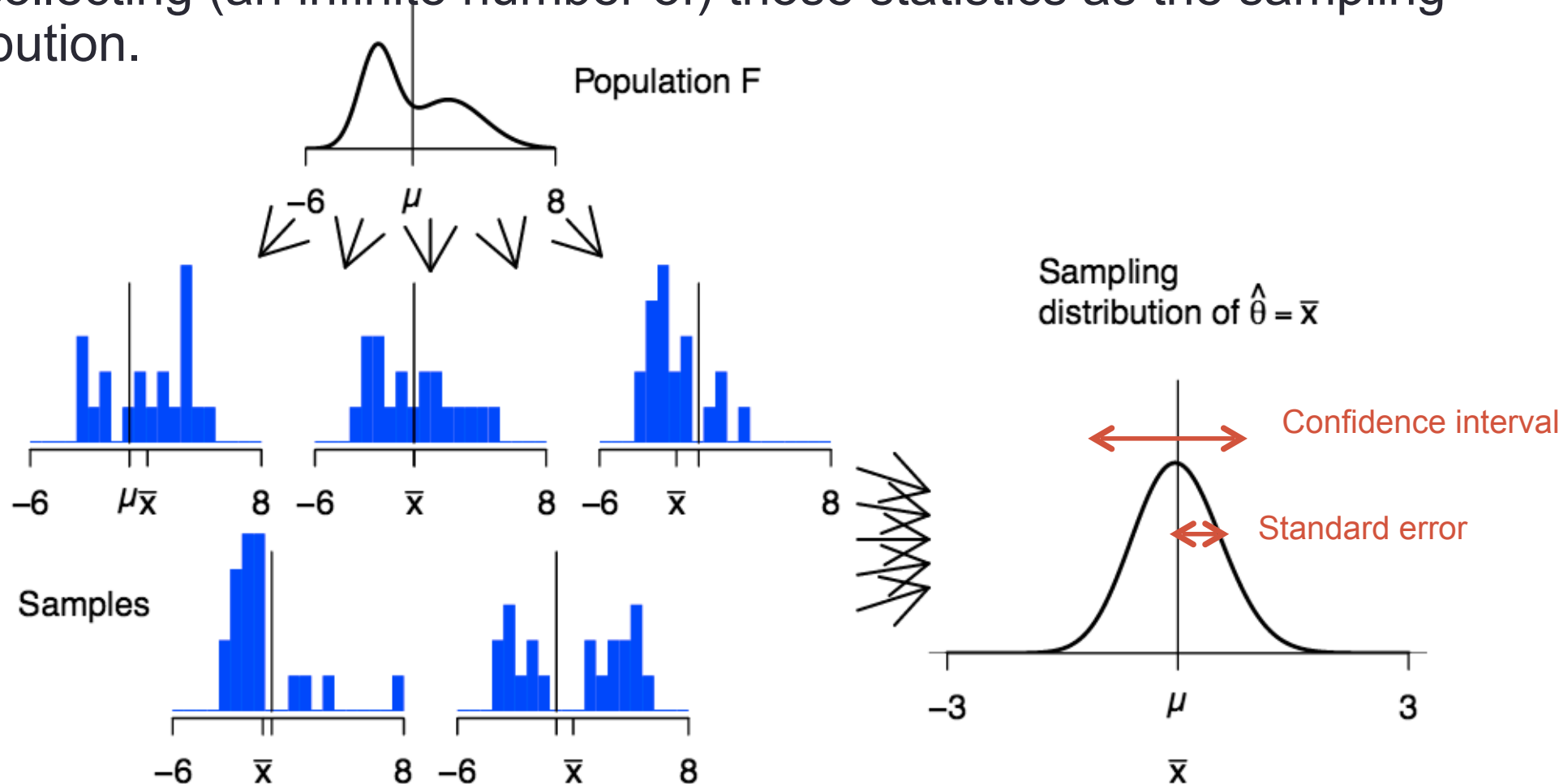


DATA ANALYSIS
Year 2019–2020

#4 Statistical Methods

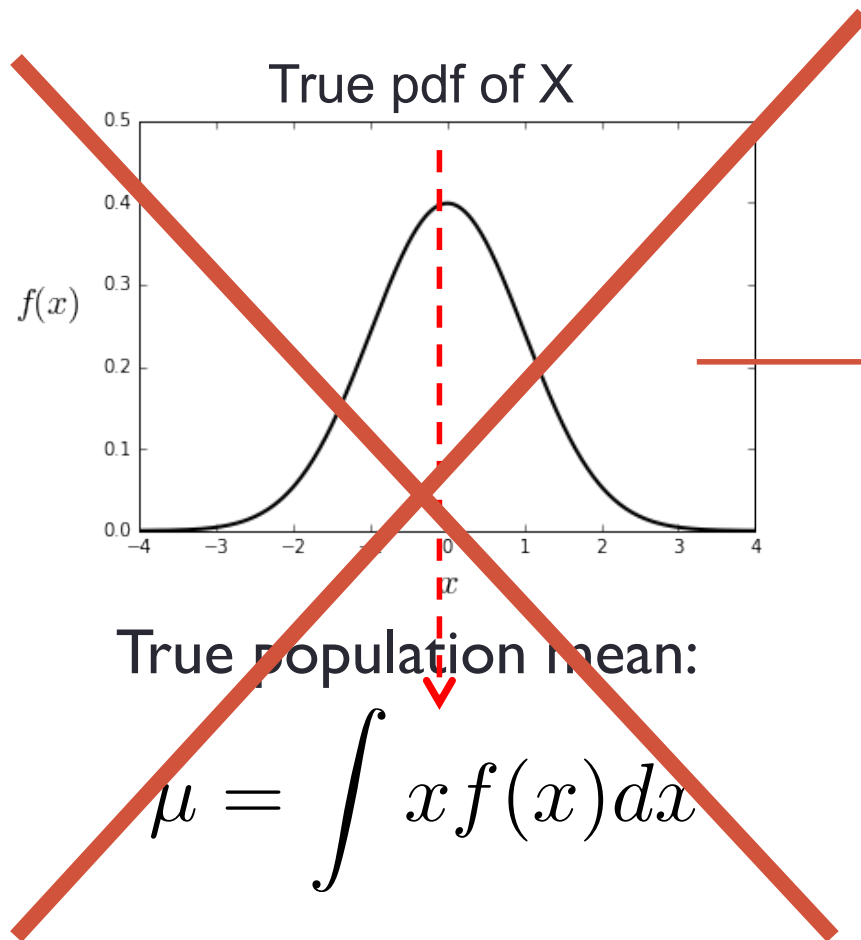
Estimating standard errors and confidence intervals

