

Inf2 - Foundations of Data Science: Estimation - Principle of confidence intervals



THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Announcements

InfPALS - especially LaTeX

Workshop

Lab

Project ideas - due Friday!

Last Lecture

1. Parameter

- value of a statistic (e.g. mean or max) in population
- parameter in distribution (e.g. mean, variance of normal)

2. Point estimator

- Method of converting sample into estimate of parameter
- E.g. Mean of sample (\bar{x}) estimates mean of population μ

$$\hat{\mu}$$

3. Point estimator is random variable

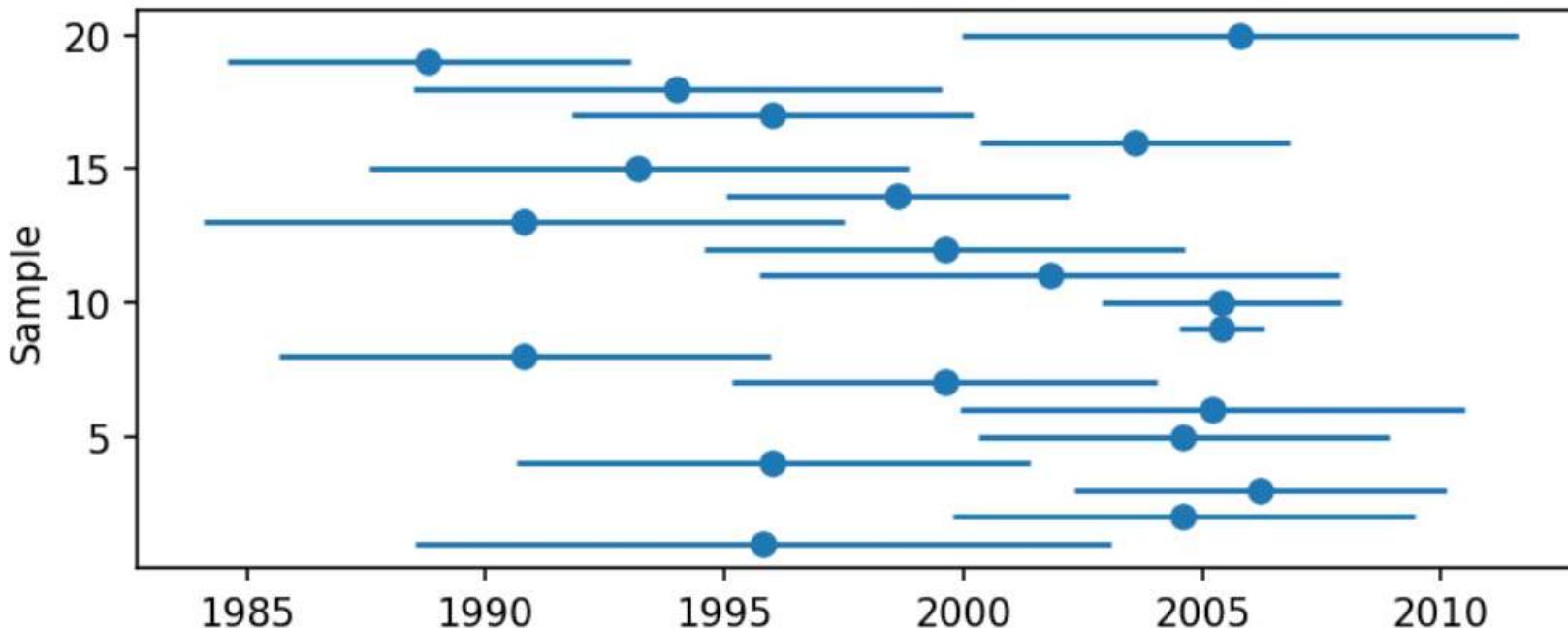
- a different random sample from population => different value of point estimator
- But we only have one sample, so only one value

4. Idea of confidence intervals for estimator

- based on sample standard error of estimator

Confidence intervals from different samples

Estimates $\hat{\mu}$ of the mean year of a coin μ
and CIs formed by $+\/- 1 \text{ SEM}$



How do we calibrate the width of the confidence interval so that there is a specified chance that it encloses the true value?

Today

1. How to convert inferred sampling distribution of estimator into a confidence interval with a specified chance of enclosing true value
2. How to compute a confidence interval for mean of large sample
 - z distribution
3. Choosing confidence levels and how much data to collect
4. Confidence intervals of parameters other than the mean
 - Bootstrap
5. How to calculate a confidence interval for mean of a small sample - t distribution

Theory reminder: Standard error of mean for known distribution variance σ^2

Standard deviation
 $\sigma = 1$

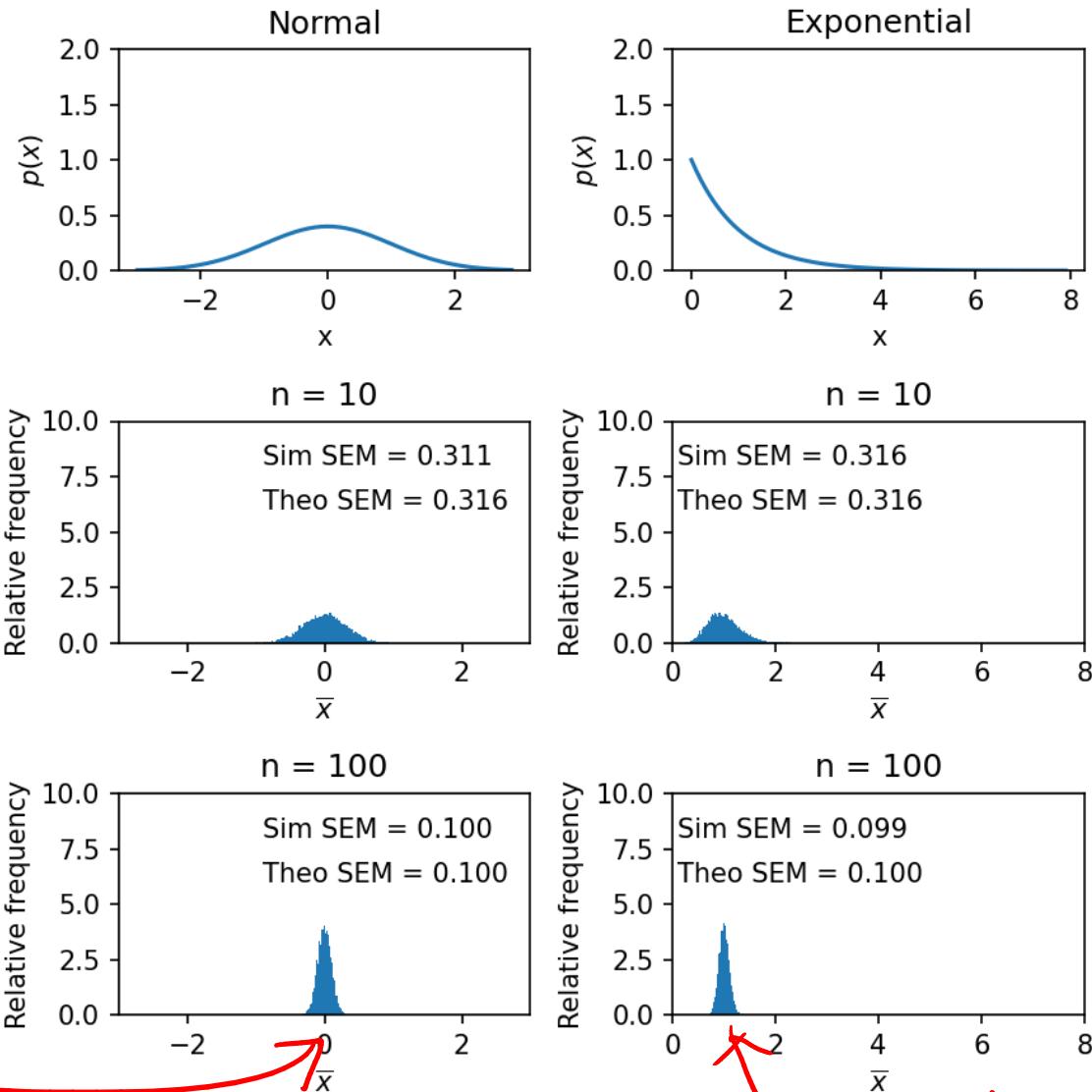
SEM

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{1}{\sqrt{10}} = 0.316$$

$$\hat{\sigma}_{\bar{x}} = \frac{1}{\sqrt{100}}$$

CLT

Normal



Estimated standard error for distribution with unknown variance σ

What if we don't know σ ?

Estimated S.E.

