

Inferential Statistics II

Jafar Namdar

Postdoctoral Associate

MIT SCM & MIT Digital Supply Chain Transformation

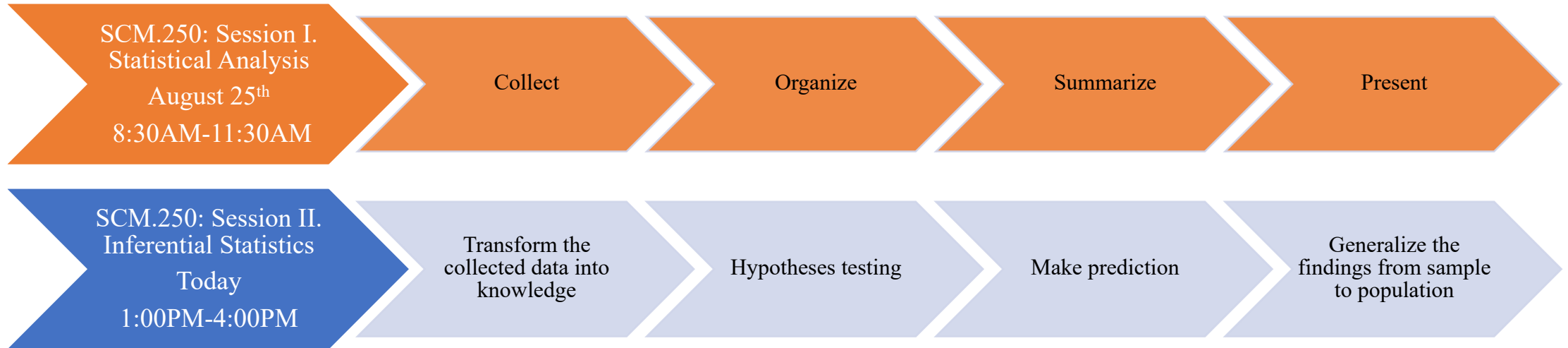


Supply Chain
MANAGEMENT

Review

- Assignments are graded.

Two Sessions

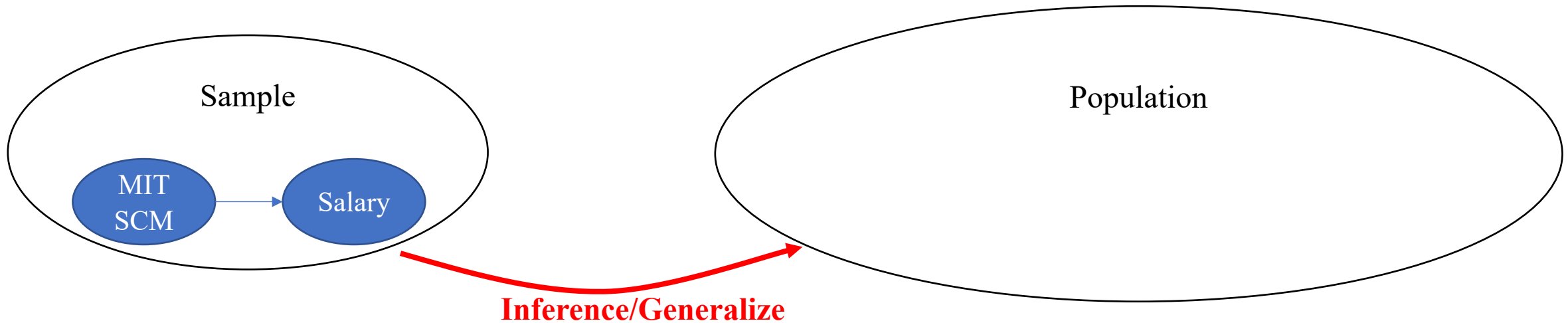


Outline

- Population vs Sample
- Sampling Distribution
- Hypothesis Testing
- Multiple Linear Regression
- ANOVA (if time permits)

Inferential Statistics

- Inferential statistics consists of a set of procedures that enables us **to draw conclusions** about the characteristics of a whole population by studying the properties of a sample of a population.



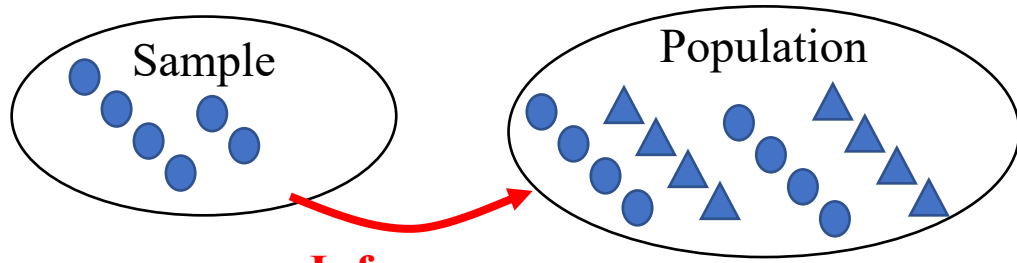
Theory of Estimation (Estimate Procedure)

1. Estimation of some characteristics of the entire population
2. Point estimate
3. Confidence Intervals

Hypothesis Testing

1. Testing the samples of the entire population meet the average.
2. Null hypothesis
3. Alternative hypothesis

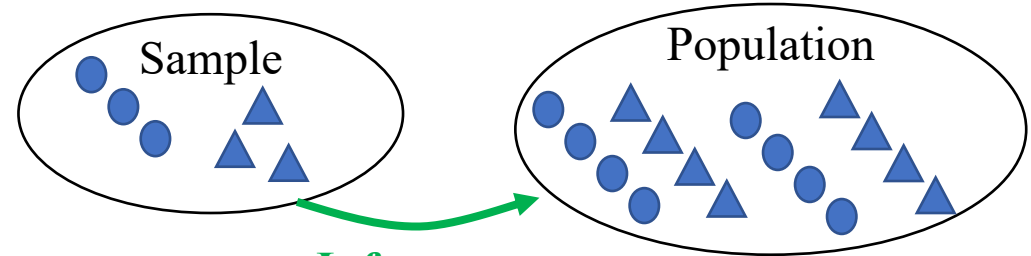
Population vs. Sample Size



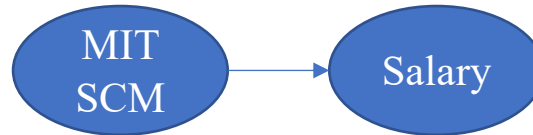
Inference



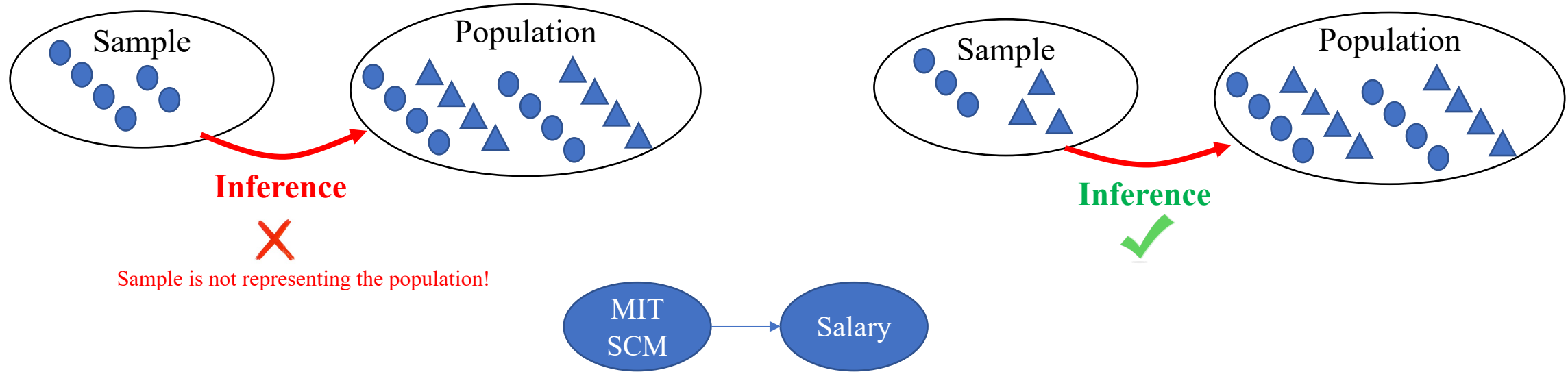
Sample not representing the population!



