# BOOTSTRAP*

HOW TO EVALUATE ACCURACY OF **THINGS**\*\* ABOUT LOTS OF **STUFF**\*\*\*
WHEN ONLY HAVING **NOT SO MUCH STUFF**\*\*\*\*

\* A statistical algorithm extremely useful when only a small set of data is available. It gives a precise idea of the accuracy of the estimate.

\*\* Estimate of statistical parameters
\*\*\* A huge population
\*\*\*\* A small sample

## START

We want to estimate the mean height of all the blueberry-loving people in Milan, but we only know the height of 10 of them. **How do we get there?**
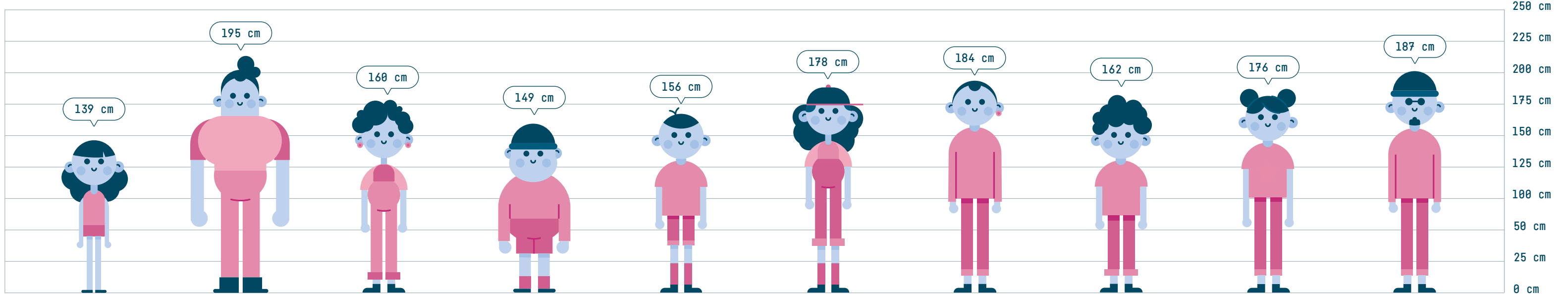
## STEP 01
## SAMPLE

### WE TAKE THE DATA FROM THE SAMPLE

As we all know, blueberry-loving people are very *very* shy, so we only managed to tackle 10 of them for our super-duper important research about the link between height and blueberryness. We'll call them our **sample**.

**The bigger the obtainable sample, the better the estimate**. In our case 10 people will do just fine.

▼ **Law of Large Numbers**
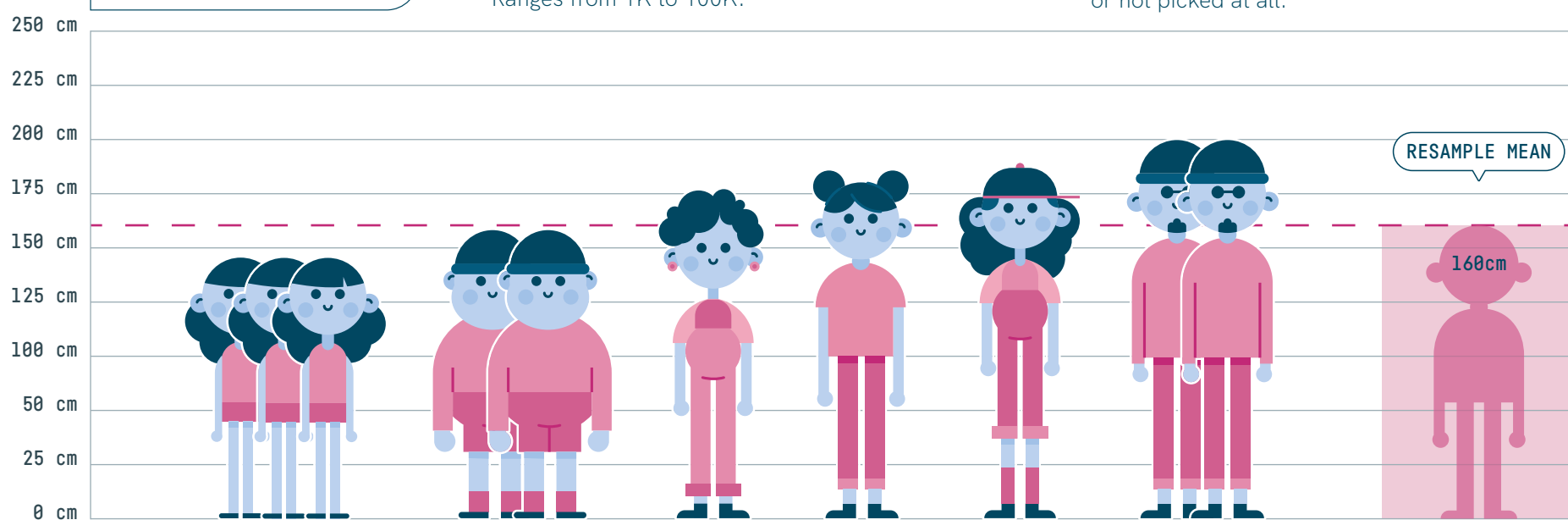As the sample size increases, its ability to give an accurate representation of the whole population increases too.

