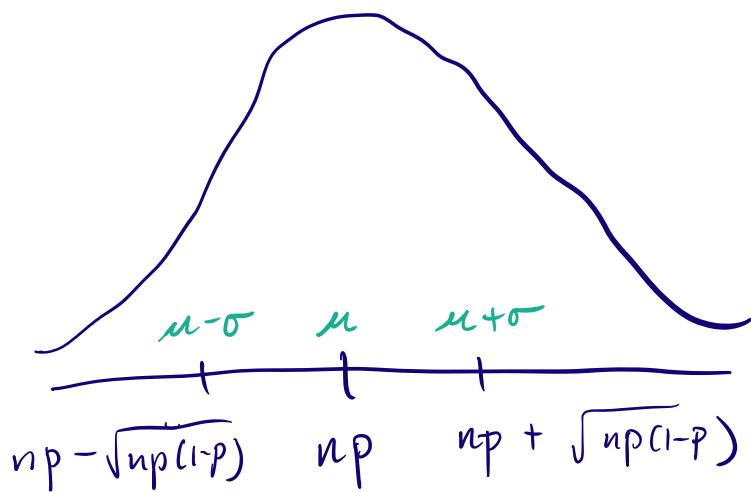
BIAS: Inappropriate Sampling Techniques, such as Convenience Sampling

\* Whenever Simple Random Sample is implemented observations from our sample are considered independent

## Normal Approximation to the Binomial Distn:

- The Binomial distribution with probability of success  $p$  is approximately normal whenever our sample size " $n$ " is sufficiently large.
- We consider sufficiently large whenever  $np \geq 10$  and  $np \underbrace{(1-p)}_{\text{failure}} \geq 10$

If the above conditions hold then the above Binomial distn will approximately be normal with parameters  $\mu = np$  and  $\sigma = \sqrt{np(1-p)}$



It turns out this is only one small piece of an Important Theorem which provides a foundation for most of statistics.

## Central Limit Theorem

- recall, the Law of Large Numbers which states as we increase the number of experiments we perform the sample mean will converge to the theoretical mean.

