

D1EAD – Análise Estatística para Ciência de Dados

2021.1



Data Distributions (Part 1)

Prof. Ricardo Sovat

sovat@ifsp.edu.br

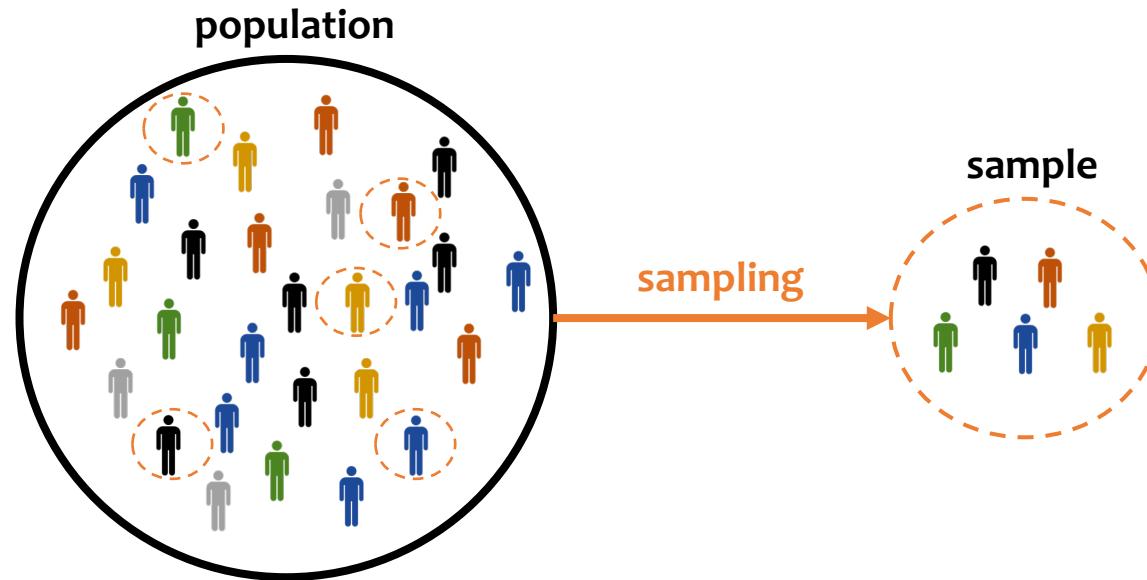
Prof. Samuel Martins (Samuka)

samuel.martins@ifsp.edu.br



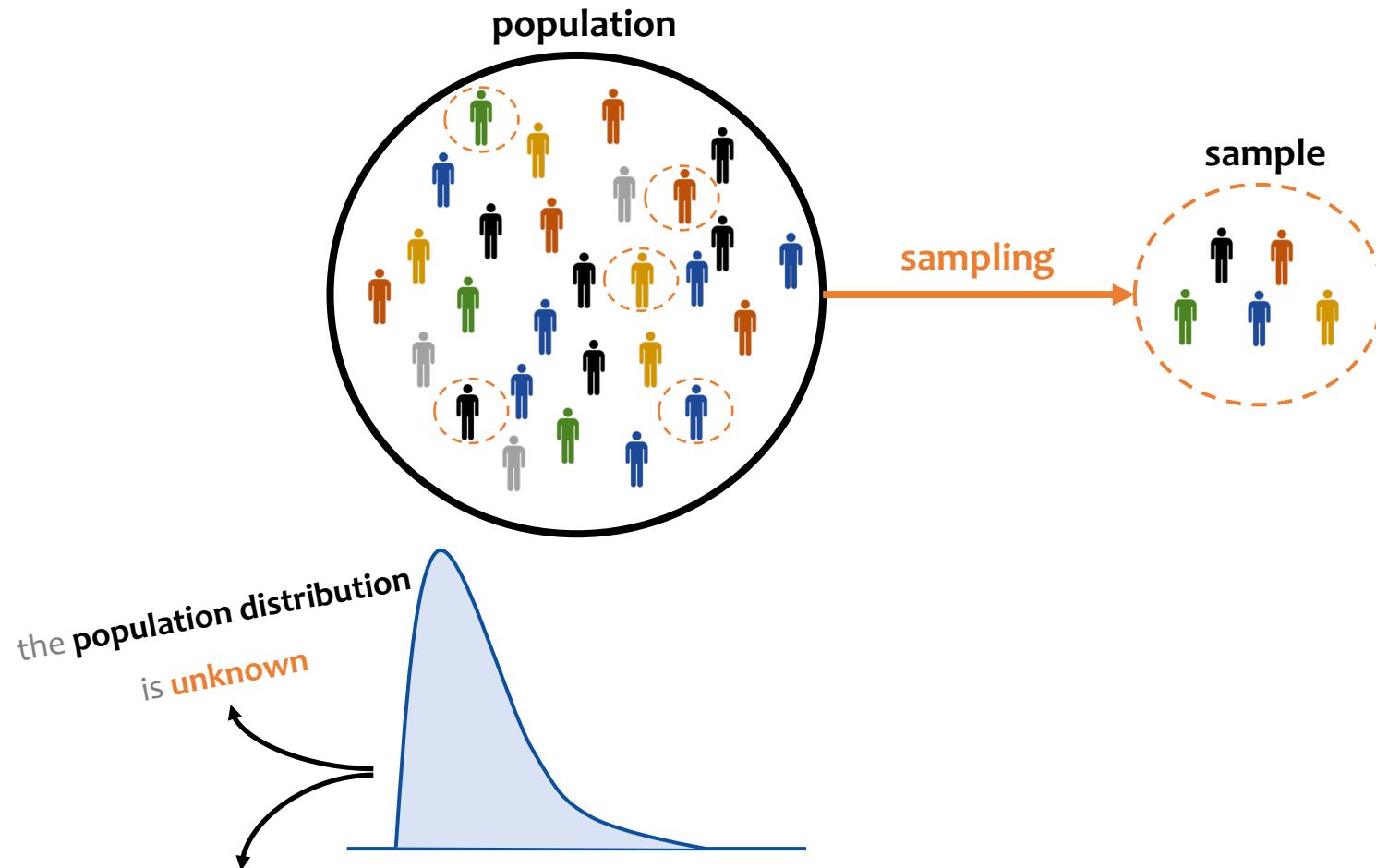
Population vs Sample

Even in the era of **big data**, **sampling** keeps being important and relevant



Population vs Sample

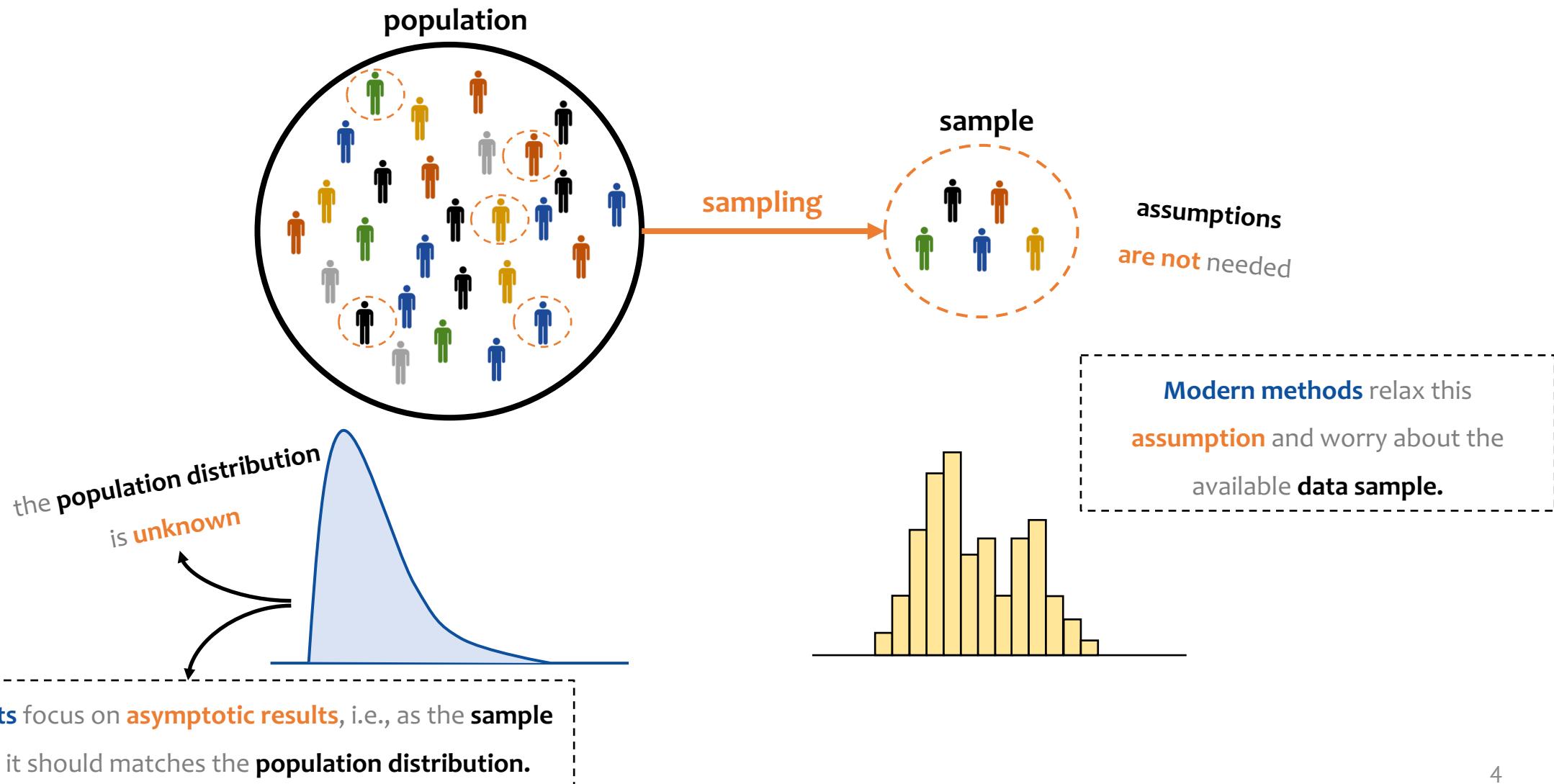
Even in the era of **big data**, **sampling** keeps being important and relevant



