

Inf2 - Foundations of Data Science:

Estimation -

Point estimation and confidence intervals



THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

Plan for statistical inference

- ✓ 1. Randomness, sampling and simulations (S2 Week 1)
- 2. Estimation, including confidence intervals (S2 Week 2)
- 3. Hypothesis testing (S2 Week 3)
- 4. Logistic regression (S2 Week 3)
- 5. A/B testing (S2 Week 4)

Last lecture...

1. Sampling

- random
- non-random

2. Inference on testing the hypothesis that the coin is biased

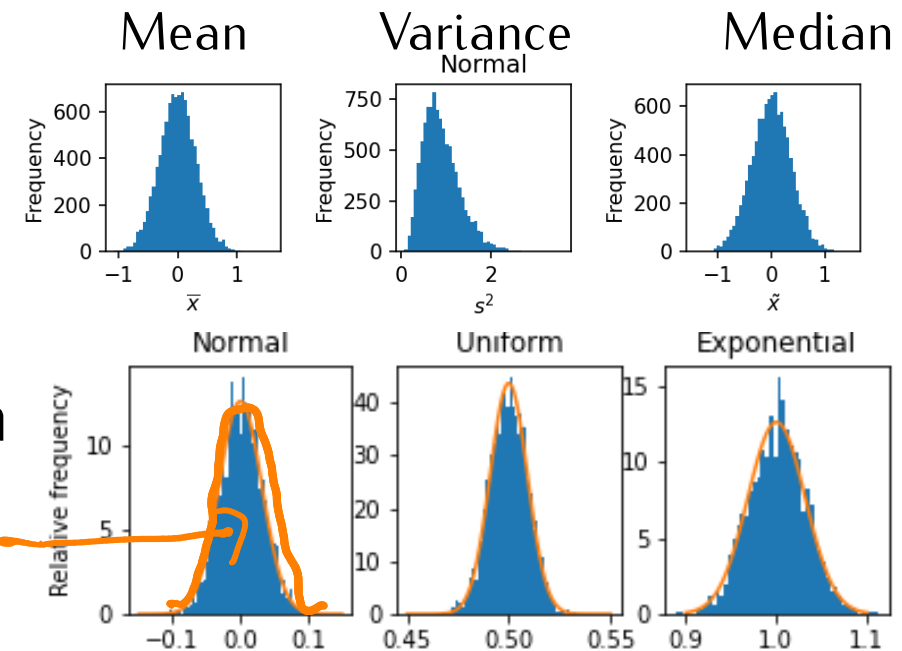
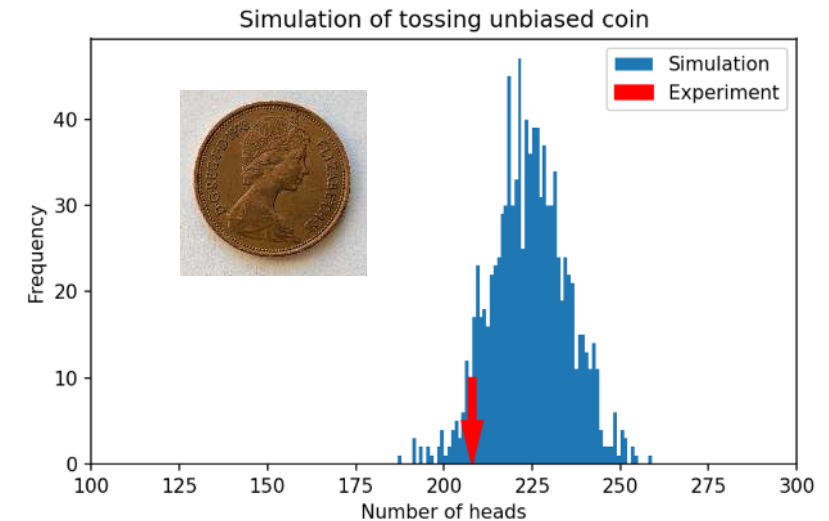
- Statistical simulations

3. Sampling distributions of statistics

- mean, variance, median

4. Sampling distribution of the mean in large samples

- Central Limit Theorem



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

Today

- Big idea: method to determine how precise our estimate of the average age of 2p coin is
- Steps:
 1. Concepts of parameters and estimators
 2. Sampling distribution of the estimator gives indication of uncertainty in estimate
 3. First attempt at theoretical method of getting confidence interval
 - worked example
 4. Theory on bias and variance of estimators

