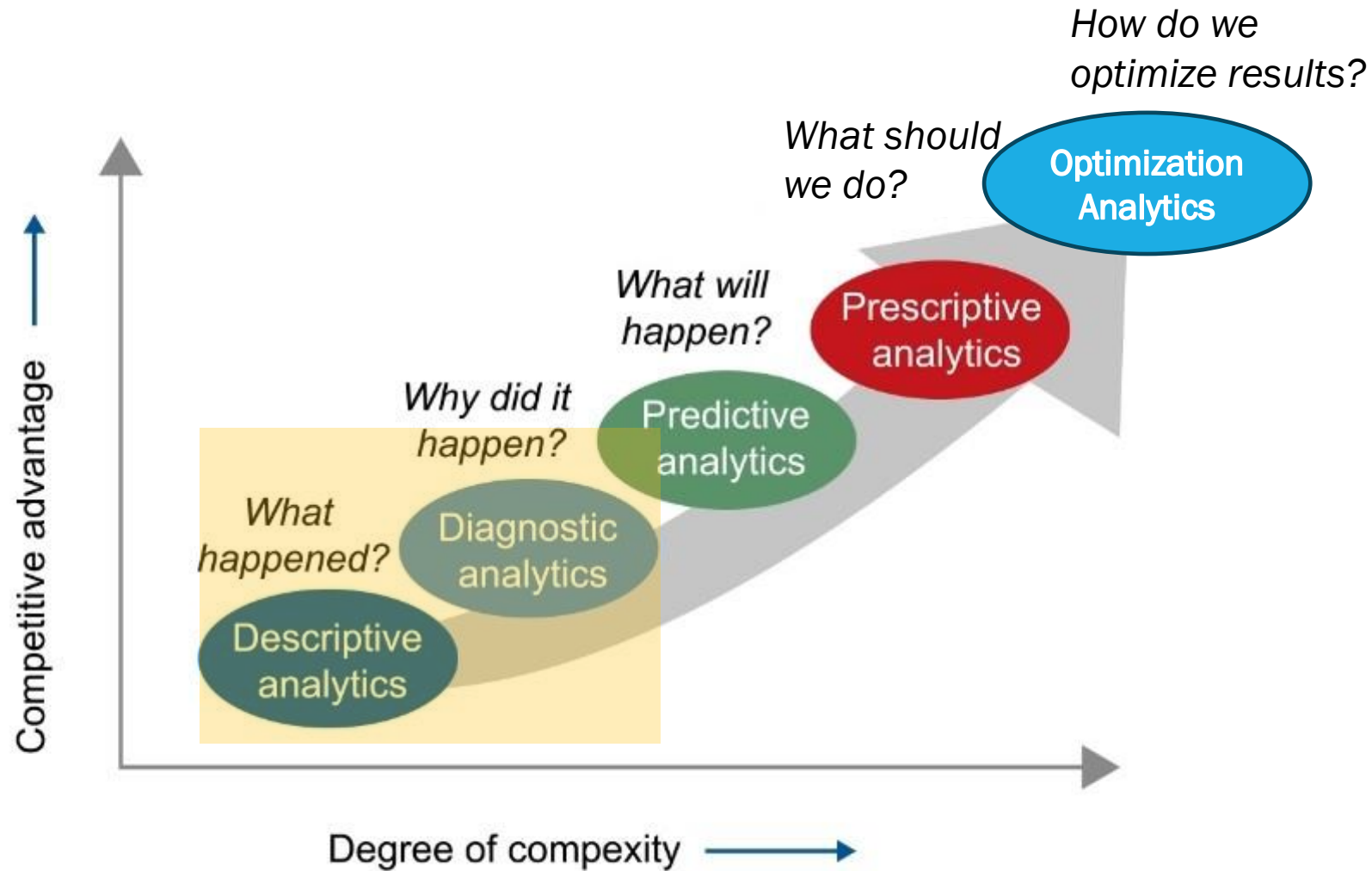


A professional team of three individuals—two women and one man—are working in a modern industrial facility, likely a food processing plant. They are dressed in business-casual attire: light blue shirts and grey blouses for the women, and a grey suit with a blue tie for the man. They are positioned along a conveyor belt that is filled with a variety of colorful, small candies or confections. The woman on the left, wearing glasses, is holding a blue clipboard and looking down at it. The woman in the middle is also holding a blue clipboard and a pen, looking at the candies. The man on the right is holding a blue folder or book and looking at the candies. In the background, other workers in similar attire are visible, working at different stations. The facility has large windows and a clean, organized appearance. The text "DATA SAMPLING" is overlaid in the center of the image.