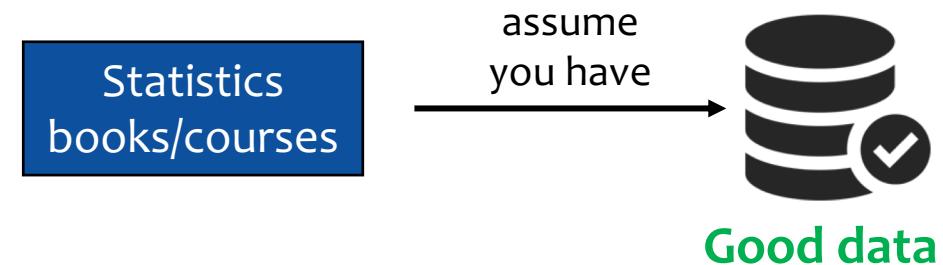
Frequentist stats focus on **asymptotic results**, i.e., as the **sample size** increases it should matches the **population distribution**.

Population vs Sample

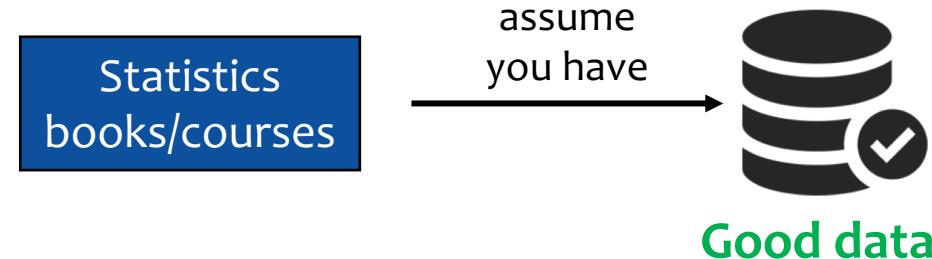
Even in the era of **big data**, **sampling** keeps being important and relevant



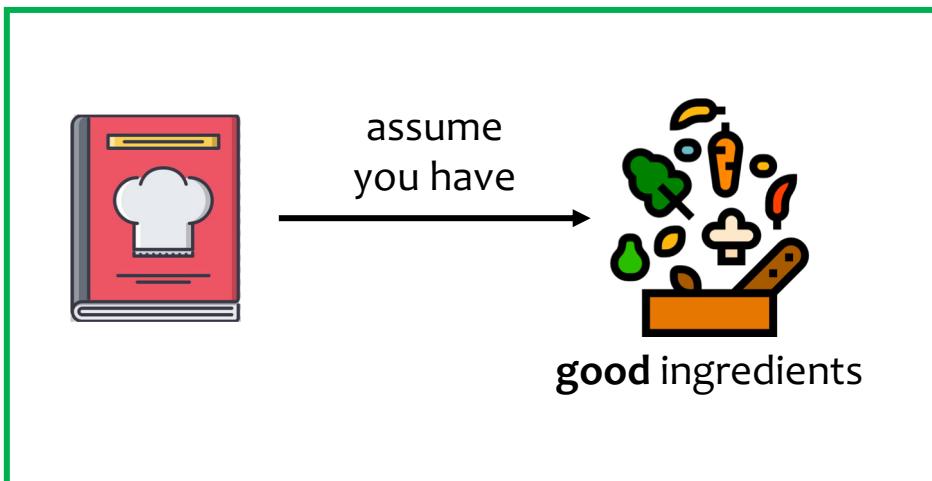
The Importance of Data



The Importance of Data

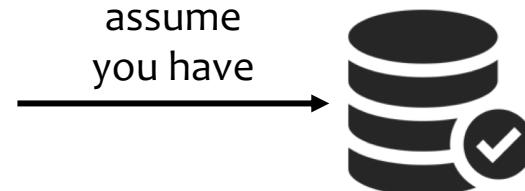


it's like



The Importance of Data

Statistics
books/courses



Good data

it's like



assume
you have →



good ingredients

BUT



+

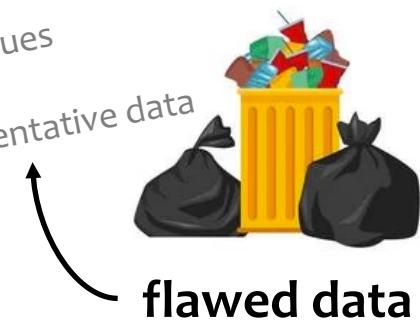


spoiled ingredients



The Importance of Data

- Noise
- Missing values
- Bias
- Unrepresentative data



+

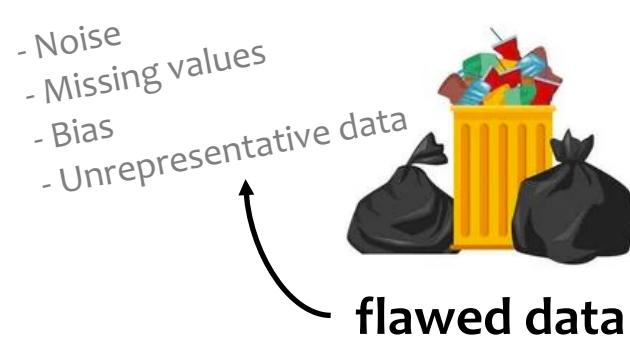


many fancy
analysis



flawed
results/conclusions

The Importance of Data



+



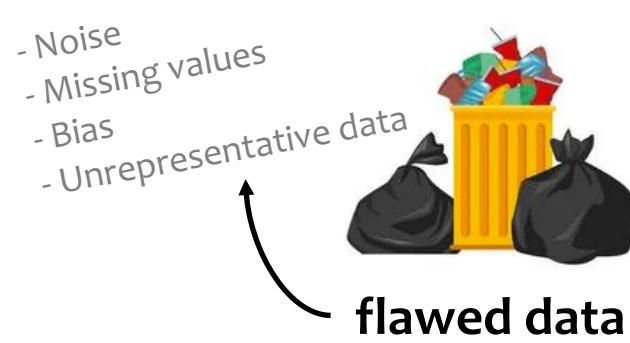
**many fancy
analysis**



**flawed
results/conclusions**

Garbage In → **Garbage Out**

The Importance of Data



+



many fancy
analysis



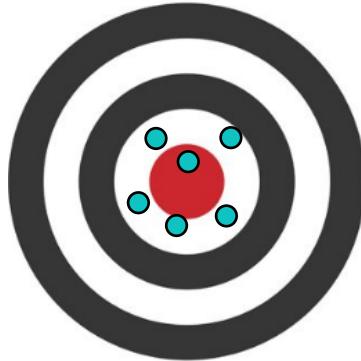
Garbage In → Garbage Out

Data quality often matters more than data quantity when making an estimate or a model based on a sample.