Inference



Population vs. Sample Size



- ✓ A population is the set or collection of all items which are analyzed in a specified purpose.
- ✓ A sample is a subset of population.
- ✓ A reasoning/**inference** from a sample to a population.
- ✓ Numbers that characterize a population are called population **parameters** (e.g., $\mu, \sigma, \beta, Y, \epsilon$, etc.). We will never know the exact or true values of the population parameters, we estimate them by using the sample.
- ✓ Numbers that characterize a sample of a population is called **statistics** (e.g., commonly used symbols $\bar{X}, S, \hat{\beta}, \hat{Y}, e$, etc.).

Why do we use Sample instead of Population?

There are various reasons for taking a sample rather than the entire population:

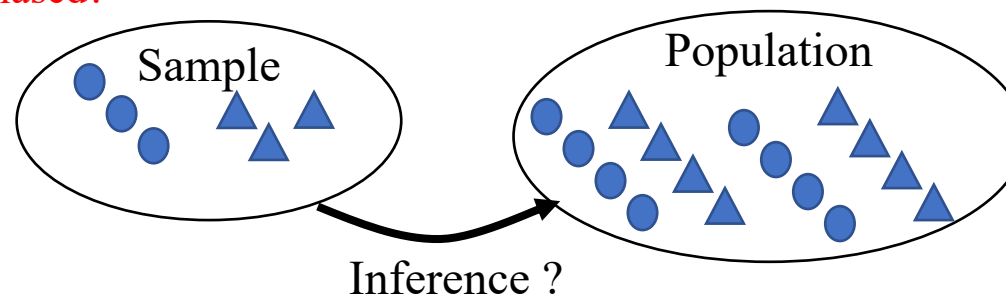
- 1) Expense
- 2) Response time (speed)
- 3) Large population
- 4) Manageability

Population vs. Sample Size

- A sample may contain a number of **atypical** elements - which will ultimately produce inaccurate estimates of the population parameters.
- **Sampling bias** – a sample is biased if it is obtained by a method that favors the selection of individuals having particular characteristics.

- Example

- Consider a situation where I have access to the MIT-SCM alumni file. In that file, 1,000 applicants voluntarily disclosed their starting salaries. The ratio of males to females is approximately one (see the following figure). **Does my sample represent the population or it is biased?**



Sampling Distribution

- To mitigate against bias, samples should be collected randomly selected using a probability distribution.
- Random variables are assumed to **independent**, and every variable has the same probability distribution:

$$X_1, X_2, \dots, X_n$$

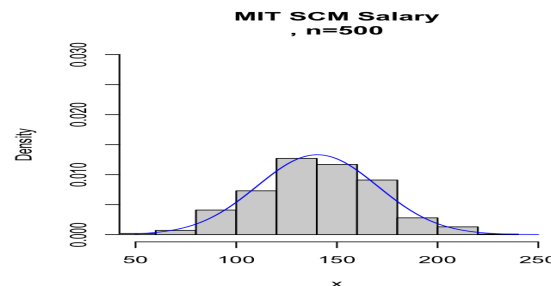
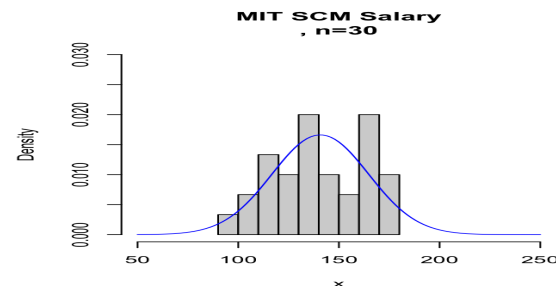
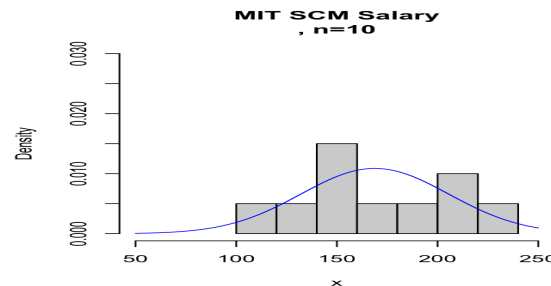
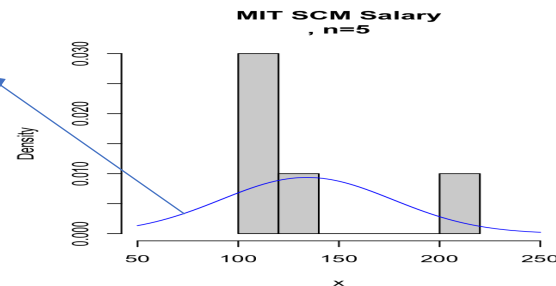
- The probability distribution of a statistic is called **sampling distribution**.

Population Parameters vs. Sample Statistics

- In addition to sample selection method, sampling distribution relies on the **sample size** and the population distribution.

Example	Population Parameters	Sample Statistics Estimated by
Mean	True mean, μ	$\bar{X} = \frac{\sum_i X_i}{n}$
Standard Deviation	True Standard Deviation, σ	$S = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n - 1}}$

Blue Graph:
Population
Distribution,
God Knows this
but we will
never!



R Code

```
set.seed(7)
par(mfrow=c(2,2))
x <- rnorm(500, mean = 140, sd = 30) # your df$col
mean_x <- mean(x)
sd_x <- sd(x)
hist(x, freq = FALSE, xlim = c(50,250), ylim = c(0,.03), main = 'MIT SCM Salary\n, n=500')
curve(dnorm(x, mean = mean_x, sd = sd_x), add = TRUE, col = "blue")
```

Central Limit Theorem

- Definition: If X_1, X_2, \dots, X_n are random samples drawn from a population with overall mean μ and finite variance σ^2 , and if \bar{X}_n is the sample mean of the first n samples, then the limiting form of the distribution, $Z = \lim_{n \rightarrow \infty} \left(\frac{\bar{X}_n - \mu}{\sigma_{\bar{X}}} \right)$, with $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$, is a standard normal distribution.*

In a simple language, the average of randomly selected sample follows normal distribution with the average of \bar{X}_n and standard deviation of $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$

$$\bar{X}_n \sim \text{Normal} \left(\bar{X}_n, \frac{\sigma}{\sqrt{n}} \right)$$

Or

$$Z = \frac{\bar{X}_n - \mu}{\sigma_{\bar{X}}} \sim \text{Normal (0,1) aka Standard Normal}$$

*Montgomery, Douglas C.; Runger, George C. (2014). *Applied Statistics and Probability for Engineers* (6th ed.). Wiley. p. 241. [ISBN 9781118539712](https://doi.org/10.1002/9781118539712).

Hypothesis Testing

- A statistical hypothesis is a statement about the parameters of one or more populations.
- **Null hypothesis (H_0)**: A claim about a population characteristics that is initially assumed to be true.
 - It must always includes equality signs ($=, \geq, \text{or } \leq$)
- **Alternative hypothesis (H_1)**: The competing claim which we favor to prove.
 - Basically anything against H_0 ; typical signs ($\neq, >, <$)
- There are two possible outcomes of hypothesis testing:
 - 1) **Reject H_0** in favor of H_1
 - 2) **Fail to reject** the null hypothesis H_0

The choice of language matters here! Thus, **never** say we **accept** H_0 or H_1 !