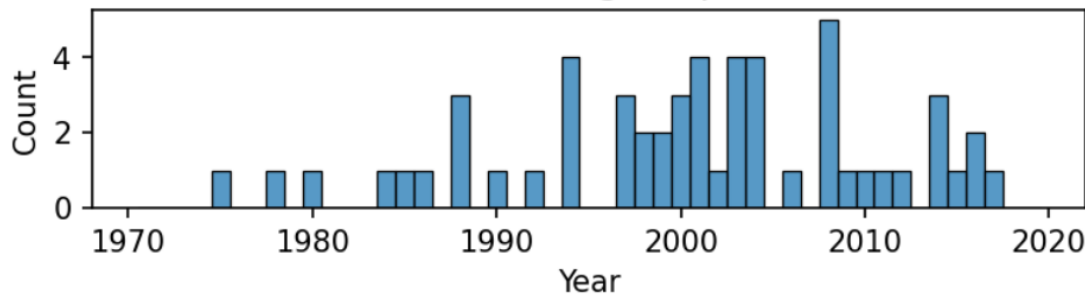
Overview

Sample



$$n = 56$$

Year of minting of 2p coins



$$\bar{x} = 2000.8 \text{ years}$$

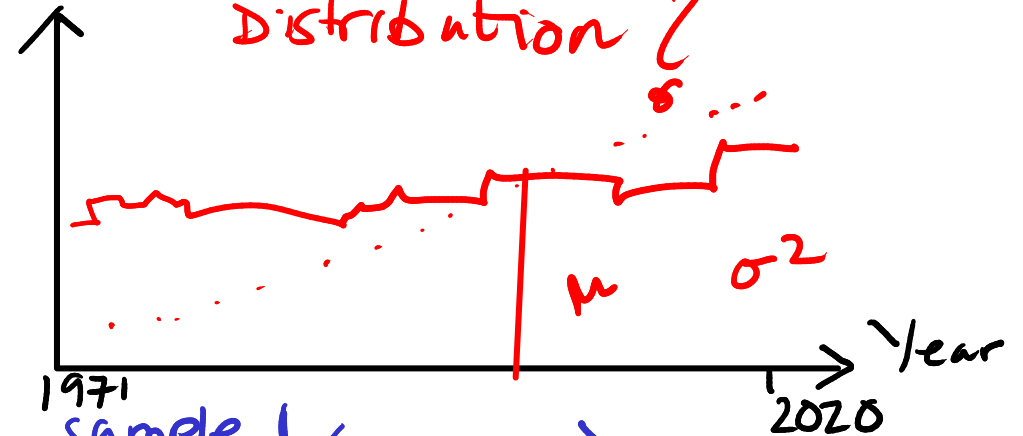
$$s = 10.4 \text{ years}$$

Population

$$N \sim 1 \times 10^9$$

Frequency

Distribution?



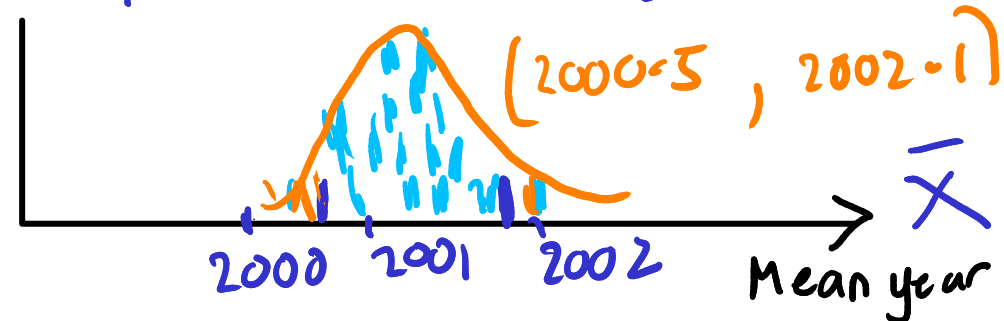
$$n = 56 \times$$



$$n = 56 \times$$



$$\bar{x}_1 = 2000.8 \quad \bar{x}_2 = 2001.8$$

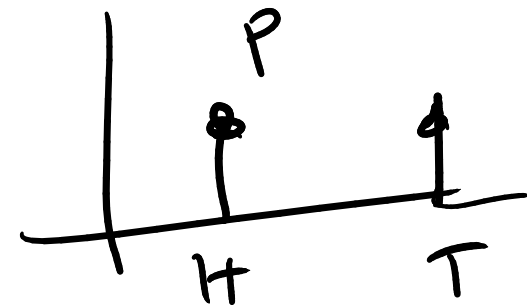
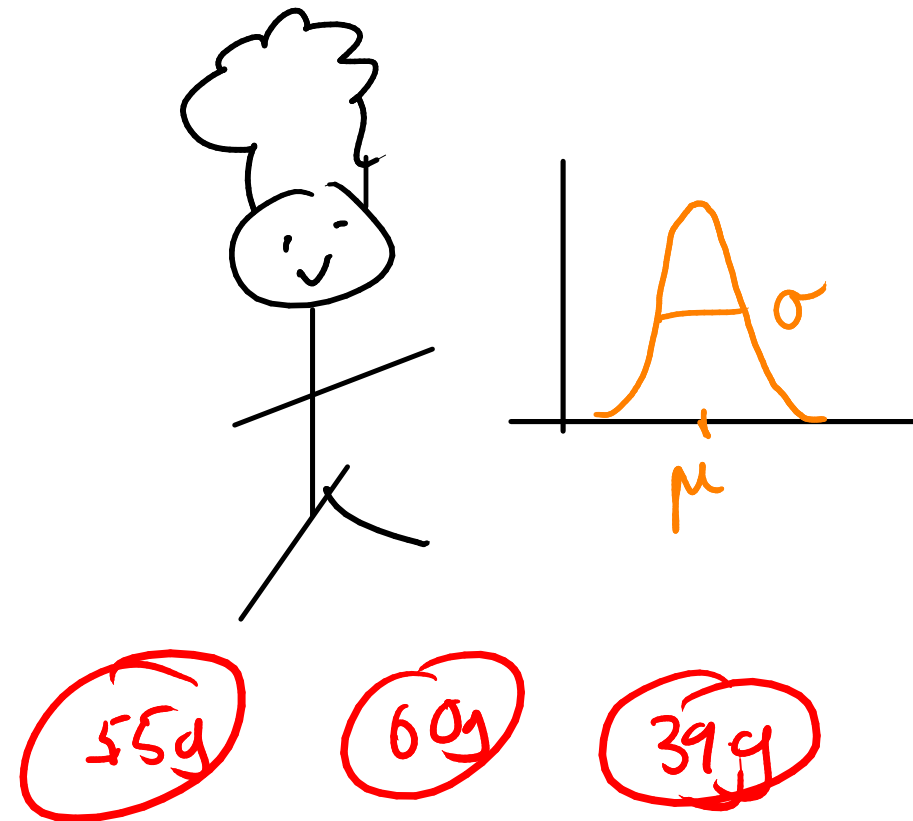