STARTING SAMPLE

139 cm · 195 cm · 160 cm · 149 cm · 156 cm · 178 cm · 184 cm · 162 cm · 176 cm · 187 cm

250 cm / 225 cm / 200 cm / 175 cm / 150 cm / 125 cm / 100 cm / 50 cm / 25 cm / 0 cm

## STEP 02
## RESAMPLE

RESAMPLE #0001

◄ **Recommended number of resamples**
Ranges from 1K to 100K.

▼ **Resampling with replacement**
Replacement means that repetition is allowed, each value can be picked more than once or not picked at all.

250 cm / 225 cm / 200 cm / 175 cm / 150 cm / 125 cm / 100 cm / 50 cm / 25 cm / 0 cm

RESAMPLE MEAN

160cm

### WE RESAMPLE OUR DATA AND CALCULATE THE MEAN *MANY MANY MANY* TIMES

We could just calculate the mean of the starting sample, but **we wouldn't have any information about the accuracy of the estimate**.
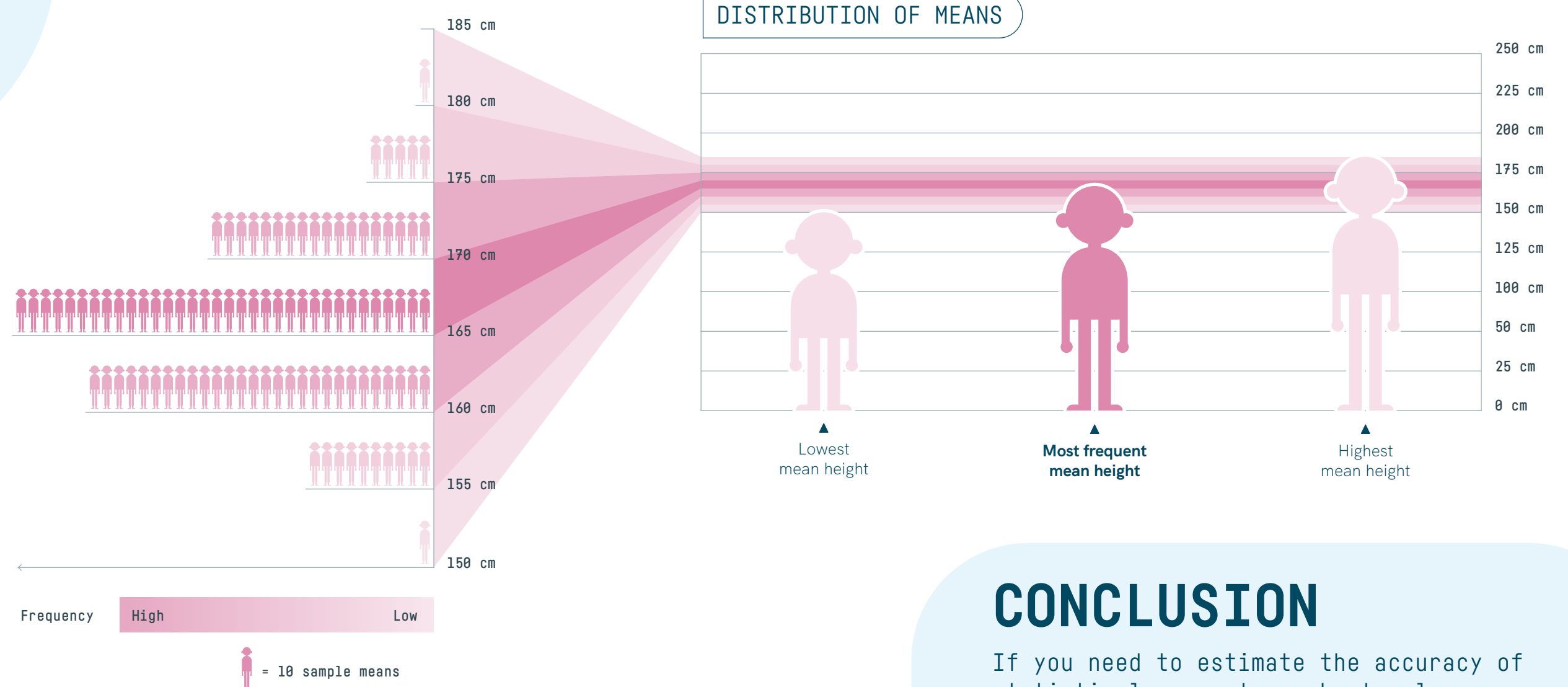
Instead, we **resample with replacement** many (*many!*) times: each resample is made of 10 height values, randomly picked from the original 10.