# DATA SAMPLING

# HIERARCHY OF DATA ANALYSIS TYPES





# POPULATION VS SAMPLE

- Ideally, analyze all data for insights.
- Examples:
  - A mobile company wants to assess all potential customers.
  - A government must consider all citizens' needs for a new service.
- Full data collection is often impractical due to:
  - High costs of data acquisition
  - Time constraints
  - Computational and storage limitations
  - Increased complexity in processing large datasets
- Solution: Select a representative, make inference



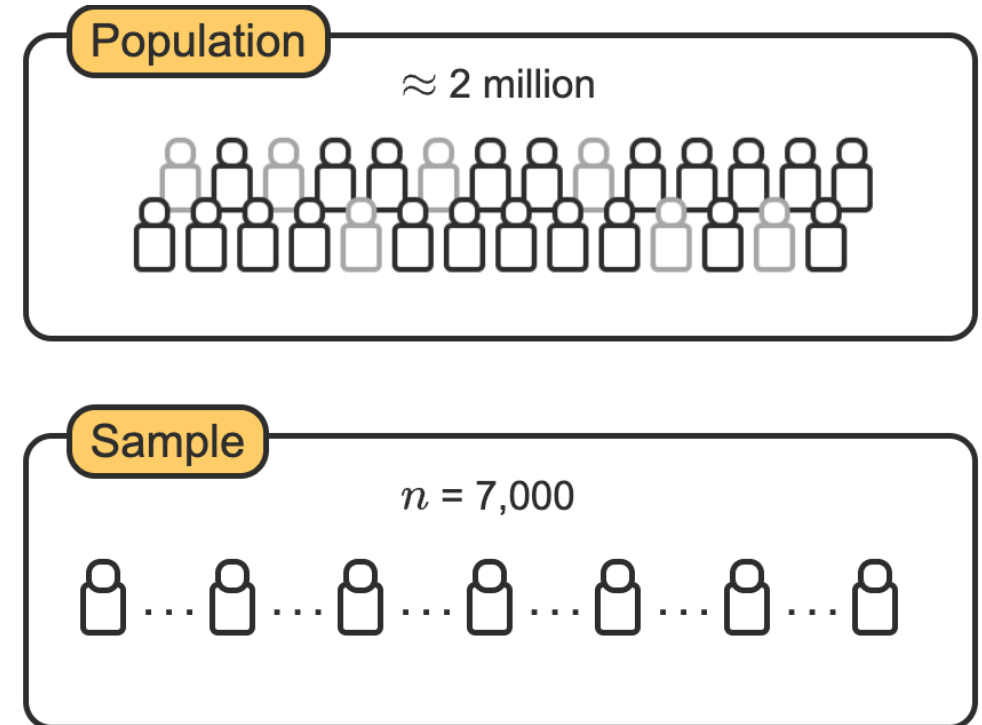
# DATA SAMPLING

A **Sampling Method** is a process to select a subset of observations from the entire population. Ideally, the observations in the sample are representative of the population. Common methods include:

- **Random Sampling:** Each subset of  $n$  units is equally likely to be chosen.
- **Stratified Sampling:** The population is divided into meaningful groups (strata), and samples are drawn from each.
- **Cluster Sampling:** The population is divided into clusters (unrelated to key study features), and some clusters are randomly selected.
- **Systematic Sampling:** Every  $k$ th observation is selected from a random starting point, where  $k \approx (\text{population size}) / n$ .
- **Convenience Sampling:** Easily accessible observations are selected (non-random).