Ideal world: Sampling distributions are obtained by drawing repeated samples from the population, computing the statistic of interest for each, and collecting (an infinite number of) those statistics as the sampling distribution.

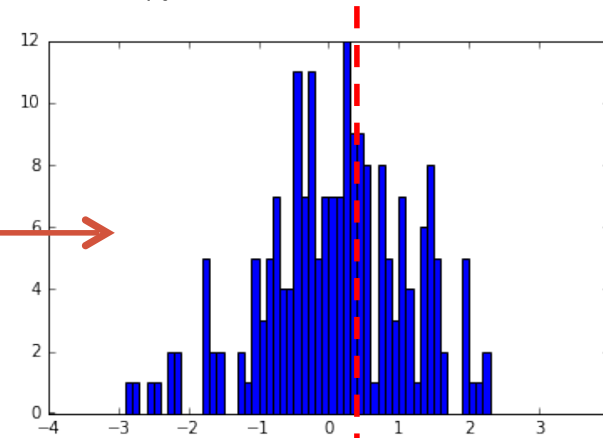


Estimating standard errors and confidence intervals

Reality: The problem is that we have only ONE sample



Sample with N values
 x_k for $k = 1..N$



Mean estimator (= sample mean)

$$\bar{x} = \frac{1}{N} \sum_k x_k$$

For example if you want to know the mean (works the same for any other statistics), you can only compute one sample mean (a-priori different than the true mean).

Estimating standard errors and confidence intervals

Reality: The problem is that we have only ONE sample

For example if you want to know the mean, you cannot do better than the sample mean. This is the best estimator.

But you also want to be able to quantify “***how far your sample mean is from the true mean?***”

This is what a **standard error** and a **confidence interval** will tell you.

The confidence interval defines the degree of certainty that a given quantity θ (mean or any other statistics) will fall between specified lower and upper bounds θ_L and θ_U :

$$P(\theta_L < \theta < \theta_U) = 1 - \alpha$$

Where α is the **level of significance (or confidence level)**, and $100(1 - \alpha)$ is the percent significance level.

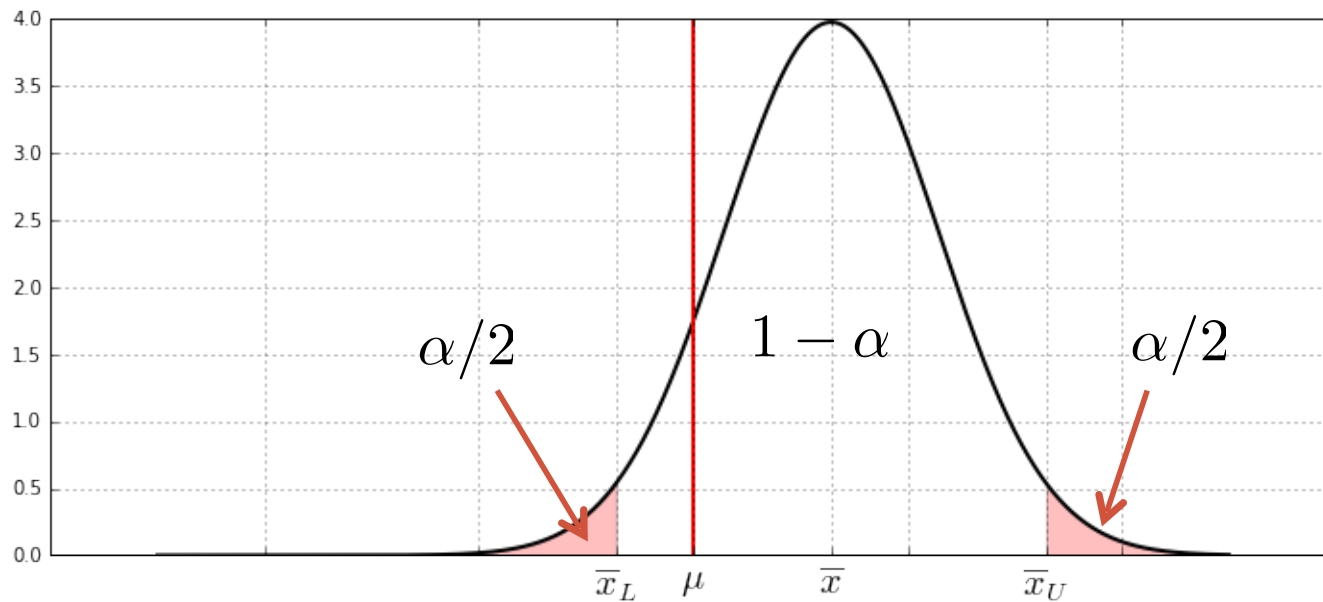
A typical value is $\alpha = 0.05$, which means $100(1 - \alpha) = 95\%$, and corresponds to the 95 % confidence interval.

