# The Standard Error

## Probability and Statistics for Data Science

Carlos Fernandez-Granda

These slides are based on the book Probability and Statistics for Data Science by Carlos Fernandez-Granda, available for purchase here. A free preprint, videos, code, slides and solutions to exercises are available at https://www.ps4ds.net

# Estimation of population parameters

Simple idea: Choose a random subset of the population

# Estimating a population mean

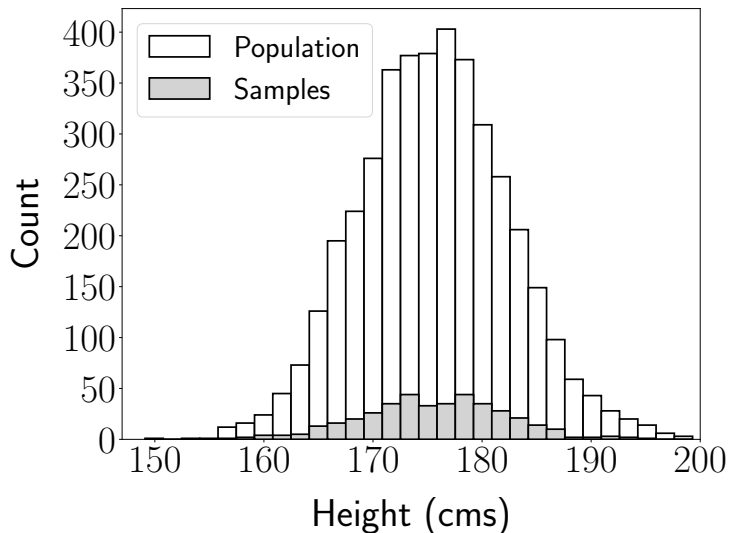Controlled scenario: True population with $N := 4{,}082$ individuals

Heights: $h_1$, $h_2$, $\ldots$, $h_N$

Population mean:

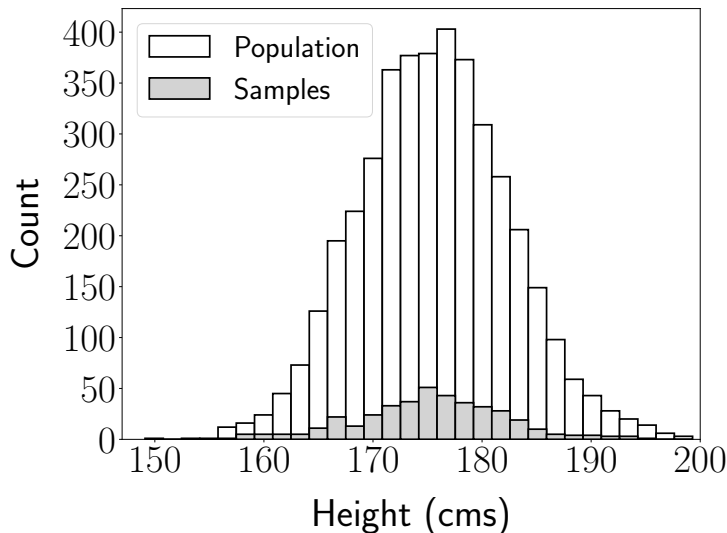$$\mu_{\text{pop}} := \frac{1}{N} \sum_{i=1}^{N} h_i = 175.6$$

# 400 random samples

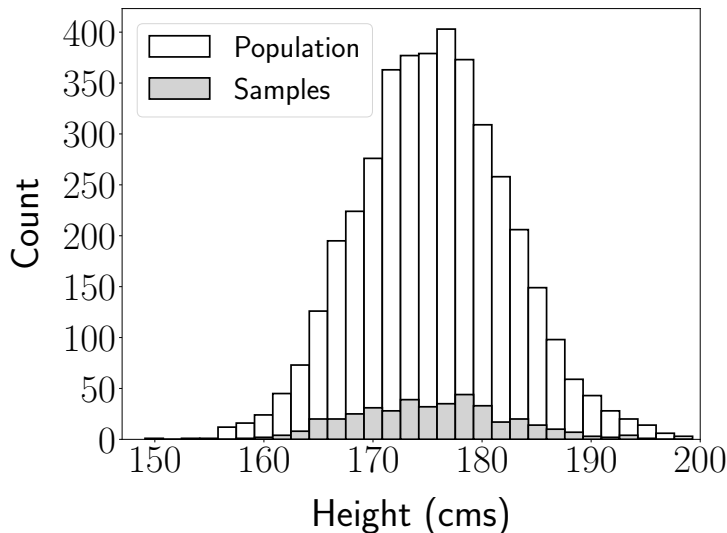Sample mean = 175.5 ($\mu_{\text{pop}} = 175.6$)

