

STATISTICS

ENG 3120

2025 - 2026 Spring Semester

 Assoc. Prof. Dr. Bora CANBULA

 www.canbula.com

 github.com/canbula/Statistics

 w520xjh

Statistics

Instructor

Assoc. Prof. Dr.
Bora CANBULA

Phone

0 (236) 201 21 08

Email

bora.canbula@cbu.edu.tr

Office Location

Dept. of CENG

Office C233

Office Hours

1 pm – 2 pm, Tuesdays

Course Overview

Statistics (Teams Code: w520xjh)

We are going to learn both the mathematical foundations and real-world application of the statistics and the probability in this course. Focus of this course will be to provide the required background for a data science / machine learning course. Python is preferred as the programming language for the applications of this course.

Required Text

Probability And Statistics for Computer Scientists, CRC Press, *Michael Baron*

Introduction to Probability and Statistics, Elsevier, *Sheldon M. Ross*

Probability and Statistics for Engineers and Scientists, Brooks/Cole, *A.J. Hayter*

Course Materials

- Python 3.x (Anaconda is preferred)
- Jupyter Notebook from Anaconda
- Pycharm from JetBrains / Visual Studio Code from Microsoft

Course Schedule

Week	Subject	Week	Subject
01	Definitions of Descriptive Statistics	08	Linear Regression
02	Data, Sampling, and Variation	09	Linear Regression with Matrix Algebra
03	Visualization of Data	10	Regression with High Degree Polynomials
04	Measures of Central Tendency	11	Data Linearization and Transformation
05	Measures of Variation	12	Chi-Square and Goodness-of-Fit Tests
06	Measures for Multiple Variables	13	Central Limit Theorem
07	Box Plots and Outliers	14	Probability Distributions

Statistics is the science of collecting, organizing, analyzing, interpreting, and presenting the **data**.

Data is any kind of information.

Data are the actual values of the variable and can be categorical or numerical.

Population is the collection of people, things, or objects under study.

Sample is a subset of the population.

Statistic is a number that represents a property of the sample.

Parameter is a characteristic of the whole population that can be estimated by a statistic.

Variable is a characteristic or measurement that can be determined for each member of a population. They can be **dependent** or **independent**.

→ **Qualitative Variables** take on **Categorical** values.

→ **Quantitative Variables** take on **Numerical** values.

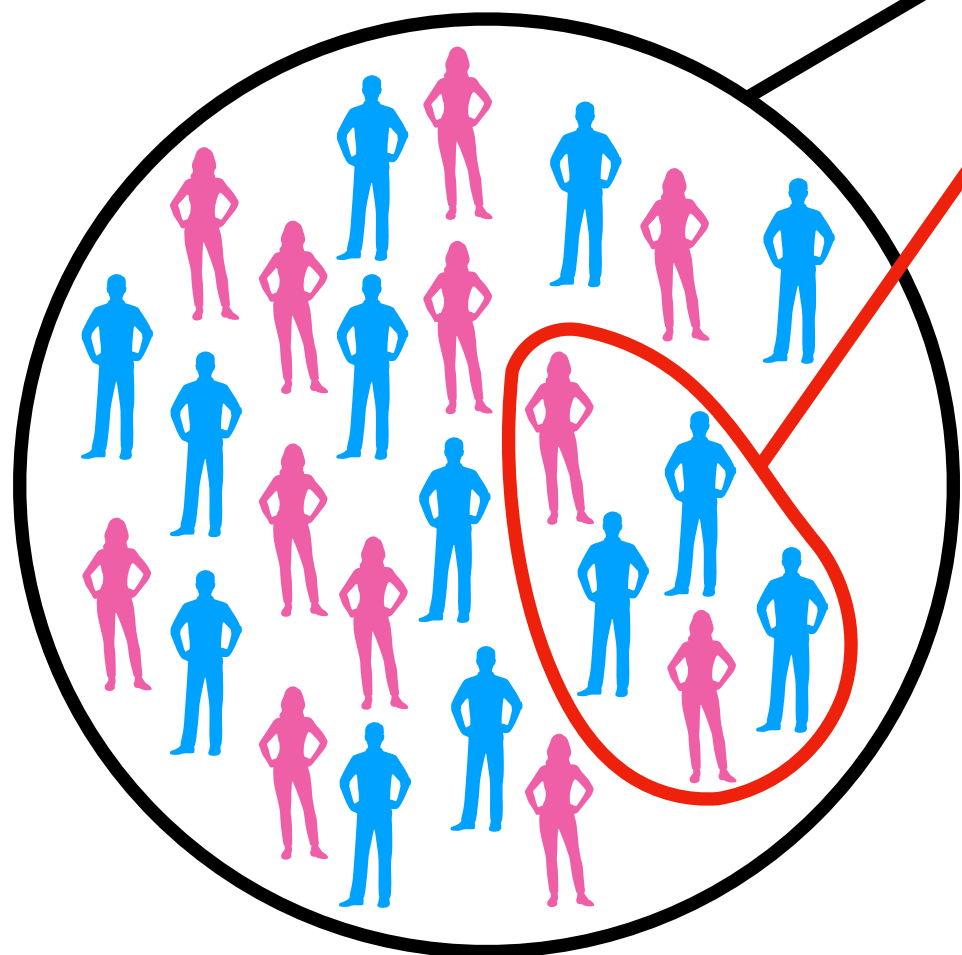
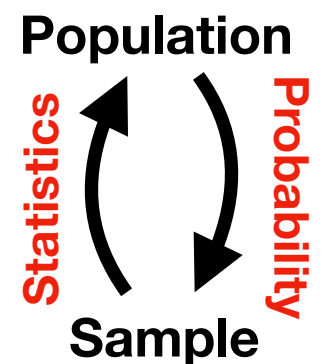
→ **Discrete Variables** take on finite number of values such as integers.