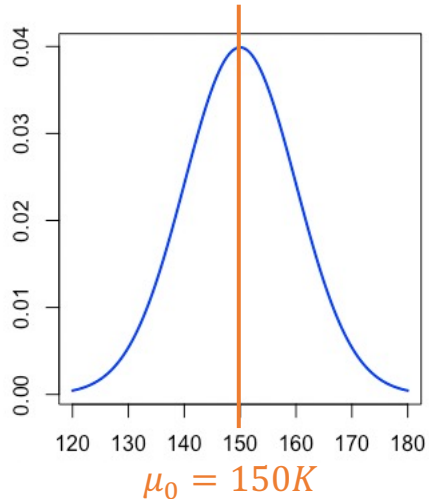
- Let's frame some hypotheses? Give me an example.

Hypothesis Testing (Type I)

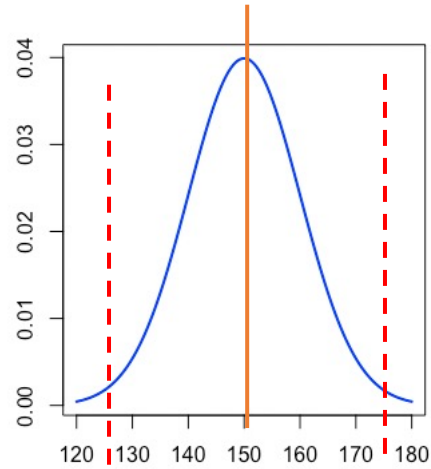
- Let's say person A says (or hypothesizes) the average salary of MIT-SCM is \$150K (This means that person is talking about the average of population or $\mu = \$150,000$). This is H_0 .
 - $H_0: \mu_0 = 150k$
- Alternatively, his/her argumentative friend argues that statement is not true. This is H_1 .
 - $H_1: \mu_0 \neq 150k$
- How to resolve this conflict?
 - Statistics gives us tools to construct the Confidence Interval (CI) when H_0 is true.

Hypothesis Testing (Type I)

- Let's say person A says (or hypothesizes) the average salary of MIT-SCM is \$150K (This means that person is talking about the average of population or $\mu = \$150,000$). **This is H_0 .**
 - $H_0: \mu_0 = 150k$
- Alternatively, his/her argumentative friend argues that statement is not true. **This is H_1 .**
 - $H_1: \mu_0 \neq 150k$
- How to resolve this conflict?
 - Statistics gives us tools to construct the **Confidence Interval (CI)** when H_0 is true.



This is what person A claims.



$(1 - \alpha)\%$ CI
Acceptance Region

Definition of $(1 - \alpha)\%$ CI: If you draw samples 100 times,
 $(1 - \alpha)\%$ times the sample mean or \bar{X} falls into this interval

This statement means, we allow error for %5 of times!

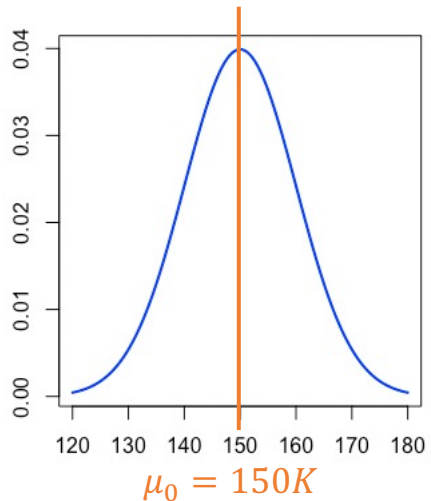
This is called α or Type I error

Typical Choice of $\alpha = 10\%, 5\%, 1\%$

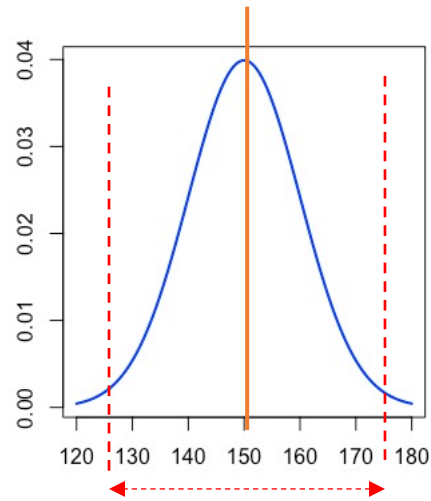
- If $\alpha = 10\%$, then $(1 - \alpha)\%$ means 90% of times our sample falls into the acceptance region (**Very Narrow Acceptance Region or CI**).
- If $\alpha = 5\%$, then $(1 - \alpha)\%$ means 95% of times our sample falls into the acceptance region. (**Commonly Used**).
- If $\alpha = 1\%$, then $(1 - \alpha)\%$ means 99% of times our sample falls into the acceptance region (**Very Wide Acceptance Region or CI**).

Hypothesis Testing (Type I)

- Let's say person A says (or hypothesizes) the average salary of MIT-SCM is \$150K (This means that person is talking about the average of population or $\mu = \$150,000$). **This is H_0 .**
 - $H_0: \mu_0 = 150k$
- Alternatively, his/her argumentative friend argues that statement is not true. **This is H_1 .**
 - $H_1: \mu_0 \neq 150k$
- How to resolve this conflict?
 - Statistics gives us tools to construct the confidence interval when H_0 is true.



This is what person A claims.



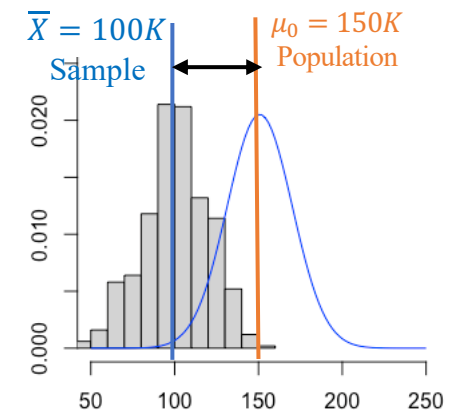
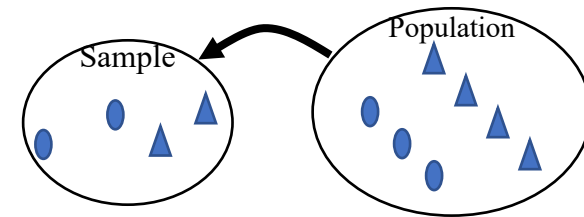
Definition of $(1 - \alpha)\% \text{ CI}$: If you draw samples 100 times, $(1 - \alpha)\%$ times the sample mean or \bar{X} falls into this interval

This statement means, we allow error for %5 of times!

This is called α or Type I error

Data collection
& Test

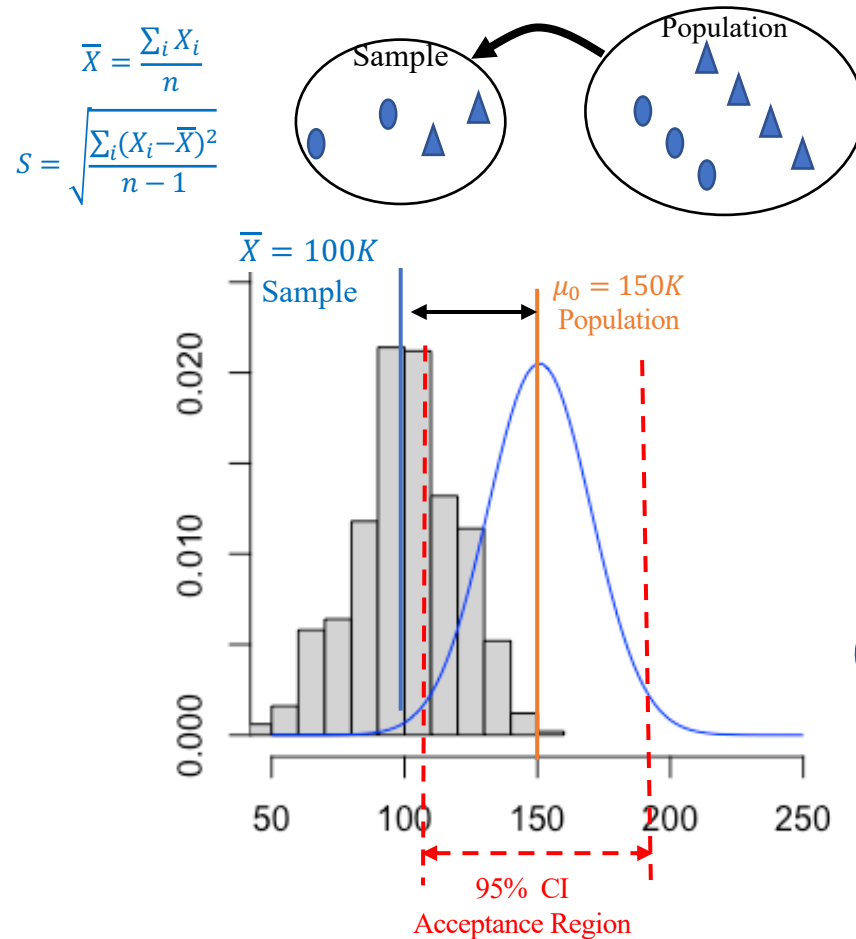
$$\bar{X} = \frac{\sum_i X_i}{n}$$
$$S = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n - 1}}$$