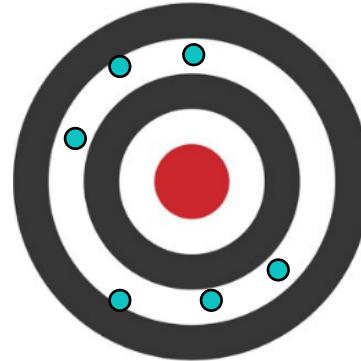
Precision vs Accuracy



✓ Precision
✗ Accuracy



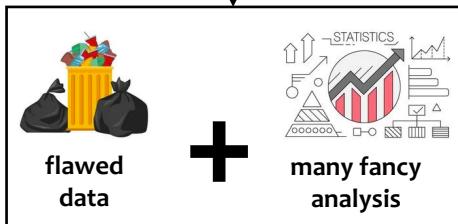
✗ Precision
✓ Accuracy



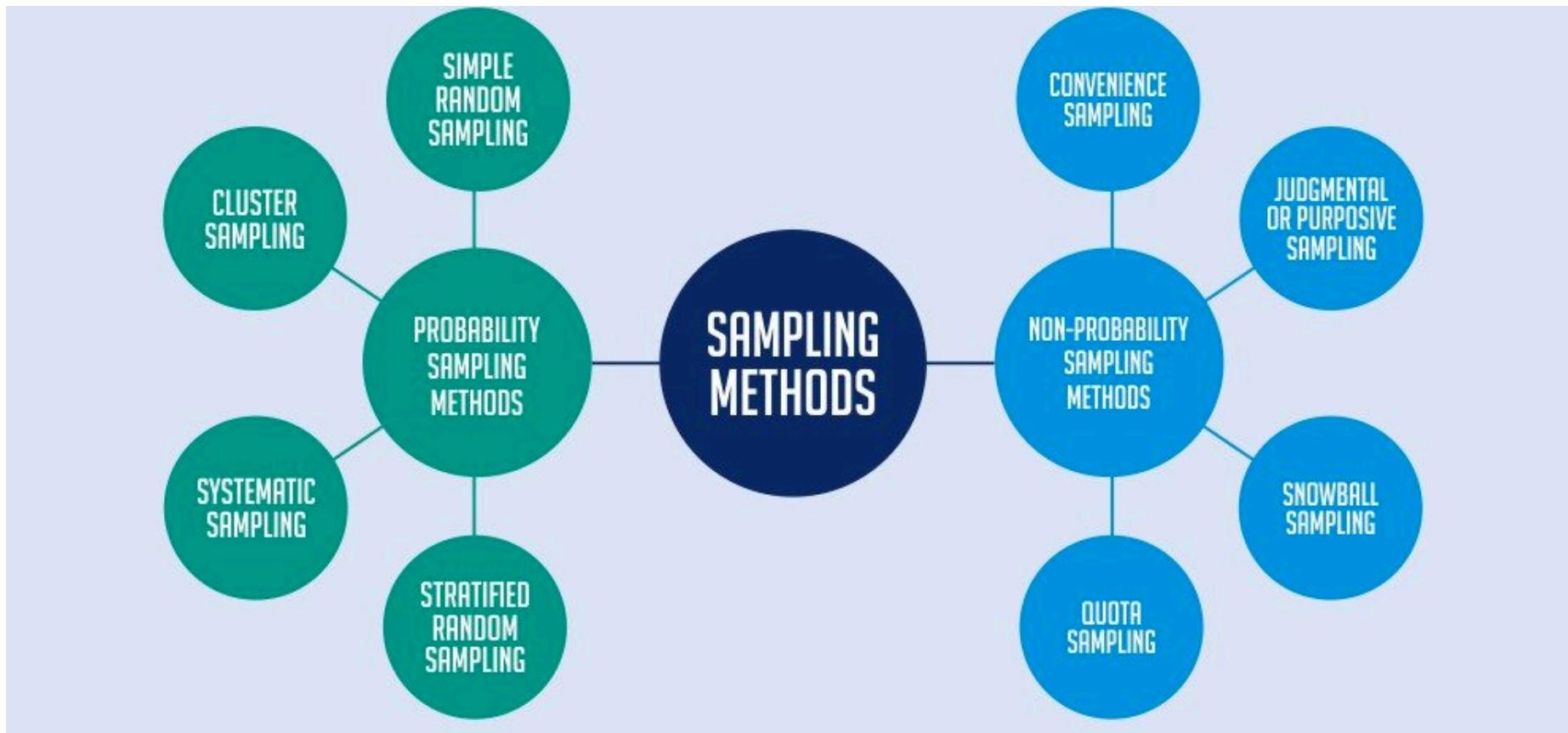
✗ Precision
✗ Accuracy

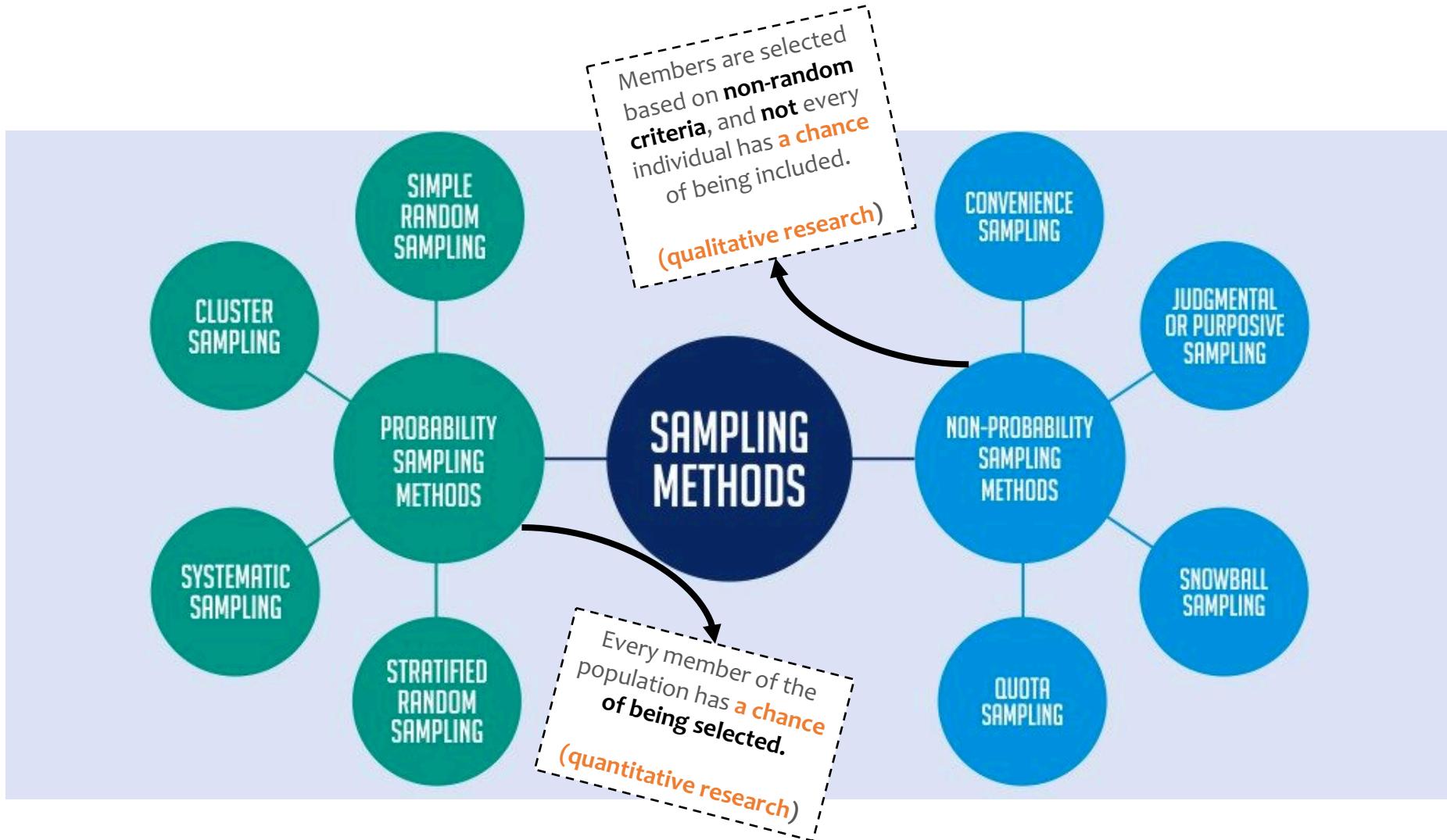


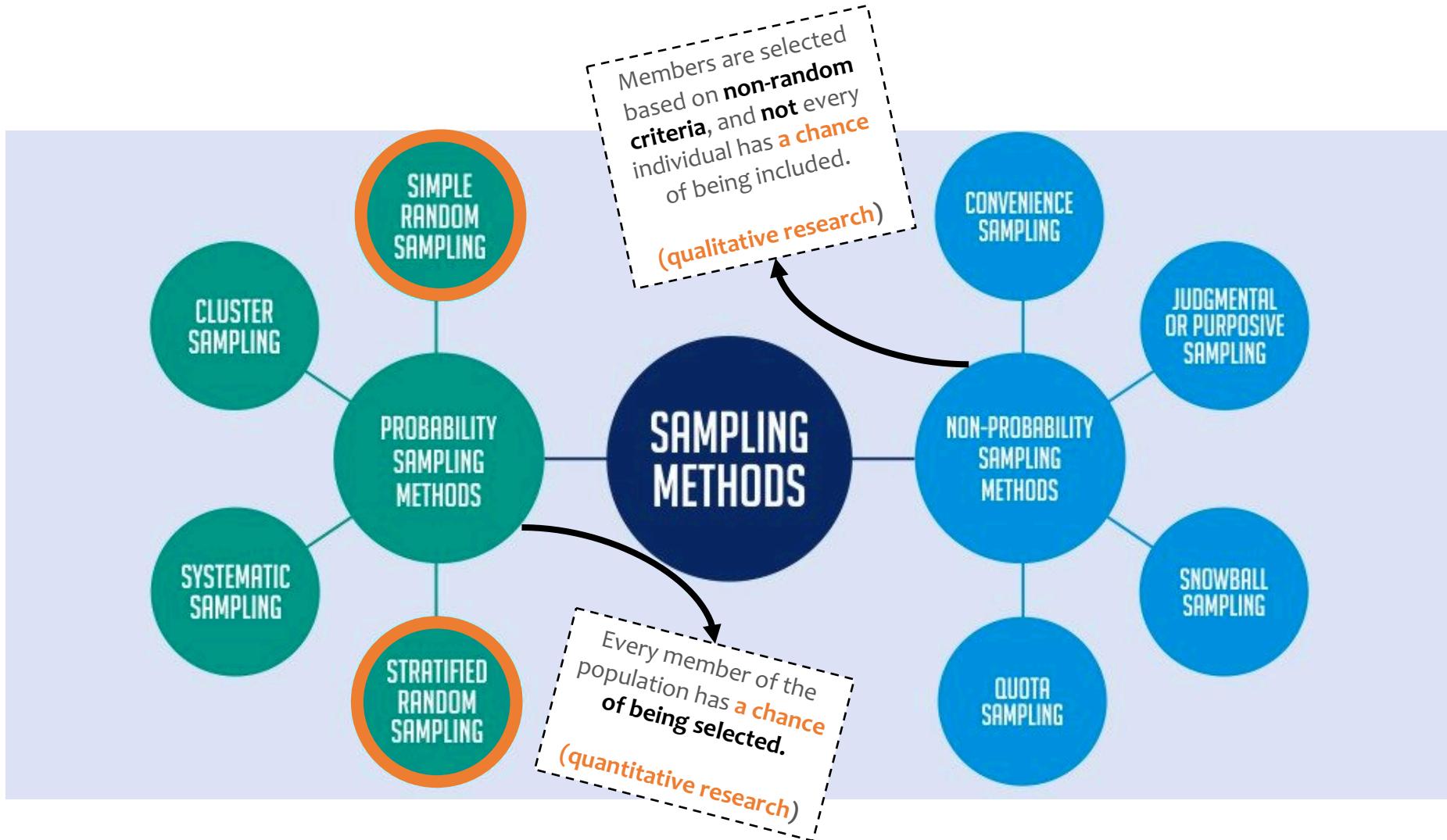
✓ Precision
✓ Accuracy