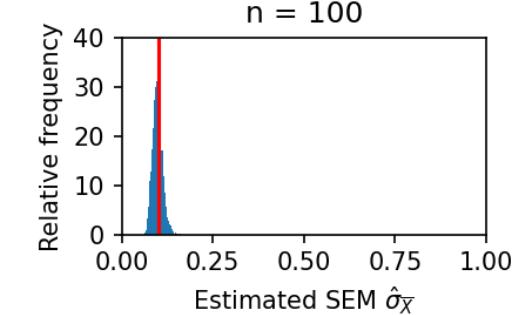
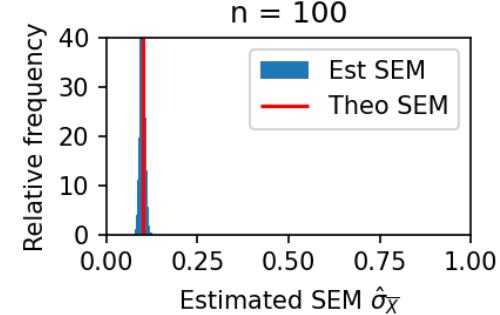
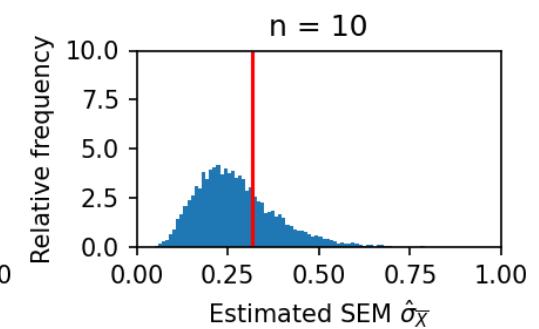
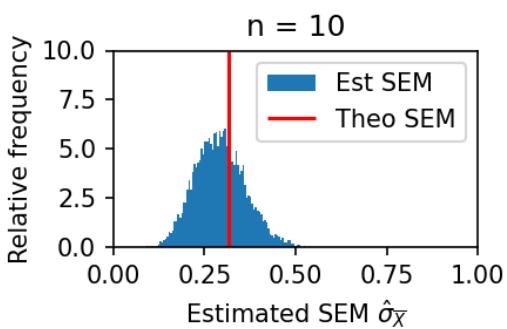
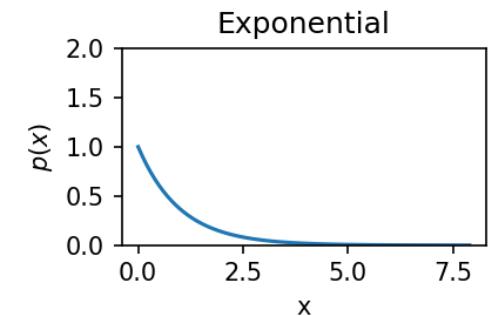
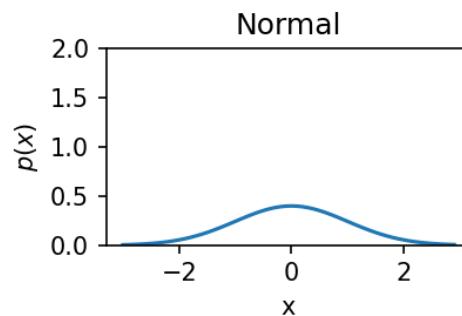
of estimator $\hat{\sigma}_{\hat{x}}$

Estimated SEM

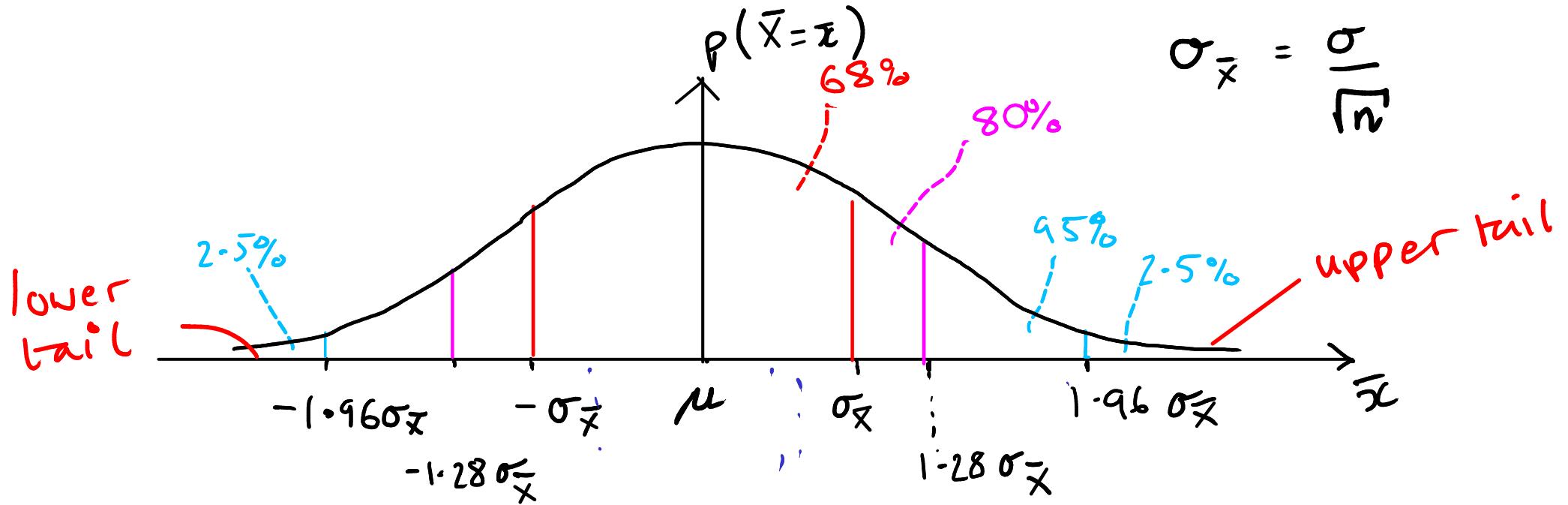
$$\hat{\sigma}_{\hat{x}} = \frac{s}{\sqrt{n}} \leftarrow r.v.$$

n large \Rightarrow

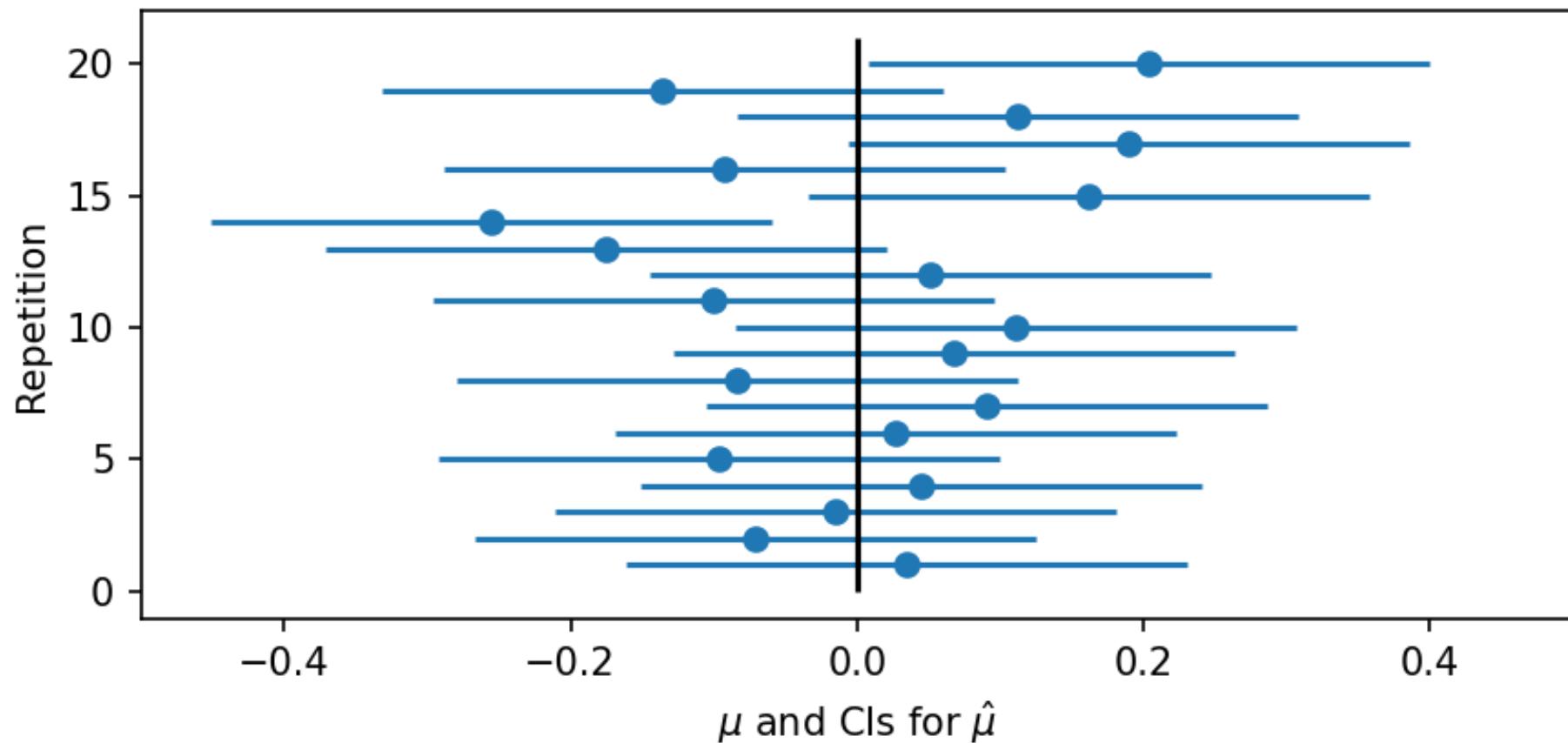
$$\hat{\sigma}_{\hat{x}} \approx \sigma_{\hat{x}}$$