Estimating standard errors and confidence intervals

For example, let's say you compute the sample mean \bar{x}

The confidence interval will correspond to the lower and upper bounds, \bar{x}_L, \bar{x}_U such that:

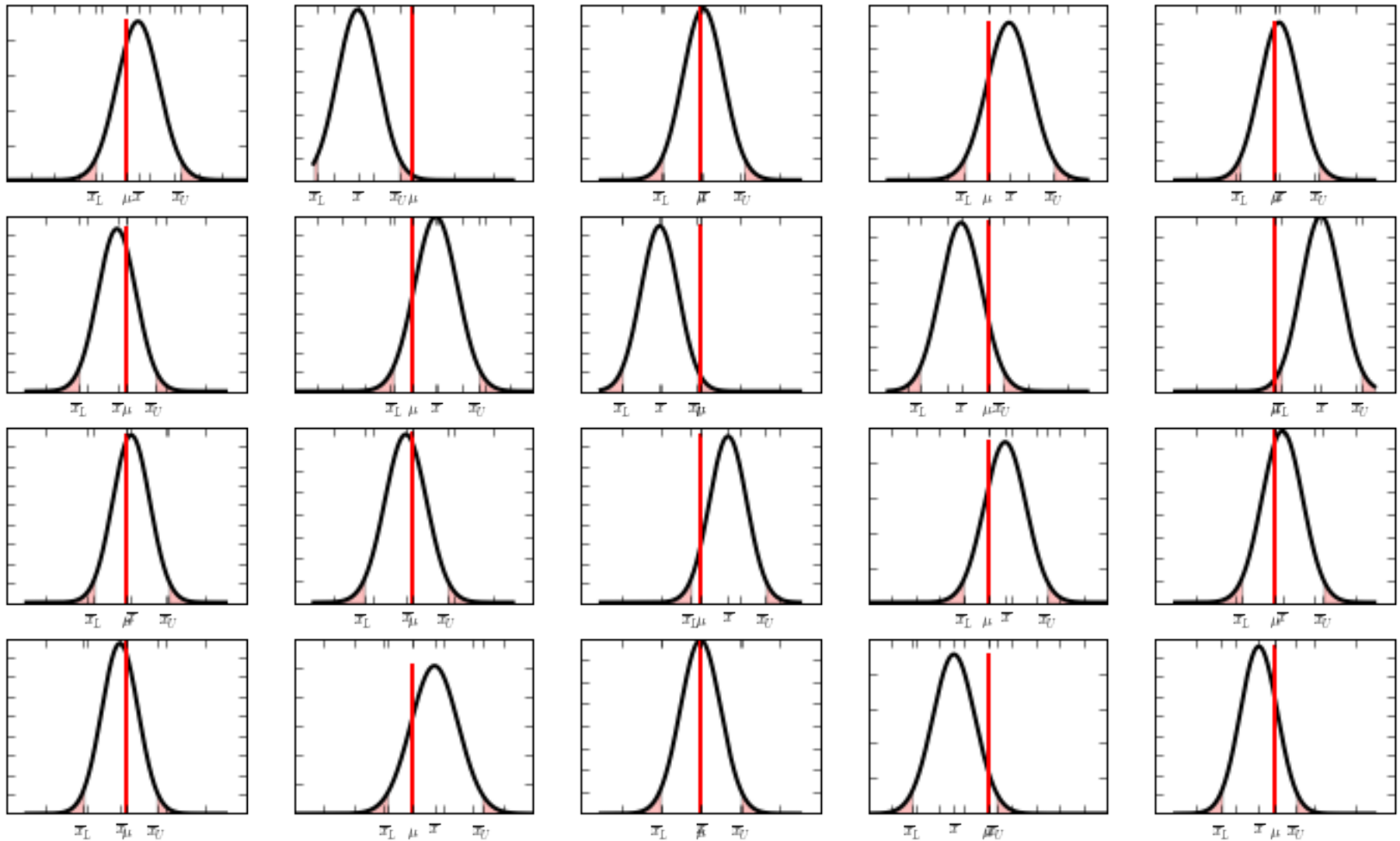
$$P(\bar{x}_L < \mu < \bar{x}_U) = 1 - \alpha$$



Which means that the interval (\bar{x}_L, \bar{x}_U) has a probability $1 - \alpha$ of containing the true mean μ .

In other words = we are 95% confident that the true value of the parameter is in the confidence interval...

Estimating standard errors and confidence intervals



If you get $N_s=100$ samples and construct $N_s=100$ confidence intervals at 95%, the true mean will be included in 95 confidence intervals out of the 100 (on average)

Estimating standard errors and confidence intervals

Question:

How do we compute the standard error and confidence interval (CI)?

There are different ways of doing that, we'll talk about 2 common ways in oceanography:

1. **Using the Central Limit Theorem** (for the mean only)
2. **Resampling methods** (bootstrap and Jackknife methods, for any statistics)

1. CI for the mean (using the CLT)

Central Limit Theorem:

Let $X_i, i = 1..N_s$ be a sequence of independent random variables (each containing N values) drawn from distributions with mean μ and variance σ^2 . Then as N_s becomes large, the distribution of the mean values X_i of each sample $\hat{\mu}_i$ approaches the normal distribution with mean μ and variance σ^2/N .

$$\hat{\mu}(x) \sim \mathcal{N}(\mu, \sigma/\sqrt{N})$$

1. CI for the mean (using the CLT)

Central Limit Theorem:

Let $X_i, i = 1..N_s$ be a sequence of independent random variables (each containing N values) drawn from distributions with mean μ and variance σ^2 . Then as N_s becomes large, the distribution of the mean values X_i of each sample $\hat{\mu}_i$ approaches the normal distribution with mean μ and variance σ^2/N .