# 400 random samples

Sample mean = 175.2 ($\mu_{\text{pop}}$ = 175.6)

# 400 random samples

Sample mean $= 176.1$ ($\mu_{\text{pop}} = 175.6$)

# Random sampling

Data: $a_1, a_2, \ldots, a_N$

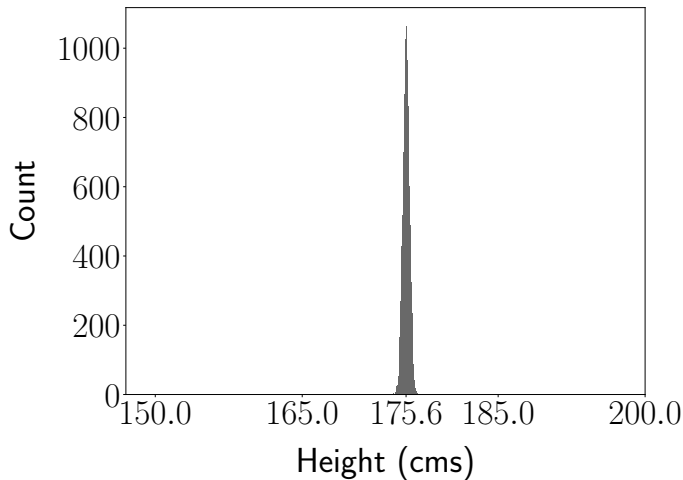Random samples: $\tilde{x}_1, \tilde{x}_2, \ldots, \tilde{x}_n$

Each $\tilde{x}_i$ is selected independently and uniformly at random with replacement

Samples are independent identically distributed (i.i.d.) random variables with pmf

$$p_{\tilde{x}_j}(a_i) = \mathrm{P}(\tilde{x}_j = a_i) = \frac{1}{N}, \qquad 1 \le i \le N,\ 1 \le j \le n$$

# Sample means of 10,000 subsets of size 400

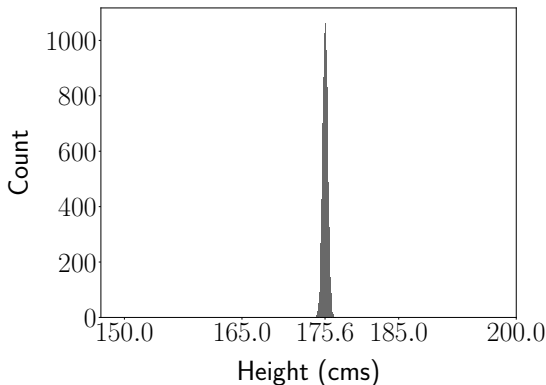Sample mean has to be analyzed probabilistically

# Sample mean is unbiased

Modeled as a random variable

$$\tilde{m} := \frac{1}{n} \sum_{i=1}^{n} \tilde{x}_i$$

$$\mathrm{E}\left[\tilde{m}\right] = \mu_{\mathsf{pop}}$$

# Standard error

Random measurements: $\tilde{x}_1$, $\tilde{x}_2$, ..., $\tilde{x}_n$

Deterministic parameter of interest: $\gamma$

Unbiased estimator: $h(\tilde{x}_1, \ldots, \tilde{x}_n)$

The standard error of the estimator is its standard deviation

$$\text{se}\left[h(\tilde{x}_1, \ldots, \tilde{x}_n)\right] := \sqrt{\text{Var}\left[h(\tilde{x}_1, \ldots, \tilde{x}_n)\right]}$$