Sampling dist of mean

Non-countable populations



Coin tossing



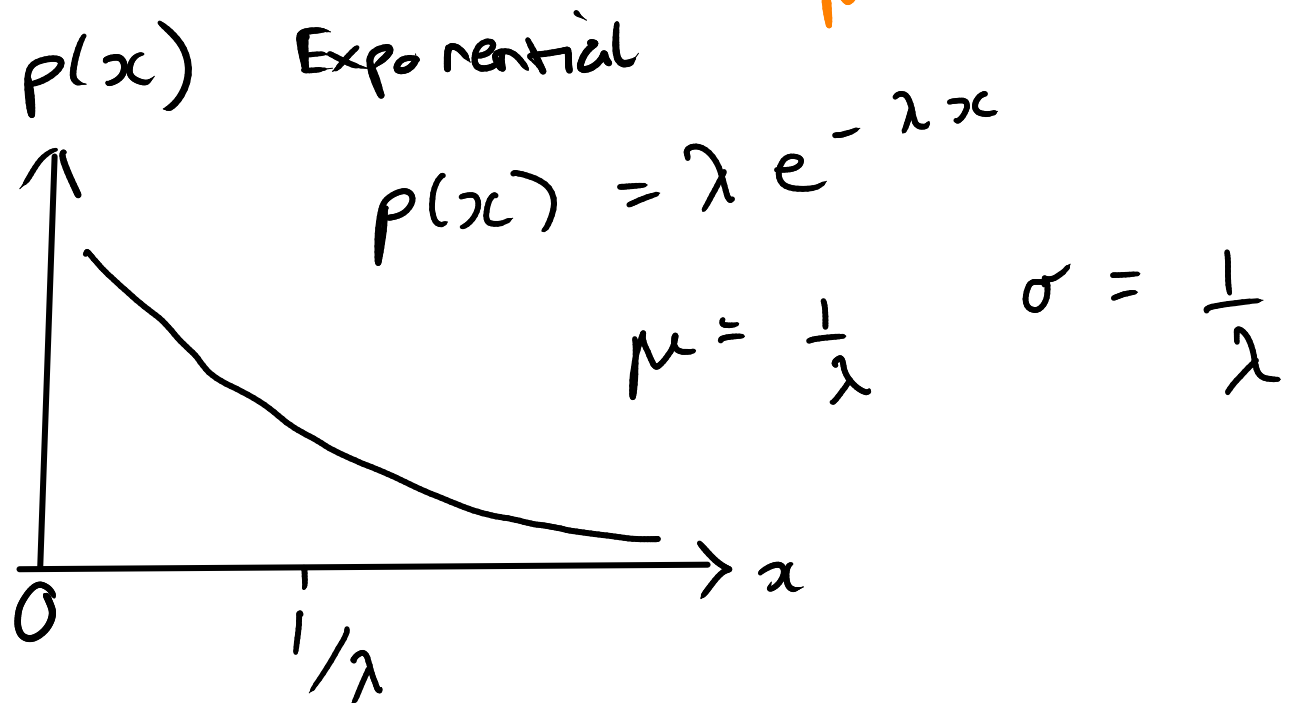
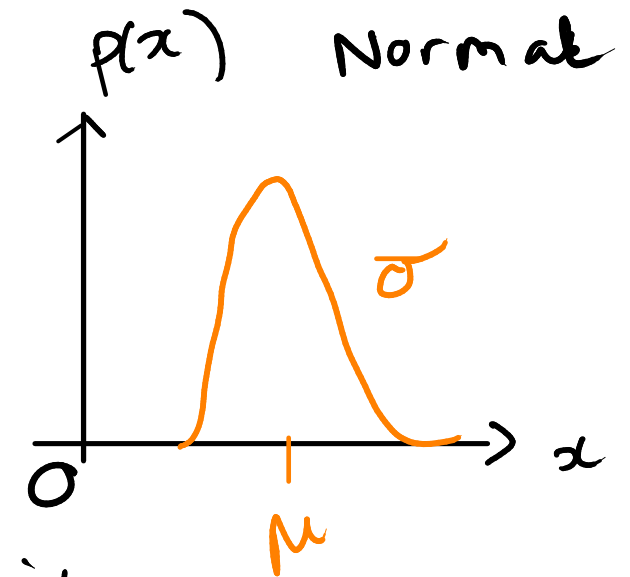
Parameters

Of a finite population



Mean μ
Variance σ^2

Of a distribution



Estimation problems

1. Construct a point estimator for parameter of population

Generic parameter:

$$\theta \begin{cases} \mu \\ \sigma^2 \\ \tilde{\mu} \end{cases}$$

Corresponding estimator:

Estimator is
function of
data

$$\hat{\theta} \begin{cases} \hat{\mu} \\ \hat{\sigma}^2 \\ \hat{\tilde{\mu}} \end{cases} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

2. Determine how accurate that estimate is using confidence intervals

Find $(u(x), v(x))$ s.t. $\underbrace{\quad}_{/95\%}$

$$P(u(x) < \theta < v(x)) = 1 - \alpha$$

Two methods to estimate confidence intervals

1. Theoretical method (start this lecture, finish next)
 - works for some parameters and estimators
2. Statistical sampling (next lecture)
 - works for some parameters and estimators that the theoretical method doesn't work for

Inf2 - Foundations of Data Science:

Estimation -

First attempt at determining confidence intervals for the mean theoretically



THE UNIVERSITY *of* EDINBURGH
informatics

F O U N D A T I O N S
O F
D A T A
S C I E N C E

Coin year point estimation and confidence interval example

What is the point estimator of the mean age of the coins?