$$\hat{\mu}(x) \sim \mathcal{N}(\mu, \sigma/\sqrt{N}) = \text{standard error}$$

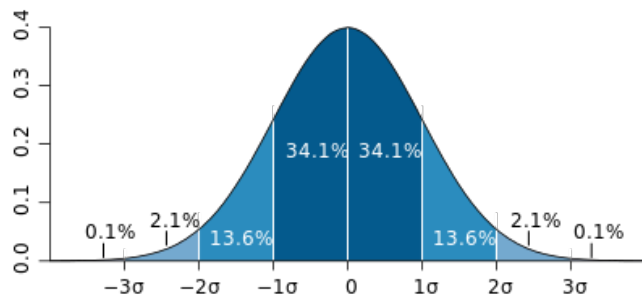
1. CI for the mean (using the CLT)

So the standard error is for the mean is $\frac{\sigma}{\sqrt{N}}$

And the $100(1 - \alpha)$ percent confidence interval for the population mean is given by:

$$\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{N}} < \mu < \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{N}}$$

Where $z_{\alpha/2}$ are the values giving the $100(1 - \alpha)$ percent confidence interval for a standard normal distribution. The $z_{\alpha/2}$ can be recomputed using the theoretical normal distribution function or they are directly read on a table



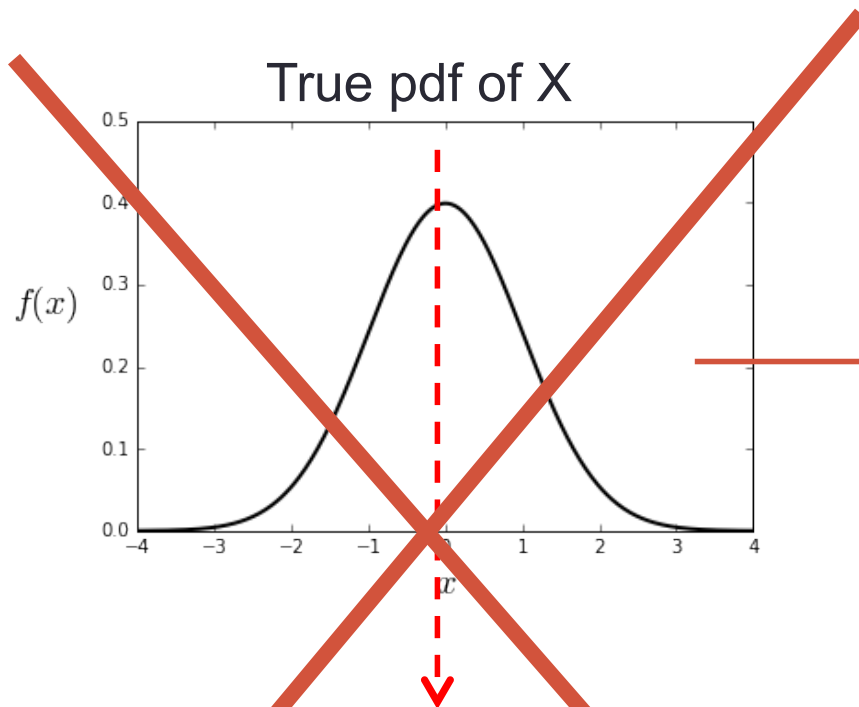
$100(1 - \alpha)$	$z_{\alpha/2}$
99%	2.576
98%	2.326
95%	1.96
90%	1.645

1. CI for the mean (using the CLT)

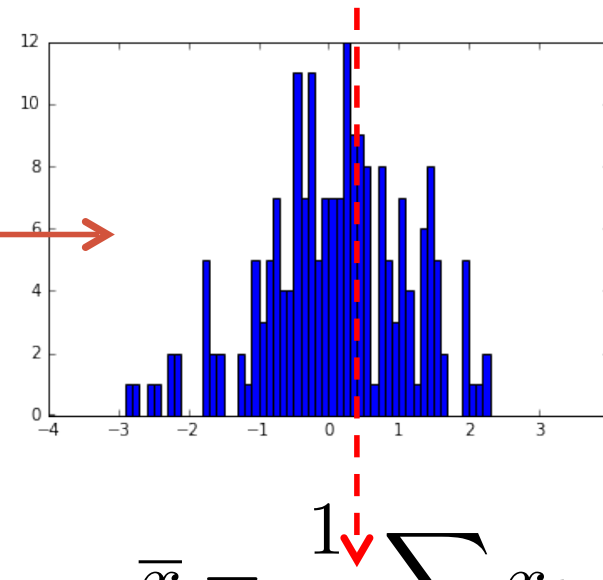
Problem: Most of the time we don't know the true standard deviation σ !!!
We only know the standard deviation estimated from our sample = s .

Sample with N values

x_k for $k = 1..N$



$$\mu = \int x f(x) dx$$
$$\sigma^2 = \int (x - \mu)^2 f(x) dx$$



$$\bar{x} = \frac{1}{N} \sum_k x_k$$

$$s^2 = \frac{1}{N-1} \sum_k (x_k - \hat{\mu})^2$$

1. CI for the mean (using the CLT)

If we don't know the true standard deviation we can estimate the standard error of the mean using: $\frac{s}{\sqrt{N}}$

(Where σ has been replaced by the sample standard deviation s)

And the $100(1 - \alpha)$ percent confidence interval for the population mean is given by:

$$\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{N}} < \mu < \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{N}}$$

(Where σ has been replaced by the sample standard deviation s and the $z_{\alpha/2}$ have been replaced by $t_{\alpha/2}$)

Where $t_{\alpha/2}$ are the values giving the $100(1 - \alpha)$ percent confidence interval for a **student's t-distribution** with $N - 1$ degrees of freedom.