So, what do you say?
Do we Reject H_0 or we fail to reject H_0 ?

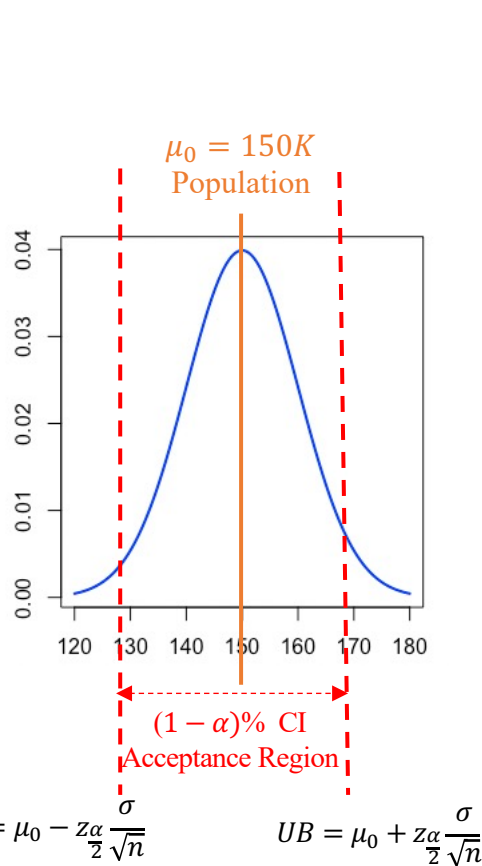
Hypothesis Testing

- To answer the Question “Do we Reject H_0 or we fail to reject H_0 ?”
 - There are two (kinda the same) methods to answer this question.
 - 1. If \bar{X} falls into the confidence interval, then we have no evidence to reject (fail to reject), otherwise we reject.

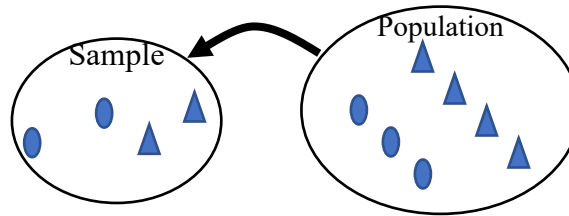


How do we calculate CI?

- To answer the Question “*Do we Reject H_0 or we fail to reject H_0 ?*”
 - There are two (kinda the same) methods to answer this question.
 1. If \bar{X} falls into the confidence interval, then we have no evidence to reject (fail to reject), otherwise we reject.



$$\bar{X} = \frac{\sum_i X_i}{n}$$
$$S = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n - 1}}$$



Ideally

← If we know σ , then we use z (normal distribution) to calculate CI.

This hardly happens in real world.

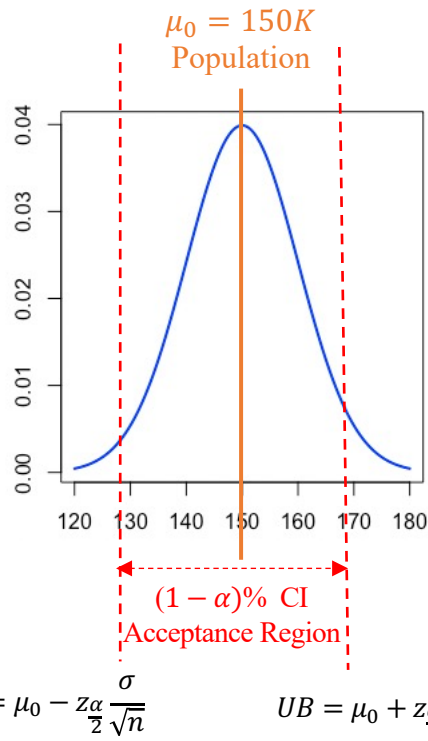
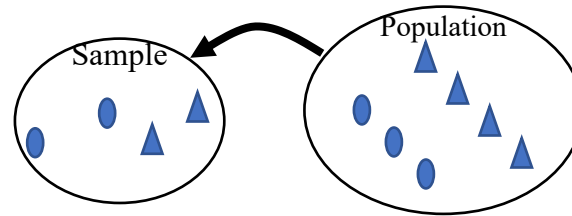
How do we calculate CI?

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

$$\bar{X} = \frac{\sum_i X_i}{n}$$

$$S = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n-1}}$$



Ideally

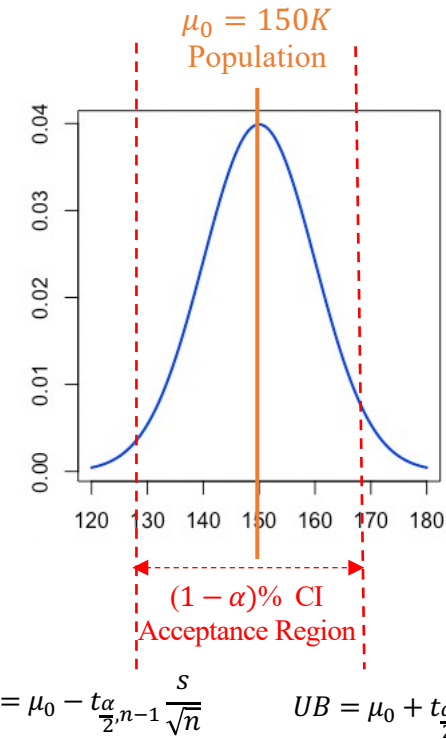
If we know σ , then we use z (normal distribution) to calculate CI.

This hardly happens in real world.

Alternatively
(In Practice)

If we end up estimating σ by using s , then we use t distribution.

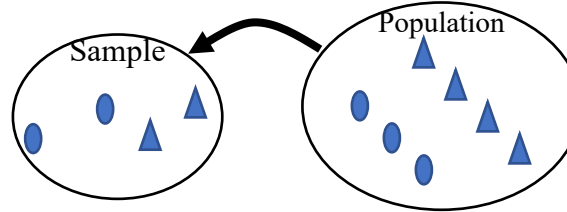
t distribution has an element of degree of freedom which is usually equals to $n - 1$



How do we calculate CI?