→ **Continuous Variables** take on infinite number of values such as real numbers.

Statistics

→ **Descriptive Statistics** is organizing and summarizing the data.

→ **Inferential Statistics** is drawing conclusions from good data.



good data == good sample

a good sample must be both random and representative

QUESTION








A study was conducted at our department to analyze the average GPA's of students who graduated last year. Match the key terms given below with the phrases that describes best.

A) Population **B)** Statistic **C)** Parameter **D)** Sample **E)** Variable **F)** Data

- ☒ **D** A group of students who graduated from our department last year
- ☒ **X** All students who attended last year
- ☒ **E** GPA of one student who graduated from our department last year
- ☒ **C** The average GPA of students who graduated from our department last year
- ☒ **A** All students who graduated from our department last year
- ☒ **F** 3.65, 2.80, 3.15, 3.90
- ☒ **B** _____

QUESTION

We plan on conducting a survey to our recent graduates to determine information on their yearly salaries. We randomly select 50 recent graduates and sent them questionnaires dealing with their present jobs. Of these 50, however, only 36 were returned. Suppose that the average of the yearly salaries reported was 415000 TL.

-  The population is: **Our all recent graduates**
-  The sample is: **36 recent graduates who returned to questionnaire**
-  The statistic is: **Yearly salary of 36 students**
-  The parameter is: **Yearly salary of our all recent graduates**
-  The variable is: **Yearly salary of one recent graduates**
-  Would we be correct in thinking that 415000 TL was a good approximation to the average salary level for all of our graduates? **No**
-  If your answer is no, can you think of any set of conditions relating to the group that returned questionnaires for which it would be a good approximation? **Suggest some questions**

QUESTION

An insurance company would like to determine the proportion of all medical doctors who have been involved in one or more malpractice lawsuits. The company selects 500 doctors at random from a professional directory and determines the number in the sample who have been in a malpractice lawsuit.

- ✍ The population is: **All medical doctors listed in the prof. directory**
- ✍ The sample is: **Selected 500 doctors**
- ✍ The statistic is: **The proportion of medical doctors in the sample**
- ✍ The parameter is: **The proportion of medical doctors in population**
- ✍ The variable is: **The number of medical doctors who have been**
- ✍ The data are: **Yes / No**