1. CI for the mean (using the CLT)

Summary

a. σ is known

The standard error is $\frac{\sigma}{\sqrt{N}}$

the $100(1 - \alpha)$ percent confidence interval for the population mean is given by:

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{N}}$$

b. σ is not known

The standard error is $\frac{s}{\sqrt{N}}$

the $100(1 - \alpha)$ percent confidence interval for the population mean is given by:

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{N}}$$

Estimating standard errors and confidence intervals

Question:

How do we compute the standard error and confidence interval (CI)?

There are different ways of doing that, we'll talk about 2 common ways in oceanography:

1. Using the Central Limit Theorem (for the mean only)

2. **Resampling methods** (bootstrap and Jackknife methods, for any statistics)

2. Bootstrapping method

The CLT is very convenient to estimate the standard error of the mean.

But in a more general case we want to compute different statistics (median, standard deviation, kurtosis, percentiles, ...) and know how much they will vary from one sample to another (standards errors and confidence intervals.)

We can use generic methods known as **resampling methods** such as jackknifing and **bootstrapping**.

2. Bootstrapping method

Bootstrapping is a very simple and efficient method to estimate the sampling distribution of almost any statistic using random sampling methods

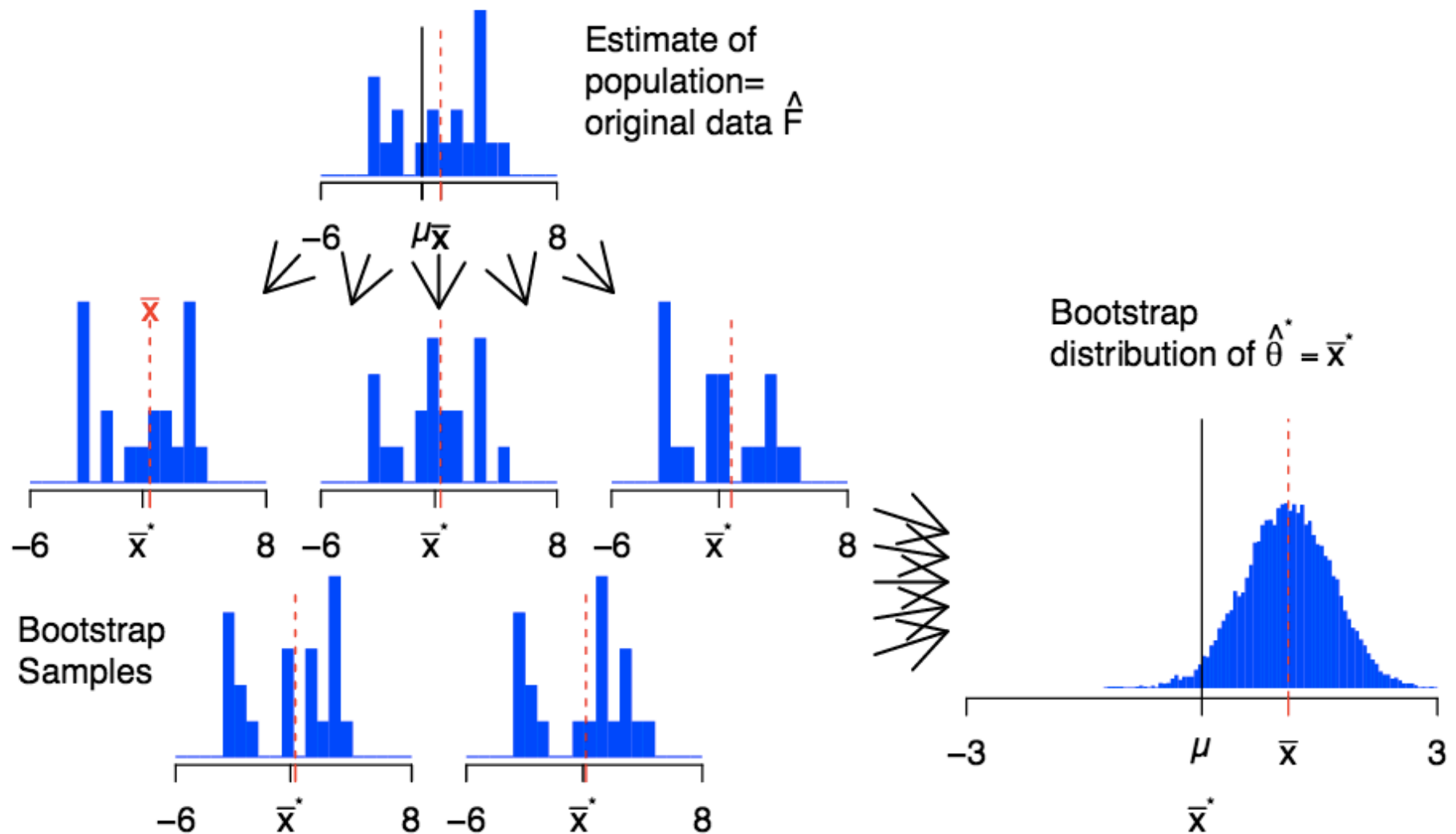
The idea is to:

- Treat the sample as the true population.
- **Sample with replacement** your actual distribution M times.
- Compute the statistic of interest on each “re-sample”.

The name bootstrap comes from the idea of “lifting yourself up by your own bootstraps” [the loop at the top of tall boots].

2. Bootstrapping method

- **Bootstrap world.** The bootstrap distribution is obtained by drawing repeated samples from an estimate of the population, computing the statistic of interest for each, and collecting those statistics.