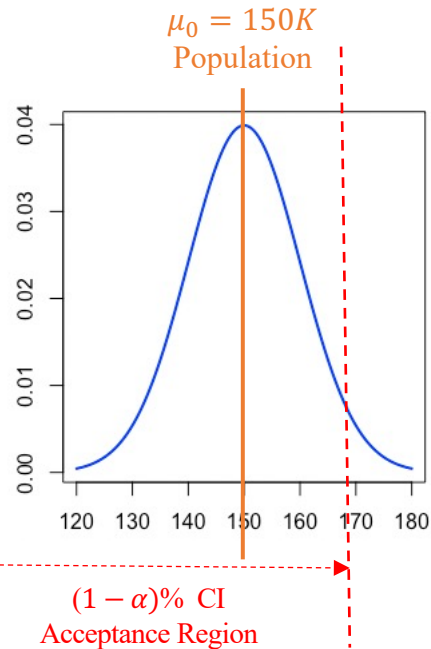
$$\bar{X} = \frac{\sum_i X_i}{n}$$

$$s = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n-1}}$$



$$H_0: \mu = \mu_0$$

$$H_1: \mu > \mu_0$$

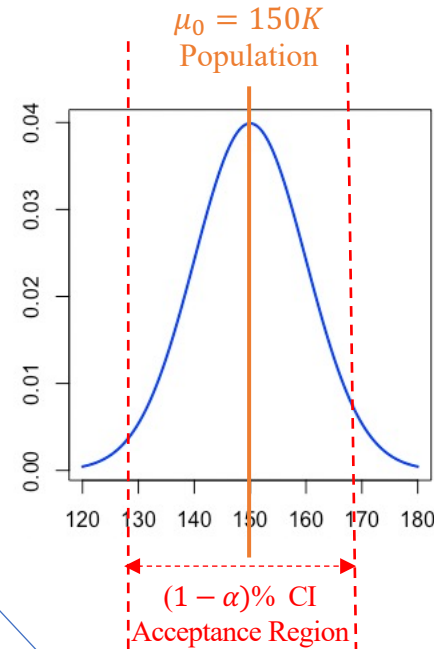


$$UB = \mu_0 + t_{\alpha, n-1} \frac{s}{\sqrt{n}}$$

One-tailed tests

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

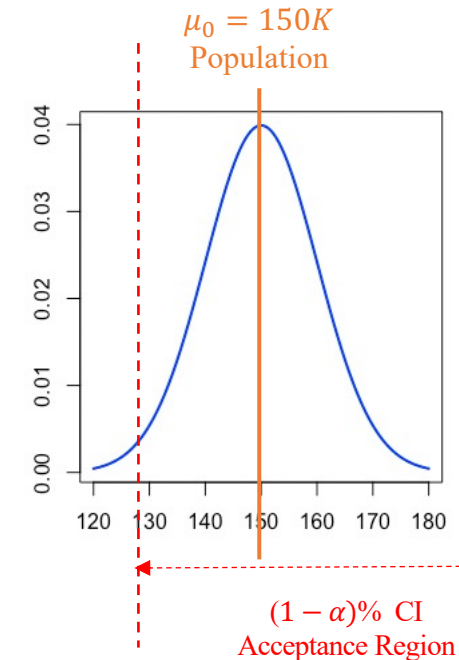


$$LB = \mu_0 - t_{\frac{\alpha}{2}, n-1} \frac{s}{\sqrt{n}} \quad UB = \mu_0 + t_{\frac{\alpha}{2}, n-1} \frac{s}{\sqrt{n}}$$

Two-tailed tests

$$H_0: \mu = \mu_0$$

$$H_1: \mu < \mu_0$$



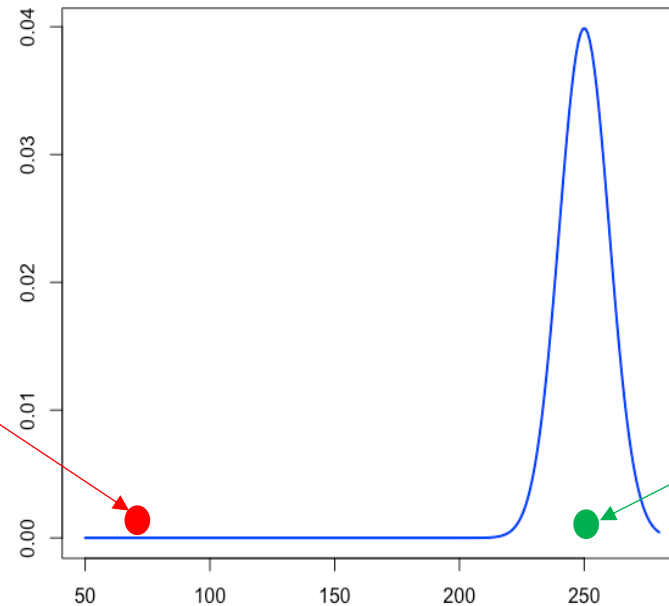
$$LB = \mu_0 - t_{\alpha, n-1} \frac{s}{\sqrt{n}}$$

One-tailed tests

Hypothesis Testing (P-value)

- To answer the Question “*Do we Reject H_0 or we fail to reject H_0 ?*”
 - There are two (kinda the same) methods to answer this question.
 1. If \bar{X} falls into the confidence interval, then we have no evidence to reject (fail to reject), otherwise we reject.
 2. If P-value is larger than α (or $\frac{\alpha}{2}$ in a two-tailed test), then we have no evidence to reject (fail to reject), otherwise we reject.

- Is there any possibility that this observation belong to this distribution? **Yes**
- What is the probability to see observation like this or even more extreme than this?
Extremely Low

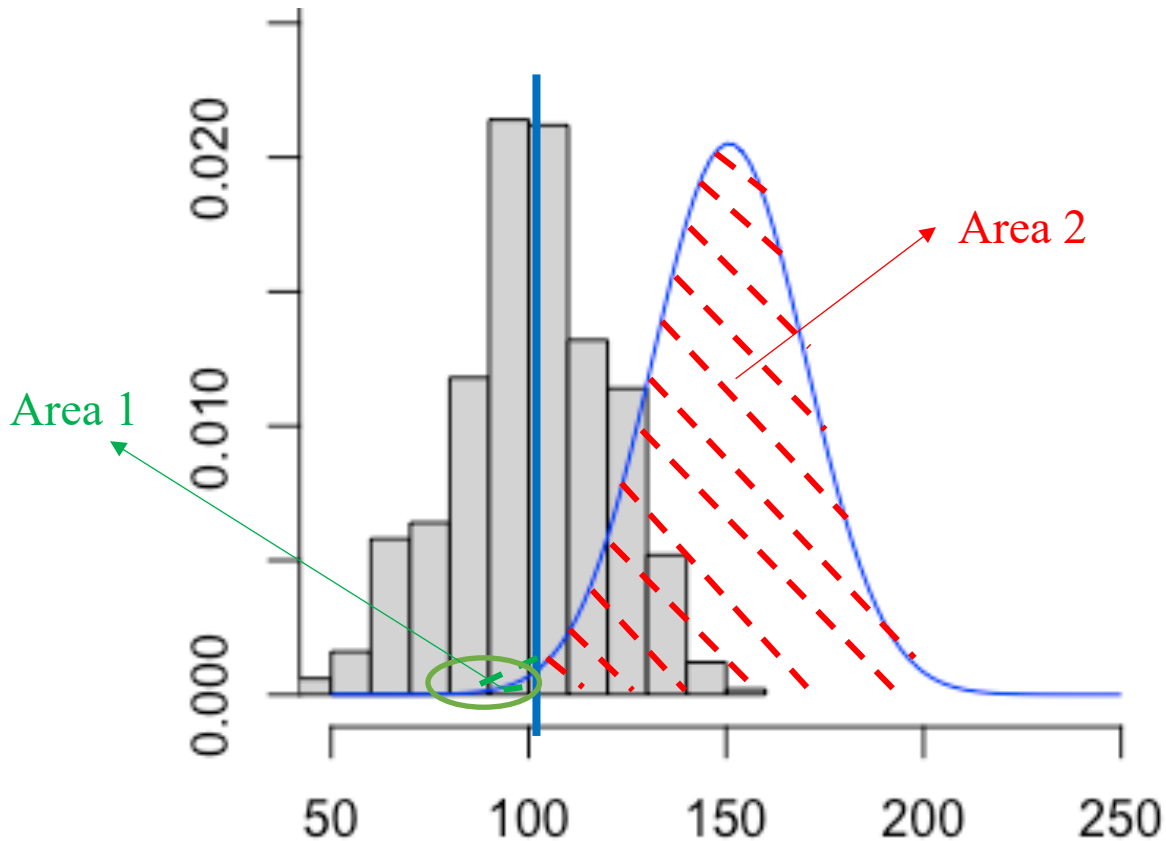


- Is there any possibility that this observation belong to this distribution? **Yes**
- What is the probability to see observation like this or even more extreme than this? **Very High**

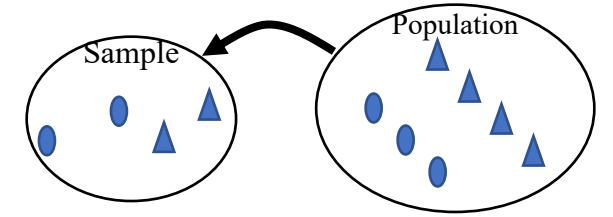
Definition: The p-value is the probability of obtaining test results at least as extreme as the result actually observed, under the assumption that the null hypothesis is correct.