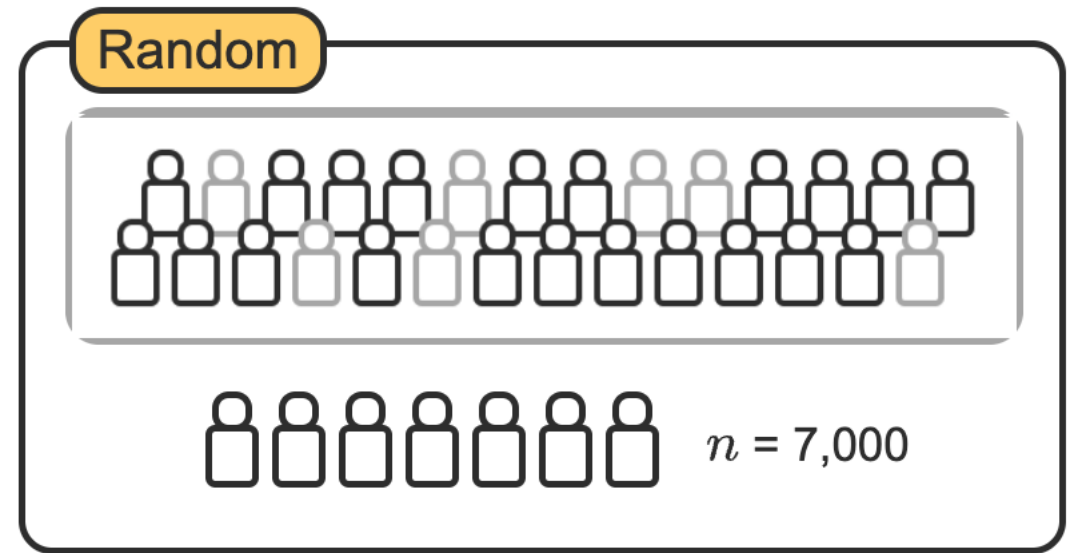
# SAMPLING SCENARIO

- A **population** is the entire set of all individuals, items, or events of interest.
- An **observational unit (aka observation)** is an individual, item, or event of the population where data is recorded.
- A **sample** is a subset of observations from the population used for analysis.
- Example: Transportation satisfaction survey across 5 cities.



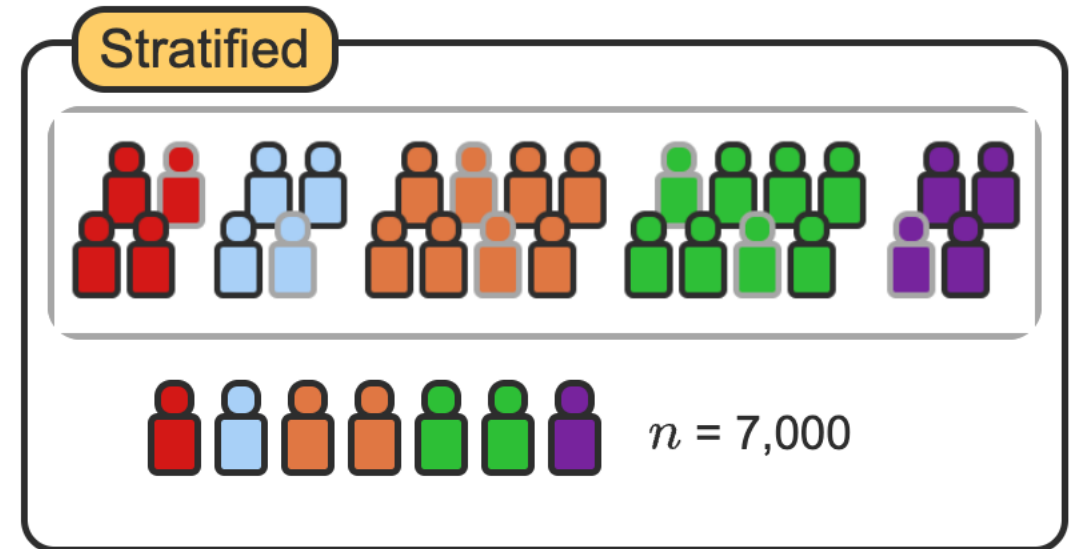
# RANDOM SAMPLING

- In random sampling, passengers are selected at random from a list of all passengers in the five cities.
- Random sampling reduces the potential for sampling bias.
- But this could result in missing important events that occur less frequently.