$$\bar{X} \xrightarrow{\text{as } n \rightarrow \infty} \mu$$

Likewise, for Bernoulli trials where we consider the proportion of successes

$$\hat{P} \xrightarrow{\text{as } n \rightarrow \infty} P$$

where  $\hat{P}$  is the proportion of successes from our sample

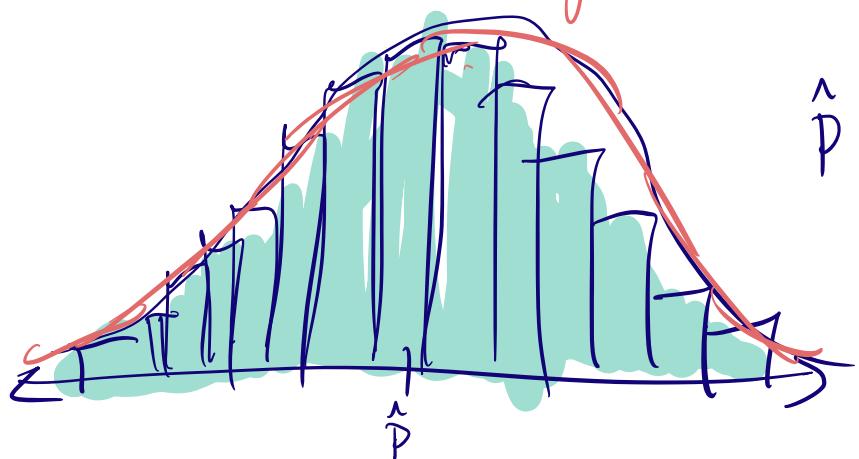
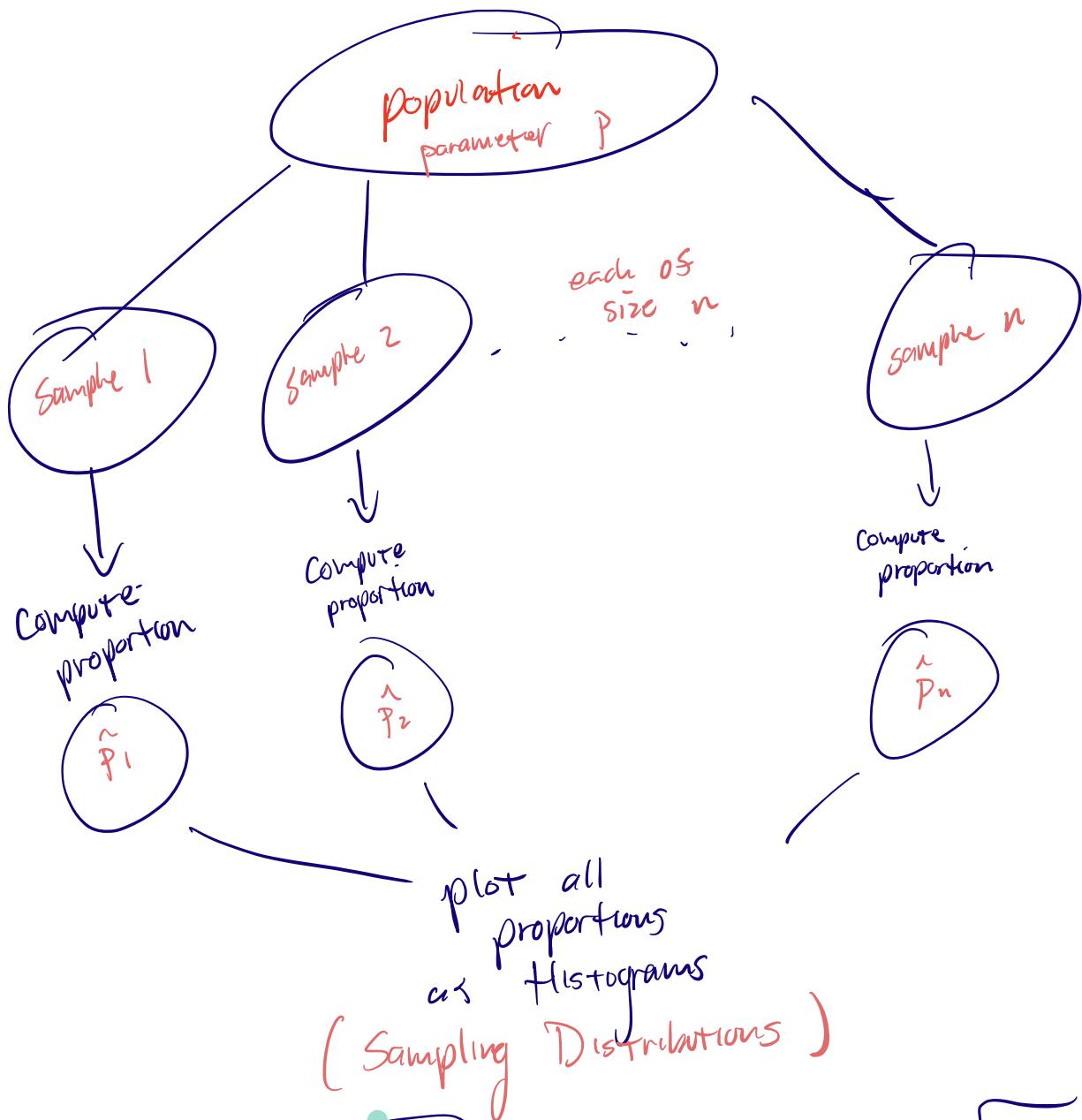
Consider proportions first, we'll do means later.

### Central Limit Theorem:

Whenever our observations are independent and the sample size is sufficiently large enough ( $np \geq 10$  and  $n(p(1-p)) \geq 10$ ) the distribution of our sample proportions  $\hat{P}$  will follow a normal distribution with

$$\mu = \hat{P} \quad \text{and} \quad SE(\hat{P}) = \sqrt{\frac{P(1-P)}{n}}$$

where  $SE(\hat{P})$  describes the variability of the distribution of our sample proportions  $\hat{P}$  [Sampling distribution]



$$\hat{p} \sim N\left(p, \sqrt{\frac{p(1-p)}{n}}\right)$$

by CLT

- \* The Center of The Sampling distribution are always at The population parameter.
- \* The Variability in our Sampling distn decreases as n increases, giving us more precise estimates
- \* In most cases, we don't know The population parameter , so  $\hat{p}$  is a reasonable substitute for p when checking assumptions such as  $[np \geq 10 \text{ and } n\hat{p}(1-\hat{p}) \geq 10]$