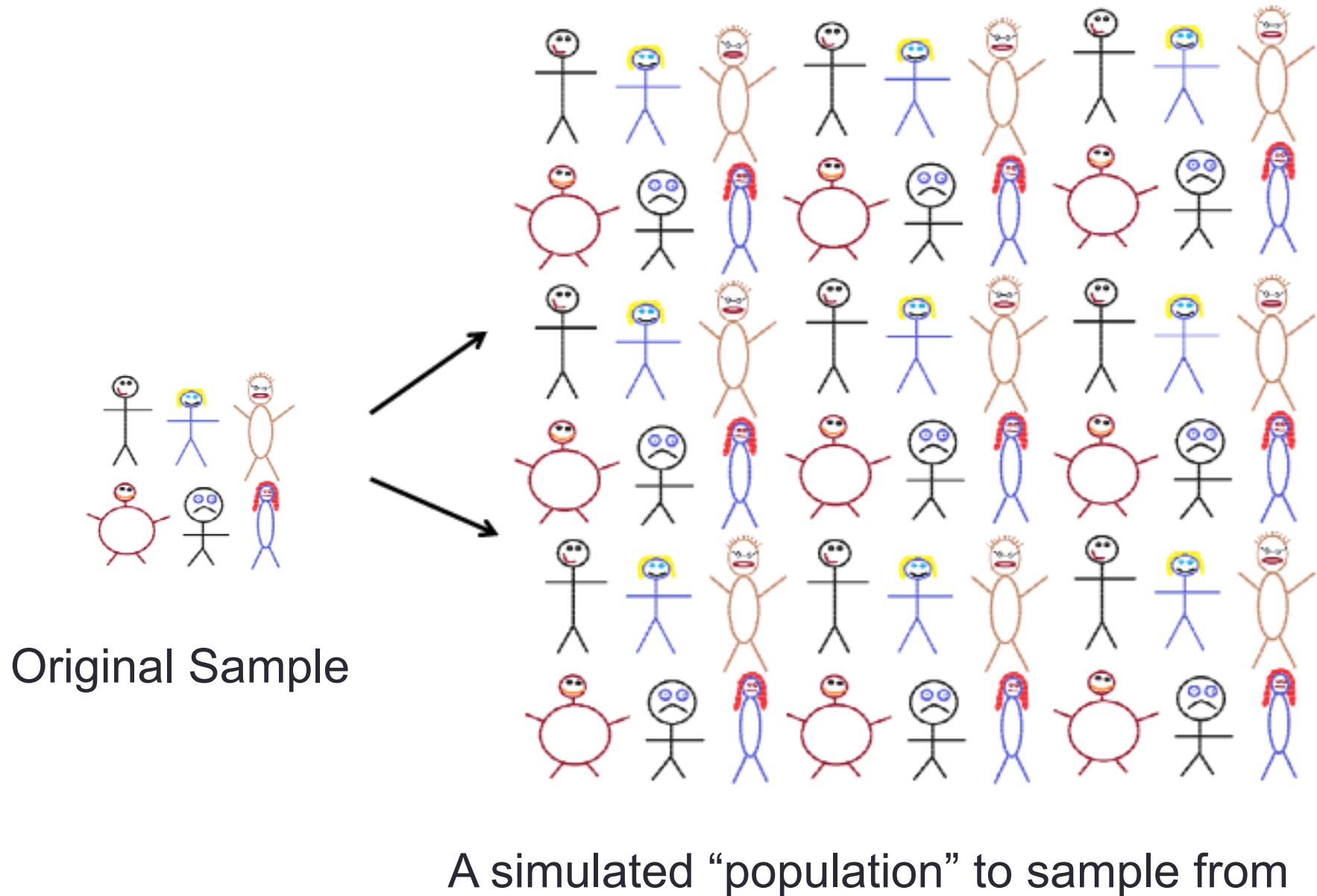
2. Bootstrapping method

Sampling with replacement:

Suppose we have a random sample of 6 people



2. Bootstrapping method



2. Bootstrapping method

Sampling with replacement:

Imagine putting papers in a hat. If you **sample with replacement**, you would choose one paper, put the paper back in the hat, and then choose another paper.

2. Bootstrapping method

Sampling with replacement:



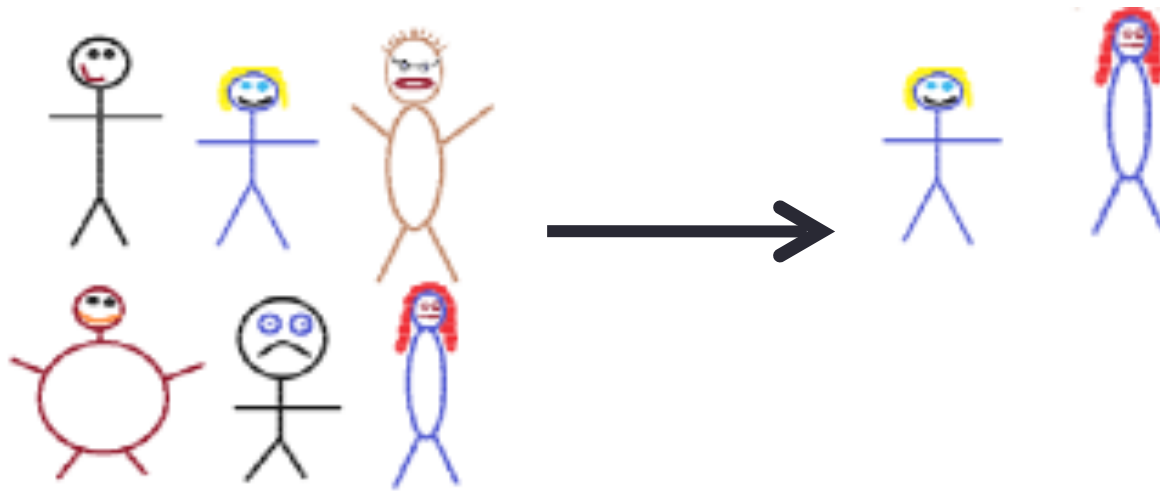
Original Sample

Bootstrap Sample

A bootstrap sample is a random sample taken with replacement from the original sample, of the same size as the original sample.