Confidence interval of the mean of a sample from a distribution with unknown mean and known variance



E.g.: 95% Confidence intervals of mean of 100 samples from normal distribution with mean 0 and variance 1



Inf2 - Foundations of Data Science: Estimation – Definition of a confidence interval



THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Definition of a confidence interval

Confidence interval: An interval

$$(\hat{\theta} - a \hat{\sigma}_{\hat{\theta}}, \hat{\theta} + b \hat{\sigma}_{\hat{\theta}})$$

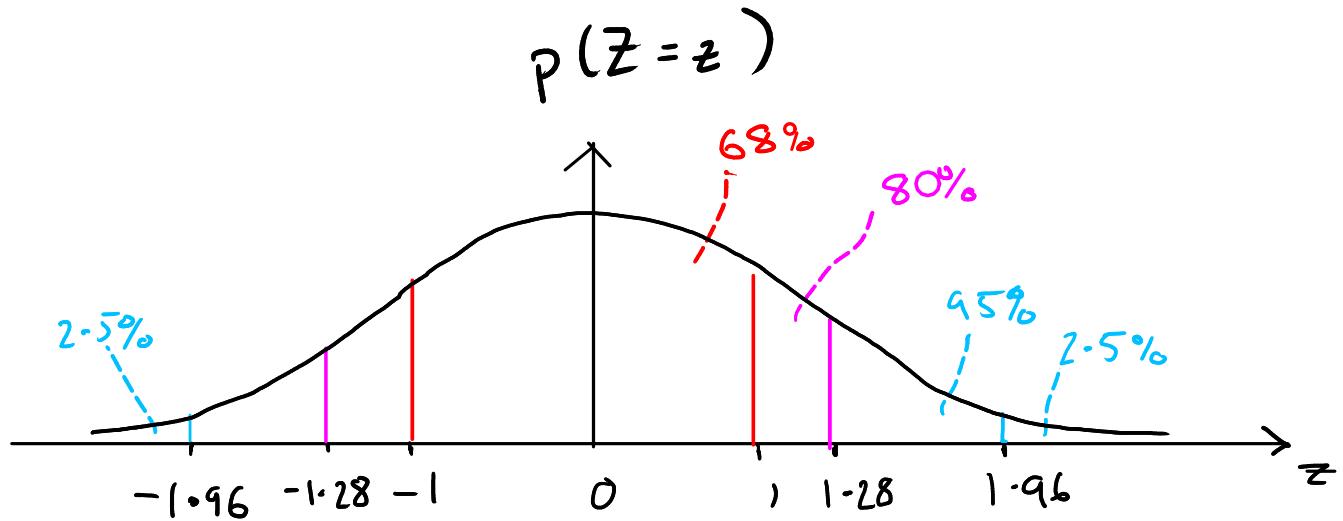
that has a specified chance $1-\alpha$ of containing the parameter θ .

e.g. $\alpha = 0.05 \Rightarrow 1 - 0.05 = 95\% \text{ C. I.}$

$$P(\hat{\theta} - a \hat{\sigma}_{\hat{\theta}} < \theta < \hat{\theta} + b \hat{\sigma}_{\hat{\theta}}) = 1 - \alpha$$

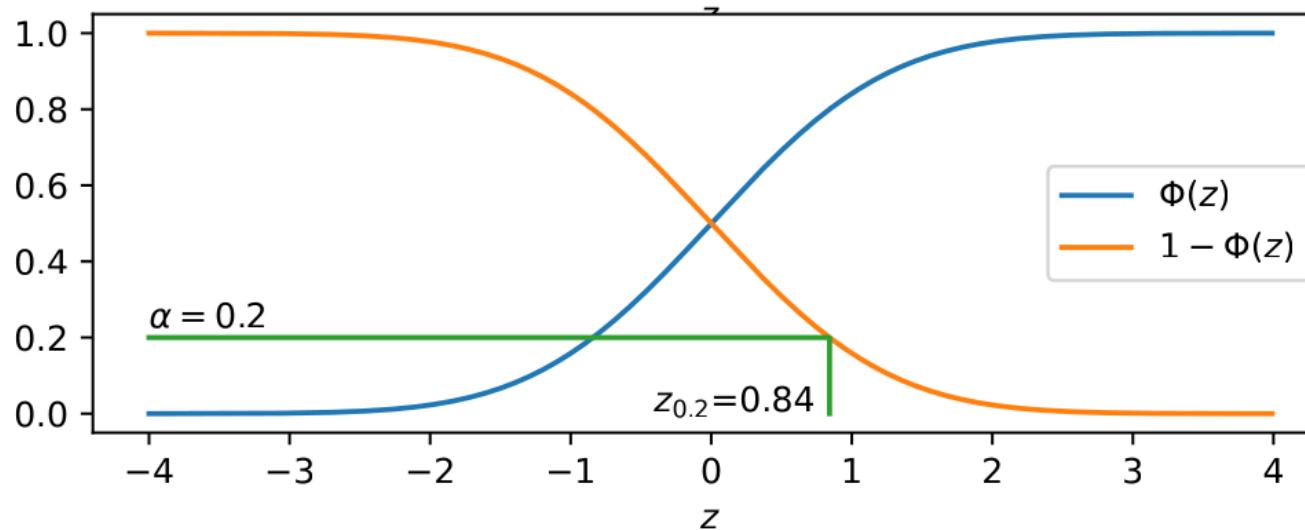
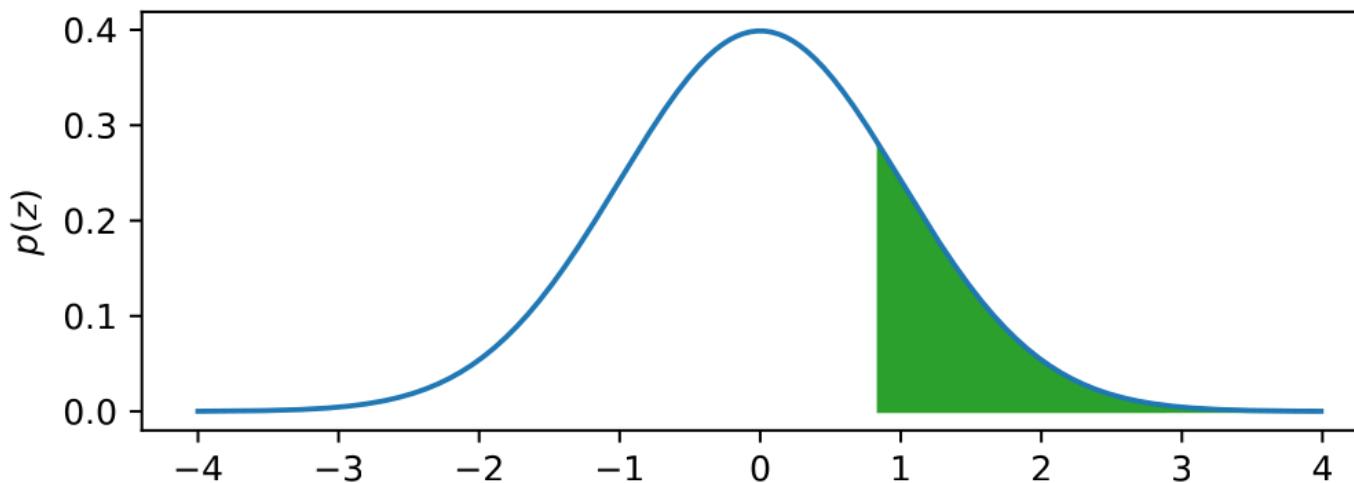
Often the interval is symmetric, ie $a = b$.

