

CDSF04, CDSF05, CDSF06

Recap II

DS Academy



TODOS OS DIREITOS RESERVADOS

2018



Don't give up

2



PYTHON
is
the
new
EXCEL



Pandas and SQL





I/O

[Files](#)

[DB's and DW's](#)

[Python Objects](#)

`pd.read_csv(), pd.read_json(), pd.read_pickle(), ...`

`pd.read_sql(), pd.read_gbq(), ...`

`pd.DataFrame.from_dict(), pd.DataFrame.from_records()`

Indexing

[Integer, Range, Slice](#)

[List](#)

[Boolean Array](#)

`df[start:end:step], df.iloc[], df.loc[]`

`df[[1,2,4,5]], s.reindex([1,2,3])`

`df[(df['col_a']=='foo') & df['col_b']=='bar']`

Transform

[Groupby](#)

[Pivot, Stack, Unstack](#)

`df.groupby('col').sum(), .apply(), .agg(), .transform(), ...`

`df.pivot(), df.stack(), df.unstack()`

Merge

[Merge](#)

[Join](#)

[Concat](#)

`df1.merge(df2, left_on='pk', right_on='fk', how='left')`

`df1.join(df2, how='inner')`

`pd.concat([df1, df2, df3])`



pandas.pydata.org/pandas-docs/stable/user_guide

Table Of Contents

- What's New in 0.25.0
- Installation
- Getting started
- User Guide
 - IO tools (text, CSV, HDF5, ...)
 - Indexing and selecting data
 - MultiIndex / advanced indexing
 - Merge, join, and concatenate
 - Reshaping and pivot tables
 - Working with text data
 - Working with missing data
 - Categorical data
 - Nullable integer data type
 - Visualization
 - Computational tools
 - Group By: split-apply-combine
 - Time series / date functionality
 - Time deltas
 - Styling
 - Options and settings
 - Enhancing performance
 - Sparse data structures
 - Frequently Asked Questions (FAQ)
 - Cookbook
- Pandas ecosystem
- API reference
- Development
- Release Notes



stackoverflow

shift
+ tab



?pd.



Probability Distributions and Frequentist Statistics

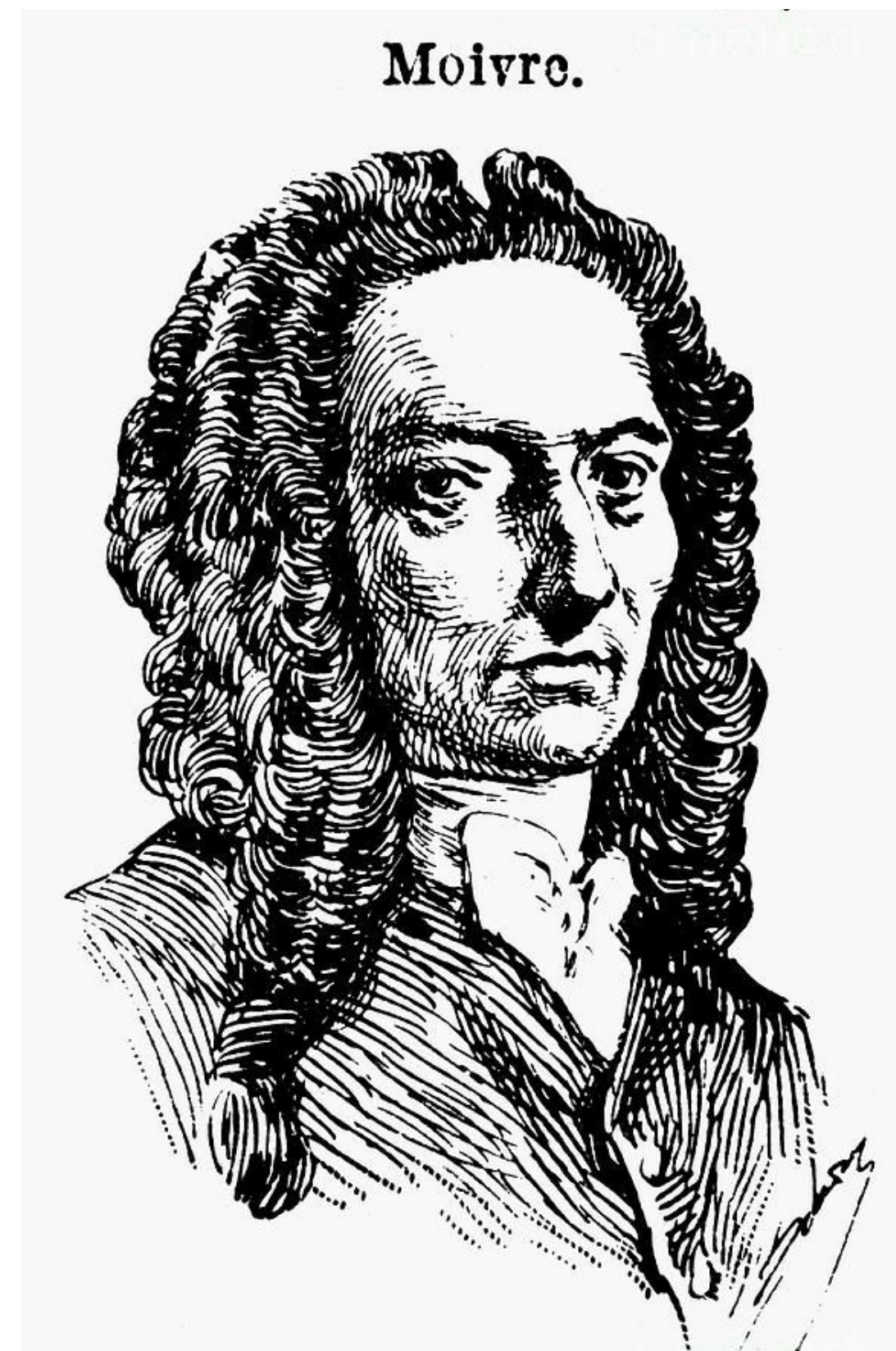
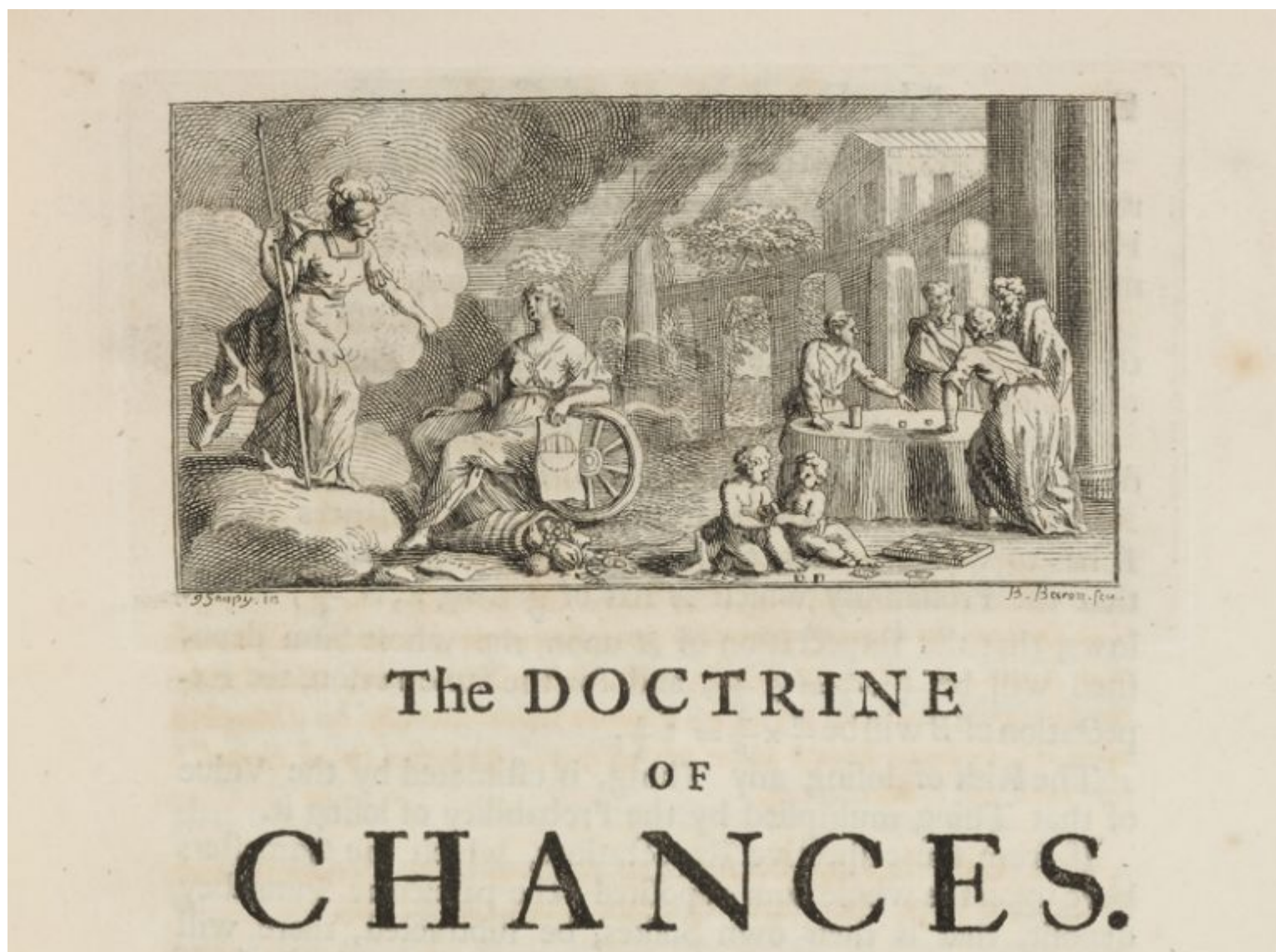




Agenda

- What is probability?
- Multiplication rule
- Independence
- Addition rule
- Mutually Exclusive
- PDF and CDF
- Mean, Median and Mode
- Standard error
- Variance
- The law of large numbers
- The Central Limit Theorem







\$

The frequency theory was originally designed to solve gambling problems





LIKELIHOOD




$$P(\text{event}) =$$

OCCURRENCES EVENT

ALL POSSIBLE OUTCOMES





$$P(\text{throw}=6) =$$

$$\frac{1}{6}$$



$$P(\text{throw not } 6) =$$

$$\frac{5}{6}$$

$$= 1 - P(\text{throw} = 6)$$



The chance
that **both** two
things will
happen

=

The chance
that the **first**
thing will
happen

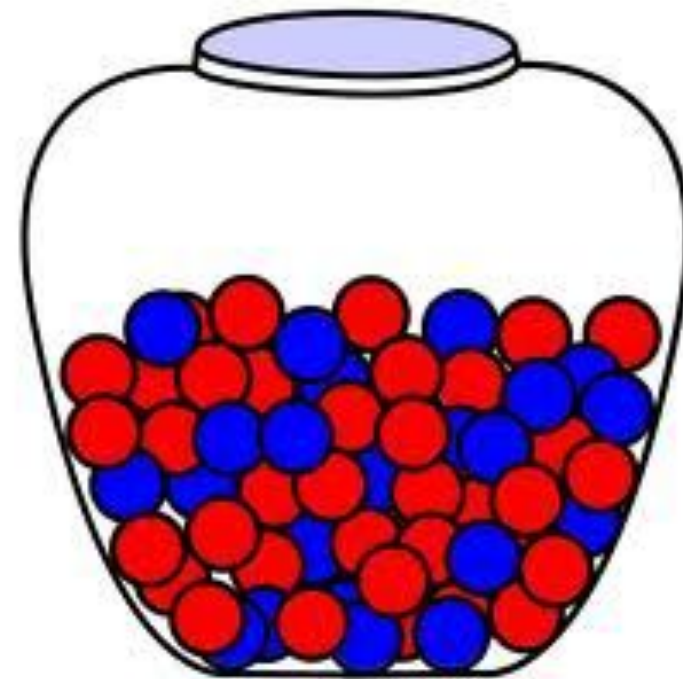
*

The chance
the **second**
will happen

Two things are **independent** if the chances for the second, given the first are the same, no matter how the first one turns out. Otherwise, the two things are **dependent**.