Types of Sampling



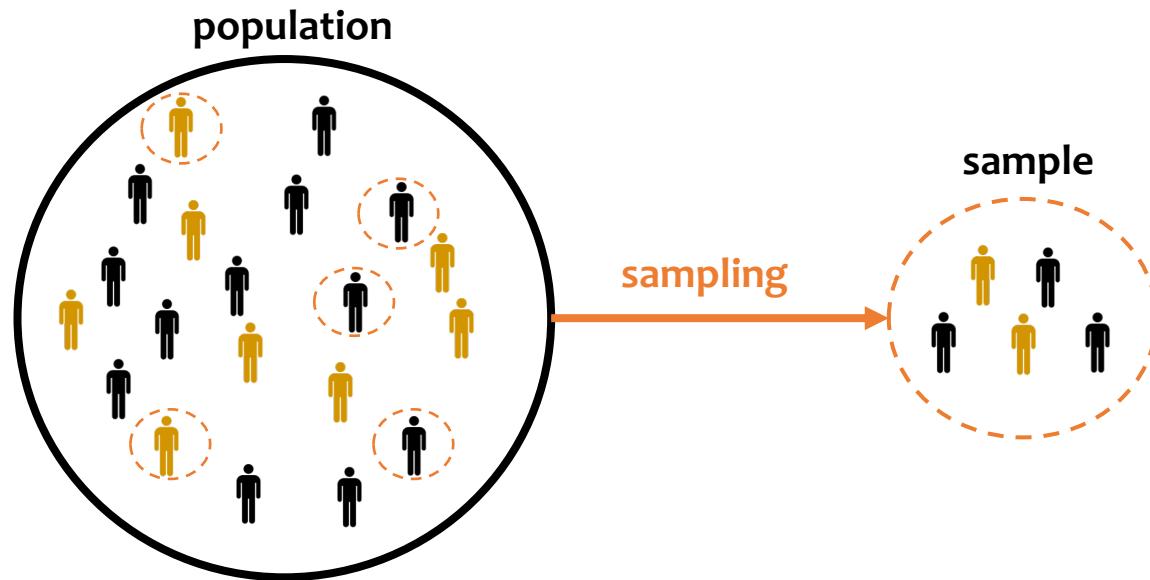




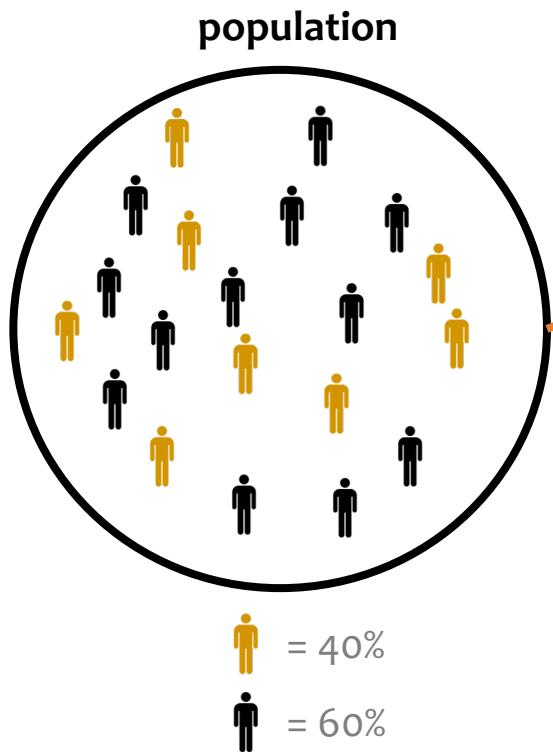
Random Sampling

The common way for **sampling** a population and to **avoid selection bias** (we'll see soon!).

Each **member/observation** has **an equal chance** of being chosen for the **sample** at each draw.