The distribution of the standardised sample mean of a large sample – the z-distribution

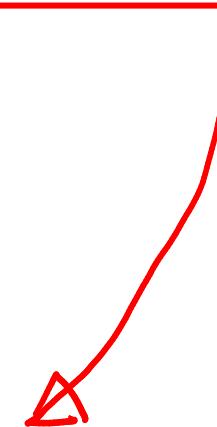


Normal dist $\mu=0$, $\sigma^2=1$

z-critical values



α	0.100	0.050	0.025	0.010	0.005	0.001
v						
1	3.078	6.314	12.706	31.821	63.657	318.309
2	1.886	2.920	4.303	6.965	9.925	22.327
3	1.638	2.353	3.182	4.541	5.841	10.215
4	1.533	2.132	2.776	3.747	4.604	7.173
5	1.476	2.015	2.571	3.365	4.032	5.893
6	1.440	1.943	2.447	3.143	3.707	5.208
7	1.415	1.895	2.365	2.998	3.499	4.785
8	1.397	1.860	2.306	2.896	3.355	4.501
9	1.383	1.833	2.262	2.821	3.250	4.297
10	1.372	1.812	2.228	2.764	3.169	4.144
20	1.325	1.725	2.086	2.528	2.845	3.552
30	1.310	1.697	2.042	2.457	2.750	3.385
40	1.303	1.684	2.021	2.423	2.704	3.307
∞	1.282	1.645	1.960	2.326	2.576	3.090



α	0.100	0.050	0.025	0.010	0.005	0.001
----------	-------	-------	-------	-------	-------	-------

∞	1.282	1.645	1.960	2.326	2.576	3.090
----------	-------	-------	-------	-------	-------	-------

Inf2 - Foundations of Data Science: Estimation -

Theoretical method of estimating the confidence interval of the mean of a large sample



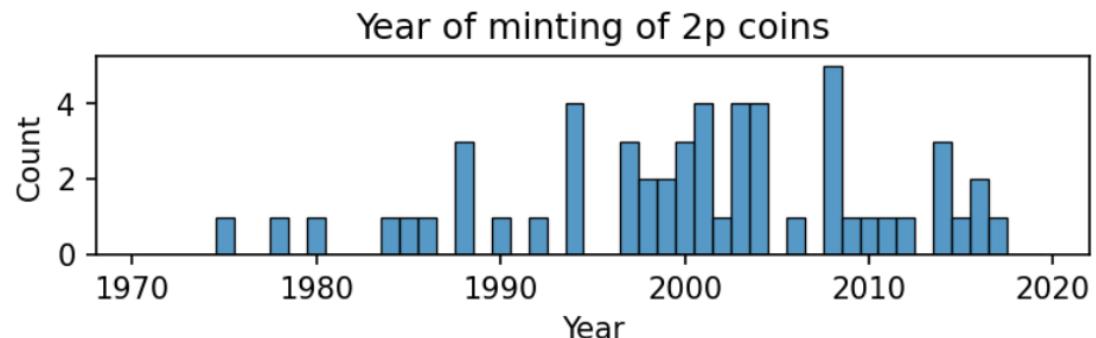
THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Confidence interval for the year of a 2p coin



$n = 56$



In a sample of 56 2p coins, the mean year of minting of the 2p coins is 2000.8 and the sample standard deviation is 10.4. Give a 95% confidence interval for the mean year of minting in the population of all 2p coins.

Practice more in this week's workshop sheet

Solution

Reporting confidence intervals

(1998 , 2004)

$$M = 2001, CI = 1998 - 2004 \quad (95\% CI)$$

$$\hat{\mu} = 2001 \pm 3 \quad (95 \% CI)$$

$$\hat{\mu} = 2006.1 \pm 1.4 \quad (\text{Mean} \pm 1. SEM)$$

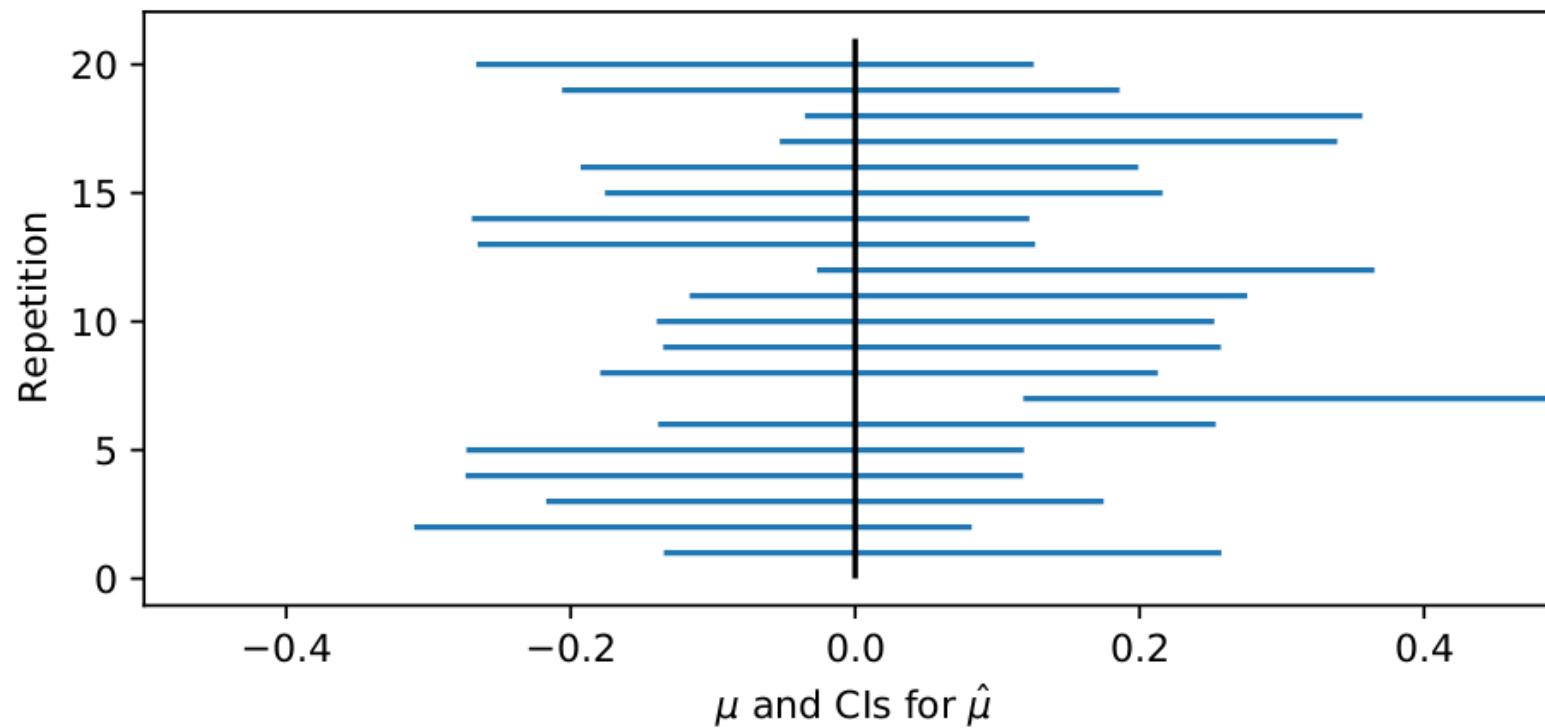
Inf2 - Foundations of Data Science: Estimation - Interpretation of confidence intervals