Hypothesis Testing (P-value)

- To answer the Question “Do we Reject H_0 or we fail to reject H_0 ?”
 - There are two (kinda the same) methods to answer this question.
 1. If \bar{X} falls into the confidence interval, then we have no evidence to reject (fail to reject), otherwise we reject.
 2. If P-value is larger than α (or $\frac{\alpha}{2}$ in a two-tailed test), then we have no evidence to reject (fail to reject), otherwise we reject.



$$\bar{X} = \frac{\sum_i X_i}{n}$$
$$s = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n - 1}}$$



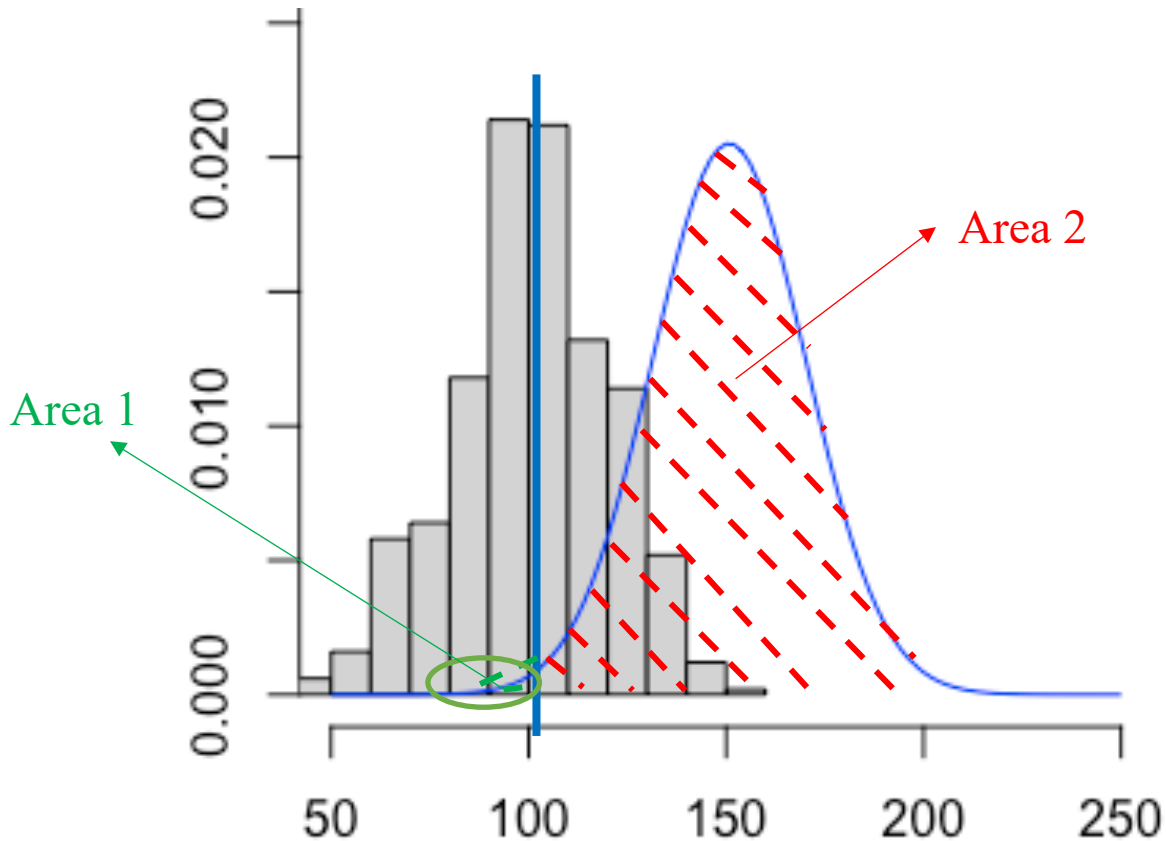
$$\text{P-value} = \min(\text{Area 1}, \text{Area 2}) = \text{Area 1}$$

The p-value is the probability of obtaining test results at least as extreme as the result actually observed, under the assumption that the null hypothesis is correct.

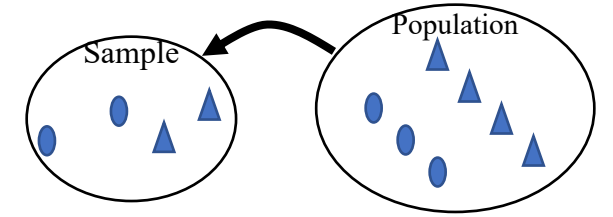
In a simple language, p-value means to what extent the observed data belong to the population.

Hypothesis Testing (P-value)

- To answer the Question “Do we Reject H_0 or we fail to reject H_0 ?”
 - There are two (kinda the same) methods to answer this question.
 - If \bar{X} falls into the confidence interval, then we have no evidence to reject (fail to reject), otherwise we reject.
 - If P-value is larger than α (or $\frac{\alpha}{2}$ in a two-tailed test), then we have no evidence to reject (fail to reject), otherwise we reject.



$$\bar{X} = \frac{\sum_i X_i}{n}$$
$$s = \sqrt{\frac{\sum_i (X_i - \bar{X})^2}{n-1}}$$



One-tailed

P-value = min(Area 1, Area 2) = Area 1

$P_value < \alpha$ Reject

$P_value > \alpha$ Fail to reject

Two-tailed

P-value = min(Area 1, Area 2) = Area 1

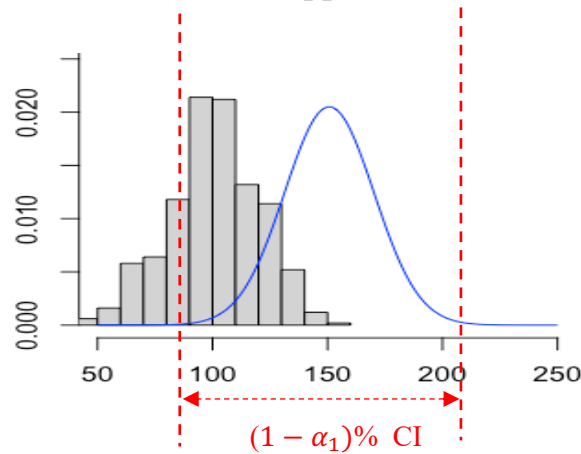
$P_value < \alpha/2$ Reject

$P_value > \alpha/2$ Fail to reject

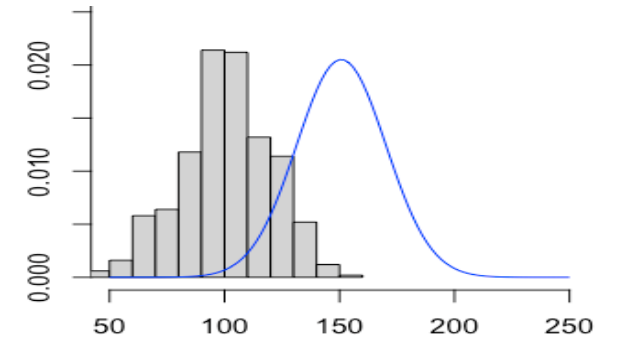
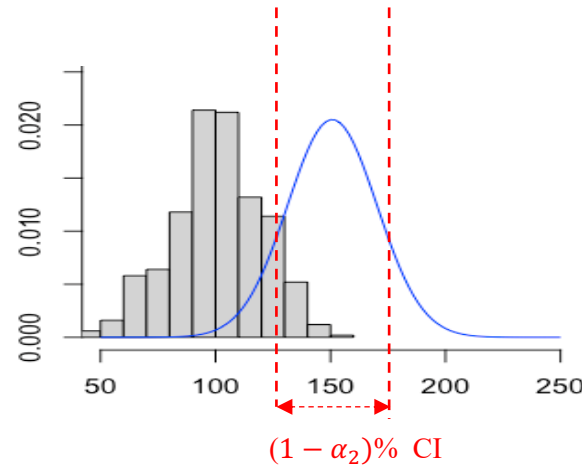
Hypothesis Testing (Type I & II errors)