Random Sampling



sampling

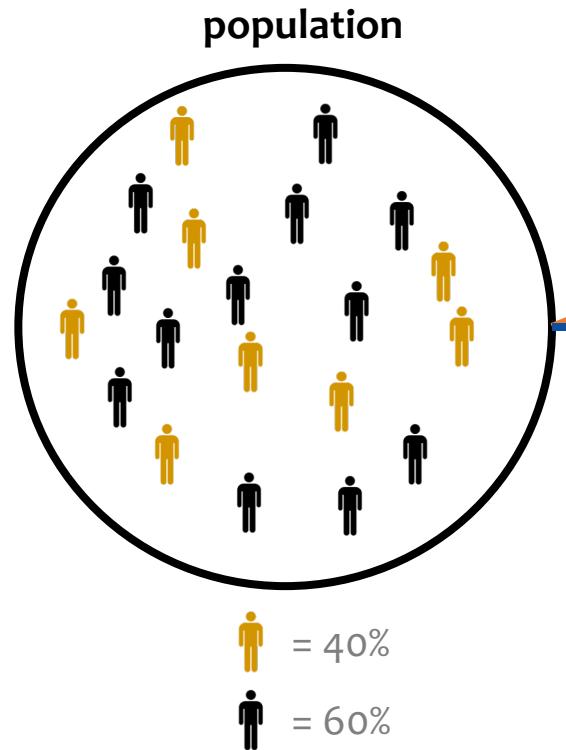
sample #1



Same proportion from
the population
distribution, but no
guarantees for that.

■ = 40%
■ = 60%

Random Sampling



sampling

sampling

Same proportion from the population distribution, but no guarantees for that.

Yellow icon = 40%

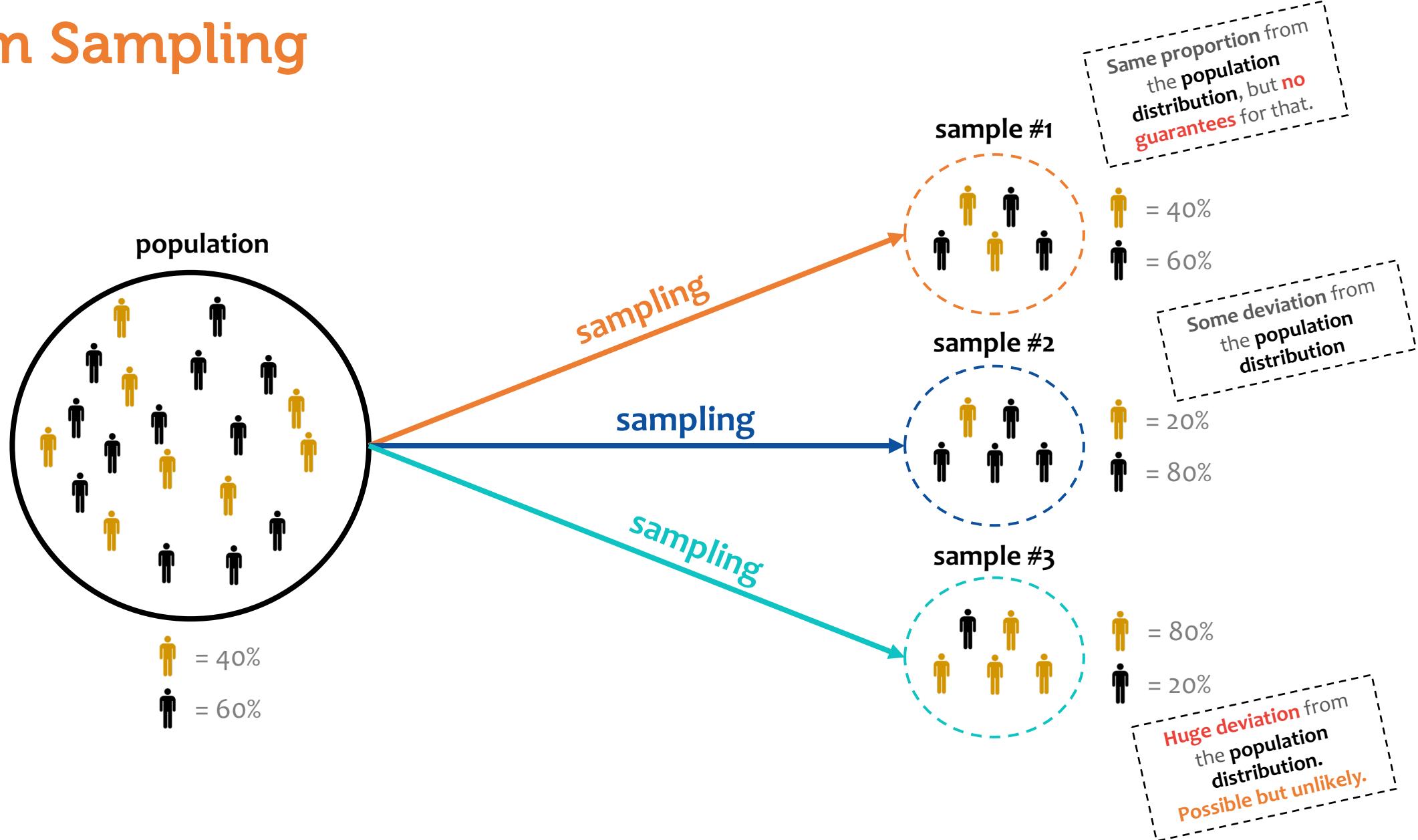
Black icon = 60%

Some deviation from the population distribution

Yellow icon = 20%

Black icon = 80%

Random Sampling



Random Sampling – Pandas

	App	Rating	Type
Duolingo: Learn Languages	Free	4.7	Free
Z City	4.3	Free	
EF Coach	4.8	Free	
Dating Network	4.0	Free	
Chess of Blades (BL/Yaoi Game) (No VA)	4.8	Paid	
I am rich	3.8	Paid	
Gangster Town	4.1	Free	
Safest Call Blocker	4.4	Free	
BJ's Bingo & Gaming Casino	4.5	Free	
Chess School for Beginners	4.3	Free	

Free: 10039 (92.62%)

Paid: 800 (7.38%)

Total: 10839 observations

random sampling

```
sample = population.sample(100)
```

	App	Rating	Type
Badoo - Free Chat & Dating App	4.3	Free	
What was I in my Past Life	3.7	Free	
Diabetes & Diet Tracker	4.6	Paid	
ez Share Android app	3.3	Free	
IP address BW	NaN	Free	
Carousell: Snap-Sell, Chat-Buy	4.3	Free	
Simpli CT	NaN	Free	
FH WiFiCam	2.6	Free	
Muscle Premium - Human Anatomy, Kinesiology, B...	4.2	Paid	
My Teacher - Classroom Play	4.0	Free	

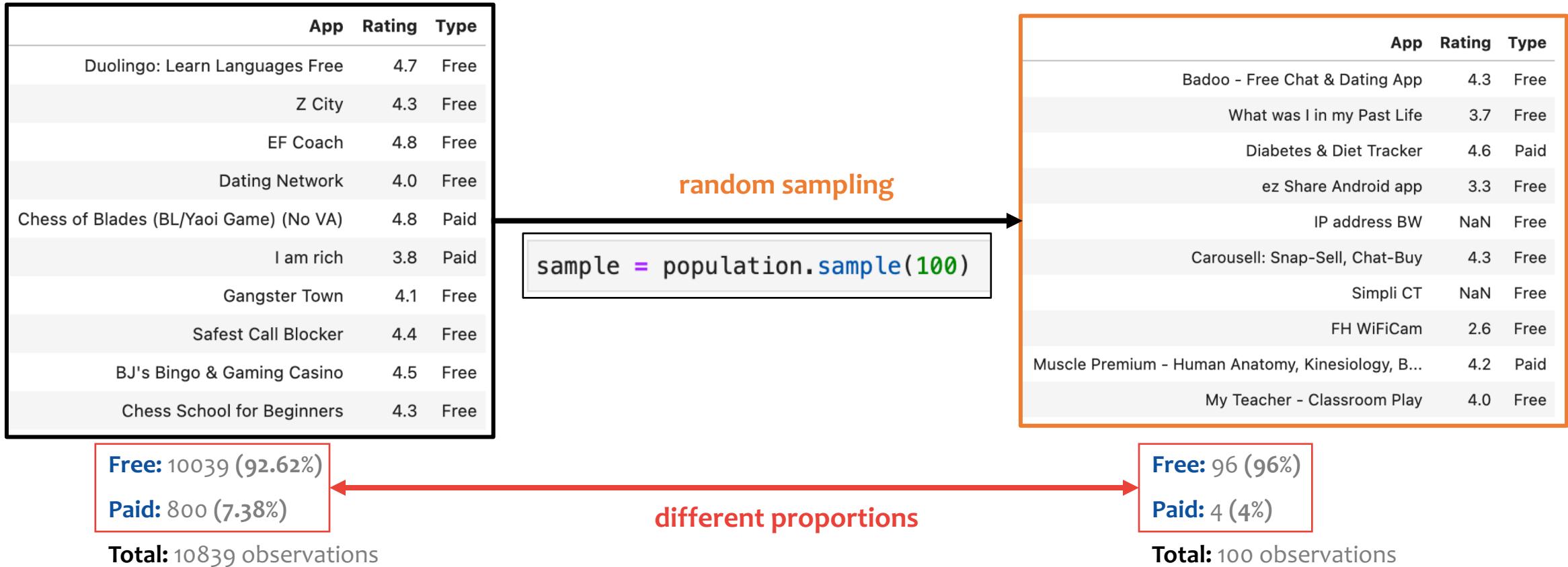
Free: 96 (96%)