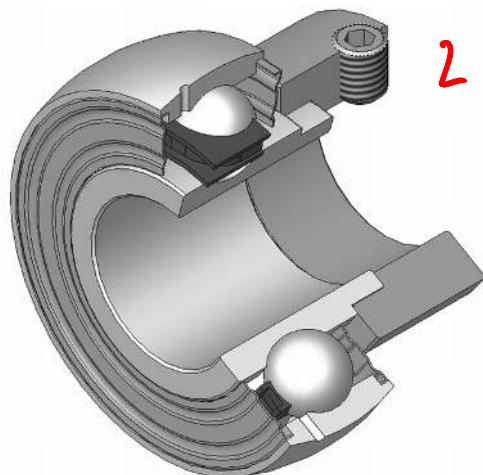
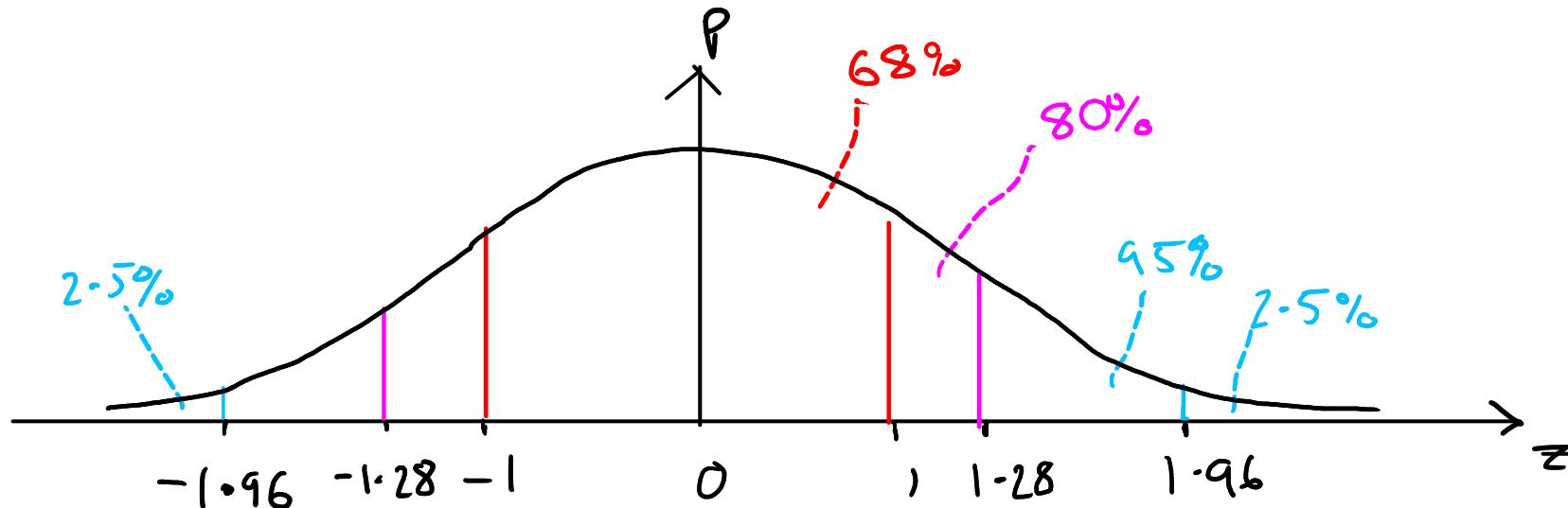
THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Confidence intervals are a random interval



What level of confidence should we choose?



Wikimedia commons, Silberwolf, CC BY 2.5

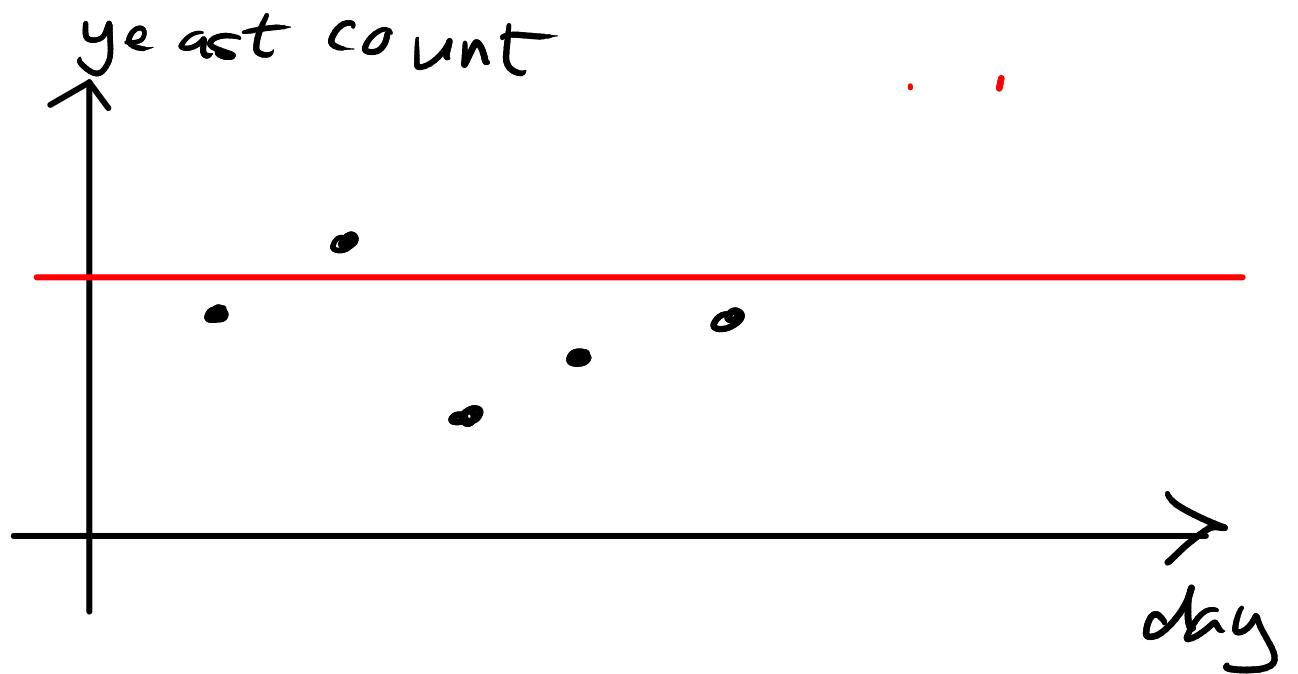


Wikimedia commons, Kiefer, CC BY SA 2.0



How much data do we collect?

A question inspired by the work of "STUDENT" (aka W. S. Gosset) in a brewery



By Satirdan kahraman - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=153514719>

Suppose we only want 1% of beer to be too bitter
What level of confidence should we have?
How many samples per day should we make?

Inf2 - Foundations of Data Science: Estimation - Bootstrapping



THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Principle of bootstrapping



- Treat the sample like a population
- Resample estimator from it to get sampling distribution
- Sample is similar to population for a large sample