WITH
REPLACEMENT

INDEPENDENT



WITHOUT
REPLACEMENT

DEPENDENT

INDEPENDENT

$$P(A \text{ and } B) = P(A) * P(B)$$

DEPENDENT

$$P(A \text{ and } B) = P(A) * P(B|A)$$

The chance that **both** two things will happen

The chance that the **first** thing will happen

The chance the **second** will happen

What is $P(\text{BBB})$?

WITH
REPLACEMENT

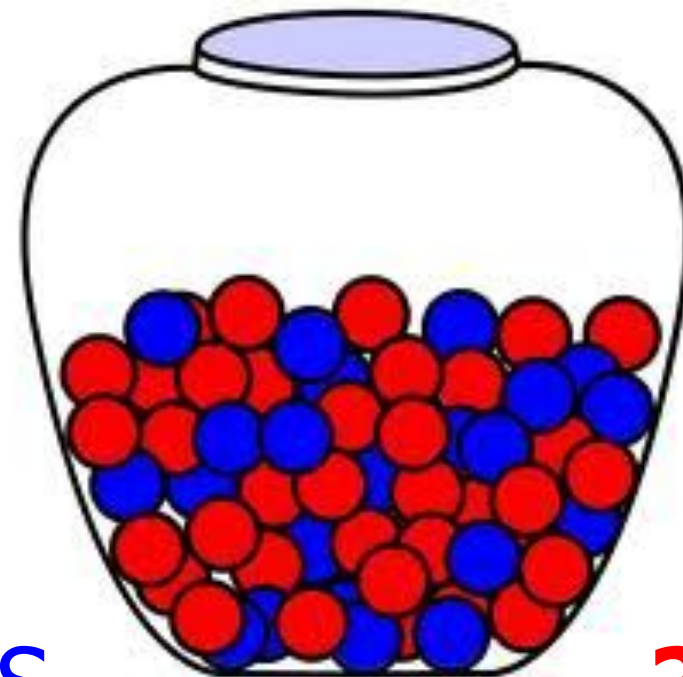
INDEPENDENT

$$30/55 * 30/55 * 30/55$$

WITHOUT
REPLACEMENT

DEPENDENT

$$30/55 * 29/54 * 28/53$$



30 BLUES

25 RED

What is $P(\text{BRR})$?

WITH
REPLACEMENT

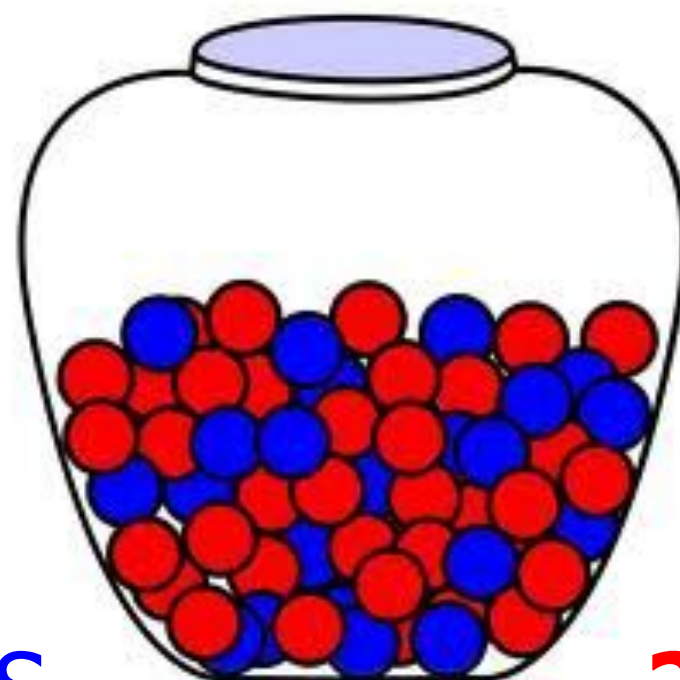
INDEPENDENT

$$30/55 * 25/55 * 25/55$$

WITHOUT
REPLACEMENT

DEPENDENT

$$30/55 * 25/54 * 24/53$$

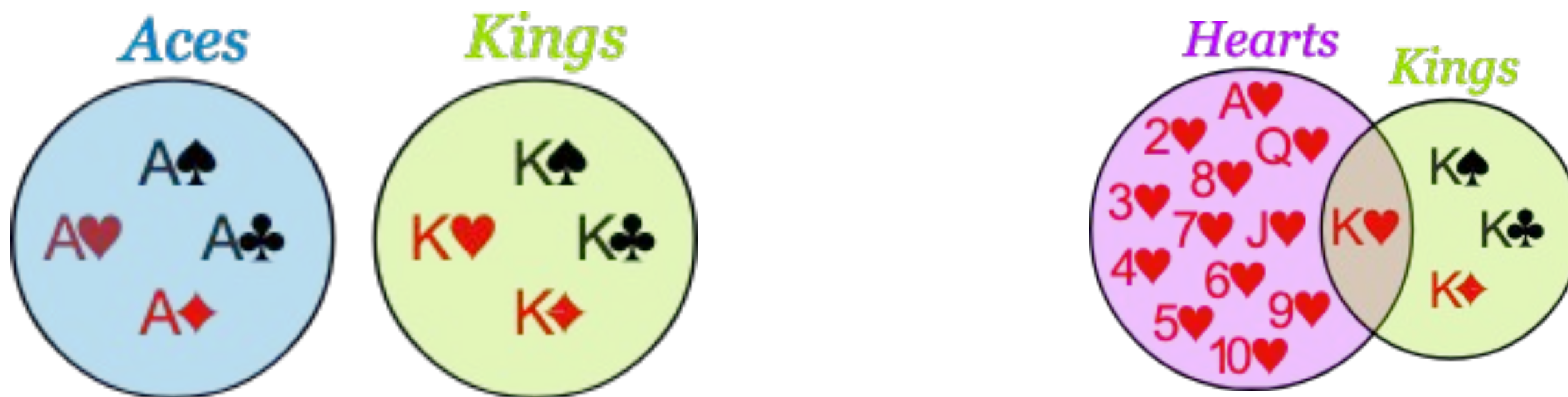


30 BLUES

25 RED



Two things are **mutually exclusive** if the occurrence of one prevents the occurrence of the other.





MUTUALLY EXCLUSIVE

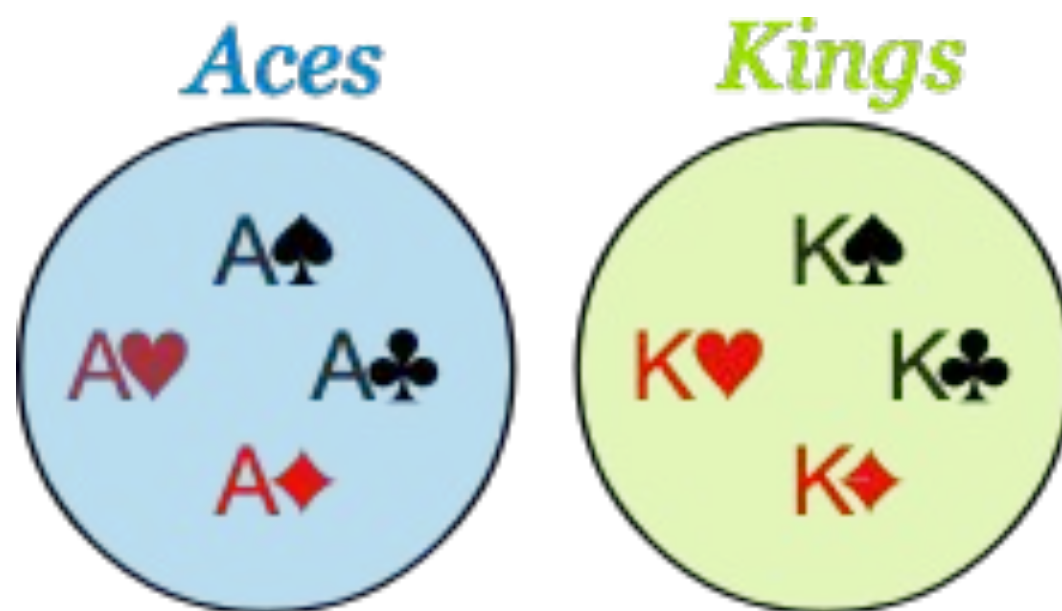
$$P(A \text{ or } B) = P(A) + P(B)$$

NOT MUTUALLY EXCLUSIVE

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

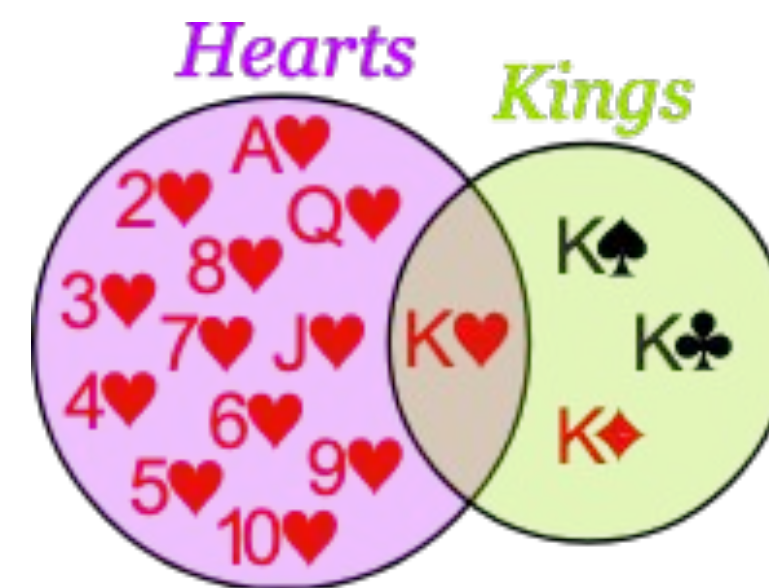
$P(\text{Aces or Kings}) =$

$$4/52 + 4/52 = 8/52$$



$P(\text{Hearts or Kings}) =$

$$4/52 + 13/52 - 1/52 = 16/52$$





$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B \mid A) = \frac{P(A \cap B)}{P(A)}$$