Parameter : mean age of coins in the pop. of all coins

$$\theta = \mu = \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad x_i - \text{year of coin } i$$

Estimator : mean of sample

$$\hat{\theta} = \hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad X_i - \text{" " " } i$$

Knowing what we learned about sampling in the last lecture, what might a good confidence interval of the mean be?

$$(\mu - a \times SEM, \mu + a \times SEM) \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$SEM = \hat{\sigma}_{\hat{\mu}} \quad \text{Estimate } \hat{\sigma}_{\hat{\mu}} = \frac{S}{\sqrt{n}}$$

$$\boxed{\text{NOT } (\mu - a \times S, \mu + a \times S)}$$

Exercise: Suppose we had only 5 samples:

1975, 1984, 1998, 2008, 2014

What is the point estimate of the mean and the SEM?

Point estimate of mean $\hat{\mu} = \bar{X}$

$$= \frac{1}{5} \sum_{i=1}^5 (1975 + 1984 + 1998 + 2008 + 2014) = 1995.8$$

$$S^2 = \frac{1}{5-1} \left[\sum_{i=1}^5 (1975^2 + 1984^2 + 1998^2 + 2008^2 + 2014^2) - 5 \times 1995.8^2 \right]$$

$$= 264 \cdot 200$$

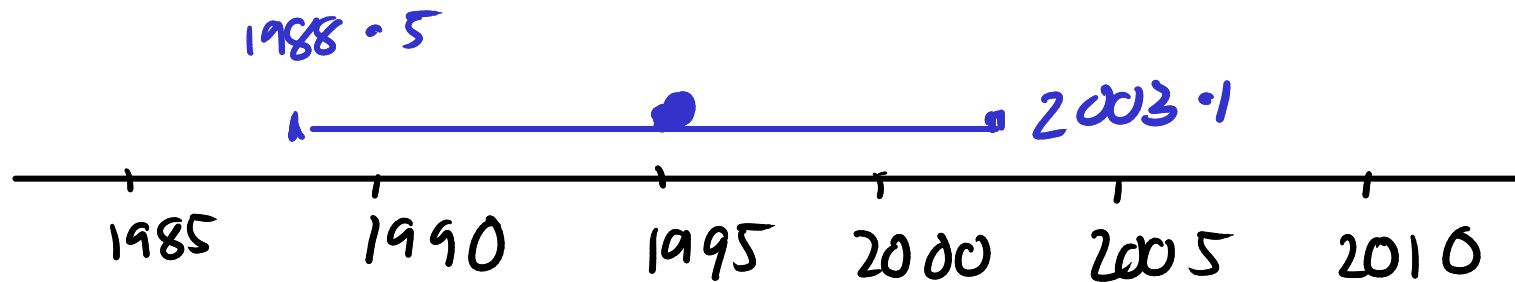
$$\Rightarrow \text{Sample st dev } S = \sqrt{264 \cdot 200} = 16 \cdot 254$$

$$\text{SEM } \frac{S}{\sqrt{n}} = \frac{S}{\sqrt{5}} = \frac{16 \cdot 254}{\sqrt{5}} = 7 \cdot 269 \simeq 7 \cdot 3$$