Similarly , we could plug in  $\hat{p}$  to our Standard Error of the Sample proportion

$$SE_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

## § 5.2 Confidence Interval for a Proportion

- Recall A point estimate is a single number estimate for the population parameter. It is more ideal to give a range of plausible values for the parameter instead with some given confidence, provided a large sample is utilized.

$$\hat{P} \rightarrow P \quad (\text{point estimate})$$

$$(L\beta, U\beta) \rightarrow P \quad (\text{Interval Estimate})$$

where  $L\beta = \text{point estimate} - \text{Margin of Error}$

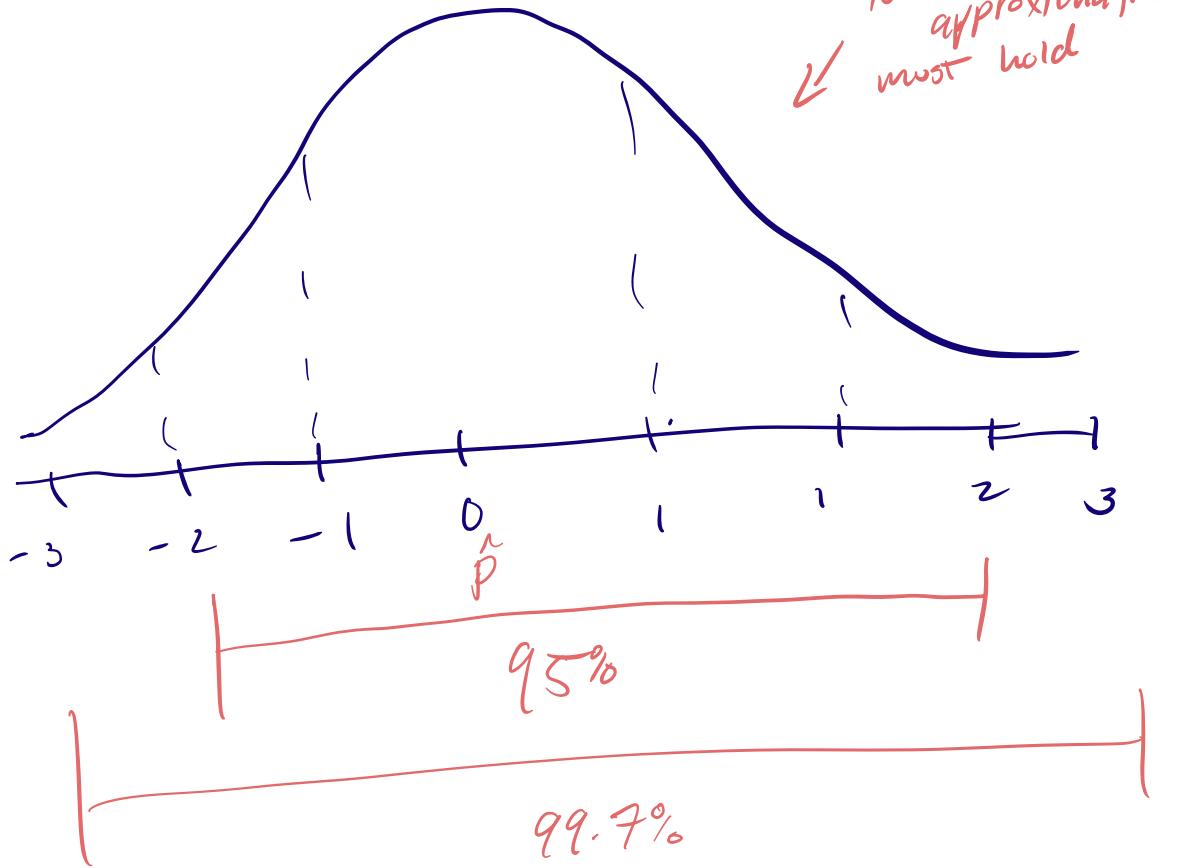
$$U\beta = \text{point estimate} + \text{Margin of Error}$$

A margin of error consists of the Standard Error and some Confidence Level

$$\text{Margin of Error} = \text{Confidence level selected} \times \text{Standard Error}$$

$$ME = z^* SE$$

Value is obtained from a Standard Normal distn given a Certain Confidence.