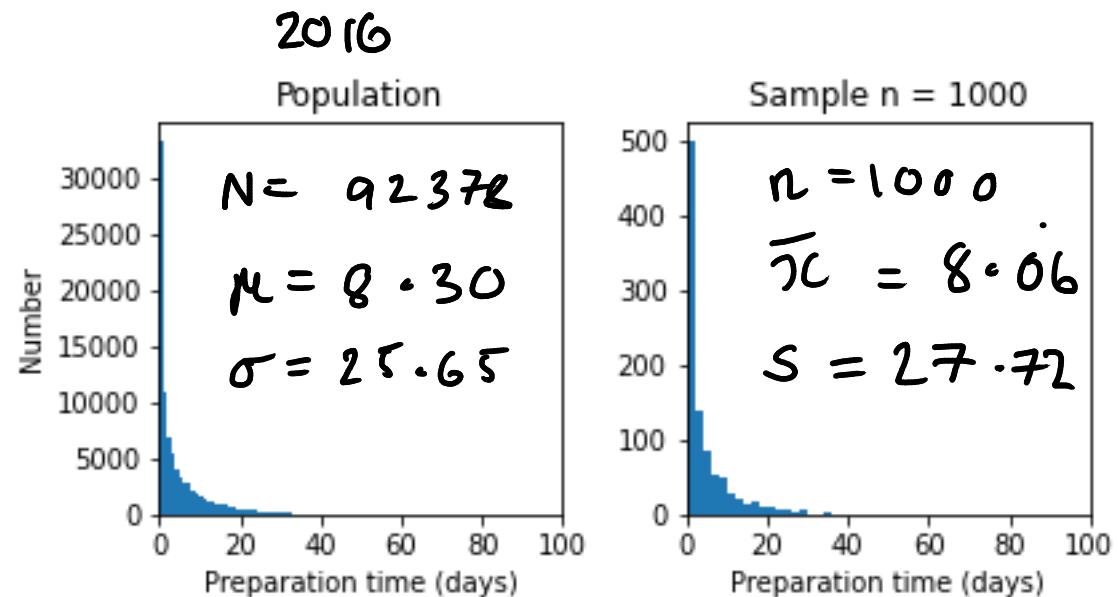
Related lab on the bootstrap

E.g. Japanese restaurant reservation times



Mstyslav Chernov, Wikimedia Commons, CC BY SA 3.0

"Preparation time"
= Time of reservation
- Time reservation made



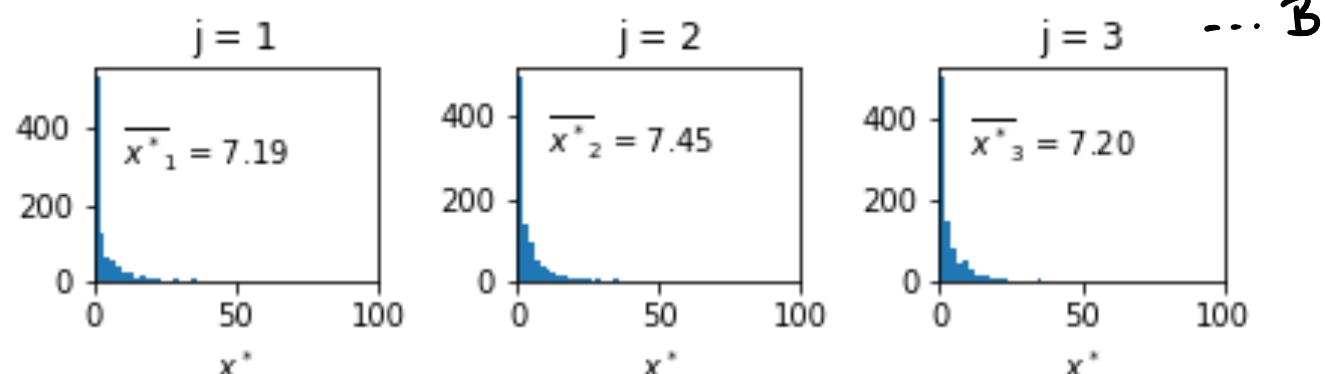
	Population	Sample
count	92378.00	1000.00
mean	8.30	8.06
std	25.65	27.72
min	0.00	0.00
25%	0.21	0.17
50%	2.08	1.96
75%	7.88	6.92
max	393.12	364.96

Bootstrap confidence interval for the mean

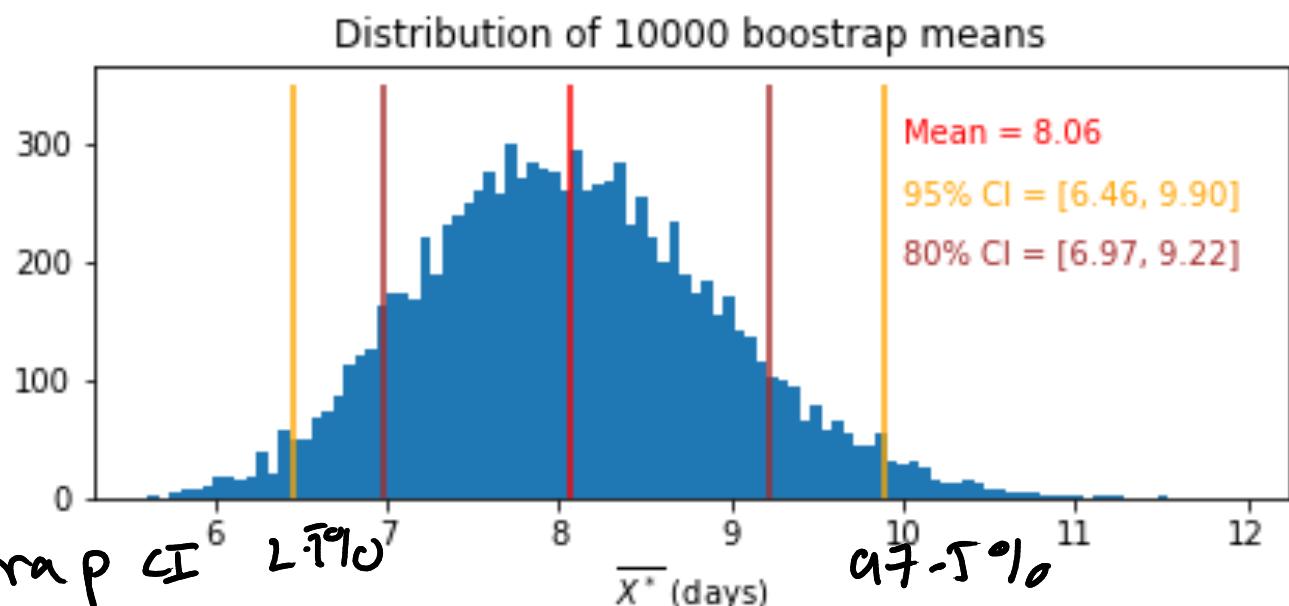
$n=1000$ x # Bootstrap samples

for $j = 1, \dots, B$
 x^* of size n
 from x
 with replacement

$$\bar{x}_j^* \leftarrow \text{mean } x^*$$



$$S_{\text{boot}}^2 = \frac{\sum_{j=1}^B (\bar{x}_j^* - \bar{x})^2}{B-1}$$



(6.46, 9.90) - Bootstrap CI 95%

(6.34, 9.78) - Normal approx

General formulation of the bootstrap

Bootstrap CI. $\hat{\theta}$ $\stackrel{\sim}{\rightarrow}$ $\begin{matrix} \hat{\mu} \\ \hat{\sigma}^2 \\ \hat{\beta} \end{matrix}$

- For $j \in 1, \dots, B$

- Sample n items from x with replacement

- Compute sample stat of the new sample $\hat{\theta}_j^*$

- Bootstrap estimator of variance of statistic

$$s_{\text{boot}}^2 = \frac{\sum_{j=1}^B (\hat{\theta}_j^* - \hat{\theta})^2}{B-1}$$

- Find CI from Bootstrap dist.