$$\text{CI: } 1995.8 \pm 7.3 \quad \text{or} \quad (1988.5, 2003.1)$$

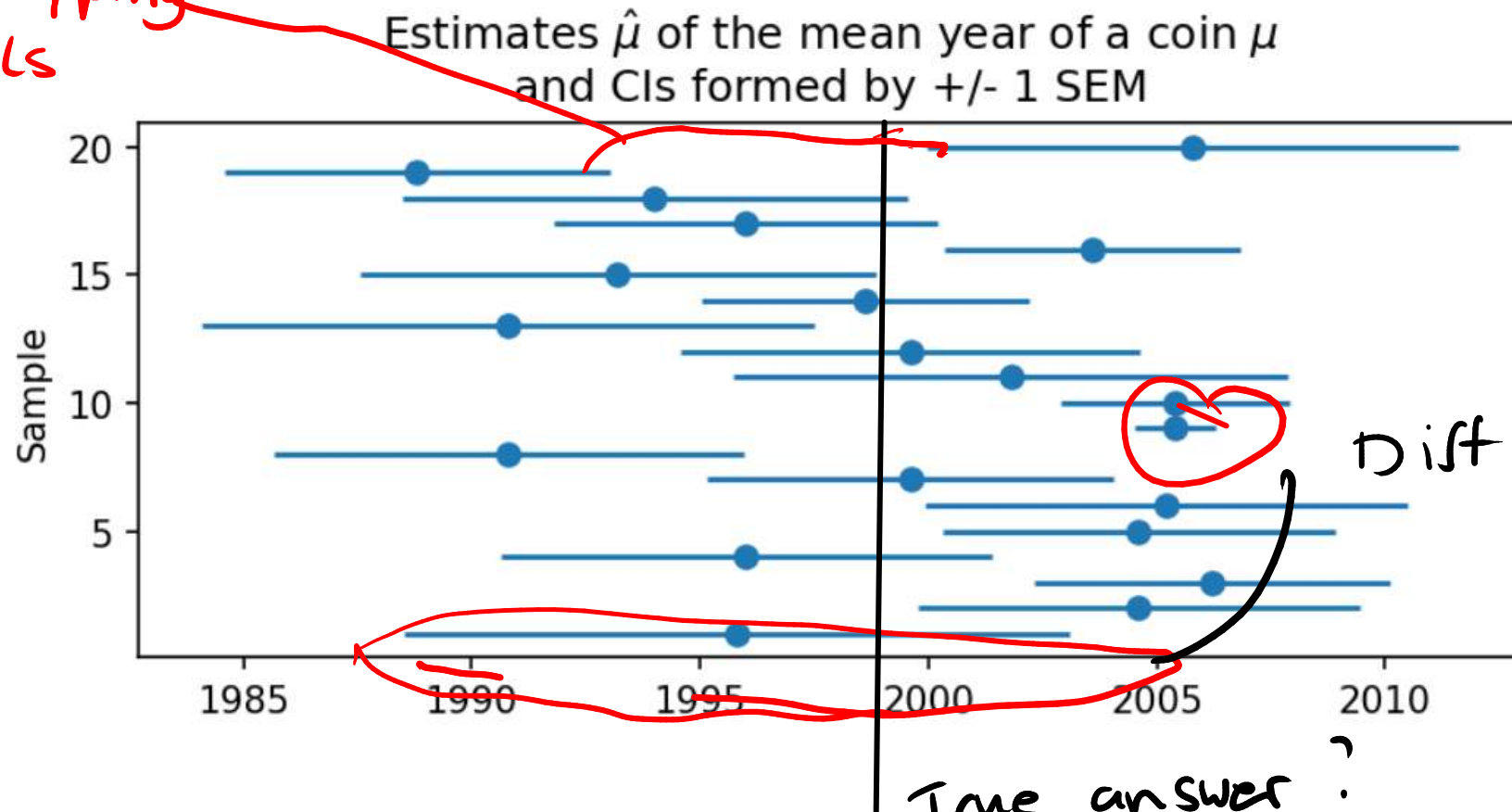
If you know how to do the above, try computing the same for:
2003, 2004, 2006, 2006, 2008

Graphical representation of confidence intervals



Confidence intervals from different samples

Non overlapping
intervals



Intervals are random variables.