# Standard error

Since the estimator is unbiased $\mathrm{E}\left[h(\tilde{x}_1, \ldots, \tilde{x}_n)\right] = \gamma$

$$
\begin{aligned}
\mathrm{se}\left[h(\tilde{x}_1, \ldots, \tilde{x}_n)\right] &:= \sqrt{\mathrm{Var}\left[h(\tilde{x}_1, \ldots, \tilde{x}_n)\right]} \\
&= \sqrt{\mathrm{E}\left[\left(h(\tilde{x}_1, \ldots, \tilde{x}_n) - \mathrm{E}\left[h(\tilde{x}_1, \ldots, \tilde{x}_n)\right]\right)^2\right]} \\
&= \sqrt{\mathrm{E}\left[\left(h(\tilde{x}_1, \ldots, \tilde{x}_n) - \gamma\right)^2\right]}
\end{aligned}
$$

# Standard error of the sample mean

$$\text{se}\,[\widetilde{m}]^2 = \text{Var}\,[\widetilde{m}] = \text{Var}\left[\frac{1}{n}\sum_{j=1}^{n}\tilde{x}_j\right]$$

$$= \frac{1}{n^2}\text{Var}\left[\sum_{j=1}^{n}\tilde{x}_j\right]$$

# Uncorrelated random variables

If $\tilde{a}$ and $\tilde{b}$ are uncorrelated

$$\mathrm{Var}[\tilde{a} + \tilde{b}] = \mathrm{Var}[\tilde{a}] + \mathrm{Var}[\tilde{b}]$$

# Sum of independent random variables

Independent random variables $\tilde{a}_1$, $\tilde{a}_2$, $\ldots$, $\tilde{a}_n$ with finite variance

$$\operatorname{Var}\left[\sum_{k=1}^{n} \tilde{a}_k\right] = \operatorname{Var}\left[\tilde{a}_1\right] + \operatorname{Var}\left[\sum_{k=2}^{n} \tilde{a}_k\right]$$

$$= \operatorname{Var}\left[\tilde{a}_1\right] + \operatorname{Var}\left[\tilde{a}_2\right] + \operatorname{Var}\left[\sum_{k=3}^{n} \tilde{a}_k\right]$$

$$= \sum_{k=1}^{n} \operatorname{Var}\left[\tilde{a}_k\right]$$

# Standard error of the sample mean

$$\text{se}\left[\widetilde{m}\right]^2 = \frac{1}{n^2}\text{Var}\left[\sum_{j=1}^{n}\tilde{x}_j\right]$$

$$= \frac{1}{n^2}\sum_{j=1}^{n}\text{Var}\left[\tilde{x}_j\right]$$

$$= \frac{\sigma_{\text{pop}}^2}{n}$$

$$\text{Var}\left[\tilde{x}_j\right] := \text{E}\left[(\tilde{x}_j - \text{E}\left[\tilde{x}_j\right])^2\right]$$

$$= \text{E}\left[(\tilde{x}_j - \mu_{\text{pop}})^2\right]$$

$$= \sum_{i=1}^{N}(a_i - \mu_{\text{pop}})^2 p_{\tilde{x}_j}(a_i)$$

$$= \frac{1}{N}\sum_{i=1}^{N}(a_i - \mu_{\text{pop}})^2 = \sigma_{\text{pop}}^2$$

# Standard error of the sample mean

$$se\left[\widetilde{m}\right] = \frac{\sigma_{\mathsf{pop}}}{\sqrt{n}}$$

No dependence on $N$!

# Height data: $n = 20$

$\mu_{\text{pop}} := 175.6$ cm, $\sigma_{\text{pop}} = 6.85$ cm
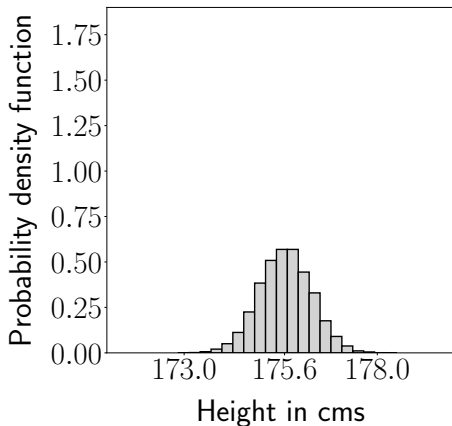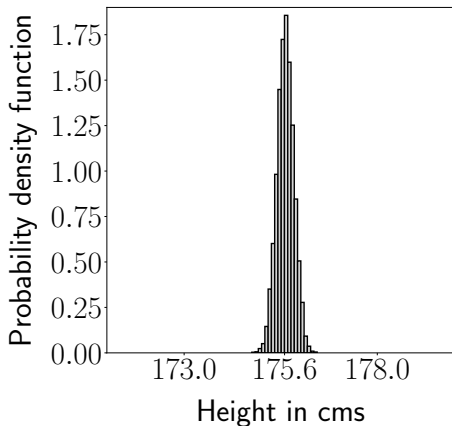
Total population $N := 4{,}082$