Paid: 4 (4%)

Total: 100 observations

Dataset: Google Play Store Apps: <https://www.kaggle.com/lava18/google-play-store-apps>

Random Sampling – Pandas



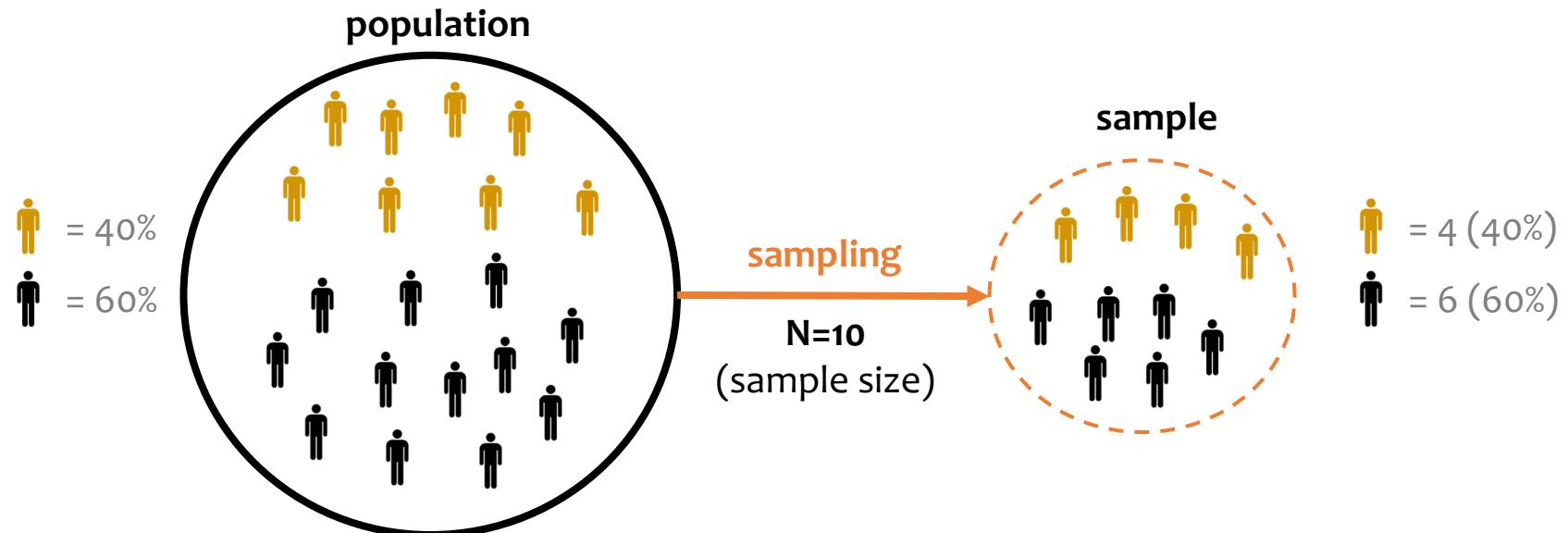
Dataset: Google Play Store Apps: <https://www.kaggle.com/lava18/google-play-store-apps>

Stratified Sampling

Divide the population into **subgroups (strata)** based on a **relevant characteristic** (e.g., gender, age range, healthy/unhealthy, ...)

Perform **random sampling** in each **subgroup**, respecting the **overall population proportion**.

Aim at drawing **more precise conclusions** by ensuring that every **subgroup/class** is properly represented in the **sample**.



Stratified Sampling – Pandas

	App	Rating	Type
Duolingo: Learn Languages	Free	4.7	Free
Z City	4.3	Free	
EF Coach	4.8	Free	
Dating Network	4.0	Free	
Chess of Blades (BL/Yaoi Game) (No VA)	4.8	Paid	
I am rich	3.8	Paid	
Gangster Town	4.1	Free	
Safest Call Blocker	4.4	Free	
BJ's Bingo & Gaming Casino	4.5	Free	
Chess School for Beginners	4.3	Free	

stratified sampling

```
sample = stratified_sampling(population,  
                             sample_size=100)
```

N=100 (sample size)

Free: 10039 (92.62%)
Paid: 800 (7.38%)

Total: 10839 observations

	App	Rating	Type
Badoo - Free Chat & Dating App	4.3	Free	
What was I in my Past Life	3.7	Free	
Diabetes & Diet Tracker	4.6	Paid	
ez Share Android app	3.3	Free	
IP address BW	NaN	Free	
Carousell: Snap-Sell, Chat-Buy	4.3	Free	
Simpli CT	NaN	Free	
FH WiFiCam	2.6	Free	
Muscle Premium - Human Anatomy, Kinesiology, B...	4.2	Paid	
My Teacher - Classroom Play	4.0	Free	

“equal” proportions

Free: 93 (93%)
Paid: 7 (7%)

Total: 100 observations

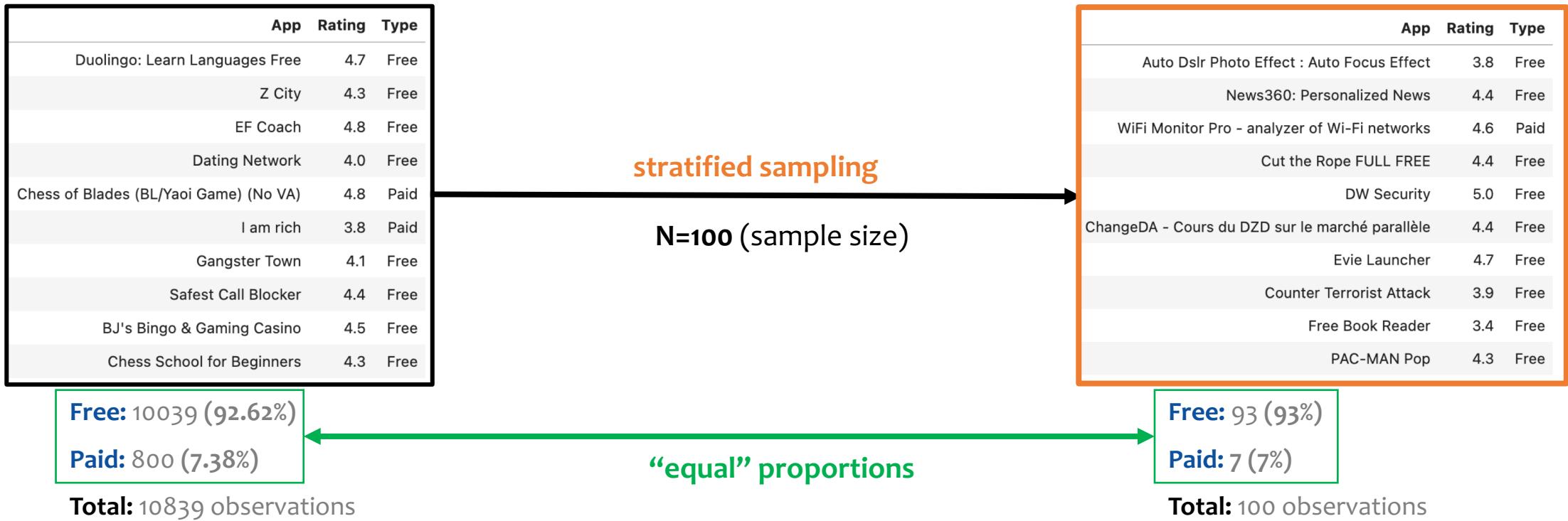
```
def stratified_sampling(population, sample_size)  
    count = population.groupby('Type').size()  
    proportion = count / population.shape[0]  
    n_obs = round(proportion * sample_size).astype('int')  
    sample = population.groupby('Type', group_keys=False)\n        .apply(lambda group: group.sample(n_obs.loc[group.name]))  
  
    return sample
```