$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B \mid A) = \frac{P(A \cap B)}{P(A)}$$



$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A) P(B \mid A) = P(A \cap B)$$



LIKELIHOOD
the probability of "B"
being TRUE given that "A" is TRUE

PRIOR
the probability of
"A" being TRUE

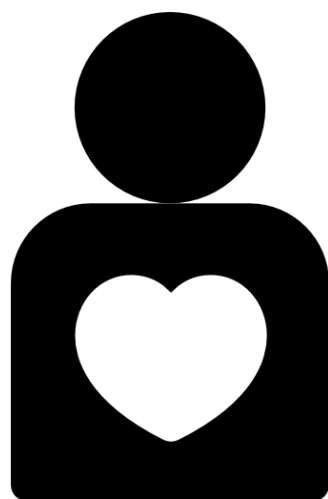
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

POSTERIOR
the probability of "A"
being TRUE given that "B" is TRUE

The probability
of "B" being
TRUE



NON
USERS



99.5%

Test/
User

+

-

+

99

1

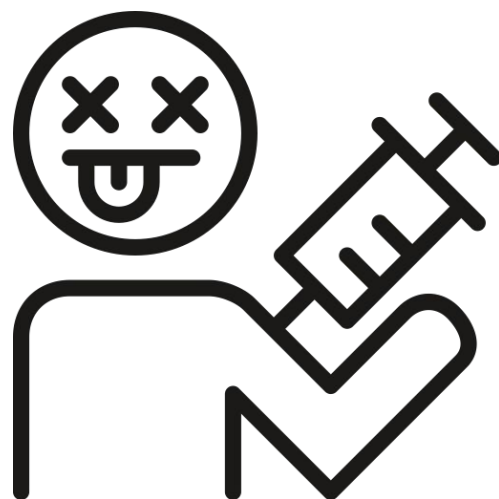
-

1

99

$$P(\text{User} \mid +) = \frac{P(+ \mid \text{User})P(\text{User})}{P(+)}$$

USERS



0.5%

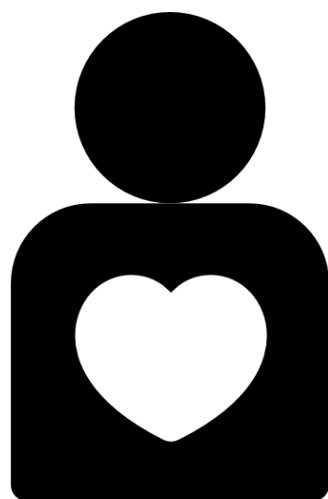
$$P(+) = P(+ \mid \text{User})P(\text{User}) + P(+ \mid \text{NotUser})P(\text{NotUser})$$

$$P(+) = .99 * .005 + .01 * .995$$

$$P(+) = 0.0149$$

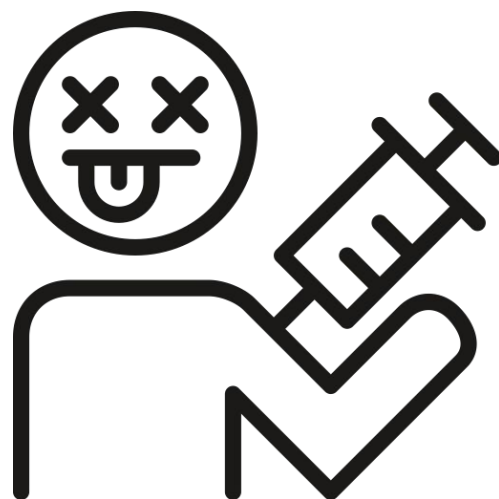


NON
USERS



99.5%

USERS



0.5%

Test/
User

+

-

+

99

1

-

1

99

$$P(\text{User} \mid +) = \frac{P(+ \mid \text{User})P(\text{User})}{P(+)}$$

$$P(\text{User} \mid +) = \frac{.99 * .005}{0.0149}$$

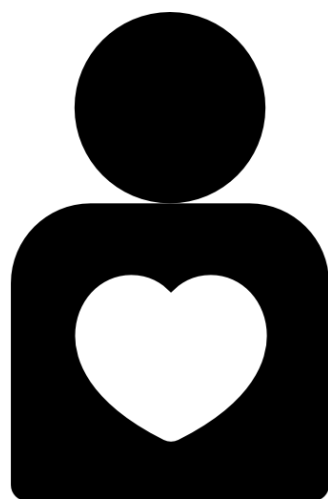


Example Bayes Theorem - Drug Testing

29

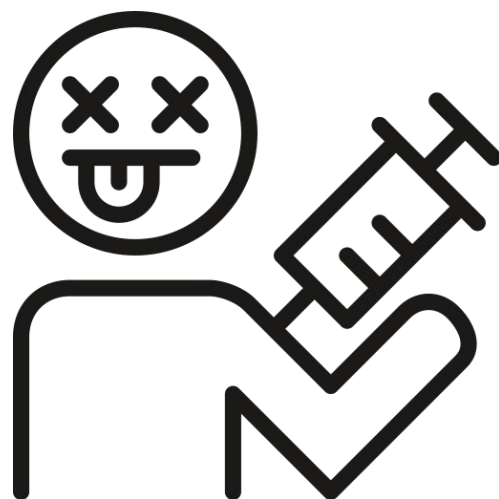


NON
USERS



99.5%

USERS



0.5%

Test/
User

+

-

+

99

1

-

1

99

$$P(\text{User} \mid +) = \frac{P(+ \mid \text{User})P(\text{User})}{P(+)}$$

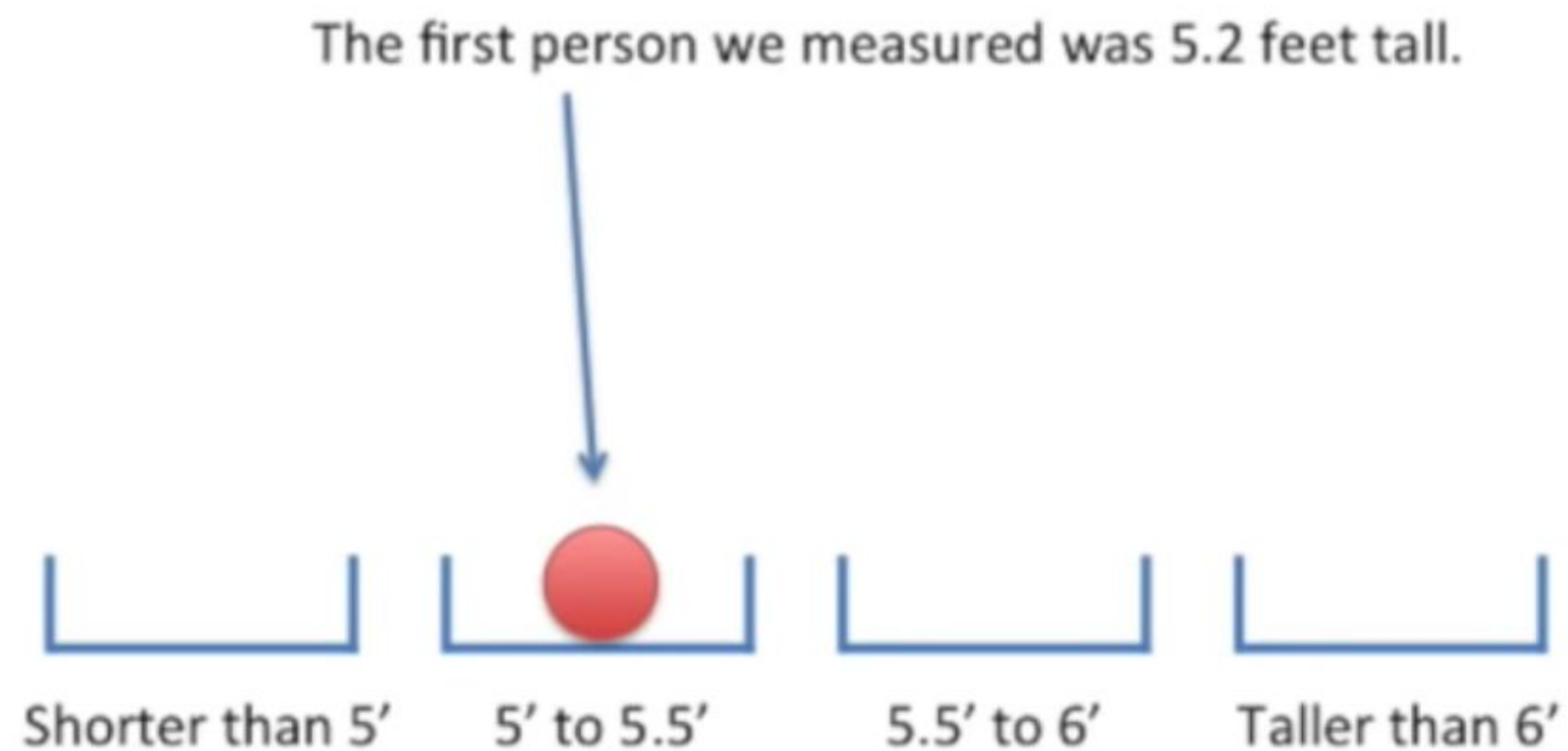
$$P(\text{User} \mid +) \approx 33.2\%$$

Event	Probability
A	$P(A) \in [0, 1]$
not A	$P(A^c) = 1 - P(A)$
A or B	$P(A \cup B) = P(A) + P(B) - P(A \cap B)$ $P(A \cup B) = P(A) + P(B) \quad \text{if A and B are mutually exclusive}$
A and B	$P(A \cap B) = P(A B)P(B) = P(B A)P(A)$ $P(A \cap B) = P(A)P(B) \quad \text{if A and B are independent}$
A given B	$P(A B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B A)P(A)}{P(B)}$