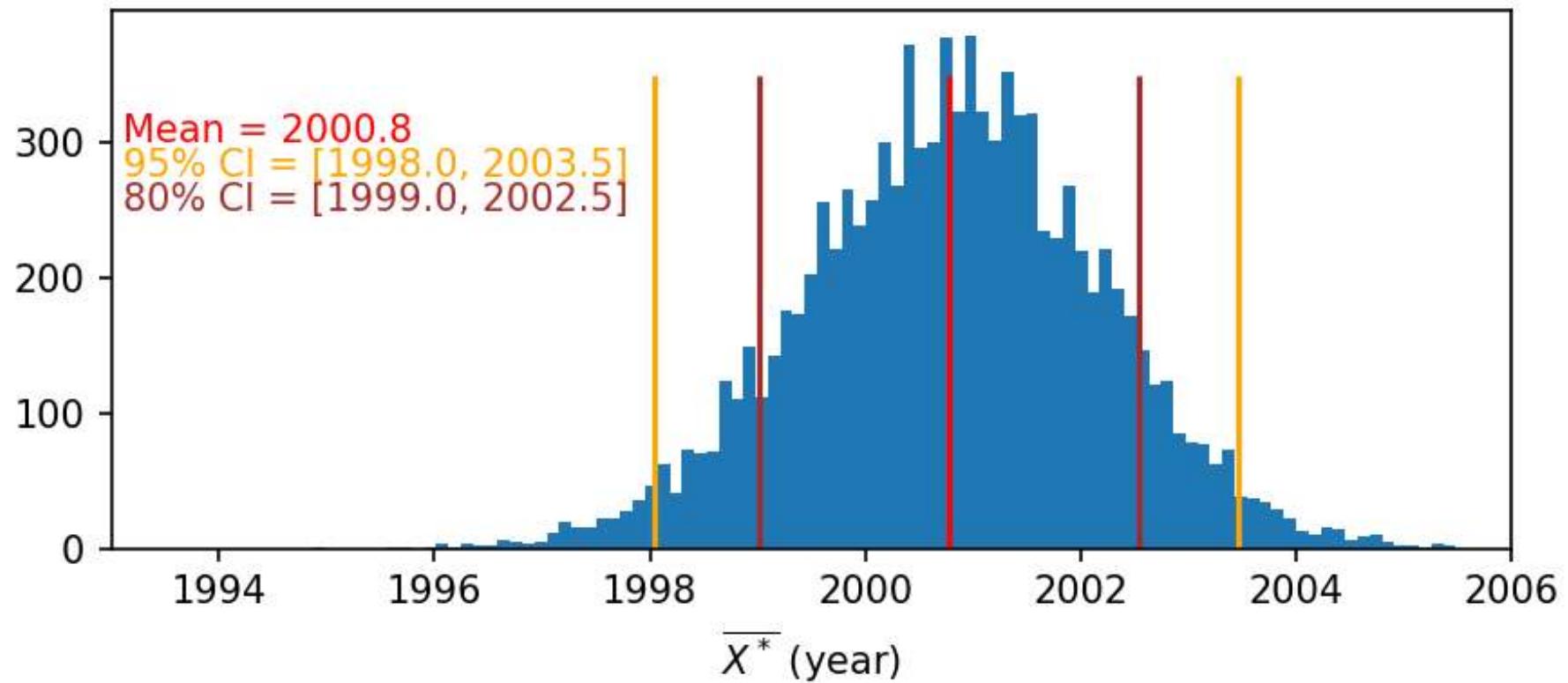
✓ Centrality - median mean

✗ Extremes - max or min

Bootstrap mean coin year



Distribution of 10000 bootstrap means

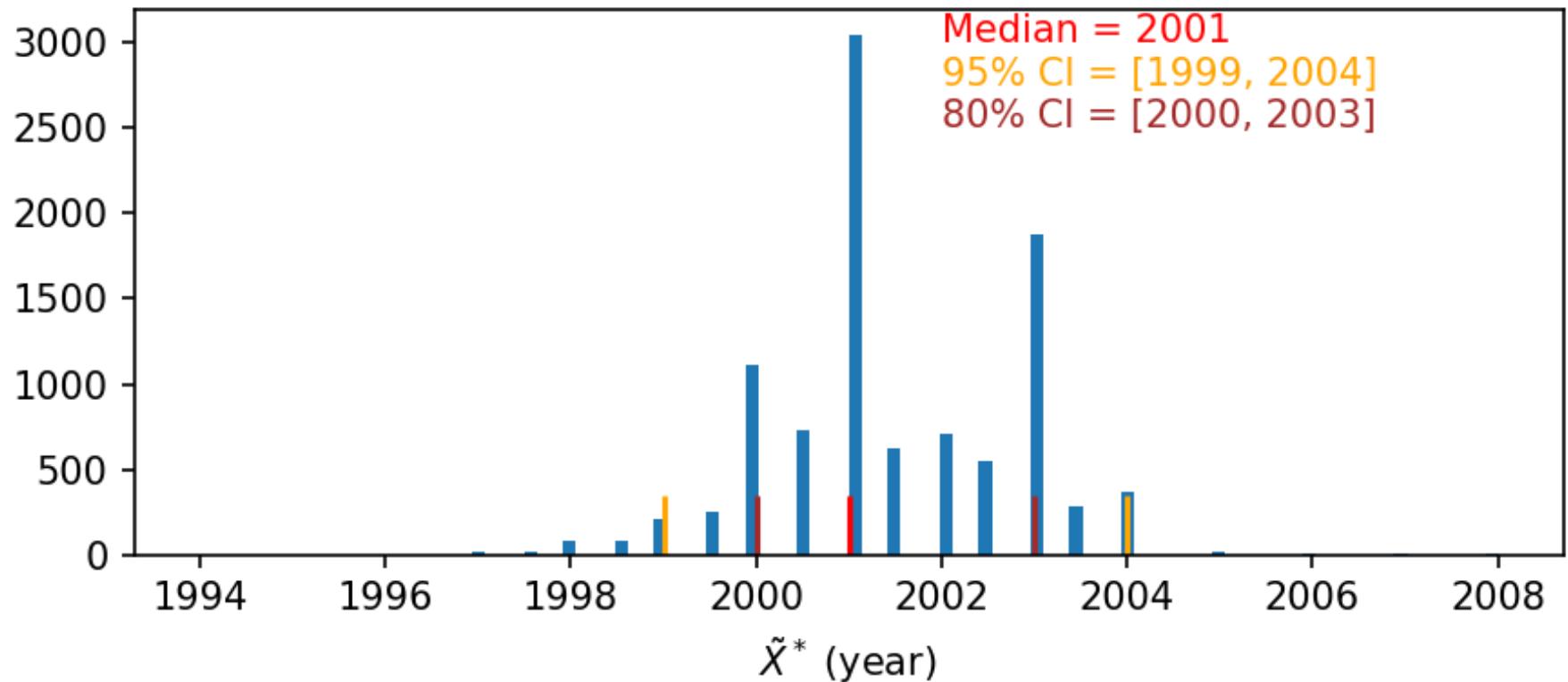


C.f theoretical estimate
(1998.1, 2003.5)

Bootstrap median coin year



Distribution of 10000 bootstrap means



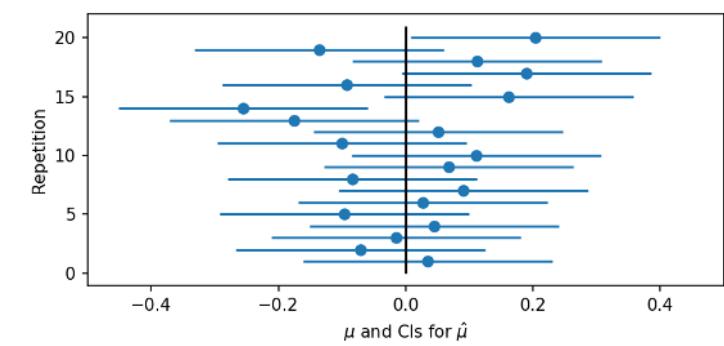
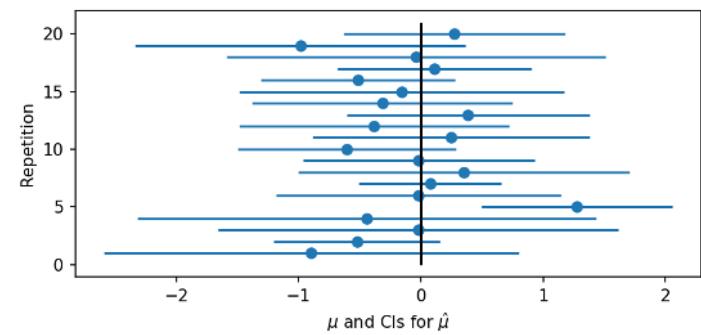
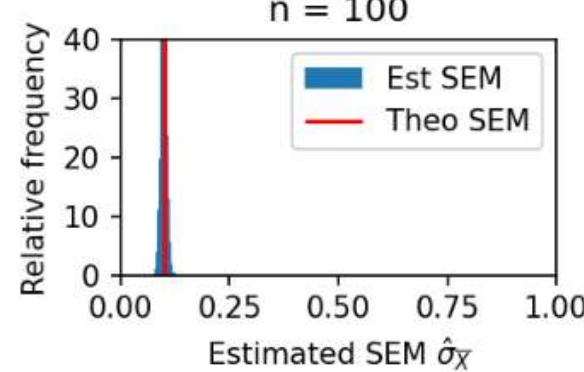
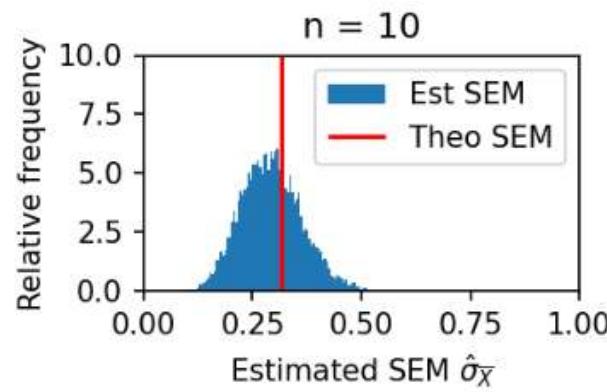
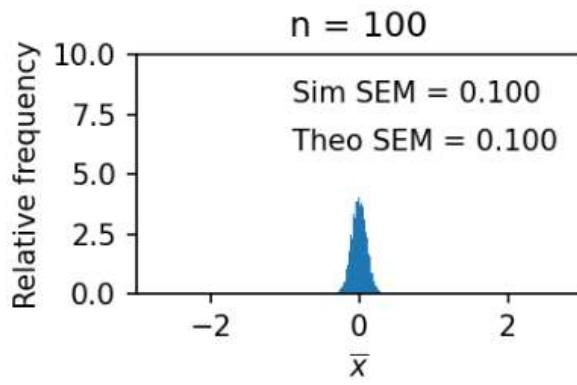
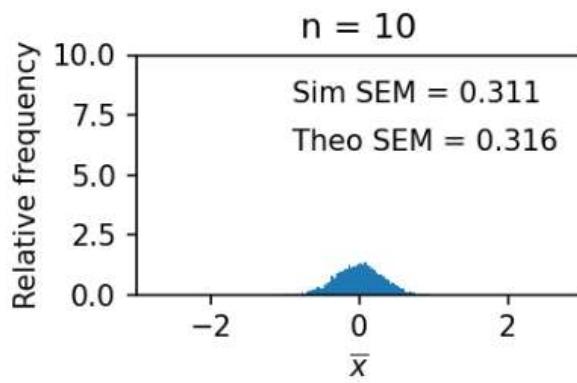
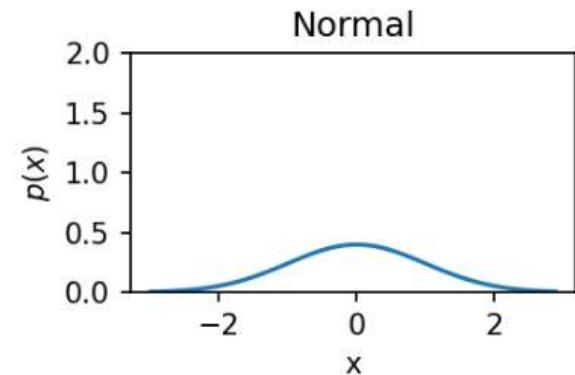
Inf2 - Foundations of Data Science: Estimation - Confidence intervals for the mean for small samples