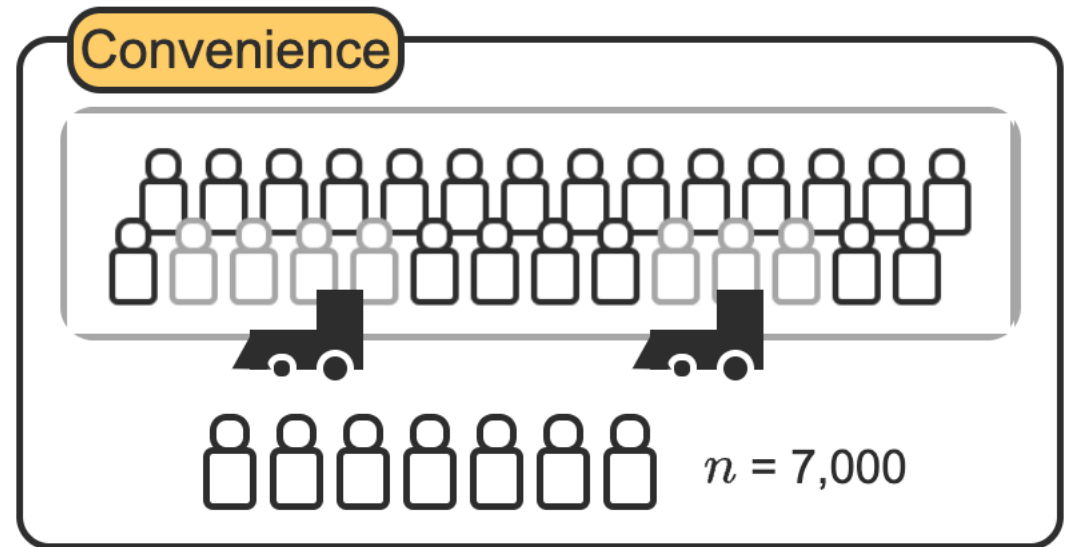
# STRATIFIED SAMPLING

- Passengers are first divided into groups based on city.
- Then from each group, passengers are selected at random.
- Unlike pure random sampling, stratified sampling ensures adequate representation from each city.
- This is especially important when working with data that includes events that are relatively rare (e.g., customer churn, network intrusion, cancer cell detection, etc.)



# CONVENIENCE SAMPLING

- Select passengers waiting in the train stations uses convenience sampling.
- This method is easy and quick, but the sample is not likely representative of all train passengers.





# SYSTEMATIC SAMPLING

- Every 286<sup>th</sup> passenger from a list of all 2 million potential passengers is selected for the sample.
- Population / sample size = selection criteria
- Depending on ordering of the list, this could be close to random, or highly biased.

Select every  
286<sup>th</sup> person

Population

≈ 2 million



Sample

$n = 7,000$



# SAMPLING IN PYTHON

- The pandas method `DataFrame.sample(n=None, frac=None, replace=False, random_state=None)` returns a random sample of items from a dataframe.
- `n` or `frac` specify the number, or fraction, of items to be returned in the sample.
- `replace=` parameter specifies whether sampling is done with (True) or without (False) replacement.
- `random_state=` parameter optionally sets the random number generator seed for reproducible sampling.
- `weights` controls the likelihood of each row (or column) being selected. Weights can be a list/array of values or a column name in the DataFrame. The weights do not need to sum to 1; Pandas normalizes them automatically.
- `axis=` The axis to sample
- `ignore_index=` reset the index to 0,1,2,3...

```
DataFrame.sample(n=None, frac=None, replace=False, weights=None,  
random_state=None, axis=None, ignore_index=False)
```