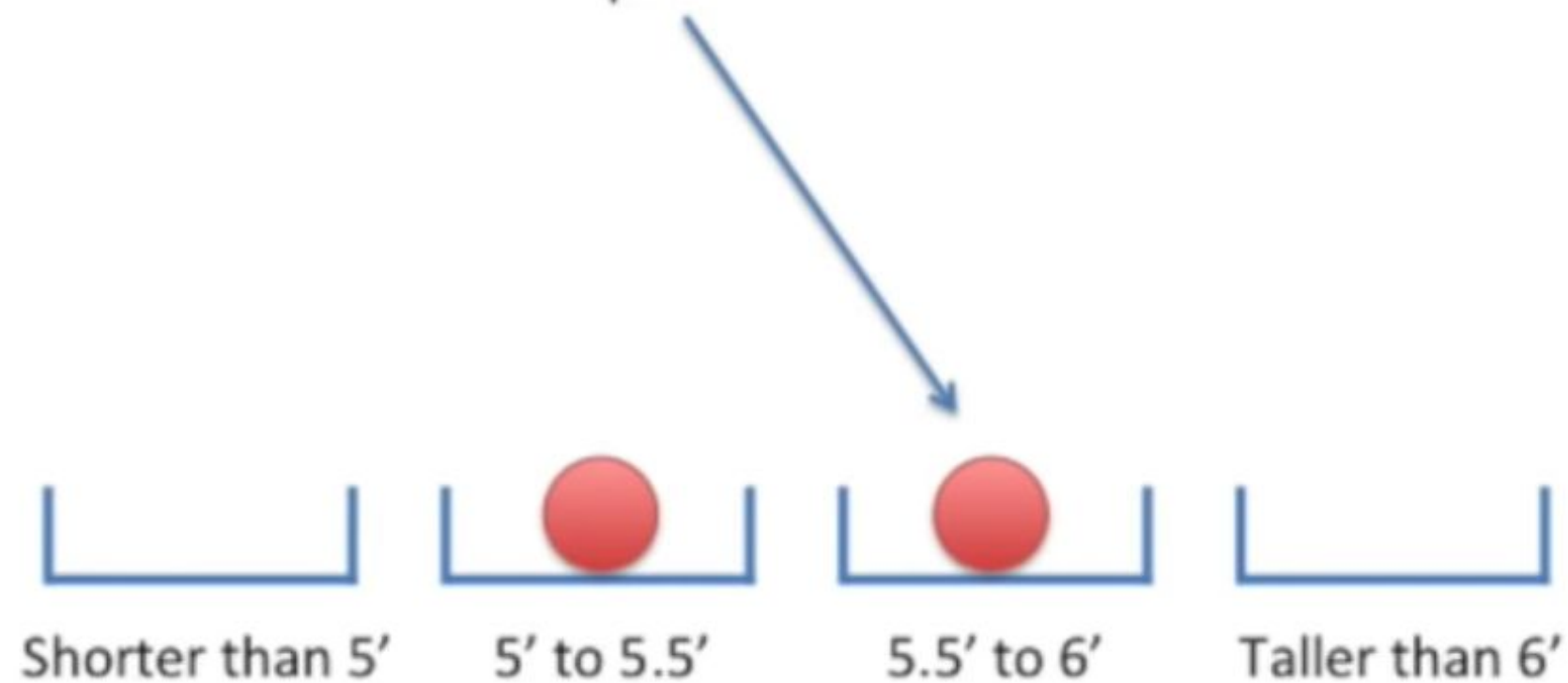
Imagine we measured the height of a lot of people.





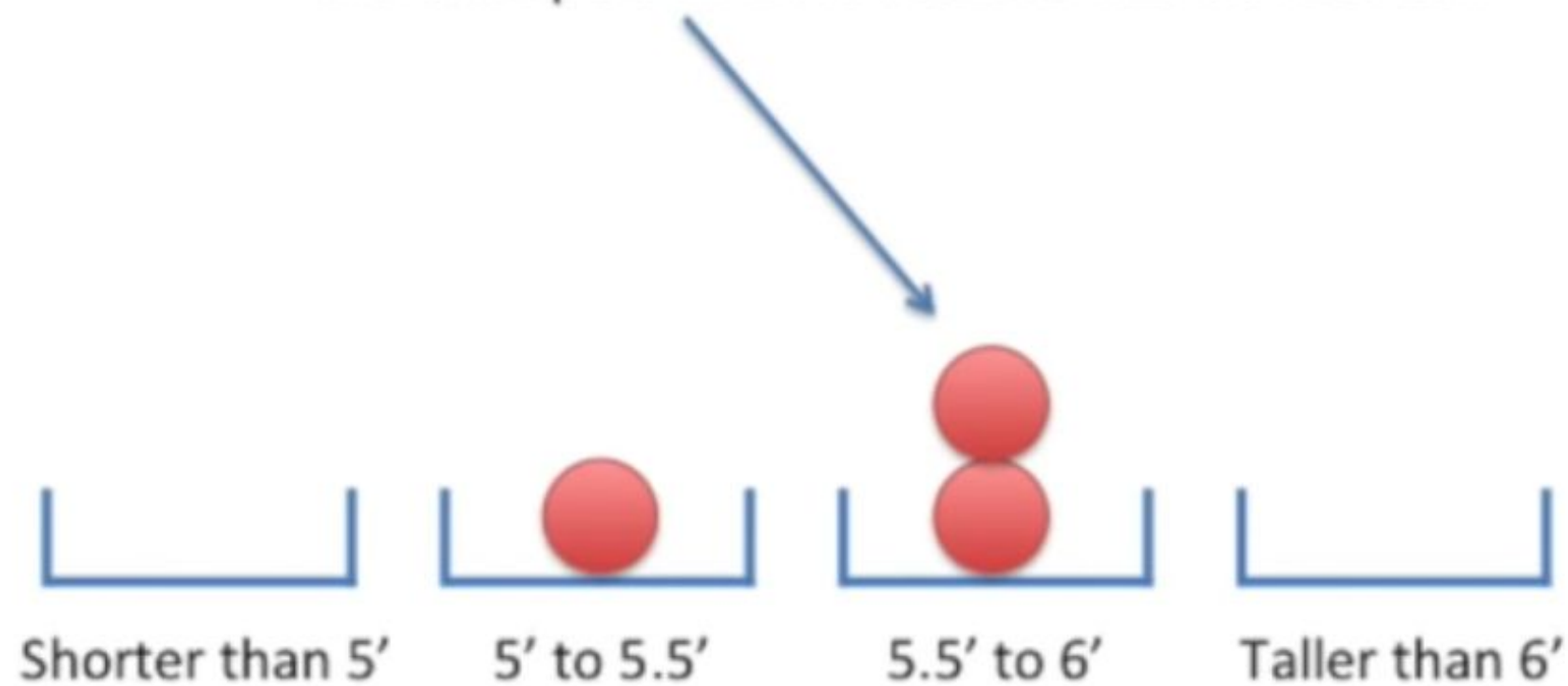


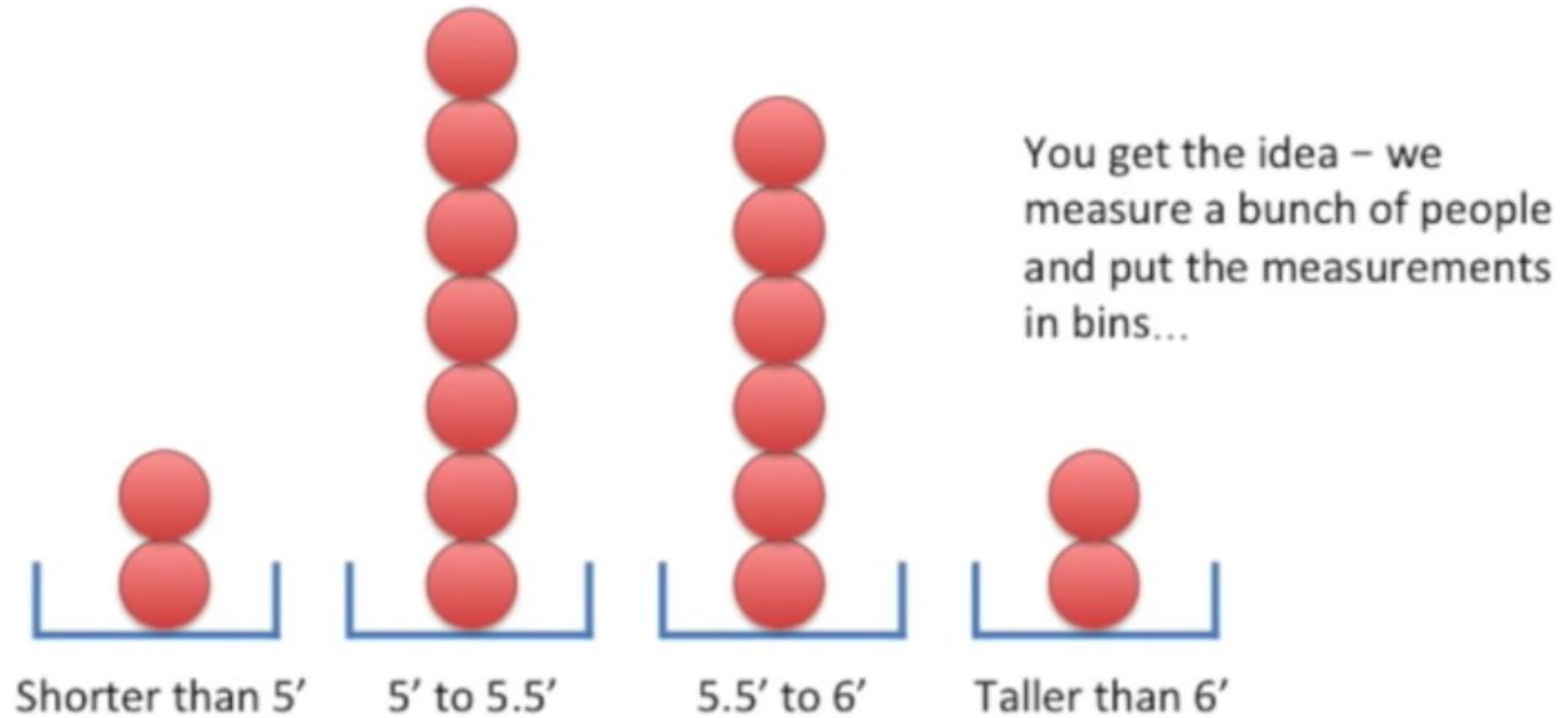
The second person we measured was 5.8 feet tall.

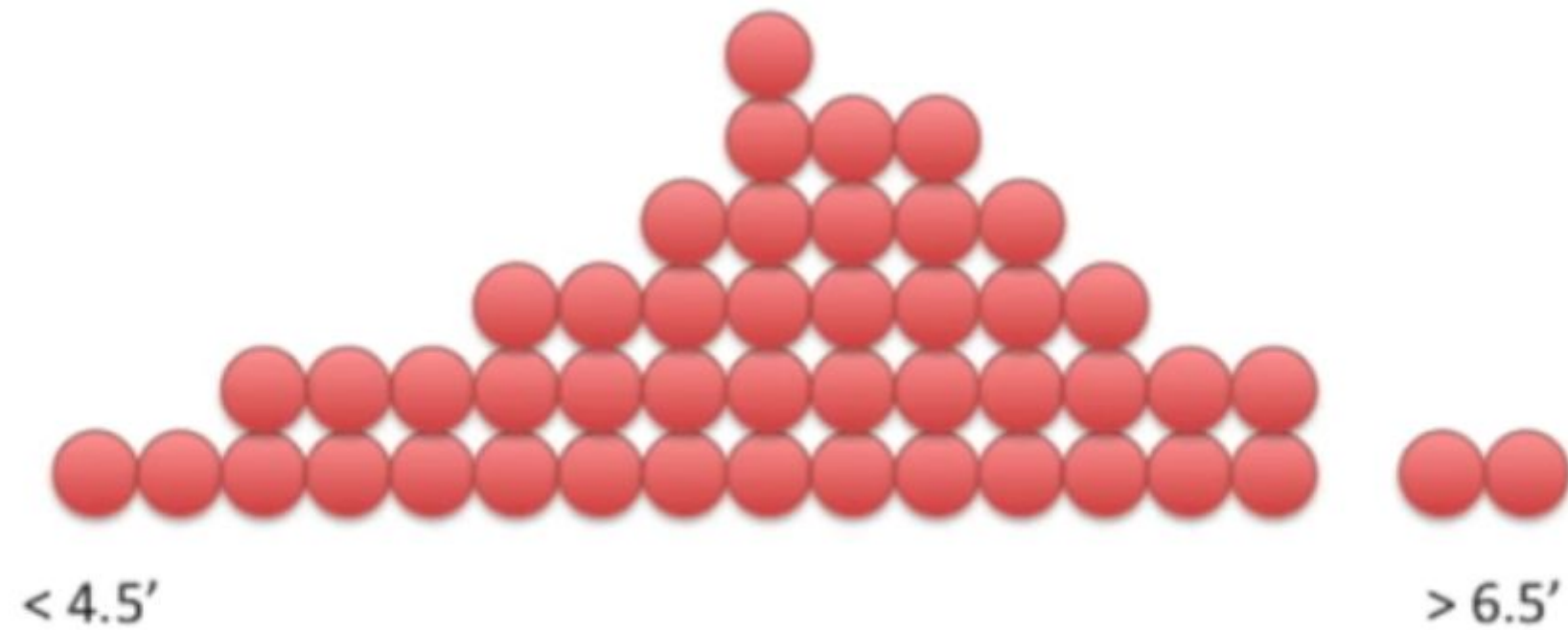




The third person we measured was 5.6 feet tall.