QUESTION

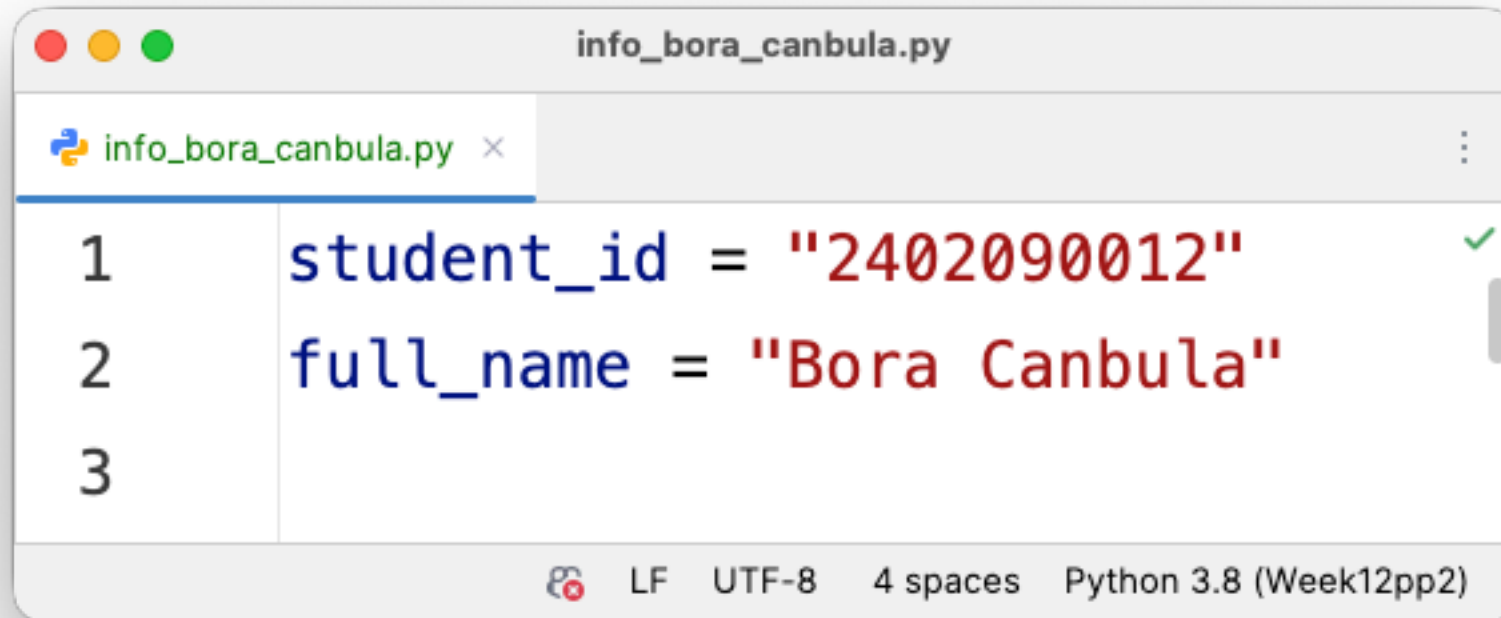
Determine the correct data type for the variables given below. Indicate whether quantitative data are continuous or discrete.

A) Numerical and discrete **B)** Numerical and continuous **C)** Categorical

- ☒ **A** The number of pairs of shoes you own
- ☒ **C** Gender
- ☒ **B** The distance from your home to university
- ☒ **A** The number of courses you take this semester
- ☒ **C** The brand of your mobile phone
- ☒ **B** Your weight
- ☒ **A** Number of correct answers on a quiz
- ☐ **?** Age

Week01/info_firstname_lastname.py

- ✓ A string variable with the name `student_id` that contains your student id.
- ✓ A string variable with the name `full_name` that contains your full name.



```
info_bora_canbula.py
1 student_id = "2402090012"
2 full_name = "Bora Canbula"
3
```

LF UTF-8 4 spaces Python 3.8 (Week12pp2)

Sampling

Population (N) \longrightarrow Sample (n)

BEST**HARDEST****Simple Random Sampling**

Every member of the population has an equal chance to be in the sample.

Stratified Sampling

The population is split into non-overlapping groups, which is called strata, then simple random sampling is applied to each group.

Cluster Sampling

The population is divided into groups (clusters), then some of the clusters are randomly selected.

Systematic Sampling

The sample is constructed with every n^{th} individual from the population.

Convenience Sampling

The sample is constructed with easily obtained members of the population.

WORST**EASIEST**

🔄 Determine the type of sampling used in the following examples:

🕒 A soccer coach selects six players from a group of boys aged eight to ten, seven players from a group of boys aged 11 to 12, and three players from a group of boys aged 13 to 14 to form a recreational soccer team.