95% extends to 1.96 and -1.96

## 95% Confidence Interval

$$(LB, UB) = (\hat{p} - ME, \hat{p} + ME)$$

where  $ME = z^* \sqrt{\frac{p(1-p)}{n}}$  and  $z^* = 1.96$

All together,

A 95% CI is

$$\left( \hat{p} - 1.96 \sqrt{\frac{p(1-p)}{n}}, \hat{p} + 1.96 \sqrt{\frac{p(1-p)}{n}} \right)$$

$\text{LB}$      $\text{UB}$

\* We are 95% Confident The population proportion of "1" is between LB and UB  
statement about population parameter

- obtain Point Estimate, check normal Assumptions
- calculate ME
- calculate LB/UB
- Interpret your Interval Estimate

Show  
Simulation!

90% CI: use  $Z^* = 1.65$

95% CI: use  $Z^* = 1.96$

99% CI: use  $Z^* = 2.58$

Q: What happens to the length of the interval if we increase our confidence?

ex: You work at Sunblock Brand A, and your main competitor is Brand B. You TEST both brands of sunblock on a SRS of 28 subjects. Brand A provides better UV protection for 18 subjects, but Brand B provides better protection for the other 10 subjects.

(a): Show that the assumptions for a large-sample Confidence Interval are satisfied.

- We have SRS, we can assume subjects effectiveness for each subject are identically independently distributed (iid)

-  $\hat{P} = \frac{18}{28}$ , we can substitute  $P$  with  $\hat{P}$

and verify  $n\hat{P} = 18$  and  $n(1-\hat{P}) = 10$

Since both of these are at least 10, by CLT

$$\hat{P} \sim N(P, \sqrt{\frac{P(1-P)}{n}})$$

(b) Construct and Interpret a 95% CI for the proportion of people who get better UV protection from brand A.

- our sample proportion is  $\hat{P} = \frac{18}{28}$  (point estimate)

-  $\hat{P} \sim N(P, \sqrt{\frac{P(1-P)}{n}})$ .

$$SE_{\hat{P}} = \sqrt{\frac{\hat{P}(1-\hat{P})}{n}} = \sqrt{\frac{(0.643)(1-0.643)}{28}} = 0.091$$

- with 95% CI level, we have

$z^* = 1.96$  since normal approximation holds.

our Margin of Error is

$$\begin{aligned} ME &= z^* \bar{S}E_{\bar{p}} \\ &= 1.96 (0.091) \\ &= 0.177 \end{aligned}$$

A 95% CI =  $(LB, UB)$

where  $LB = \hat{p} - ME = 0.643 - 0.177$

$$UB = \hat{p} + ME = 0.643 + 0.177$$

therefore, 95% CI =  $(0.465, 0.820)$

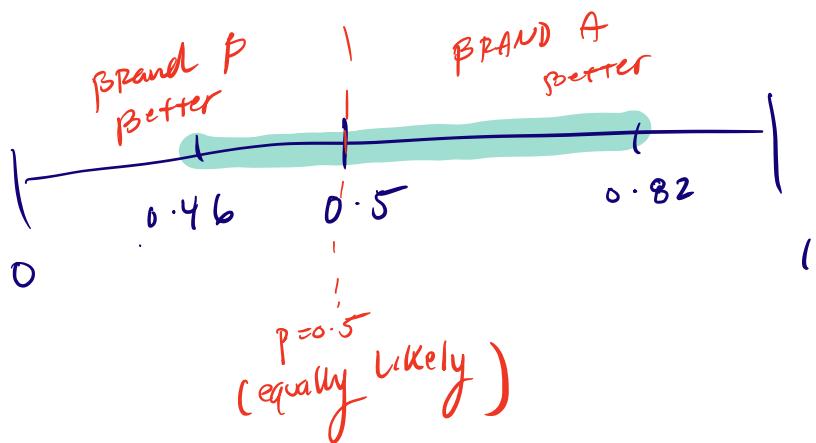
\* We are 95% CI that between 46.5% and 82% of people get better UV protection

from Brand A

(c). Given the above Confidence Interval, can we claim that Brand A is better than Brand B?