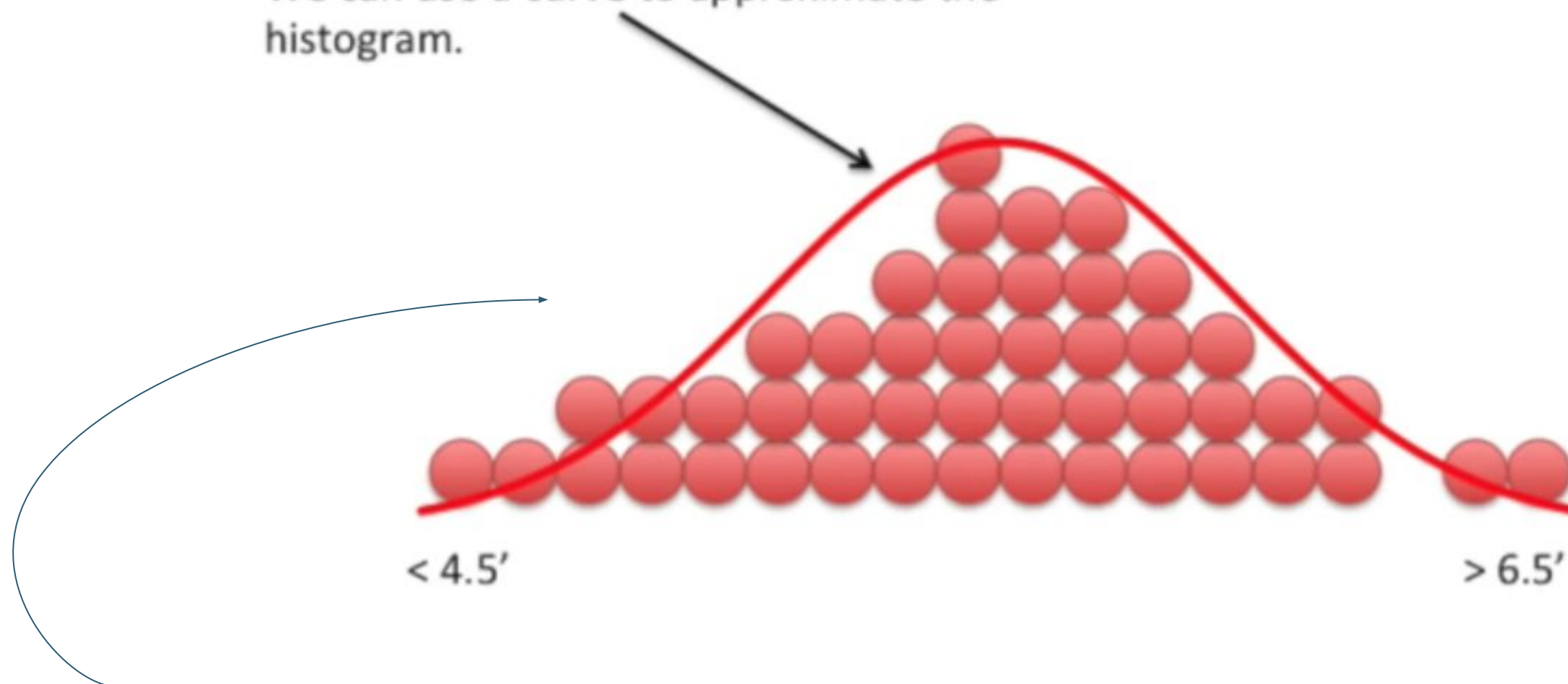




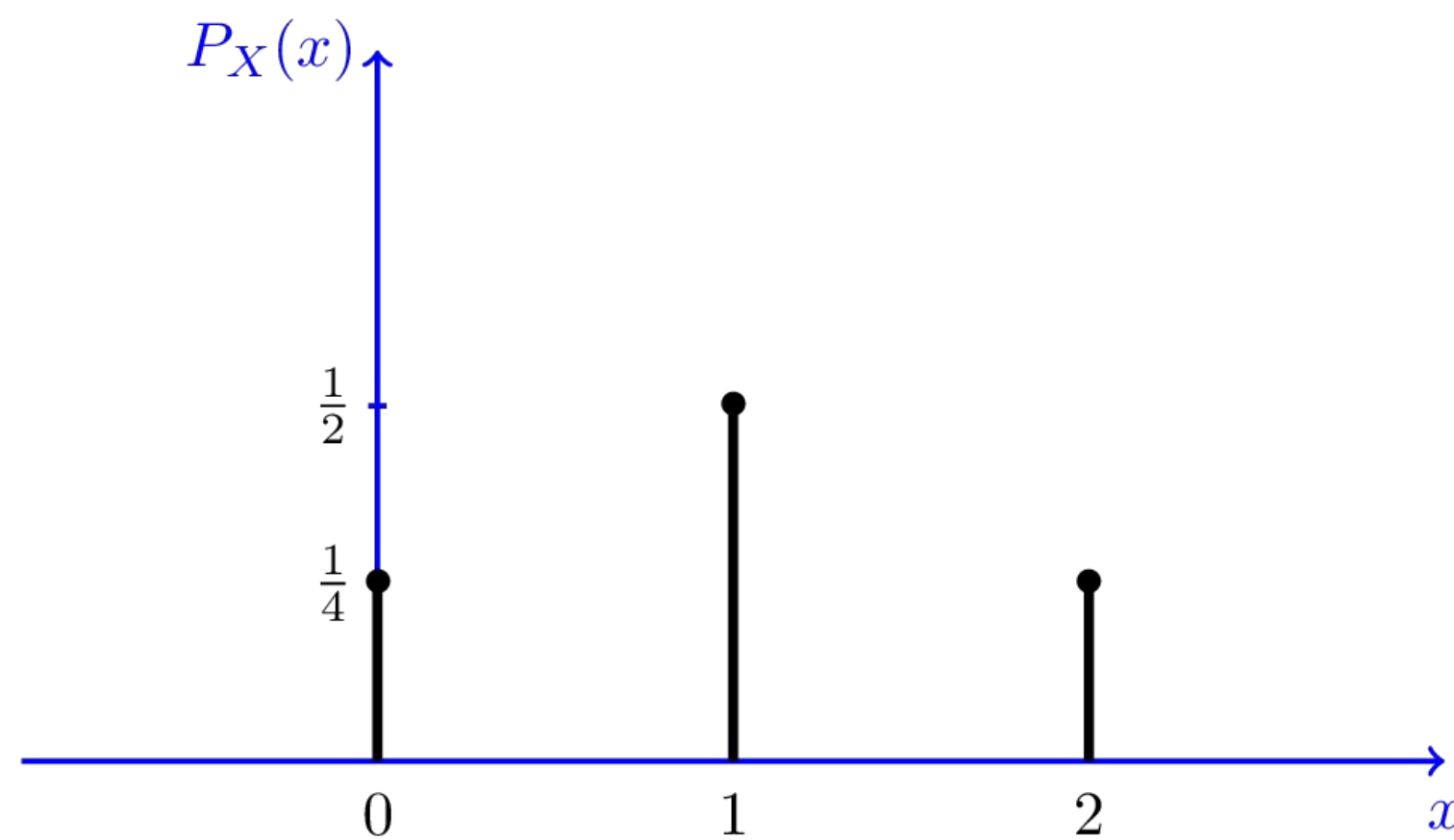
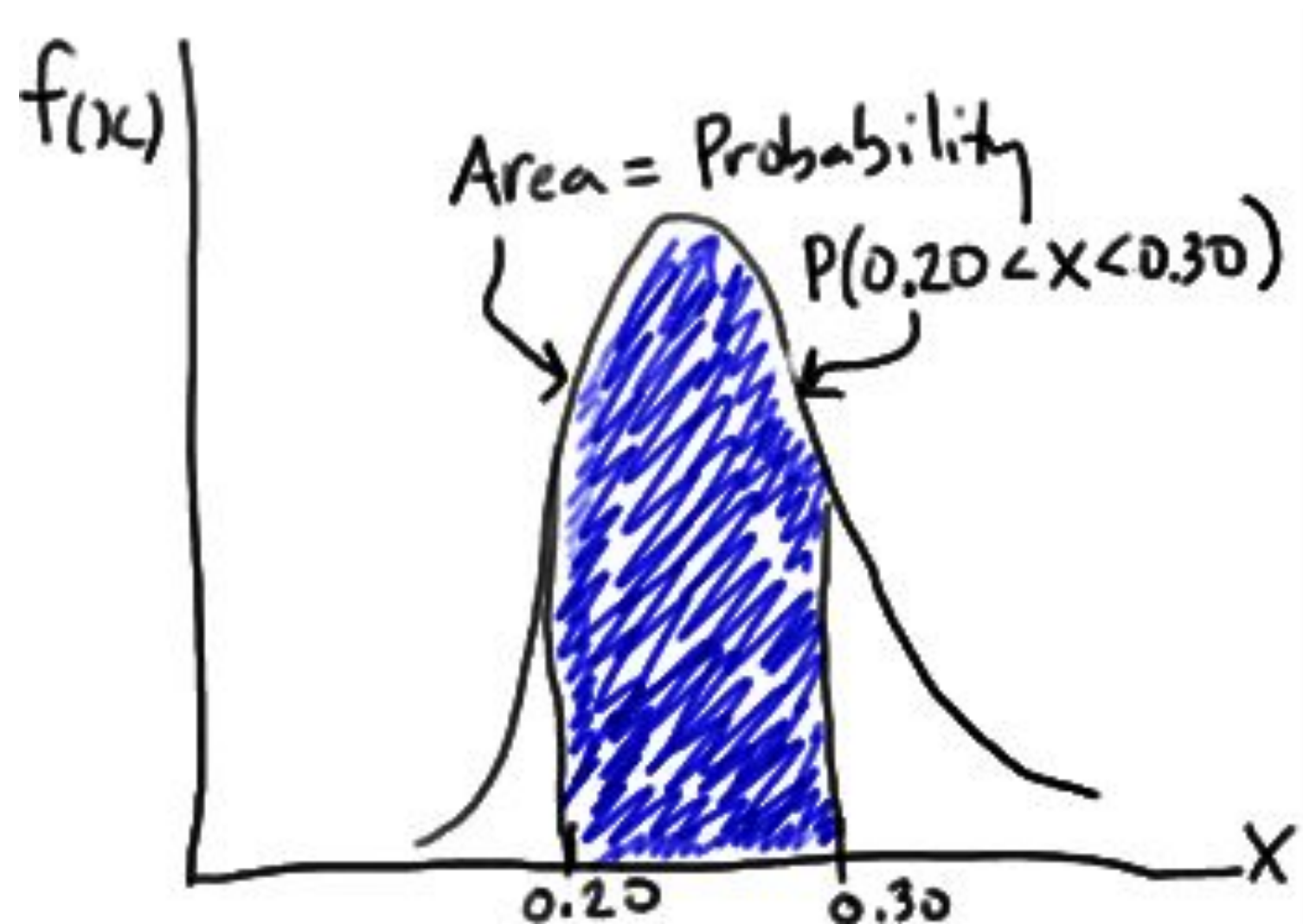