Stratified Sampling

🕒 A pollster interviews all human resource personnel in five different high tech companies.

Cluster Sampling

🕒 A high school educational researcher interviews 50 high school female teachers and 50 high school male teachers.

Stratified Sampling

🕒 A medical researcher interviews every third cancer patient from a list of cancer patients at a local hospital.

Systematic Sampling

🕒 A high school counselor uses a computer to generate 50 random numbers and then picks students whose names correspond to the numbers.

Simple Random Sampling

🕒 A student interviews classmates in his algebra class to determine how many pairs of jeans a student owns, on the average.

Convenience Sampling

Sampling

Population (N) → Sample (n)

BEST

HARDEST

Simple Random Sampling

Every member of the population has an equal chance to be in the sample.

```
cities = [  
    "Adana",  
    "Adıyaman",  
    "Afyonkarahisar",  
    "Ağrı",  
    # ...  
    "Kilis",  
    "Osmaniye",  
    "Düzce",  
]
```

```
import sys  
  
sys.path.append(".")  
  
from Week02 import data  
  
print(data.cities)
```

import random

```
[  
    "betavariate",  
    "binomialvariate",  
    "choice",  
    "choices",  
    "expovariate",  
    "gammavariate",  
    "gauss",  
    "getrandbits",  
    "getstate",  
    "lognormvariate",  
    "normalvariate",  
    "paretovariate",  
    "randbytes",  
    "randint",  
    "random",  
    "randrange",  
    "sample",  
    "seed",  
    "setstate",  
    "shuffle",  
    "triangular",  
    "uniform",  
    "vonmisesvariate",  
    "weibullvariate",  
]
```

Help on method sample in module random:

`sample(population, k, *, counts=None)`
Chooses k unique random elements from the population sequence or set.

Returns a new list leaving the original population sequence or set unchanged. This also works for population being a set.

Members of the population are chosen without replacement. This is equivalent to selection in the sample without replacement.

Repeated elements counts parameter.

`sample(['red', 'blue', 'green'], 2)`
is equivalent to:

`sample(['red', 'blue', 'green'], 2, counts=[1, 1, 1])`

To choose a sample of size k from a population of size n, you can use `sample(range(n), k)`.

import sys

sys.path.append(".")

```
from Week02 import data  
import random
```

```
def simple_random_sampling(data, n):  
    return random.sample(data, n)
```

```
sample = simple_random_sampling(data.cities, 10)  
print(sample)
```

WORST

EASIEST

Population (N) \longrightarrow Sample (n)

Stratified Sampling

The population is split into non-overlapping groups, which is called strata, then simple random sampling is applied to each group.

```
def stratified_sampling(data, n, strata):
    sample = []
    for key in strata:
        sample += random.sample(data[key], n)
    return sample
```

► **Fix this!**

```
sample = stratified_sampling(data.cities_by_region, 3, data.regions)
print(sample)
```

```
cities_by_region = {
    "Marmara": ["Edirne", "Kırklareli", "Tekirdağ", "Trakya"],
    "İç Anadolu": ["Aksaray", "Ankara", "Çankırı", "Konya", "Sakarya"],
    "Ege": ["İzmir", "Manisa", "Aydın", "Denizli", "Muğla"],
    "Akdeniz": ["Adana", "Osmaniye", "Antalya", "Hatay", "Mersin"],
    "Karadeniz": ["Rize", "Trabzon", "Artvin", "Giresun", "Samsun"],
    "Doğu Anadolu": ["Ağrı", "Ardahan", "Bingöl", "Erzurum", "Van"],
    "Güneydoğu Anadolu": ["Adıyaman", "Batman", "Diyarbakır", "Gaziantep", "Mardin", "Siirt", "Şanlıurfa", "Tunceli"]
}
```

HARDEST

WORST

EASIEST

Population (N) \longrightarrow Sample (n)

Cluster Sampling

The population is divided into groups (clusters), then some of the clusters are randomly selected.

```
def stratified_sampling(data, n, strata):
    sample = []
    for key in strata:
        sample += random.sample(data[key], n)
    return sample
```