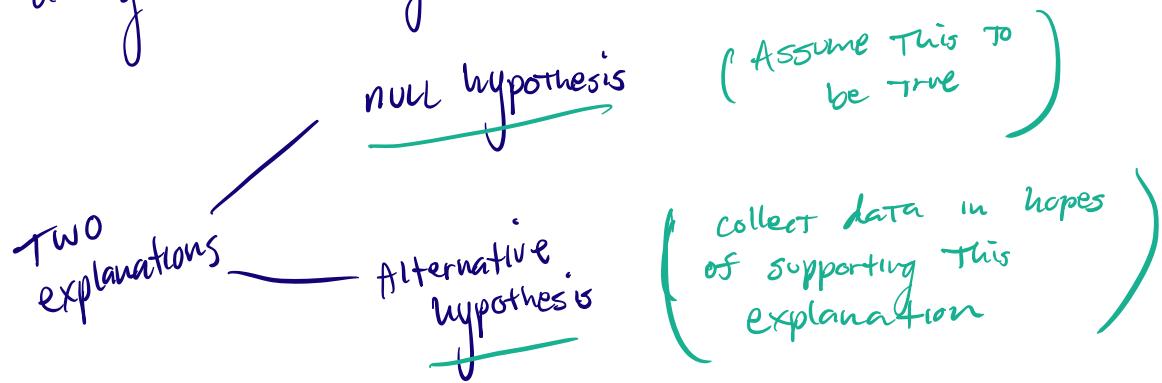
Nope! If  $p > 0.5$  Then brand A proves better protection. However our interval contains values where  $p < 0.5$ . Therefore it is plausible for brand B to be better. Similarly our interval also contains  $p = 0.5$  (Both are the same)

\* Inconclusive Confidence Interval.



### § 5.3 Hypothesis Test for Proportions

we want to test two competing explanations  
during an investigation



**null hypothesis**: explanation that claims nothing is happening (no difference)

**Alternative hypothesis**: explanation that claims something is happening (some difference)

we are interested in trying to "prove" this is TRUE.

**Test of Significance**:  
Collect and Evaluate Evidence  
to Decide between The  
two  
Competing explanations  
(null vs Alternative )

Statement about  
a population  
parameter

- Based on our Test of significance, There are two possible outcomes

(1) we fail to Reject The null hypothesis

(2) we Reject The null hypothesis in favor of Alternative hypothesis

\* If we fail to reject our null hypothesis  
(we did not have enough evidence to support the Alternative hypothesis)

That does not imply The null hypothesis is true



Recall: you work at Sunblock Brand A.

from 28 subjects, Brand A provides better UV protection for 19 subjects, but Brand B provides better protection for the remaining 10 subjects.

Ex: Translate The Real-world claim  
 "people get better UV protection from Brand A than Brand B" into null and alternative hypotheses  
 (in words and Mathematical Sentences)

$H_0$ : Both Brand A and Brand B provide the same UV protection

$H_a$ : Brand A provides better UV protection than Brand B.