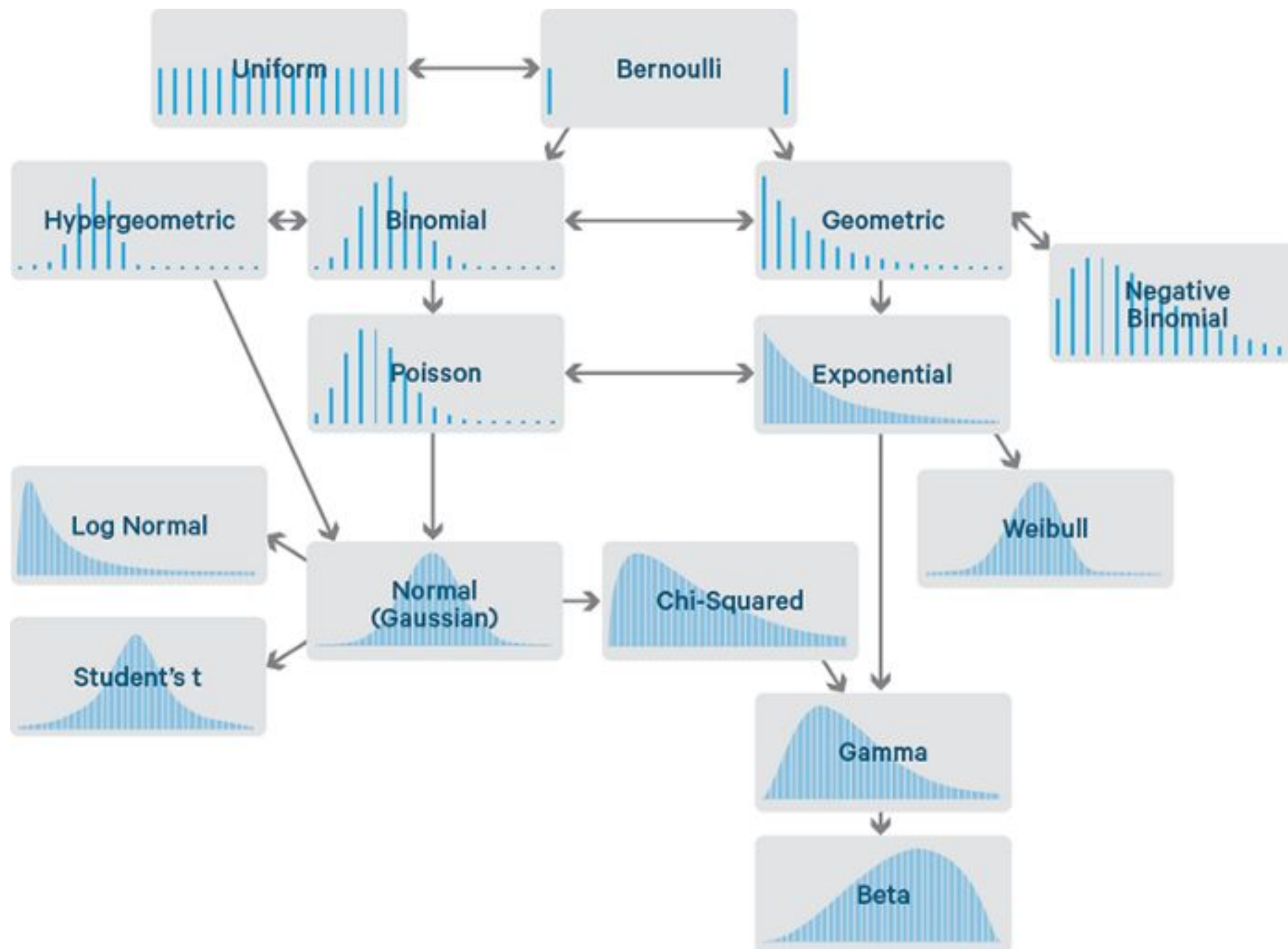
By measuring more people and using smaller bins, we get a more accurate and more precise estimate of how heights are distributed.

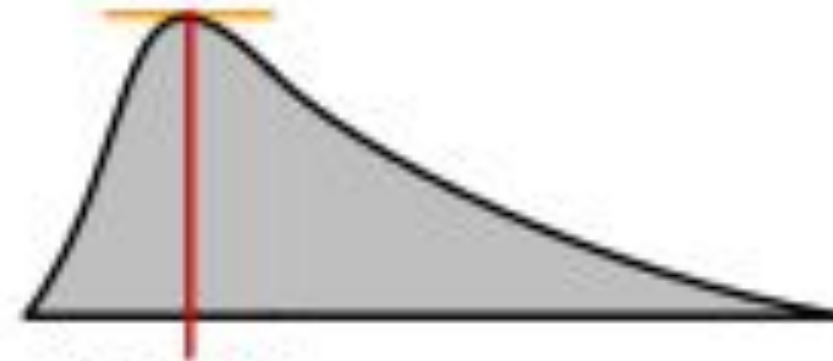
We can use a curve to approximate the histogram.



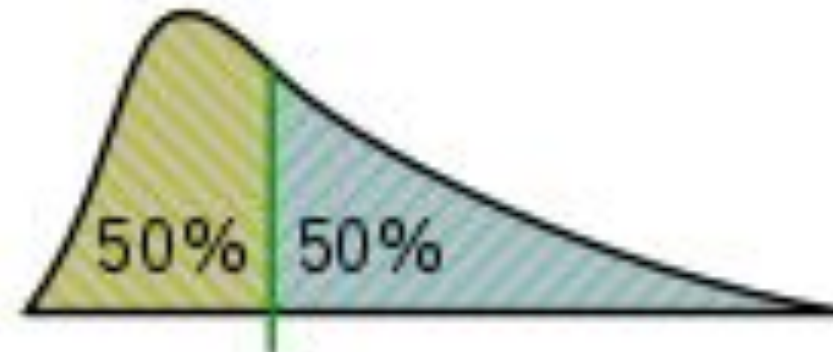
It is a function!