Fix this!

```
sample = stratified_sampling(data.cities_by_region, 3, data.regions)
print(sample)
```

```
cities_by_region = {  
    "Marmara": ["Edirne", "Kırklareli", "Tekirdağ",  
        "İç Anadolu": ["Aksaray", "Ankara", "Çankırı",  
            "Ege": ["İzmir", "Manisa", "Aydın", "Denizli",  
                "Akdeniz": ["Adana", "Osmaniye", "Antalya",  
                    "Karadeniz": ["Rize", "Trabzon", "Artvin", "Görecik",  
                        "Doğu Anadolu": ["Ağrı", "Ardahan", "Bingöl",  
                            "Güneydoğu Anadolu": ["Adıyaman", "Batman", "Diyarbakır"]  
}
```

```
def cluster_sampling(data, n, clusters):
    picked_clusters = random.sample(clusters, n)
    sample = []
    for cluster in picked_clusters:
        sample += data[cluster]
    return sample
```

Duplicates?

Duplicates?

HARDEST

EASIEST

WORST

Sampling

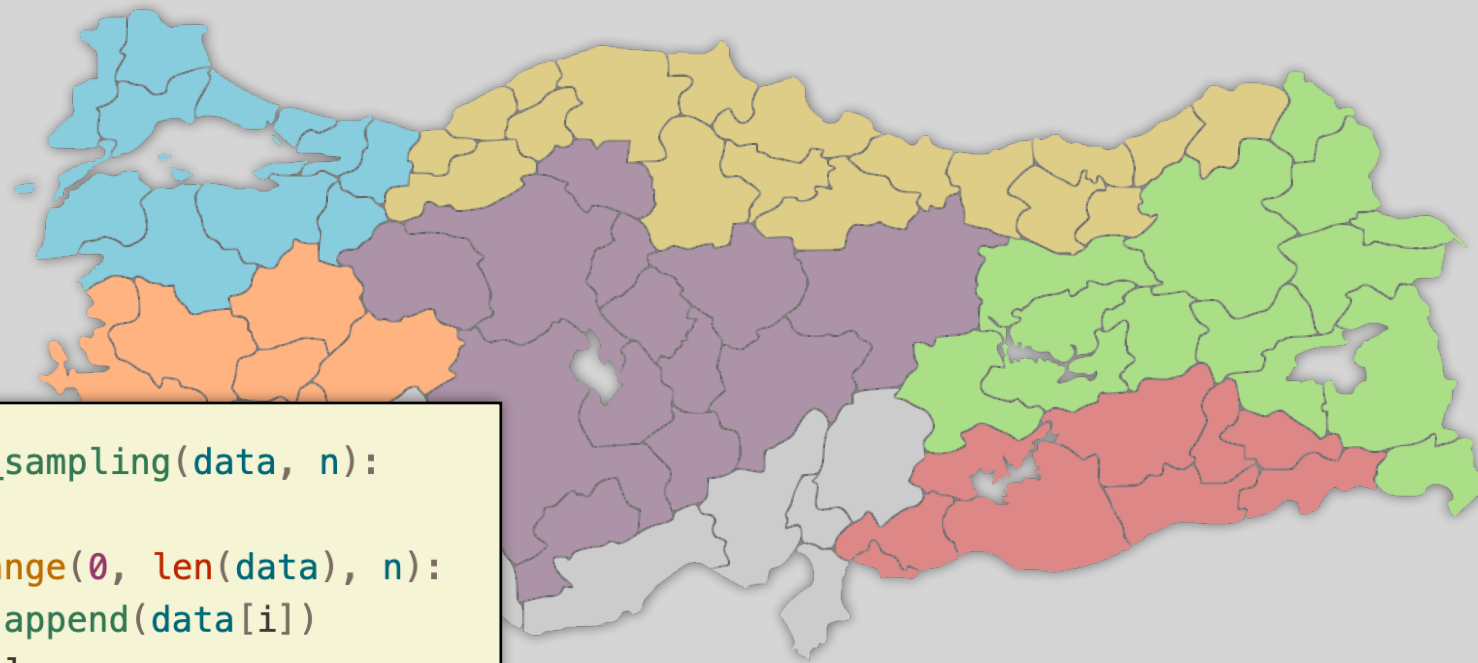
Population (N) → Sample (n)