1. Type I error or α

- $\alpha = P(\text{Type I Error}) = P(\text{Reject } H_0 \text{ when } H_0 \text{ is True})$
- What would happen if I increase α ?



$$\alpha_1 < \alpha_2$$

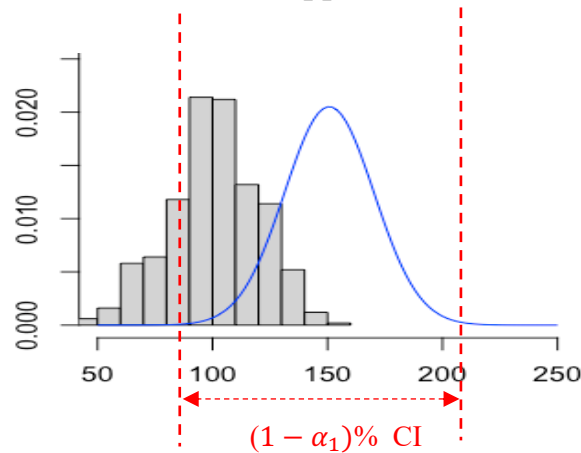


Maybe this sample was very odd!
Because we allow error for $\alpha\%$ of times.

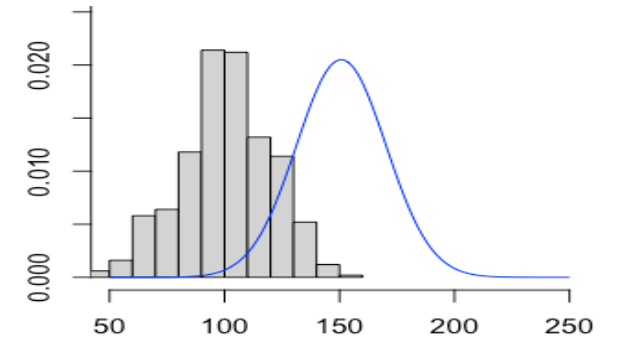
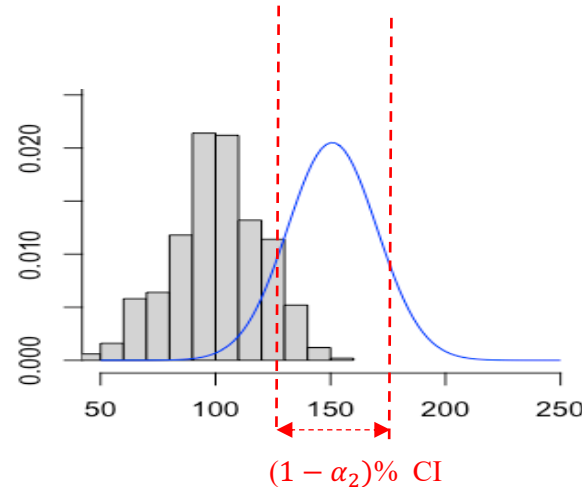
Hypothesis Testing (Type I & II errors)

1. Type I error or α

- $\alpha = P(\text{Type I Error}) = P(\text{Reject } H_0 \text{ when } H_0 \text{ is True})$
- What would happen if I increase α ?



$$\alpha_1 < \alpha_2$$



Maybe this sample was very odd!
Because we allow error for $\alpha\%$ of times.

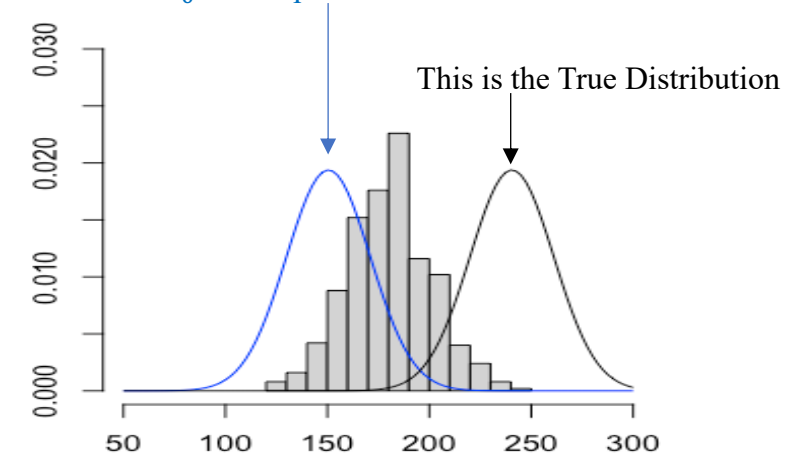
2. Type II error

- $\beta = P(\text{Type II Error}) = P(\text{Fail to reject } H_0 \text{ when } H_0 \text{ is False})$

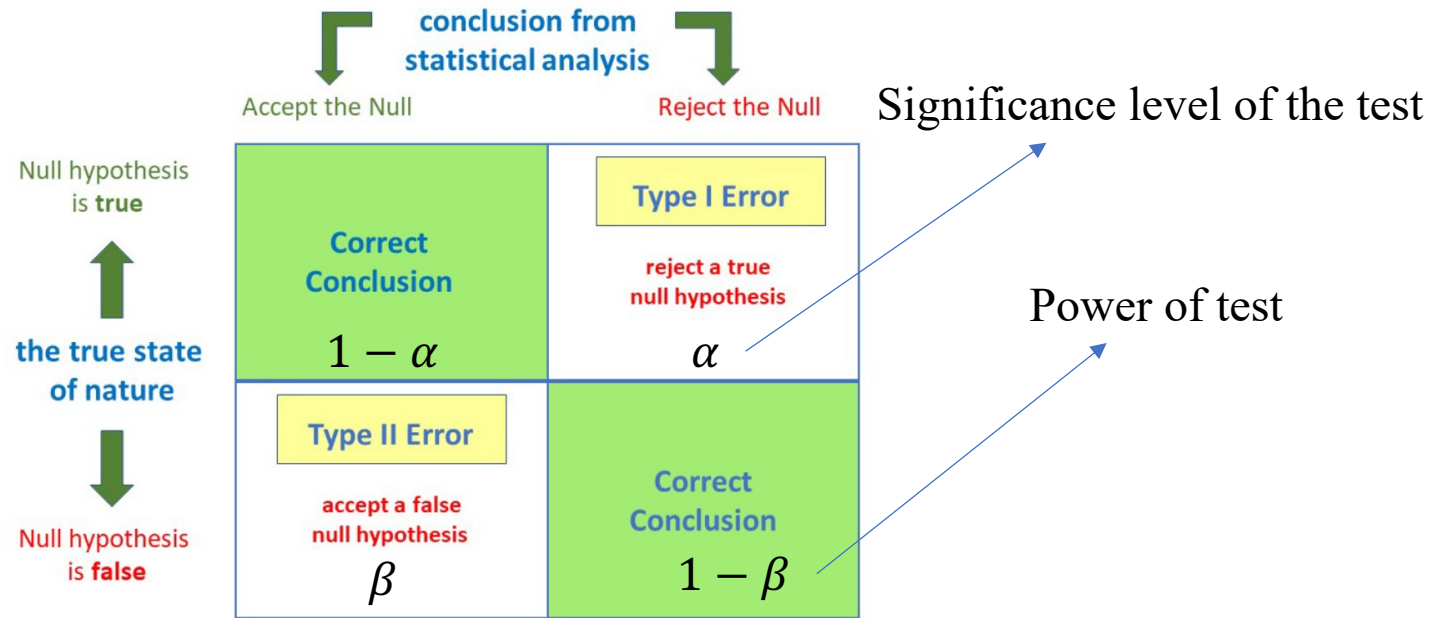
Clearly, H_0 or the claim of person A is wrong, but we collect data and unfortunately we fail to reject it since our sample overlaps a lot with the blue graph.