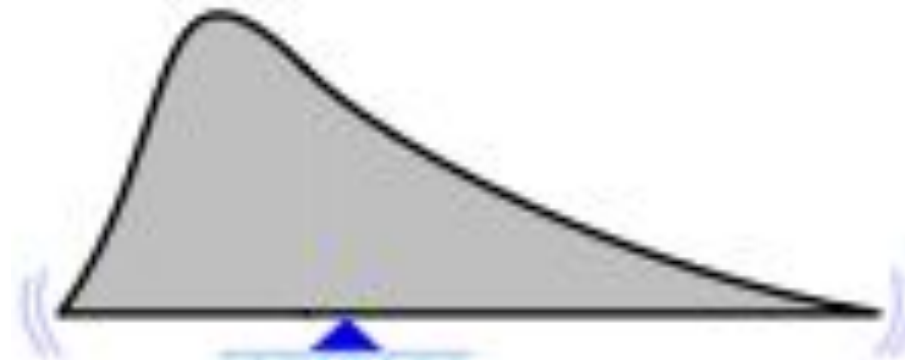




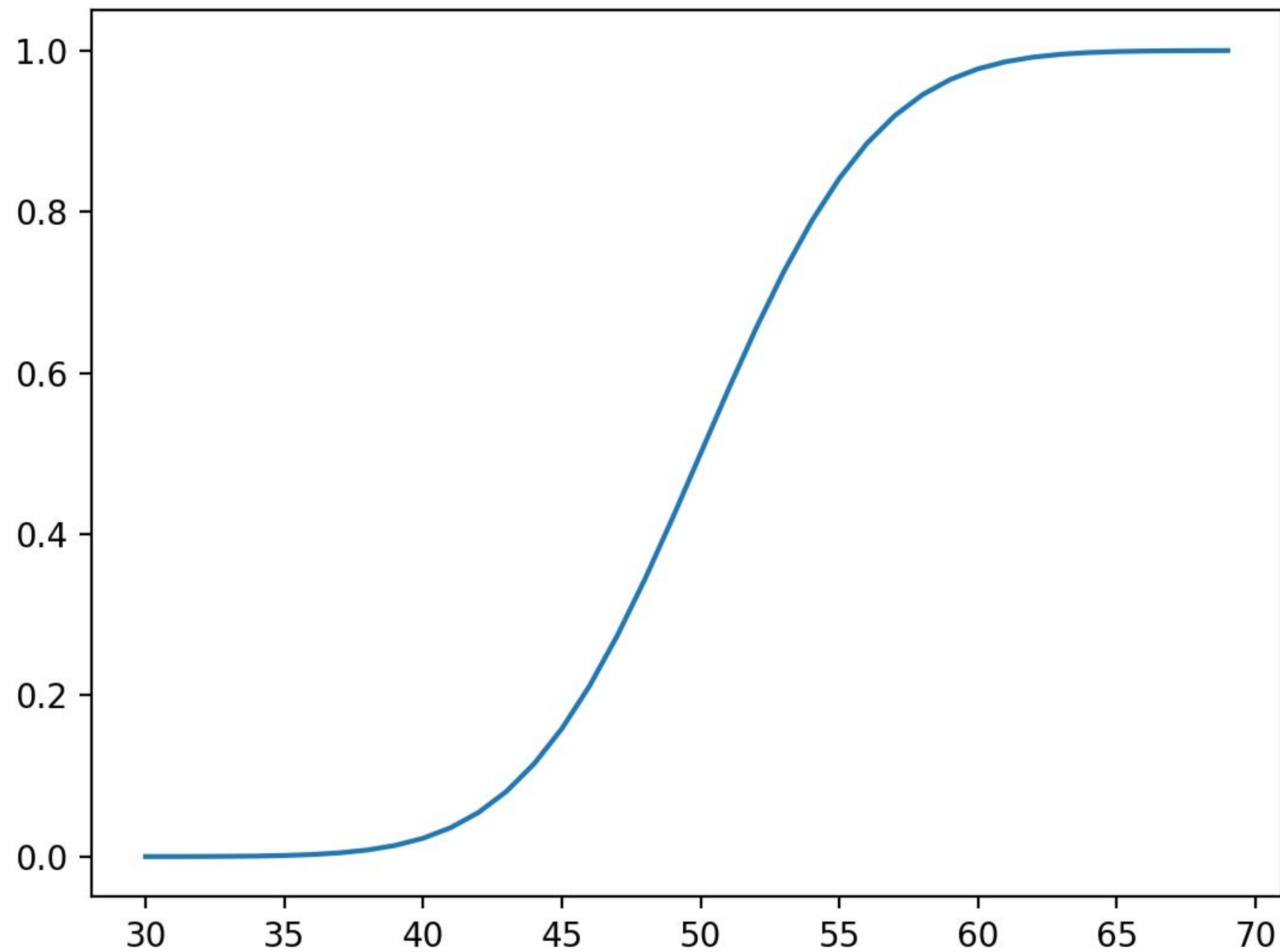
mode



median



mean



**STANDARD
DEVIATION**

$$\sigma = \sqrt{\frac{\sum (X - \bar{X})^2}{n - 1}}$$

VARIANCE

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N}$$



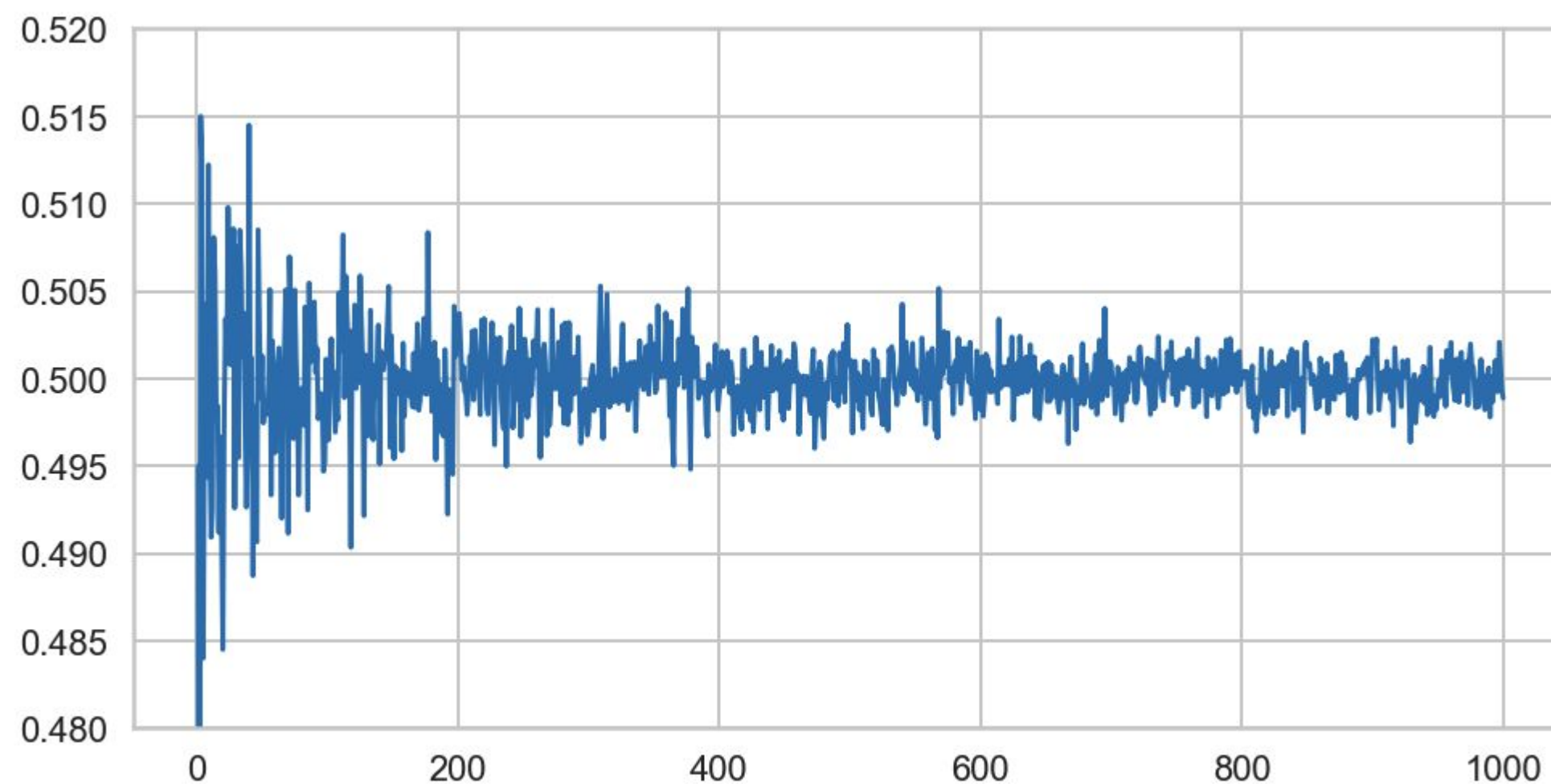
$$SE = \frac{\sigma}{\sqrt{n}}$$

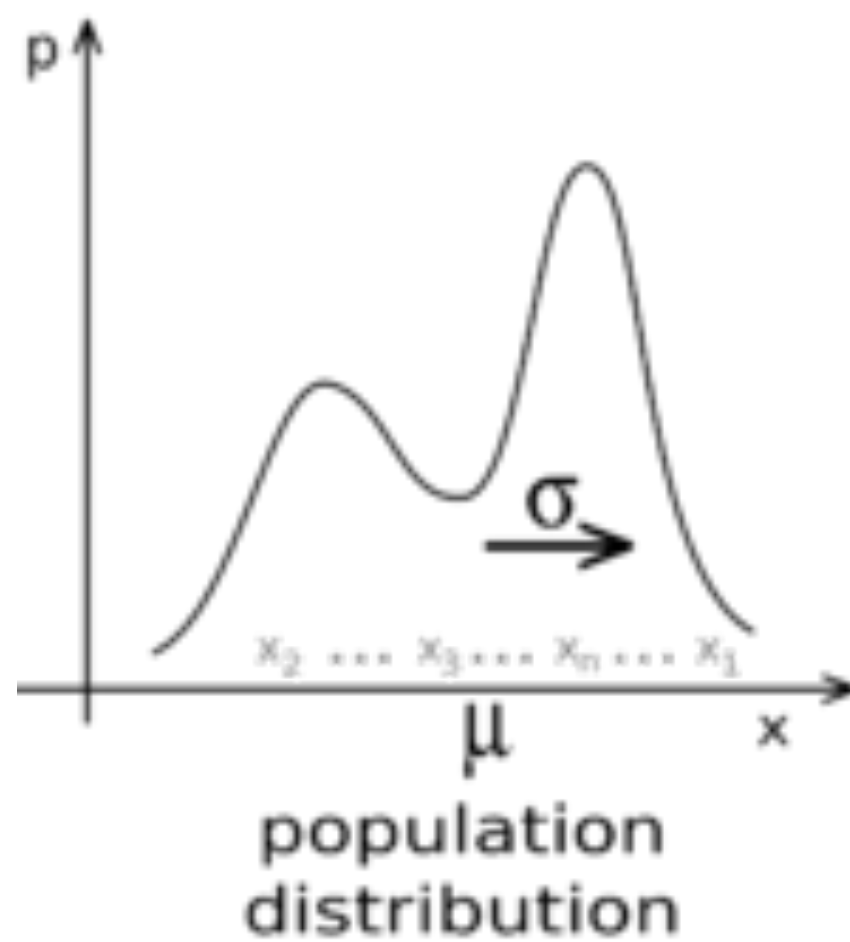
← Standard deviation

← Number of samples



The Law of Large Numbers

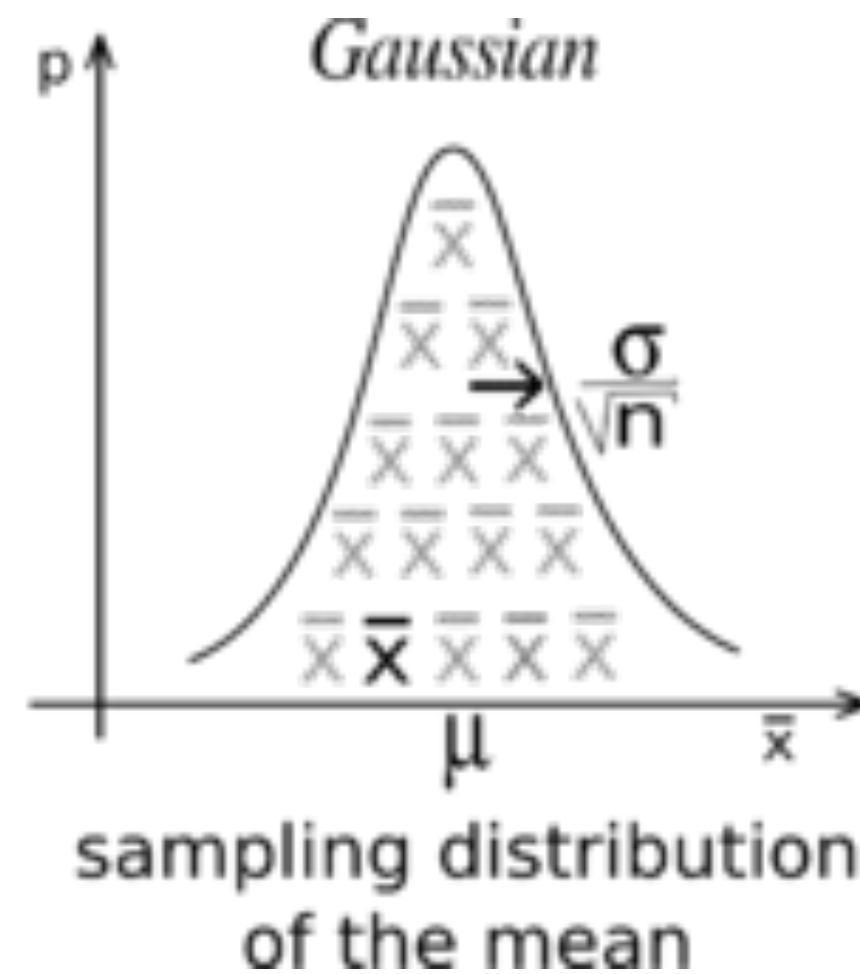




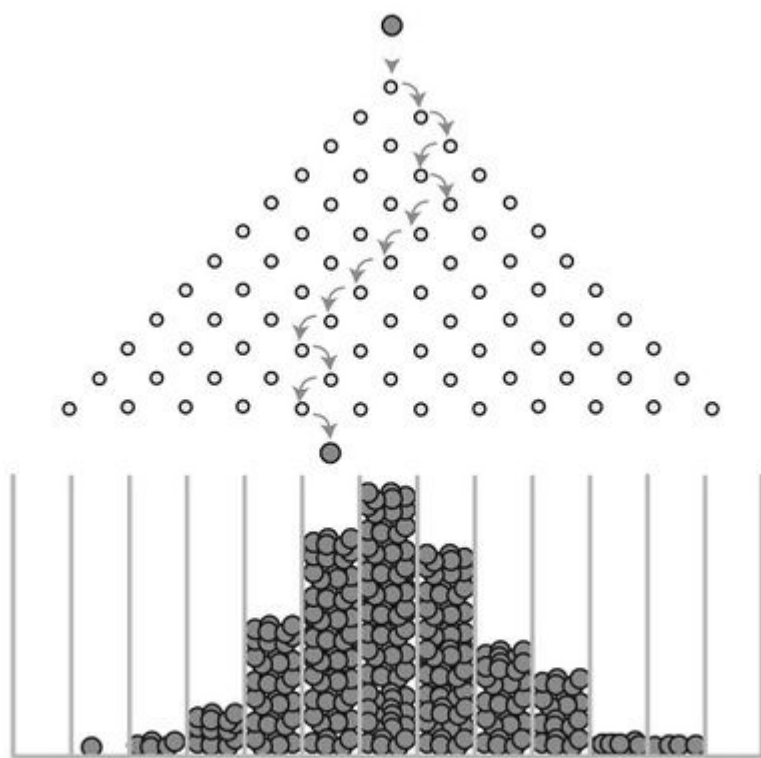
samples
of size n

\bar{x}

\bar{x}



PROOF



46



Follow

Comments v

Statistical Models



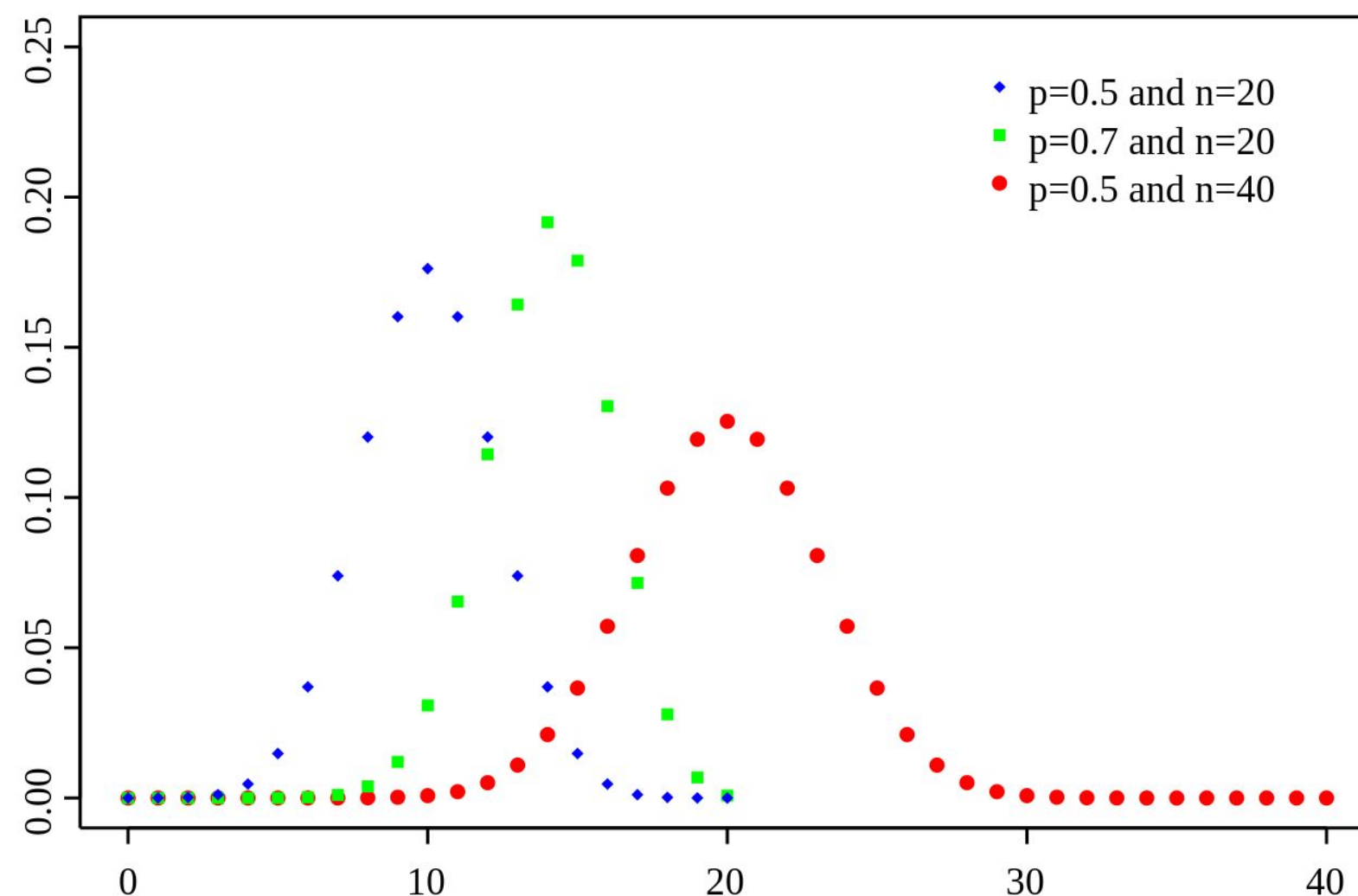