Free
Paid
10039
800

Free
Paid
0.926192
0.073808

Free
Paid
93
7

Stratified Sampling – Pandas + Scikit-learn



```
from sklearn.model_selection import train_test_split  
sample, _ = train_test_split(population, train_size=100,  
                           stratify=population['Type'])
```

sample size

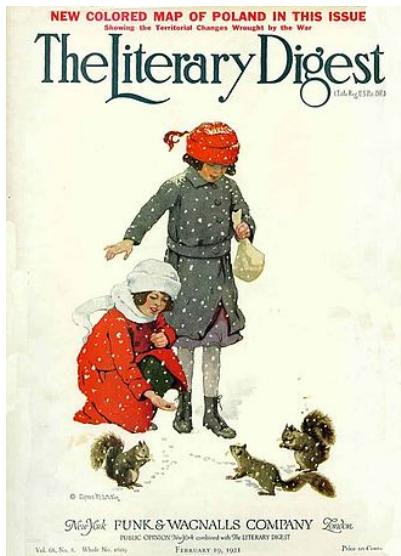
Bias

Bias

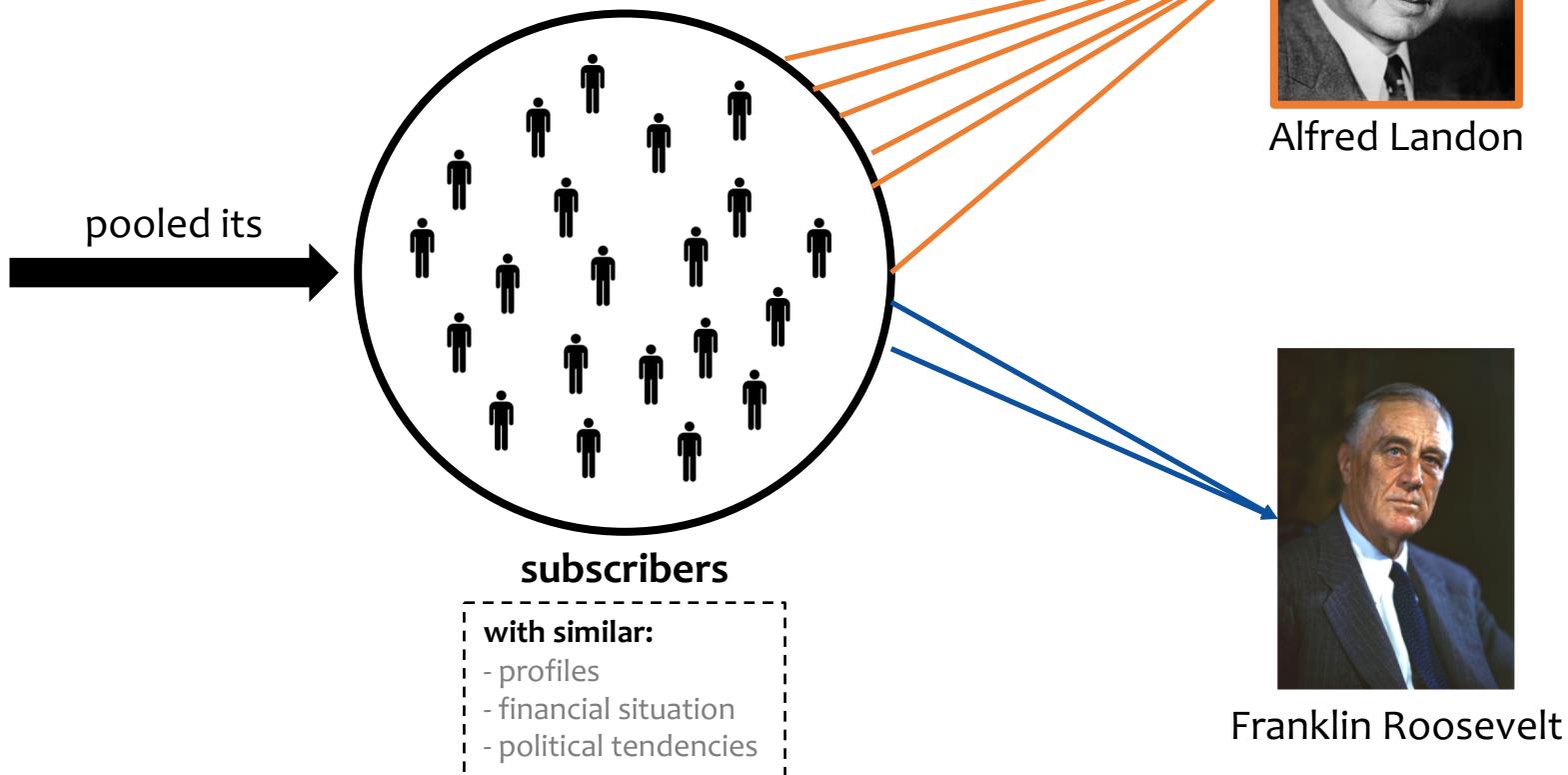
- **Measurement or sampling error** that are systemic and produced by the measurement or **sampling process**.
- **Tendency** of a statistic **overestimate** or **underestimate** a parameter.

Classical Example of Selection Bias

The Literary Digest (1936)