# PYTHON EXAMPLES, SAMPLING

- <https://colab.research.google.com/github/rhodes-byu/cs180-winter25/blob/main/notebooks/06-sampling.ipynb>

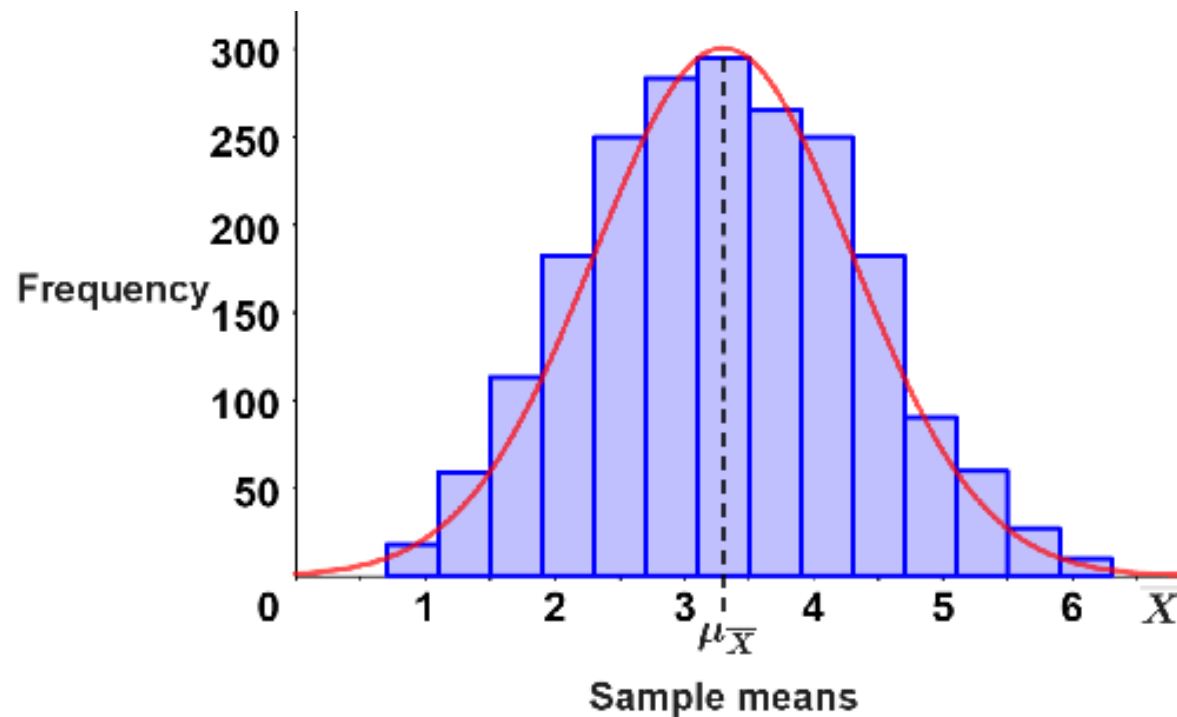
# SAMPLING DISTRIBUTION – NOT ALL SAMPLES EQUAL

A **sampling distribution** is like taking many small samples from a big jar of jellybeans and calculating the average number of red beans in each sample. If you repeat this process over and over, you'll get a bunch of different averages.

Now, imagine plotting those averages on a graph—you'll notice that most of them cluster around the true average of the whole jar. This pattern of sample averages is the **sampling distribution**. It helps us understand how much the results can vary and lets us make better guesses about the whole jar without checking every single jellybean!