We calculate the **mean height value of each resample**, and we store them away for later. (in case we're hungry).

## STEP 03
## DISTRIBUTION

### WE DRAW THE DISTRIBUTION OF THE MEAN HEIGHT VALUES

The distribution shows the **estimated mean height** and the **accuracy of the estimate**. The frequency is expected to be higher for values near the **real population mean**, while lower for further values.

185 cm / 180 cm / 175 cm / 170 cm / 165 cm / 160 cm / 155 cm / 150 cm

DISTRIBUTION OF MEANS

250 cm / 225 cm / 200 cm / 175 cm / 150 cm / 125 cm / 100 cm / 50 cm / 25 cm / 0 cm

▲ Lowest mean height

▲ **Most frequent mean height**

▲ Highest mean height

Frequency   High — Low
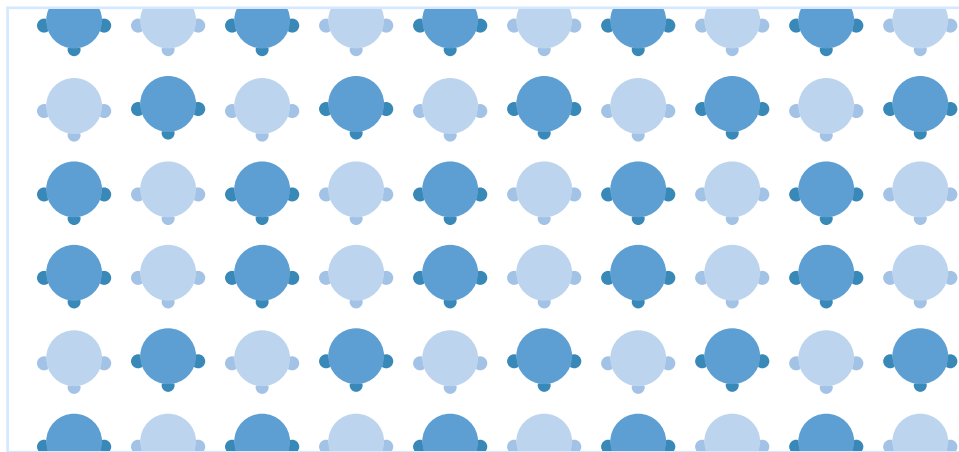
🧍 = 10 sample means

## CONCLUSION

If you need to estimate the accuracy of statistical parameters about a large population, but only have a small sample, the bootstrap algorithm is the way to go.

**CAUTION!** overconsumption of blueberries might result in your skin turning blue.

VISUAL EXPLANATIONS OF STATISTICAL METHODS

Bootstrap

AUTHORS

Daniele Dell'Orto
Martina Francella
Octavian Husoschi
Martina Melillo
Matteo Pini

Alessandro Quets
Shan Huang

FACULTY

Michele Mauri
Ángeles Briones
Gabriele Colombo
Simone Vantini
Salvatore Zingale

TEACHING ASSISTANTS

Elena Aversa
Andrea Benedetti
Tommaso Elli
Beatrice Gobbo
Anna Riboldi

RESAMPLE #N — RESAMPLE MEAN 168cm

[...]

RESAMPLE #0019 — RESAMPLE MEAN 178cm

RESAMPLE #0018 — RESAMPLE MEAN 176cm

RESAMPLE #0017 — RESAMPLE MEAN 168cm

RESAMPLE #0016 — RESAMPLE MEAN 171cm

RESAMPLE #0015 — RESAMPLE MEAN 154cm

RESAMPLE #0014 — RESAMPLE MEAN 164cm

RESAMPLE #0013 — RESAMPLE MEAN 169cm

RESAMPLE #0012 — RESAMPLE MEAN 166cm

RESAMPLE #0011 — RESAMPLE MEAN 166cm

RESAMPLE #0010 — RESAMPLE MEAN 167cm

RESAMPLE #0009 — RESAMPLE MEAN 162cm

RESAMPLE #0008 — RESAMPLE MEAN 156cm

RESAMPLE #0007 — RESAMPLE MEAN 174cm

RESAMPLE #0006 — RESAMPLE MEAN 172cm

RESAMPLE #0005 — RESAMPLE MEAN 173cm

RESAMPLE #0004 — RESAMPLE MEAN 172cm

RESAMPLE #0003 — RESAMPLE MEAN 157cm

RESAMPLE #0002 — RESAMPLE MEAN 170cm

POPULATION MEAN

169cm

## HOW ACCURATE IS IT THO?

In an ideal world, in which we have all the height values of the entire blueberry-loving population, we could just calculate the real mean height of the population.

We can compare the real mean height with the estimated distribution to verify if the algorithm did its job. (if it did, give it a cookie or something).

PULL