$10^4$ sample means

$n = 100$

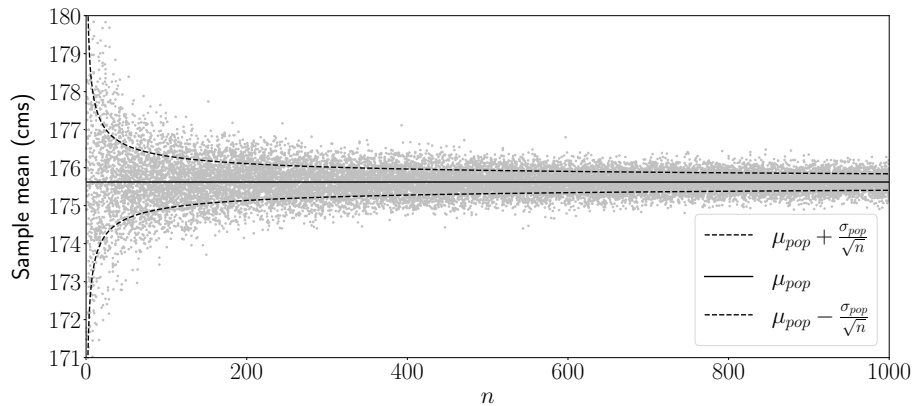$\mu_{\mathsf{pop}} := 175.6$ cm, $\sigma_{\mathsf{pop}} = 6.85$ cm

Total population $N := 4{,}082$

$10^4$ sample means

## $n = 1{,}000$

$\mu_{\mathsf{pop}} := 175.6$ cm, $\sigma_{\mathsf{pop}} = 6.85$ cm

Total population $N := 4{,}082$

$10^4$ sample means
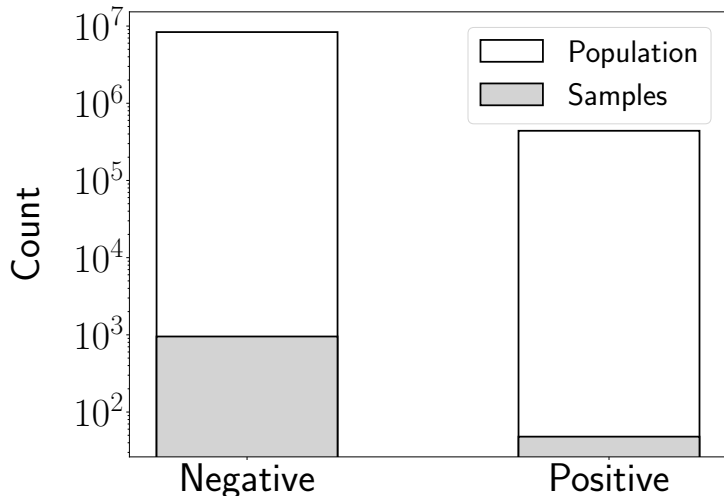
# Height data

# Estimating a population proportion

COVID-19 prevalence in New York

Population proportion:

$$\theta_{\mathsf{pop}} = 0.05$$

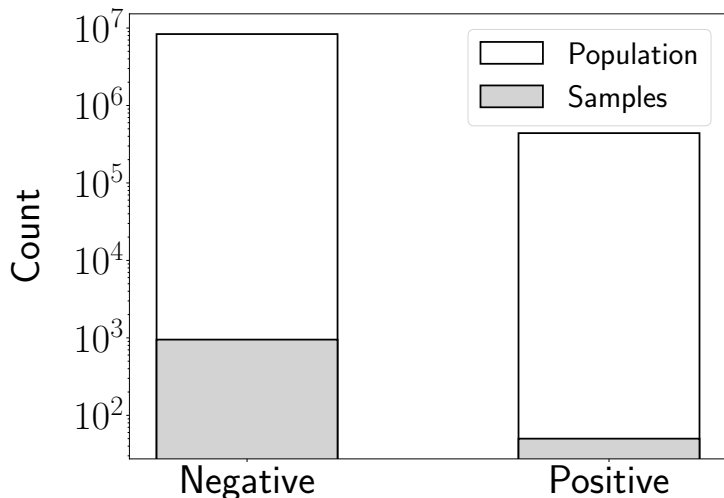# 1,000 random samples out of 8.8 million
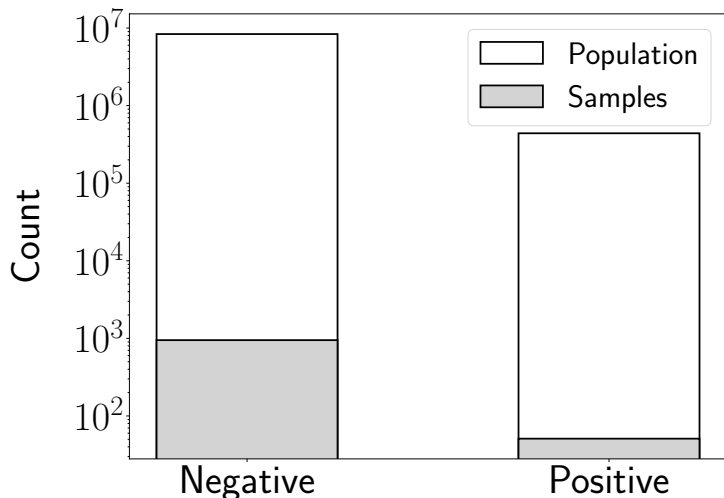
Sample proportion = 0.055 ($\theta_{pop} = 0.05$)

# 1,000 random samples out of 8.8 million

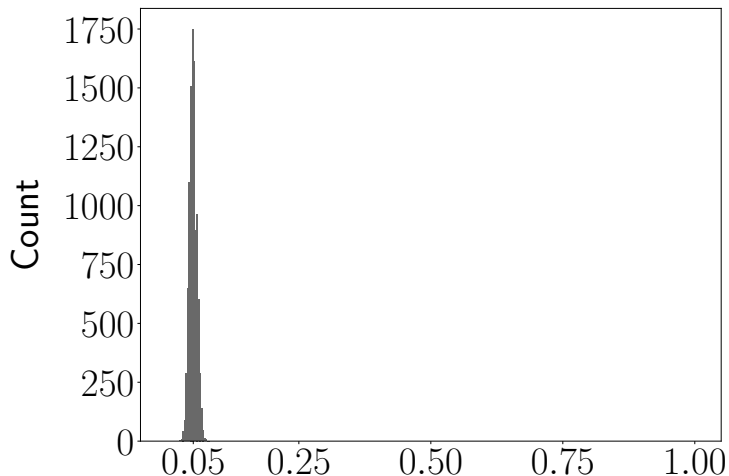Sample proportion = 0.049 ($\theta_{\mathsf{pop}} = 0.05$)

# 1,000 random samples out of 8.8 million

Sample proportion = 0.052 ($\theta_{\text{pop}} = 0.05$)

Sample proportions of 10,000 subsets of size 1,000

# Standard error of sample proportion

Data: $a_1$, $a_2$, ..., $a_N$

$a_i = 1$ if $i$th data point satisfies a certain condition

Random samples: $\tilde{x}_1$, $\tilde{x}_2$, ..., $\tilde{x}_n$

Sample proportion is sample mean $\tilde{m} := \frac{1}{n} \sum_{j=1}^{n} \tilde{x}_j$

$$\text{se}[\tilde{m}] = \frac{\sigma_{\text{pop}}}{\sqrt{n}}$$

# Population variance

$$\sigma_{\mathsf{pop}}^2 := \frac{1}{N} \sum_{i=1}^{N} (a_i - \theta_{\mathsf{pop}})^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} a_i^2 - \frac{2\theta_{\mathsf{pop}}}{N} \sum_{i=1}^{N} a_i + \frac{1}{N} \sum_{i=1}^{N} \theta_{\mathsf{pop}}^2$$