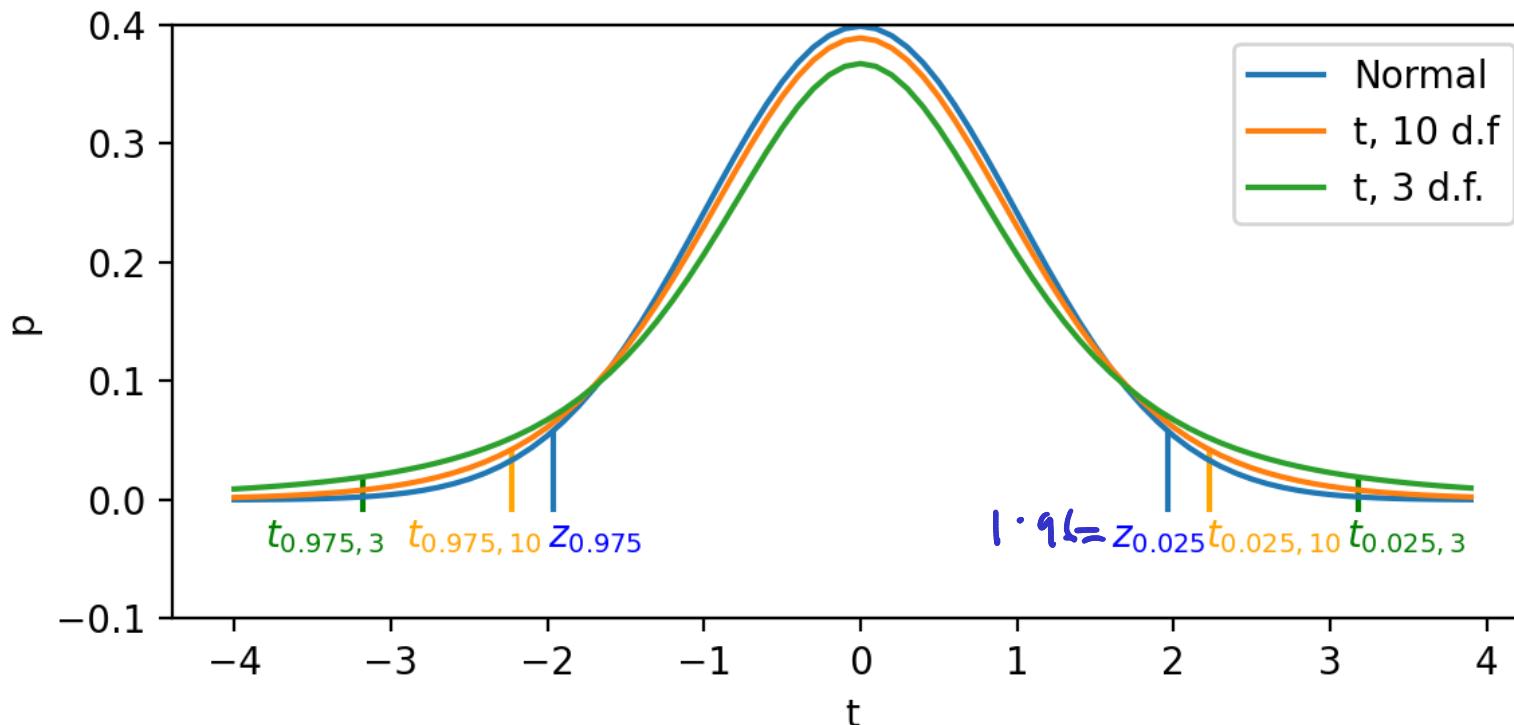
THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

Small samples



The t-distribution



$$T = \frac{\bar{x} - \mu}{\hat{\sigma}_{\bar{x}}}$$

t -critical value $t_{\alpha, \nu}$: value of t at which, in a t -distribution with ν degrees of freedom has an area under the curve of α to its right.

Small samples

$$n \leq 40$$

Number of degrees of freedom $v = n - 1$

Example

$$n = 29 \text{ coins}$$

$$\bar{x} = 2001.551 \text{ years} \quad s = 11.444 \text{ years}$$

$$\text{Estimated SEM, } \sigma_{\bar{x}} = \frac{s}{\sqrt{n}} = \frac{11.444}{\sqrt{29}} = 2.125 \text{ years}$$

$$\hat{\mu} = \bar{x}$$

$$\underline{t - \text{statistic}} \quad T = \frac{\bar{x} - \mu}{\hat{\sigma}_{\bar{x}}}$$

Using the t-distribution to calculate a confidence interval

$$95\% \text{ C.I.} \Rightarrow \alpha = 0.05$$

$$\text{Sample size } n \Rightarrow v = n-1 \text{ d.f.}$$

$$t_{\alpha/2, v} = t_{\alpha/2, n-1} \text{ t-critical value}$$

$$\bar{x} - t_{\alpha/2, n-1} \hat{\sigma}_{\bar{x}}, \bar{x} + t_{\alpha/2, n-1} \hat{\sigma}_{\bar{x}}$$

$$n=29, \alpha=0.05 \Rightarrow t_{0.025, 29-1} = t_{0.025, 28} = 2.281$$

$$t_{0.025, 28} \hat{\sigma}_{\bar{x}} = 2.281 \times 2.125 = 4.871 \text{ years}$$

$$\Rightarrow \hat{\mu} = 2001 \pm 5 \text{ years} \quad (95\% \text{ C.I.})$$

Summary

- .
- 1. Principle and meaning of confidence intervals
- 2. Confidence intervals of the mean of a large samples ($n > 40$)
computed theoretically
 - z distribution
- 3. Confidence intervals for more types of estimator
computed using the bootstrap
- 4. Confidence intervals of the mean of a small sample ($n < 40$)
computed theoretically
 - t distribution

Japanese restaurant confidence interval calculation

$$N = 92372$$

$$\mu = 8.30$$

$$\sigma = 25.65$$

$$n = 1000$$

$$\bar{x} = 8.06$$

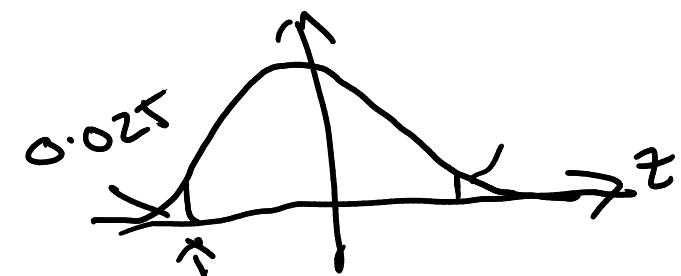
$$s = 27.72$$

Estimated SEM $\hat{\sigma}_{\bar{x}} = \frac{s}{\sqrt{n}} = \frac{27.72}{\sqrt{1000}} = 0.88$ days

Large sample \Rightarrow Normal distribution of sample mean \Rightarrow "z" distribution

$$95\% \Rightarrow \alpha = 0.05$$

$$\alpha/2 = z_{\alpha/2} = z_{0.025} = 1.96$$



$$(\bar{x} - z_{0.025} \hat{\sigma}_{\bar{x}}, \bar{x} + z_{0.025} \hat{\sigma}_{\bar{x}}) = (\underline{6.34}, \underline{9.78})$$

Reporting confidence intervals

$$\underline{(6.34, 9.78)}$$

$$M = 8.06, CI = 6.34 - 9.78 \quad (95\% CI)$$

$$\hat{\mu} = 8.06 \pm 1.72 \quad (95\% CI)$$

$$\text{with } z_{0.025} \frac{\hat{\sigma}_x}{\sqrt{n}} = 1.96 \times 0.88$$

$$\hat{\mu} = 8.06 \pm 0.88 \quad (\text{Mean} \pm 1. SEM)$$