2. Bootstrapping method

Sampling with replacement:



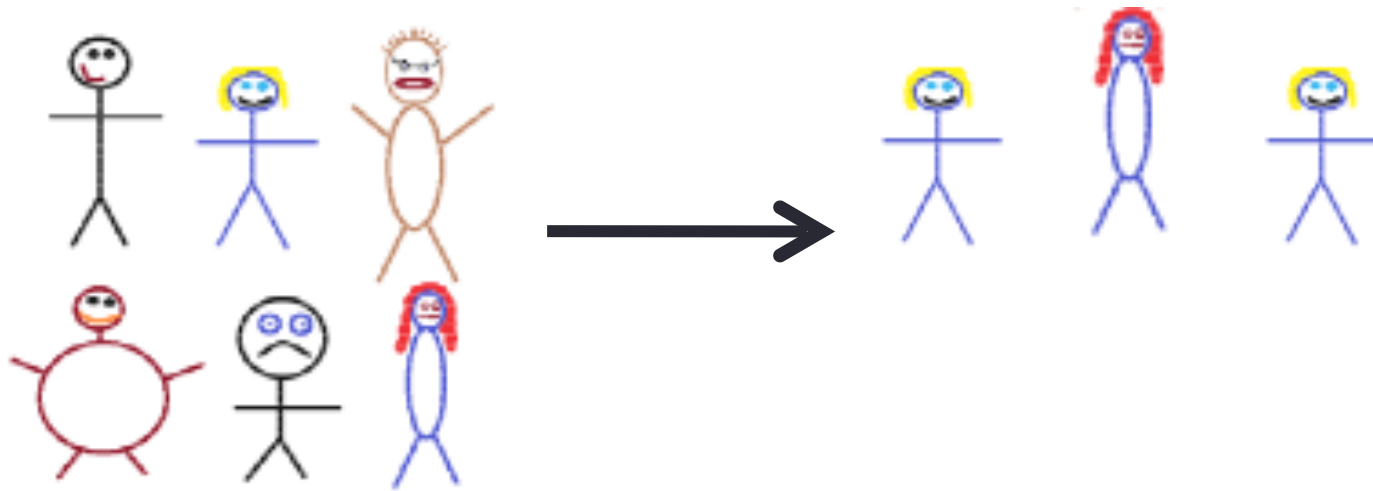
Original Sample

Bootstrap Sample

A bootstrap sample is a random sample taken with replacement from the original sample, of the same size as the original sample.

2. Bootstrapping method

Sampling with replacement:



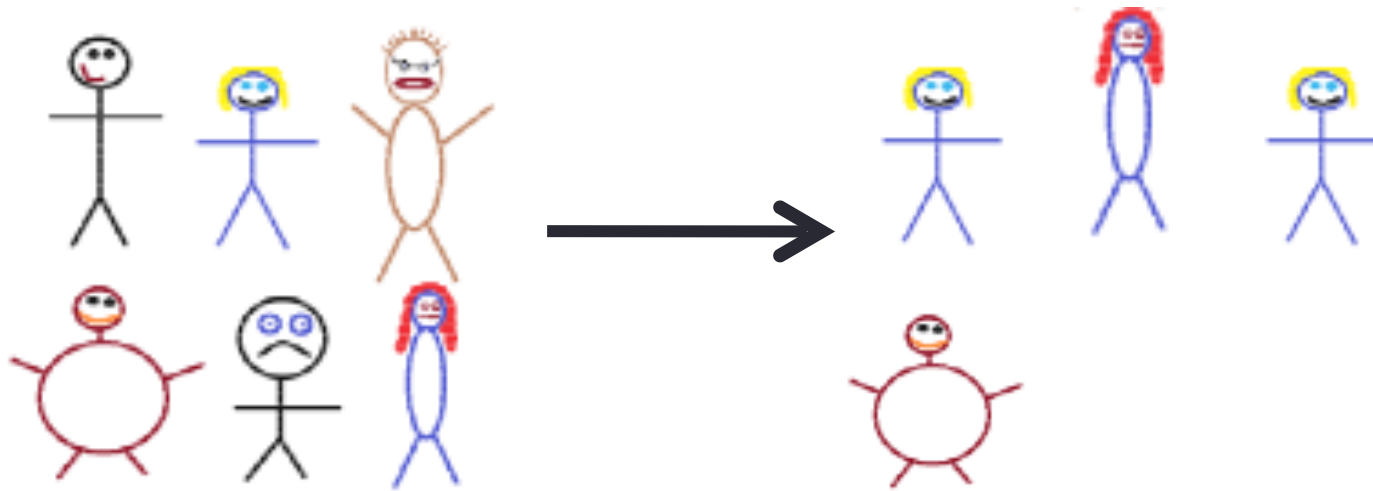
Original Sample

Bootstrap Sample

A bootstrap sample is a random sample taken with replacement from the original sample, of the same size as the original sample.