Define parameter  
 $\downarrow$

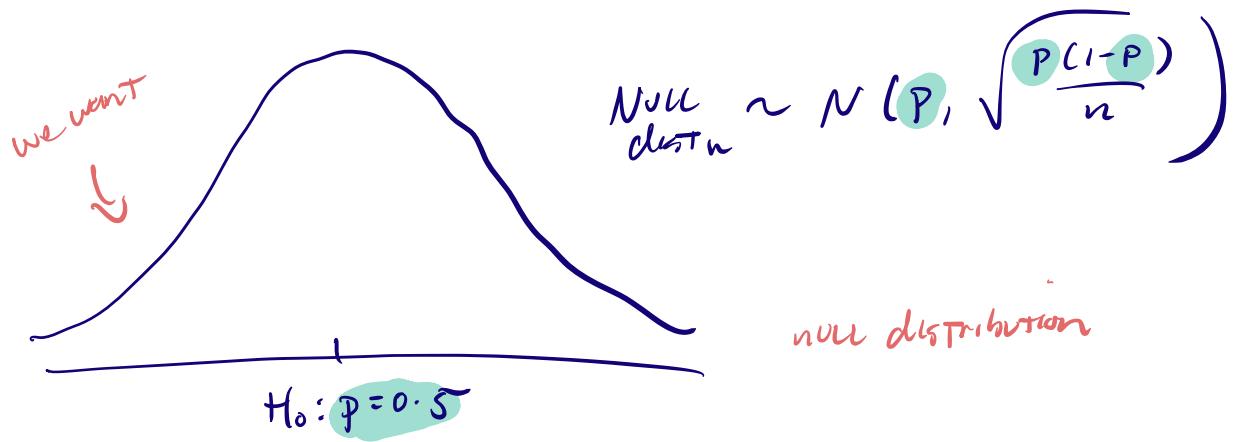
Let  $P$  = population proportion of individuals for which Brand A was better than Brand B

$H_0: P = 0.5$  (assume to be true) ] Specify null / Alternative

$H_a: P > 0.5$  (try to show its true)

## Obtain and Evaluate Evidence

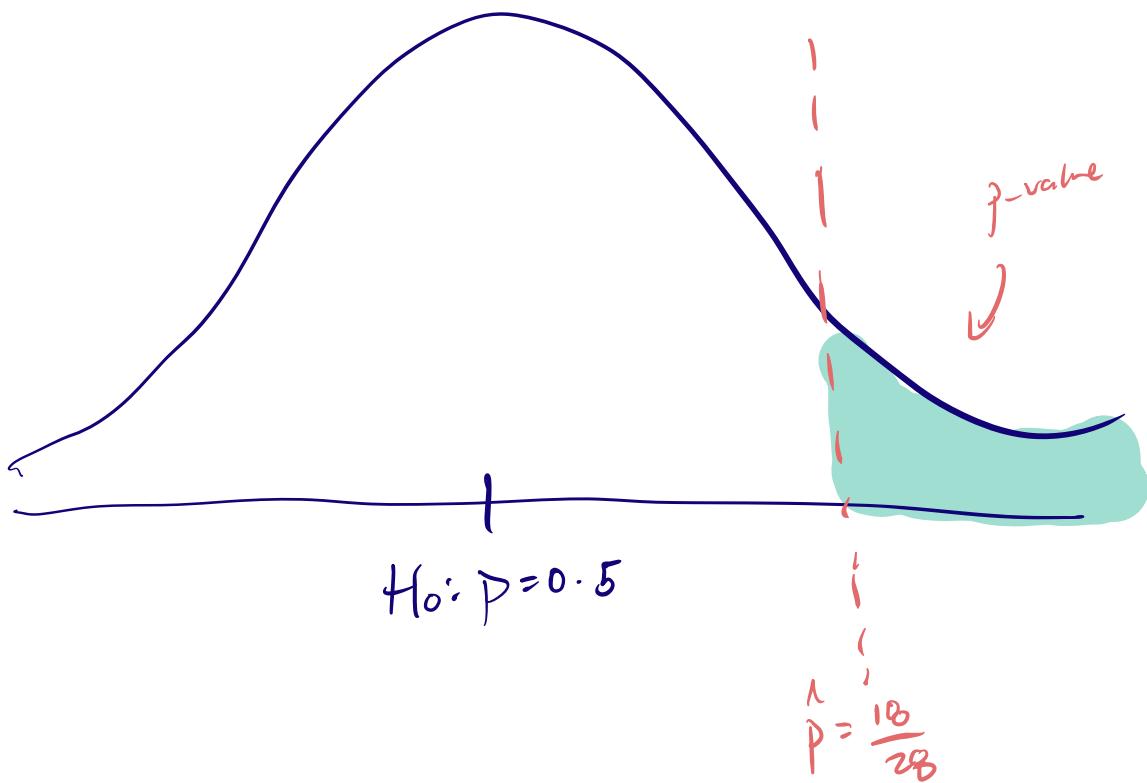
- Specify an initial significance level.  
(A threshold to determine if our results are significant)
- we use significance level =  $\alpha$ , it is common to use  $\alpha = 0.05$ , but it is not required.
- ~ check large-sample assumptions for our null distribution



$np \geq 10$  and  $n(1-p) \geq 10$  using our null value

$28(\frac{1}{2}) \geq 10$  and  $28(1-\frac{1}{2}) \geq 10$ , since it was SRS independence also holds.