$$= \theta_{\mathsf{pop}} - 2\theta_{\mathsf{pop}}^2 + \theta_{\mathsf{pop}}^2$$

$$= \theta_{\mathsf{pop}}(1 - \theta_{\mathsf{pop}})$$

# Standard error of sample proportion

Data: $a_1, a_2, \ldots, a_N$

$a_i = 1$ if $i$th data point satisfies a certain condition

Random samples: $\tilde{x}_1, \tilde{x}_2, \ldots, \tilde{x}_n$

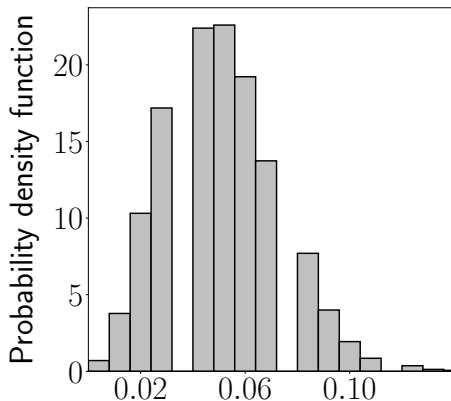Sample proportion is sample mean $\tilde{m} := \frac{1}{n} \sum_{j=1}^{n} \tilde{x}_j$

$$\begin{aligned}
\text{se}\left[\tilde{m}\right] &= \frac{\sigma_{\text{pop}}}{\sqrt{n}} \\
&= \sqrt{\frac{\theta_{\text{pop}}(1 - \theta_{\text{pop}})}{n}}
\end{aligned}$$

# COVID-19

$\theta_{\text{pop}} := 0.05$

Total population $N := 8$ million
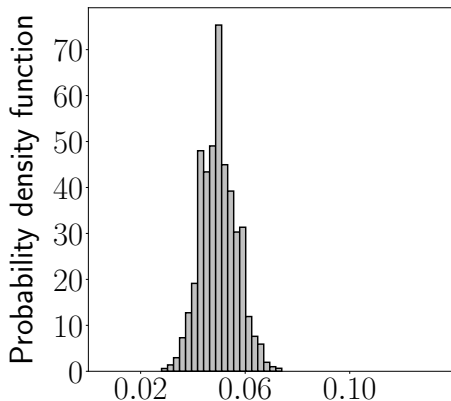
Distribution of $10^4$ sample means for $n = 100$

# COVID-19

$\theta_{\mathsf{pop}} := 0.05$

Total population $N := 8$ million

Distribution of $10^4$ sample means for $n = 1,000$

# COVID-19

$\theta_{\text{pop}} := 0.05$

Total population $N := 8$ million

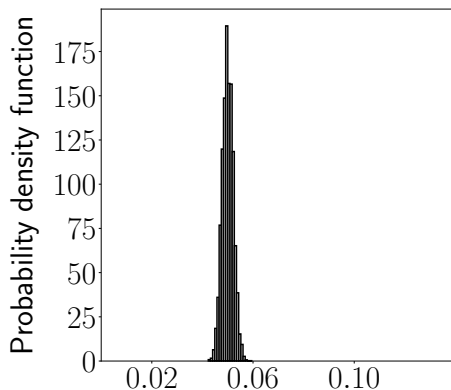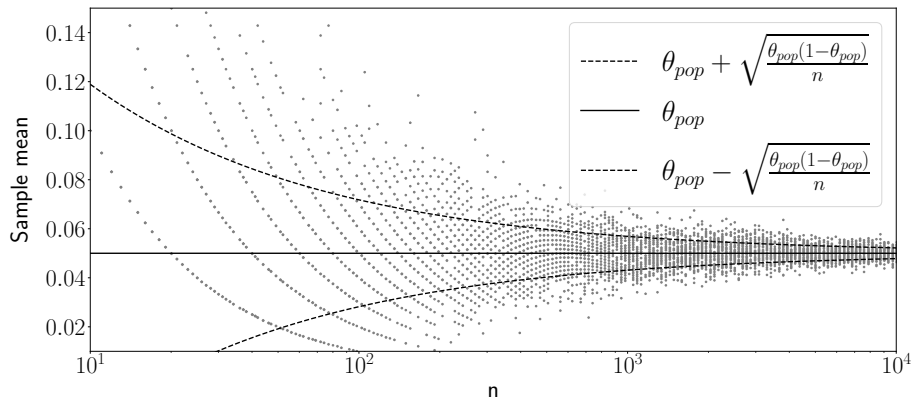Distribution of $10^4$ sample means for $n = 10{,}000$

# What have we learned

Definition of standard error

Standard error of sample mean and sample proportion

Random sampling works!