2. Bootstrapping method

Sampling with replacement:



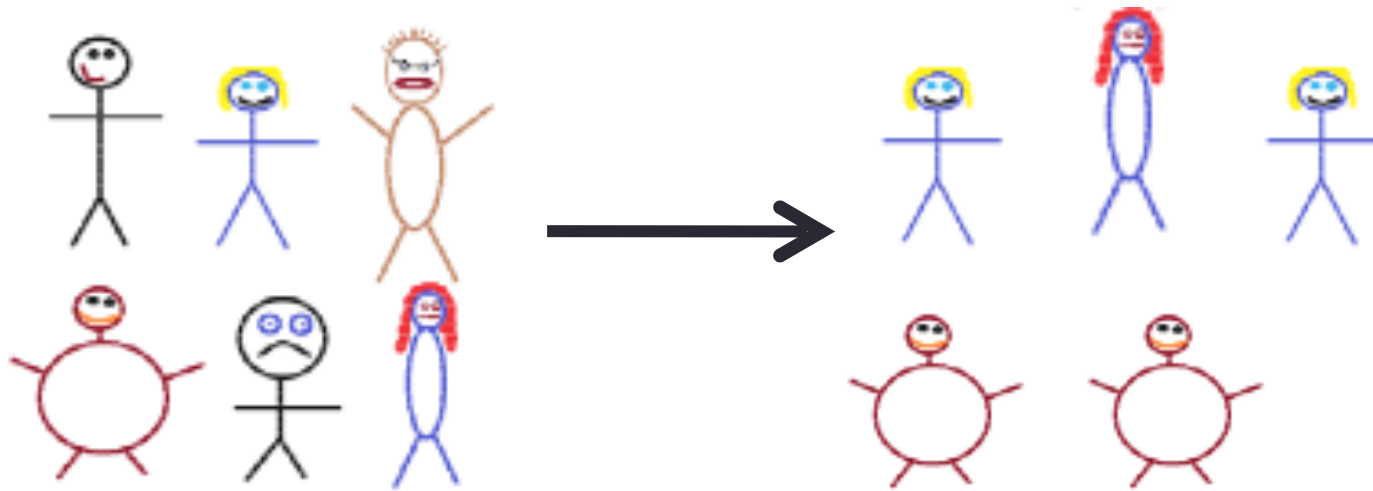
Original Sample

Bootstrap Sample

A bootstrap sample is a random sample taken with replacement from the original sample, of the same size as the original sample.

2. Bootstrapping method

Sampling with replacement:



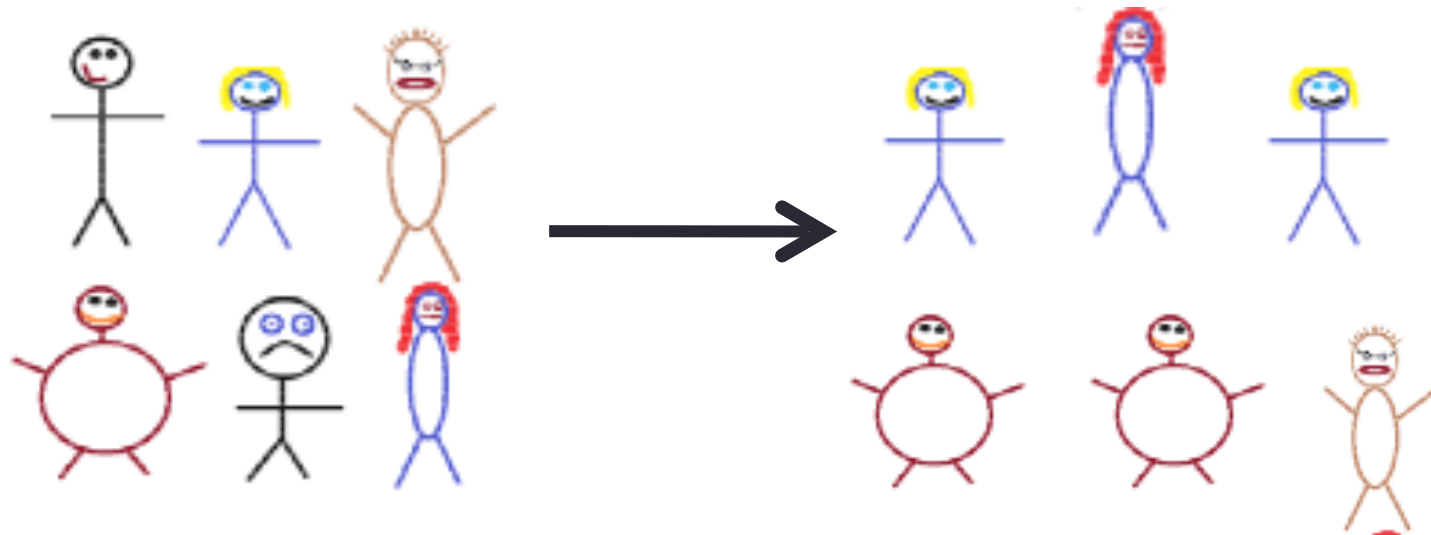
Original Sample

Bootstrap Sample

A bootstrap sample is a random sample taken with replacement from the original sample, of the same size as the original sample.

2. Bootstrapping method

Sampling with replacement:



Original Sample

Bootstrap Sample

A bootstrap sample is a random sample taken with replacement from the original sample, of the same size as the original sample.

2. Bootstrapping method

- A ***bootstrap statistic*** is the statistic computed on a bootstrap sample
- A ***bootstrap distribution*** is the distribution of many bootstrap statistics
- The variability of the bootstrap statistics is similar to the variability of the sample statistics
- The standard error of a statistic can be estimated using the standard deviation of the bootstrap distribution!
- Confidence intervals can be created using the standard error or the percentiles of a bootstrap distribution

2. Bootstrapping method

Number of Bootstrap Samples

- When using bootstrapping, you may get a slightly different confidence interval each time. This is fine!
- The more bootstrap samples you use, the more precise your answer will be.
- Increasing the number of bootstrap samples will not change the standard error or interval
- In real life, you probably want to take 10,000 or even 100,000 bootstrap samples

Estimating standard errors and confidence intervals

- See TD4 – Confidence intervals