\* If the null distribution were true,  
 Then our sampling distn should follow a  
 Normal distn with parameter  $N(p, \sqrt{\frac{p(1-p)}{n}})$   
 where  $p = 0.5$  (our null value)



- Our observed statistic will be used to test the evidence we have in favor of the alternative hypothesis
- We measure such evidence using a p-value.

Defn: A p-value is the probability of observing a test statistic giving us "as much or more evidence in favor" of the alternative hypothesis assuming the null distn is true

In our Example,

$$Z = \frac{\text{observed} - \text{null}}{SE}$$

$$Z = \frac{18/20 - 0.5}{\sqrt{\frac{0.5(0.5)}{20}}}$$

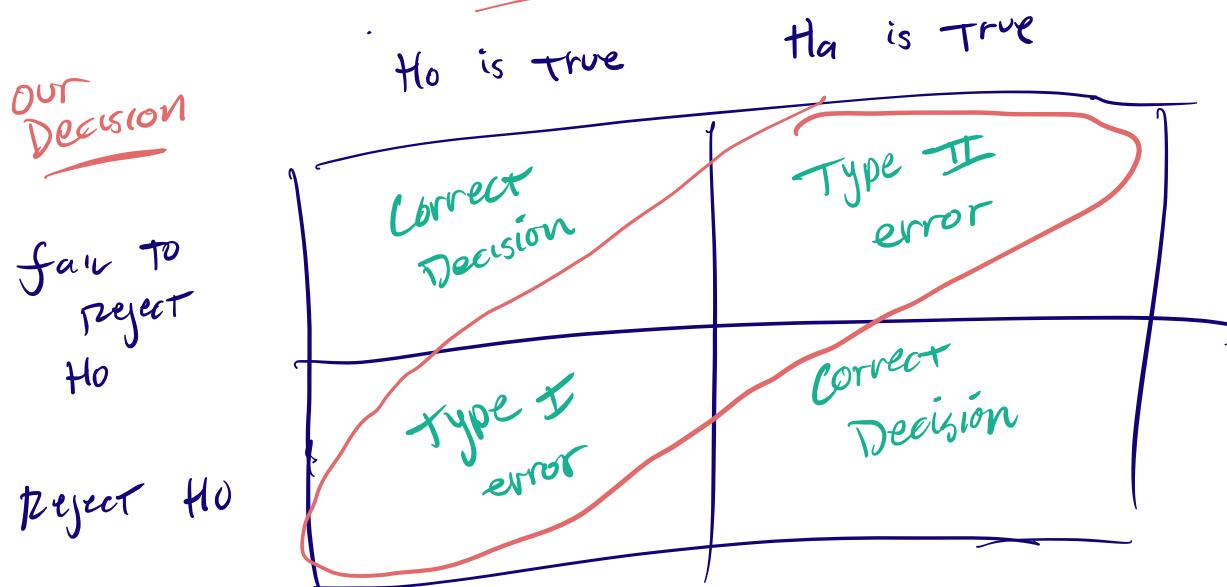
$$P(Z \geq z^*) = p\text{-value}$$

## Conclusion:

- If The p-value is less than our initial significance level , we have sufficient evidence in favor of the Alternative hypothesis and we can Reject our null hypothesis .
- Otherwise , we fail to Reject our null hypothesis if our p-value is greater than our significance level
- Lastly , state conclusions in the context of the problem .

Based on our conclusion above, we could have made the correct decision or made an error.

The Truth



- Type I error and Type II errors have an inverse relationship. If we minimize type I error we are increasing Type II error.
- Depending on the context, we have to choose which type of error we would like to minimize!