Thus, we end up not rejecting a claim which is wrong, aka Type II error.

This is the H_0 or the person A claims.

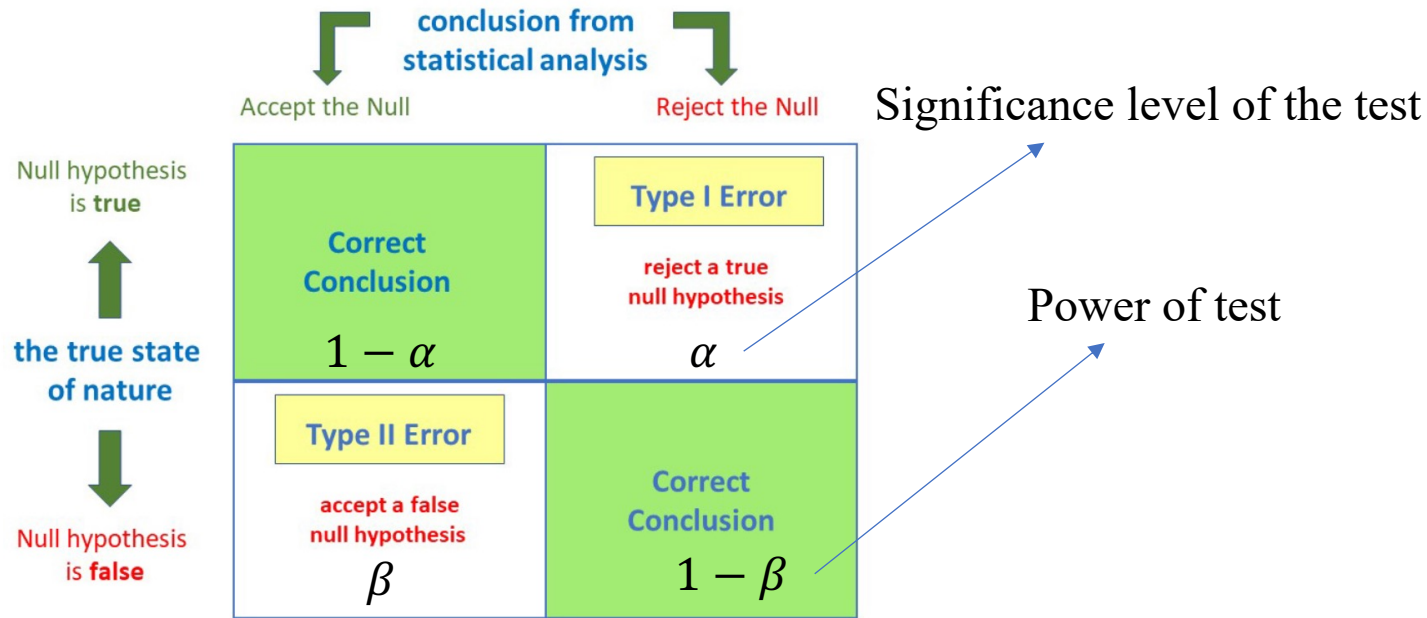


Hypothesis Testing (Type I & II errors)



https://www.simplypsychology.org/type_I_and_type_II_errors.html

Hypothesis Testing (Type I & II errors)



Sometimes we select the minimum sample size by these two measures.

1. First we set α at the maximum allowable Type I error rate (typically 1%, 5% or 10%).
2. Next, we calculate Type II error or $(1 - \beta)$.
$$1 - \beta = f(n)$$
3. Sometime $f(n)$ is too complicated, so we use simulation to find out the minimum sample size.

https://www.simplypsychology.org/type_I_and_type_II_errors.html

Example (Type I & II Errors)

- Suppose there is a test for a particular disease.
 - If the disease really exists and is diagnosed early, it can be successfully treated
 - If it is not diagnosed and treated, the person will become severely disabled
 - If a person is erroneously diagnosed as having the disease and treated, no physical damage is done.
 - First, clearly state your null hypothesis, then,
 - What is Type I error?
 - What is Type II error?
 - Which error do you want to minimize here?
- * Depending the way you frame your hypotheses, Type I and II might change.

Example (Type I & II Errors)

Suppose you are a manager of a company.

- One of your consultant tells you that you have to construct new production facility because the current production can handle up to 200% increase in demand. The consultant hypothesize that the demand is going increase 300%. The new facility cost 10 times more than annual lost sales & reputation cost.

What is null hypothesis?

Type I error?

Type II error?

Which error do you want to minimize the most?

- General Procedure in Hypothesis Tests:
 - 1) Identify parameters of interest
 - 2) State null and alternative hypotheses
 - 3) Select a significance level (α) and the critical values
 - 4) Determine the appropriate test statistics
 - 5) Calculate the test statistics
 - 6) Calculate rejection criteria and compare p-value
 - 7) Make statistical decision and interpret results

Simple Linear Regression

- A regression model is used to model and explore relationships between variables that are related in a nondeterministic manner.
- Simple Linear Regression: Only one independent variable x (regressor or predictor) and one dependent variable Y (response).

$$E(Y|x) = \mu_{Y|x} = \beta_0 + \beta_1 x$$

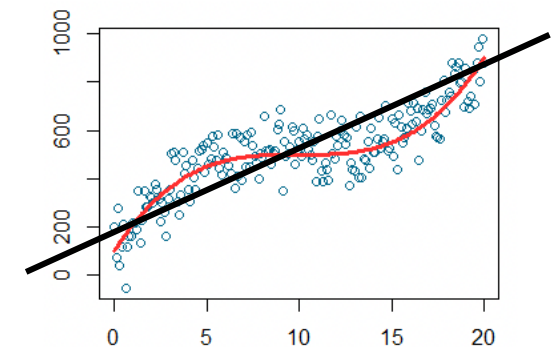
- The mean of Y is a linear function of x . However, the actual observed value y does not have exact linear relationship.
- **The fitted or estimated regression line:**

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$y = \beta + \beta_1 x + \epsilon$$

Estimated

True relationship



Significance of Regression

OLS Regression Results

```

=====
Dep. Variable:          MPG      R-squared:                0.201
Model:                  OLS      Adj. R-squared:             0.157
Method:                 Least Squares      F-statistic:           4.532
Date:                   Mon, 29 Aug 2022    Prob (F-statistic):      0.0473
Time:                   18:24:56    Log-Likelihood:         -53.272
No. Observations:       20      AIC:                   110.5
Df Residuals:           18      BIC:                   112.5
Df Model:                1
Covariance Type:        nonrobust
=====
  
```

What percentage of variation is explained by the model.

H_0 : all Coefficients are zero
 $\beta_0 = \beta_1 = 0$

H_1 : Otherwise

95% CI for each β

P-value

β_0
 β_1

	coef	std err	t	P> t	[0.025	0.975]
const	33.5348	2.614	12.829	0.000	28.043	39.027
EngDis	-0.0354	0.017	-2.129	0.047	-0.070	-0.000

$H_0: \beta_0 = 0$
 $H_1: \beta_0 \neq 0$

$H_0: \beta_1 = 0$
 $H_1: \beta_1 \neq 0$

$$y = \beta_0 + \beta_1 x = 33.53 - 0.035x$$

What would be regression if the p-value of β_1 is equal to 0.2?