BEST

HARDEST

Systematic Sampling

The sample is constructed with every n^{th} individual from the population.



```
def systematic_sampling(data, n):  
    sample = []  
    for i in range(0, len(data), n):  
        sample.append(data[i])  
    return sample
```

WORST

EASIEST

Sampling

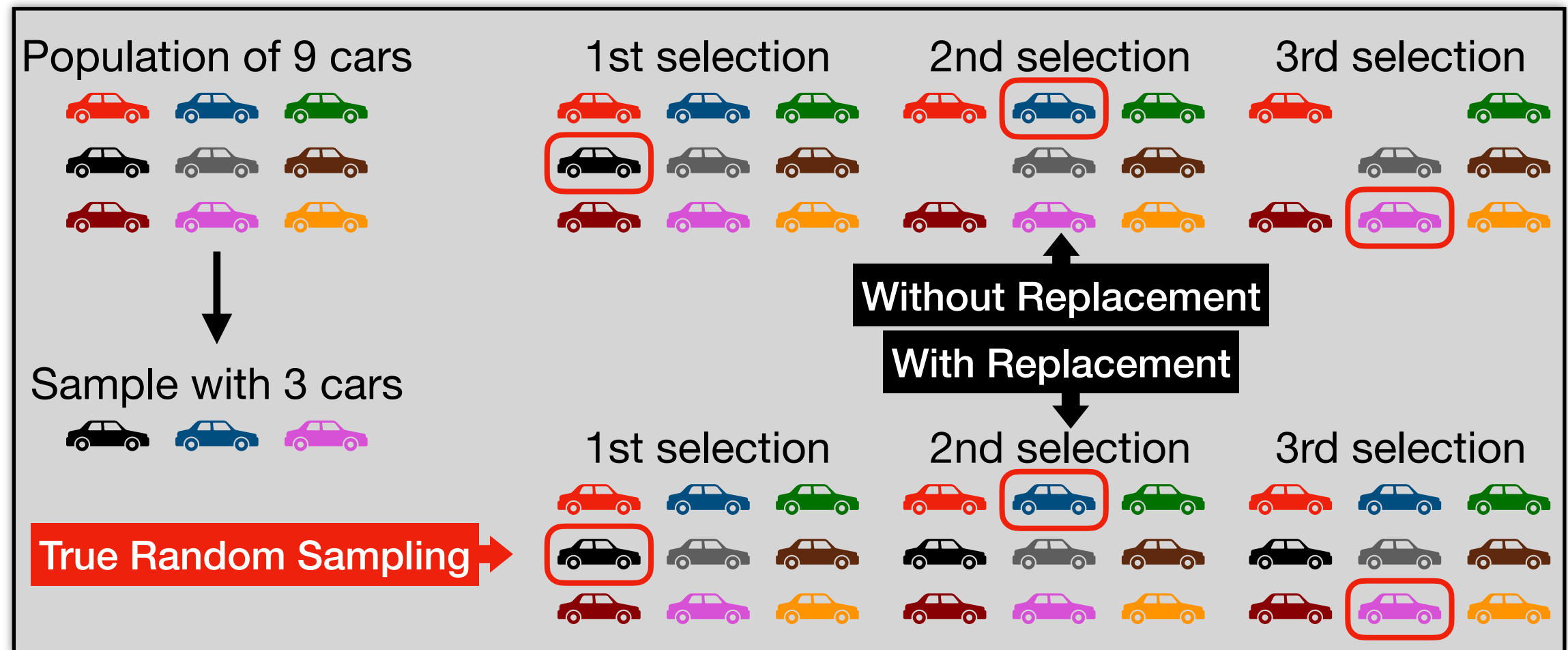
Population (N) → Sample (n)

BEST

HARDEST

Simple Random Sampling

Every member of the population has an equal chance to be in the sample.



WORST

EASIEST

IMPLEMENT THIS FEATURE AS A HOMEWORK

Weighted Simple Random Sampling w/ Replacement Support**Name of Your File**

Week02/weighted_firstname_lastname.py

Name of the Function

weighted_srs

Input Parameters

- data [list]: population
- n [int]: sample size
- weights [list]: weights for members of population
- with_replacement [bool]: flag for true random sampling

Other Rules

- Using maximum 10 lines of codes is allowed
- Using any modules, other than random, is forbidden