# THE CENTRAL LIMIT THEOREM

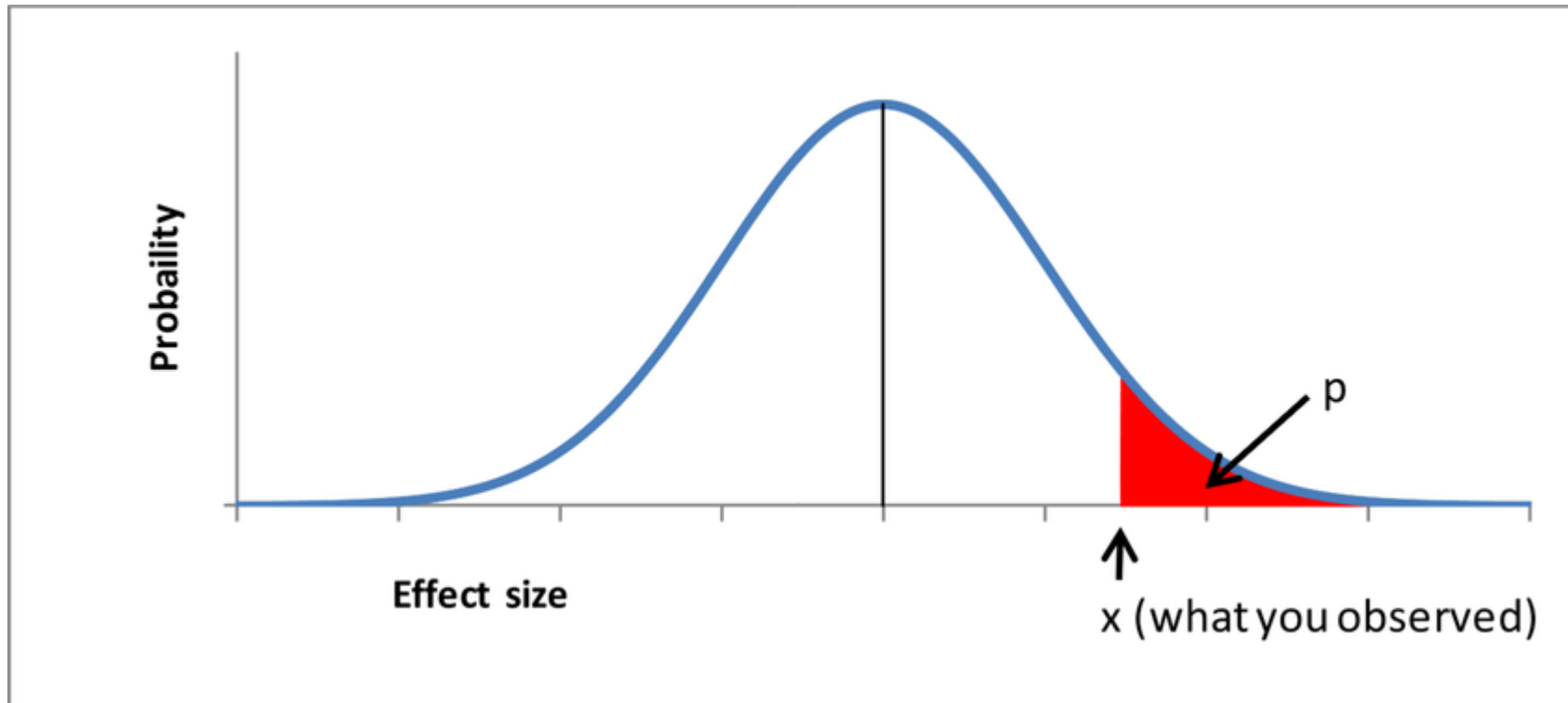
- The sampling distribution of the sample mean approaches a normal distribution as sample size increases, regardless of the population's original distribution (assuming random sampling).
- Works well for large  $n$  ( $>30$ ).





# IS OUR SAMPLE SIGNIFICANTLY DIFFERENT THAN THE POPULATION?

Proportion of red beans...



# ONE-SAMPLE T-TEST

In the case of a **one-sample t-test** (where you are comparing the sample mean against a known population mean), the equation becomes:

$$t = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}}$$

Where:

- $\bar{X}$  is the sample mean,
- $\mu$  is the population mean,
- $S$  is the sample standard deviation,
- $n$  is the sample size.

This formula tests whether the sample mean  $\bar{X}$  significantly differs from the population mean  $\mu$ .

The denominator represents the **standard error** of the mean.

# STATISTICAL SIGNIFICANCE IN PYTHON

- <https://colab.research.google.com/github/rhodes-byu/cs180-winter25/blob/main/notebooks/07-stat-significance.ipynb>