Remaining question for next lecture

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

Inf2 - Foundations of Data Science: Estimation - Bias and variance

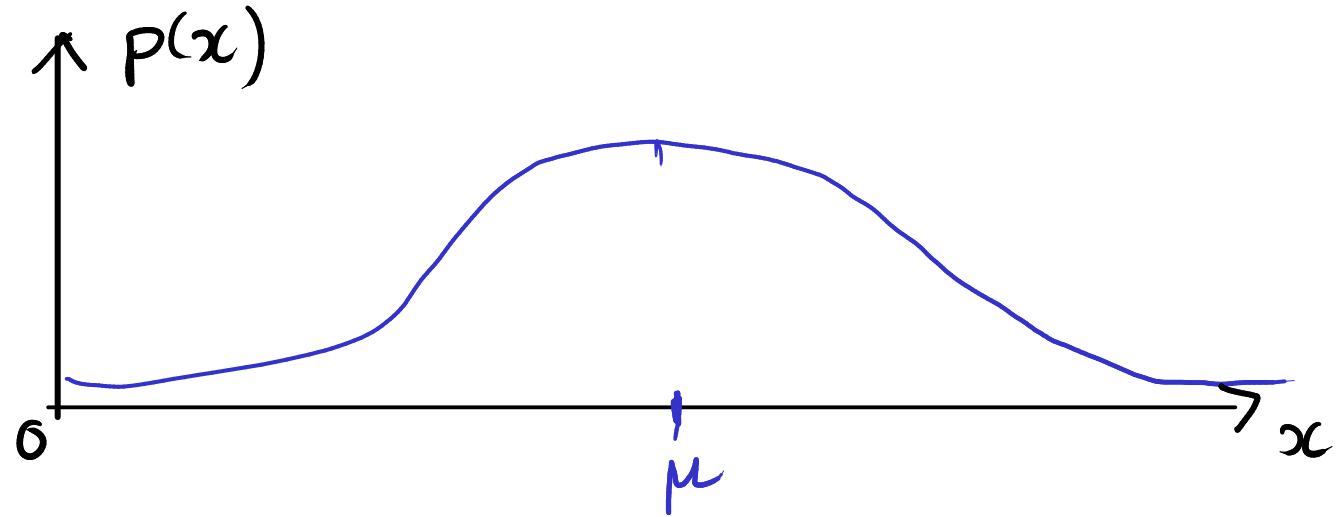


THE UNIVERSITY *of* EDINBURGH
informatics

FOUNDATIONS
OF
DATA
SCIENCE

There can be more than one estimator for a parameter

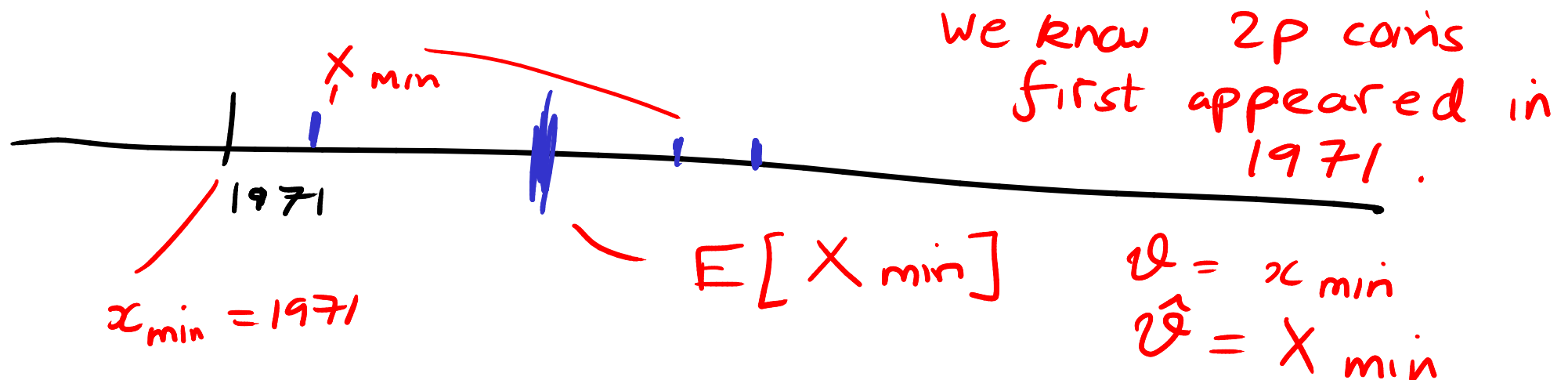
E.g. Symmetric distribution:



$$\begin{aligned}\hat{\mu} &= \bar{x} \\ &= \tilde{x}\end{aligned}$$

How would we estimate the first year 2p coins were minted (made) from the population of years of 2p coins?

- Would there be any problems with your estimate?



NOTE ADDED AFTER LECTURE :

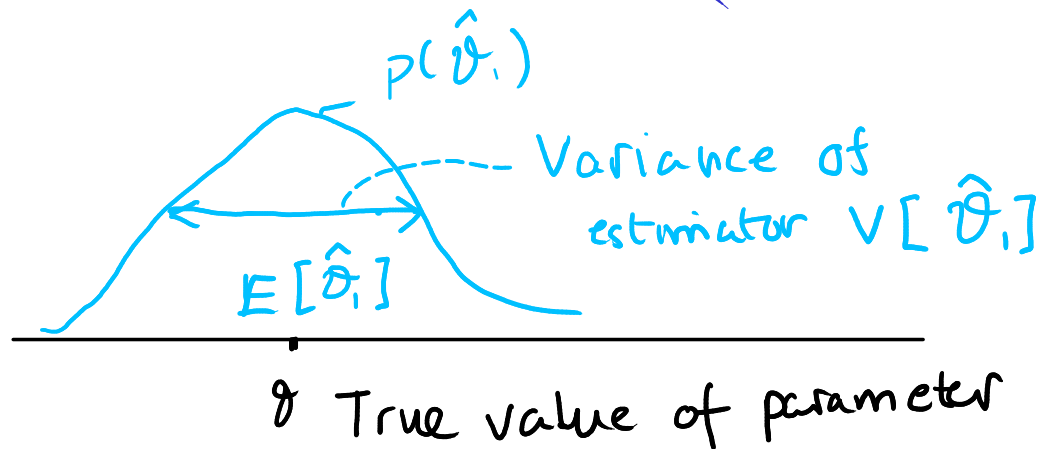
- The estimator is the minimum value of the years
- In some samples the earliest year may not appear
- E.g. we may estimate 1973 or 1980