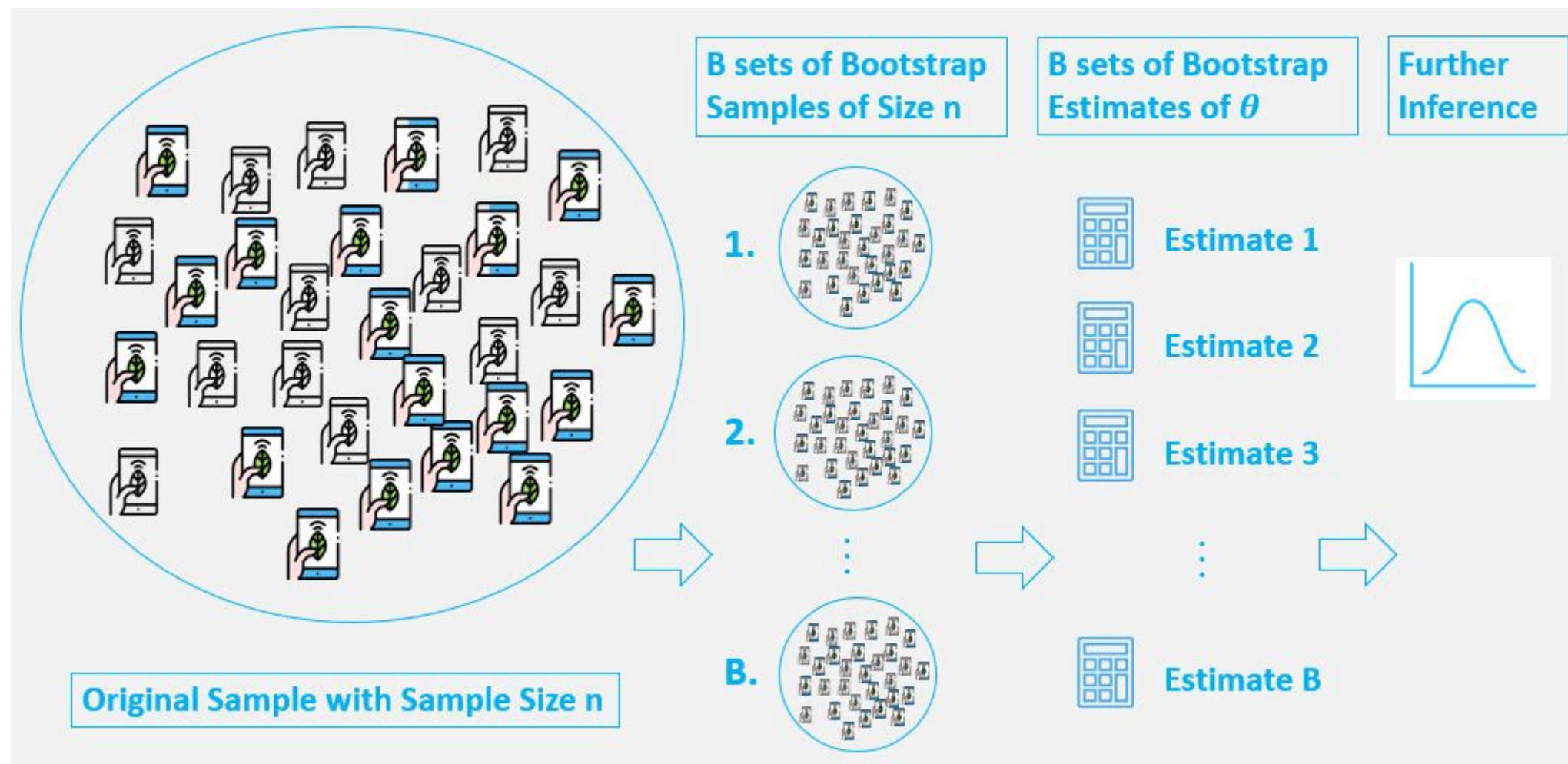
Bernoulli $n=1$

$$f(k, n, p) = \Pr(k; n, p) = \Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

for $k = 0, 1, 2, \dots, n$, where

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$





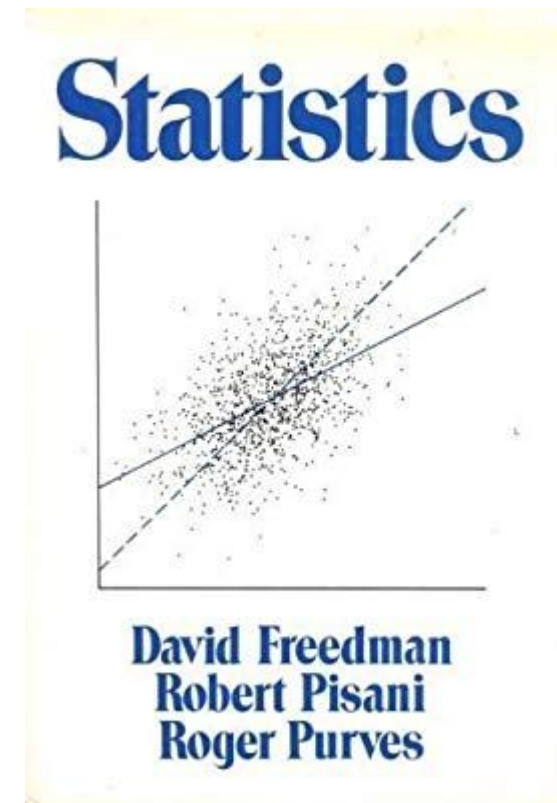


Useful Resources





Stat Quest



THANK YOU



GABRIEL MAGALHÃES

Data Scientist

gabriel.magalhaes@totvs.com.br

**Tecnologia + Conhecimento são nosso DNA.
O sucesso do cliente é o nosso sucesso.
Valorizamos gente boa que é boa gente.**

 [totvs.com](https://www.totvs.com)

 [company/totvs](https://www.linkedin.com/company/totvs)

 [@totvs](https://twitter.com/totvs)

 [fluig.com](https://www.fluig.com)

#SOMOSTOTVERS