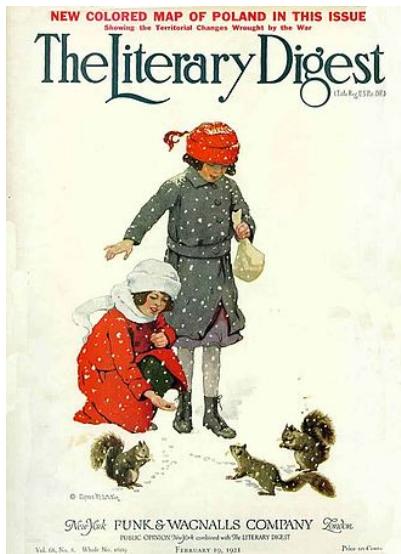


pooled its

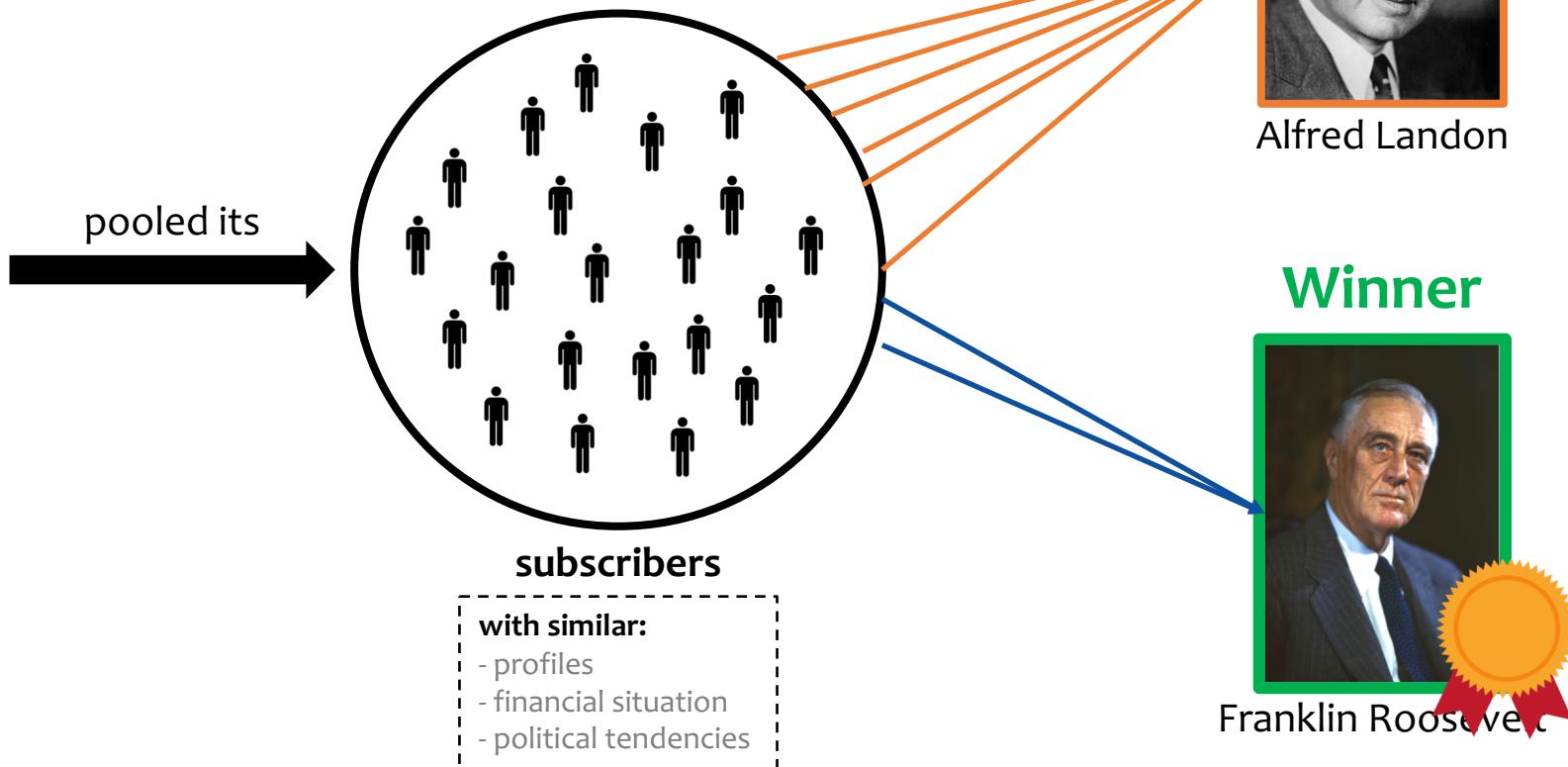


Classical Example of Selection Bias

The Literary Digest (1936)

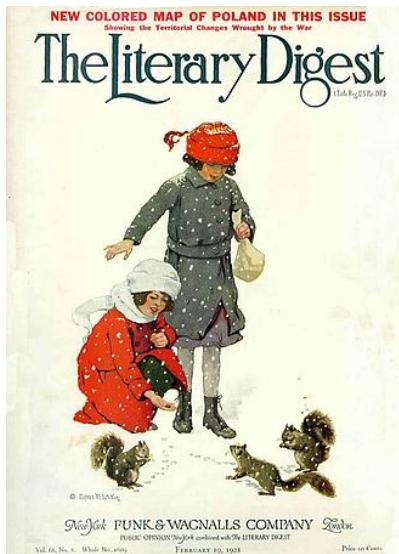


pooled its



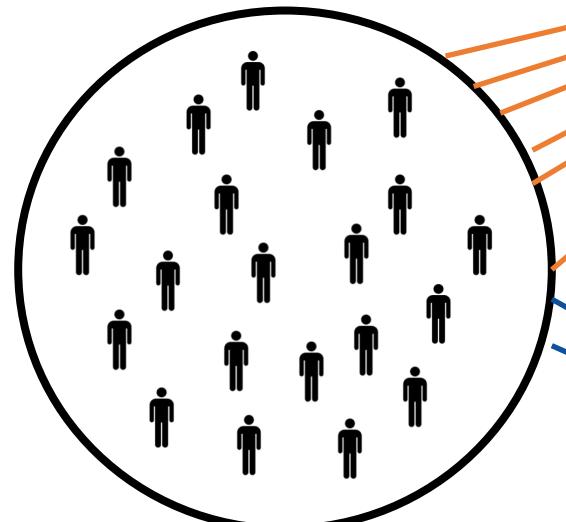
Classical Example of Selection Bias

The Literary Digest (1936)



Selection Bias

pooled its



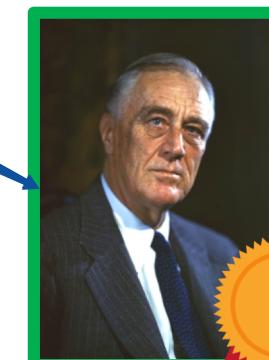
subscribers

with similar:

- profiles
- financial situation
- political tendencies



Alfred Landon



Winner
Franklin Roosevelt

Types of Sampling Biases

- Selection bias
- Self-selection bias
- Publication bias
- Recall bias
- Survivorship bias
- Healthy user bias
- ...

Types of Sampling Biases

- **Selection bias**
- **Self-selection bias**
- Publication bias
- Recall bias
- Survivorship bias
- Healthy user bias
- ...

Selection Bias

Observations or groups in a study **differ** systematically from the **population** of interest, leading to **errors** in association or outcome.

Selection Bias

Observations or groups in a study differ systematically from the population of interest, leading to errors in association or outcome.

Example: Survey at a specific neighborhood of a city to make conclusions about the city population.



fancy neighborhood

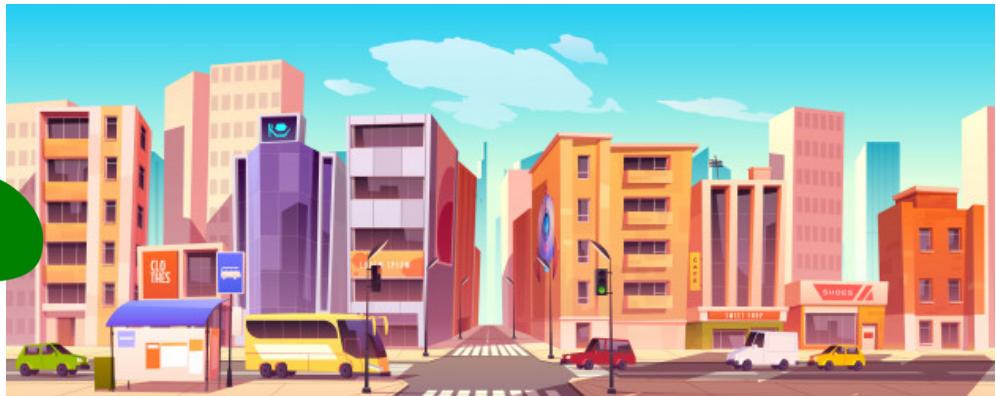


poor neighborhood

Selection Bias

Observations or groups in a study **differ** systematically from the **population** of interest, leading to **errors** in association or outcome.

Example: Survey at a **specific neighborhood** of a city to **make conclusions** about the **city population**.



fancy neighborhood



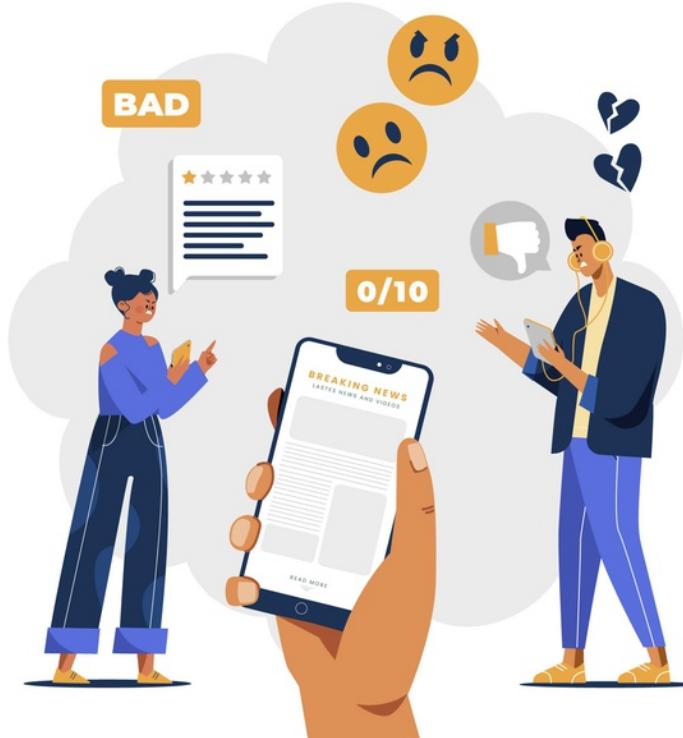
poor neighborhood

Both surveys are likely to be **biased** by the fact that people have **different characteristics**.

We can use them to investigate their **specific people/neighborhood**.

Self-selection Bias

Individuals **not randomly selected and motivated** to be part of a sample.



D1EAD – Análise Estatística para Ciência de Dados

2021.1



Data Distributions (Part 1)

Prof. Ricardo Sovat

sovat@ifsp.edu.br

Prof. Samuel Martins (Samuka)

samuel.martins@ifsp.edu.br