\Rightarrow On average $E[X_{\min}] - x_{\min} > 0$

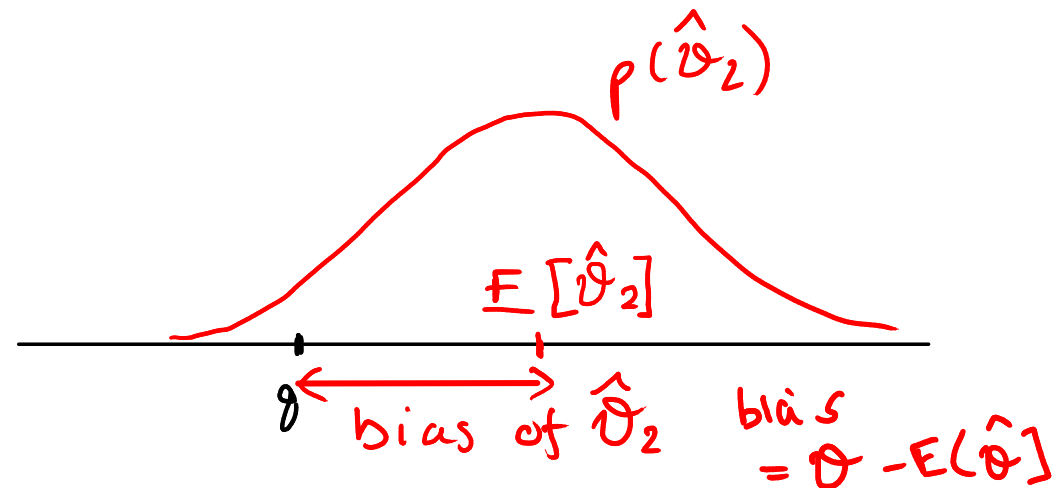
\Rightarrow Biased estimator

Estimation bias and variance

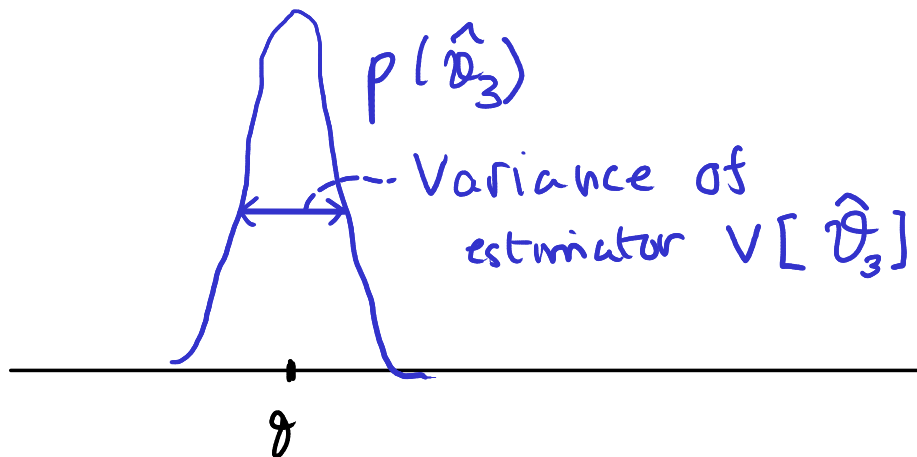
Unbiased estimator



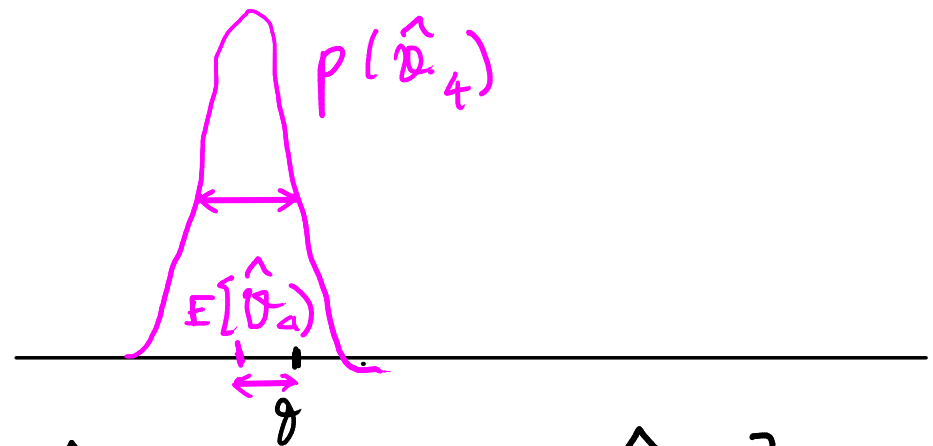
Biased estimator



Unbiased estimator with low variance



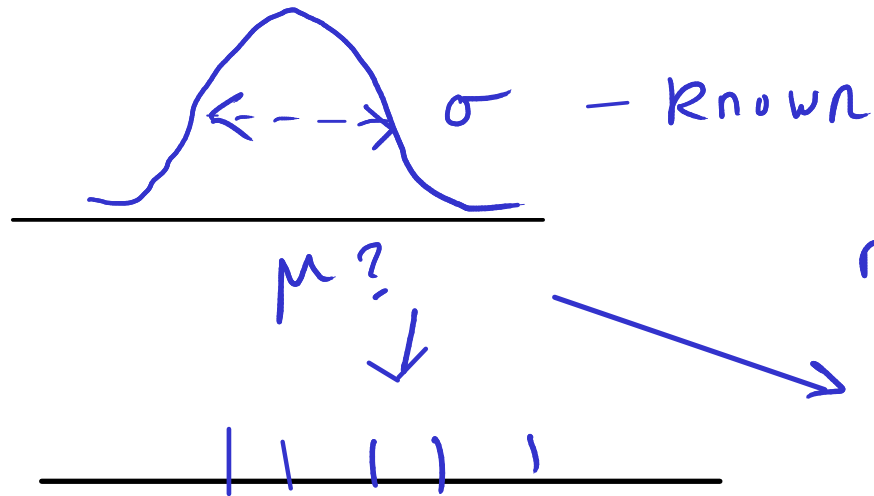
Biased estimator with low variance



$$MSE = E[(\theta - \hat{\theta})^2] = \underbrace{V[\hat{\theta}]}_{\text{variance}} + \underbrace{(\theta - E[\hat{\theta}])^2}_{\text{bias}}$$

Example: estimator of mean of normal distribution with known variance

Normal distribution



Estimator : $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

$n=5$



ADDED AFTER LECTURE :

$$E[\bar{X}] = \mu$$

$$\Rightarrow E[\bar{X}] - \mu = 0$$

$\Rightarrow \bar{X}$ is an unbiased estimator

Mean Square Error $E[(\bar{X} - \mu)^2] = V[\bar{X}] = \frac{\sigma^2}{n}$

Examples: estimator with bias

Contrived estimator: $\hat{\mu} = \bar{X} + 1$

Estimator of variance (see lecture notes):

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \Rightarrow \frac{1}{n-1}$$

See lecture notes

Bias can be helpful! See comprehension questions on bias and variance

Example from machine learning

Suppose

1. We train k-Nearest Neighbours on some training data
2. We measure the accuracy achieved on the training set
3. We use this accuracy as an estimate of accuracy on unseen data

Identify ϑ and $\hat{\vartheta}$

Is $\hat{\vartheta}$ an unbiased estimator of ϑ ?

AFTER LECTURE NOTES:

ϑ - accuracy on unseen data ("true accuracy")

$\hat{\vartheta}$ - accuracy measured on training data

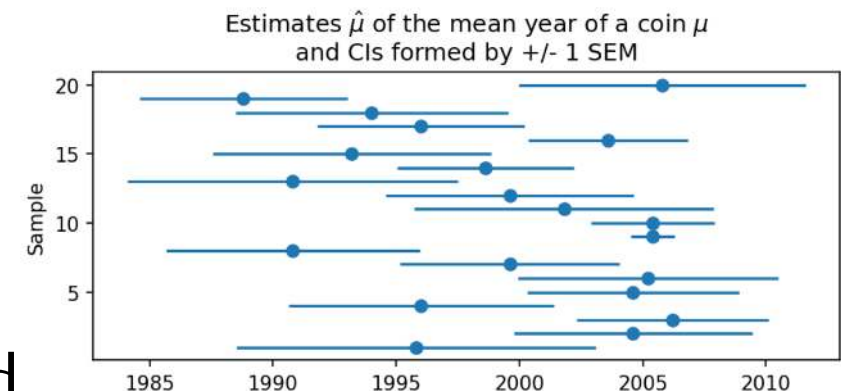
$\hat{\vartheta}$ is a biased estimator - on average it will be too optimistic

Next lecture

We have only one sample.

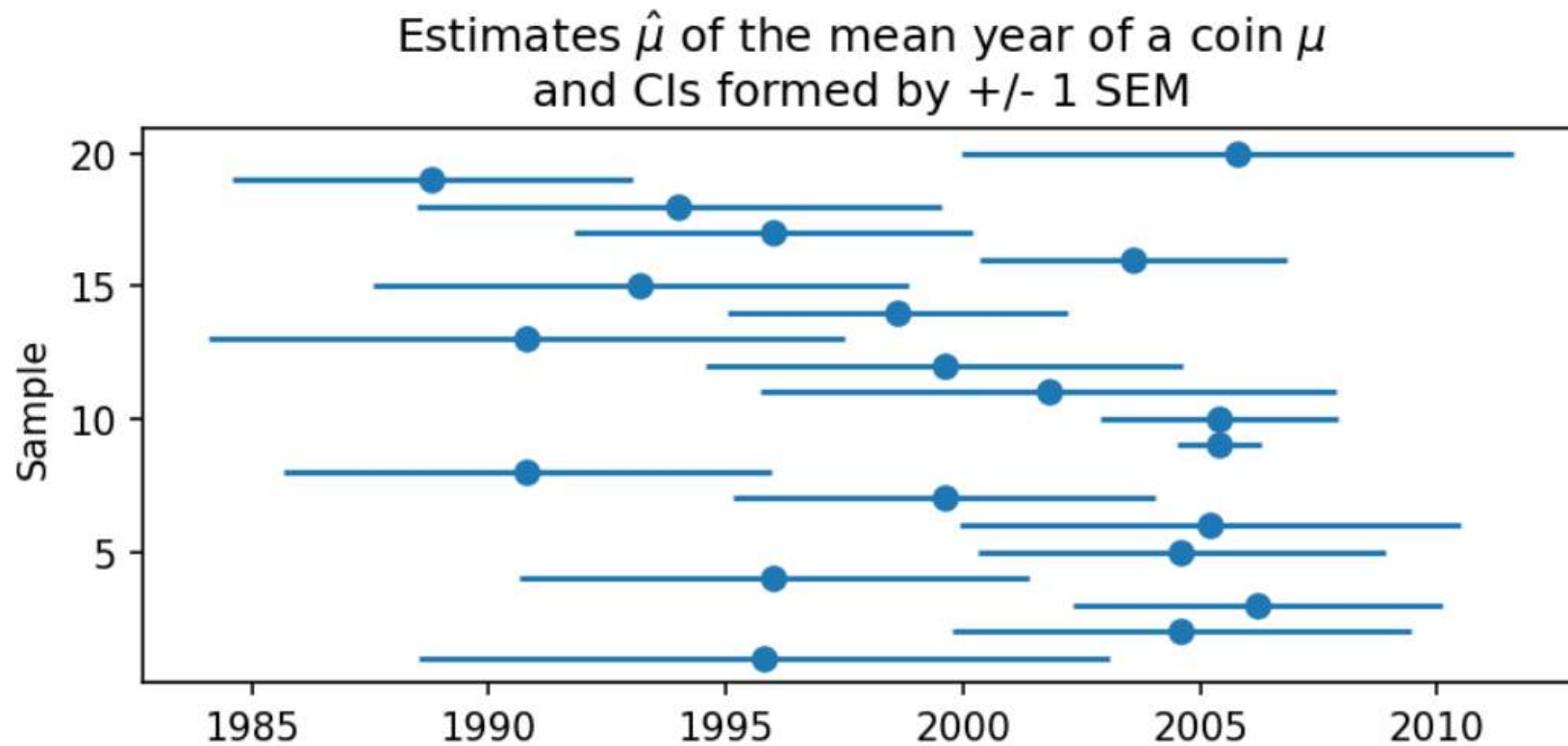
We can't resample from the population to estimate $\forall [\hat{\theta}]$ and hence get a confidence interval

1. For the mean, we can estimate the standard error of the mean using the sample variance of the sample (see above) and adjust to get desired level of confidence



2. For other estimators, we can use the bootstrap method to estimate the distribution of the estimator, and thus the standard error of the estimator

Confidence intervals from different samples



- Non-overlapping intervals
- Different sized intervals.

Summary

1. Progress on estimating the uncertainty in the estimate of the average year of a 2p coin
2. Estimators and parameters
3. Introduction to the confidence interval – theoretical method
4. Bias and variance of estimators