

Confidence Intervals

R-Training workshop, Pwani University
17TH February 2016

Dr. David Mburu (PhD)

Point Estimation:

- Provides a single value, based on observations from one sample
- Gives no information about how close the value is to the unknown population parameter
- Example: Sample mean $\bar{x} = 3$ is point estimate of unknown population mean

Interval estimation:

- Provides a range of values based on observations from one sample
- Gives information about closeness to unknown population parameter,
Stated in terms of probability,
Knowing exact closeness requires knowing unknown population parameter
- Example: The range between 50 and 70 contains the true unknown parameter value with 95% confidence]

Importance of confidence intervals

The confidence interval quantifies our uncertainty in an estimate.

A survey of haemoglobin status in children <5yrs in Kilifi district.
30 children gave a finger prick blood sample. Mean Hb was 9.6g/dl

How reliable is this estimate?

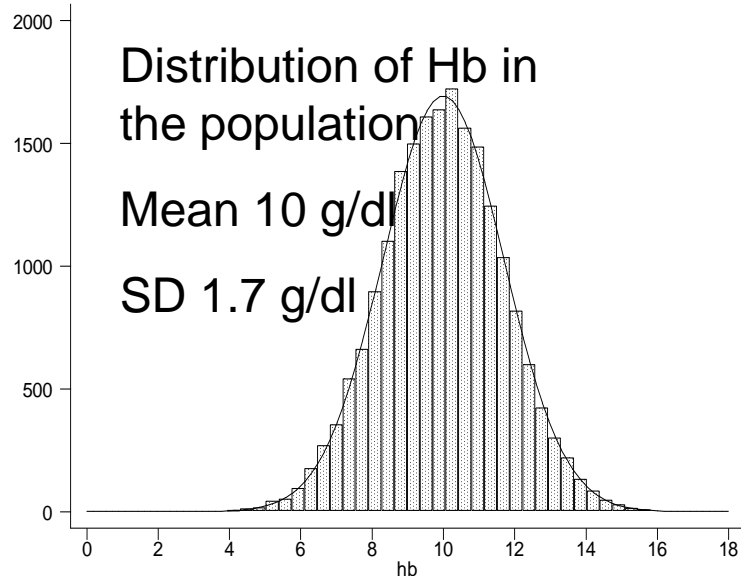
Confidence intervals quantify the magnitude of treatment benefit in a clinical trial

The confidence interval is the range of plausible values for the true vaccine efficacy

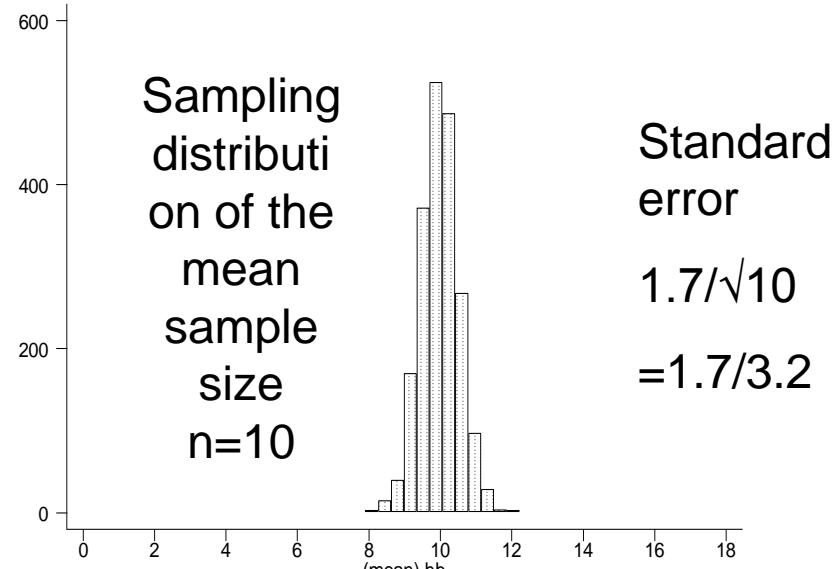
Design studies to ensure confidence intervals are narrow enough to draw conclusions. The confidence interval for vaccine efficacy is too wide to be useful. It includes useful benefit (52%) and no effect (0%) as plausible values

The **larger** the study the **narrower** the confidence interval

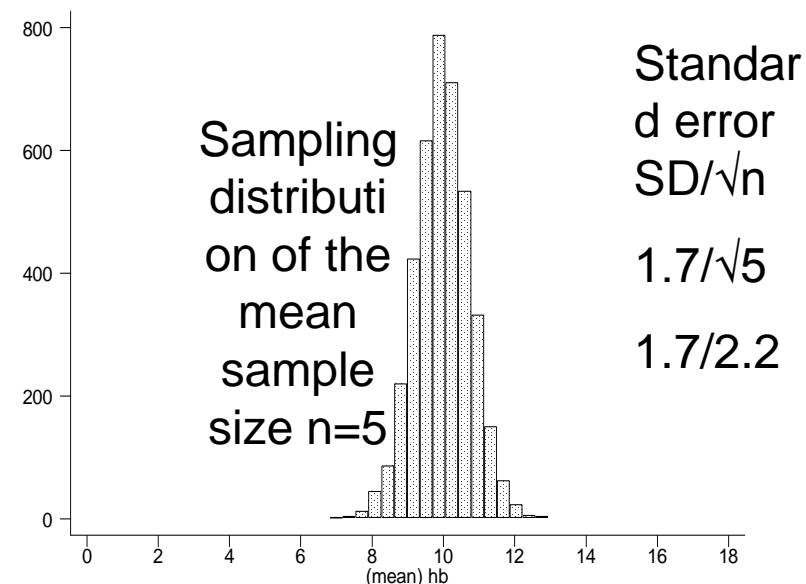
Frequency



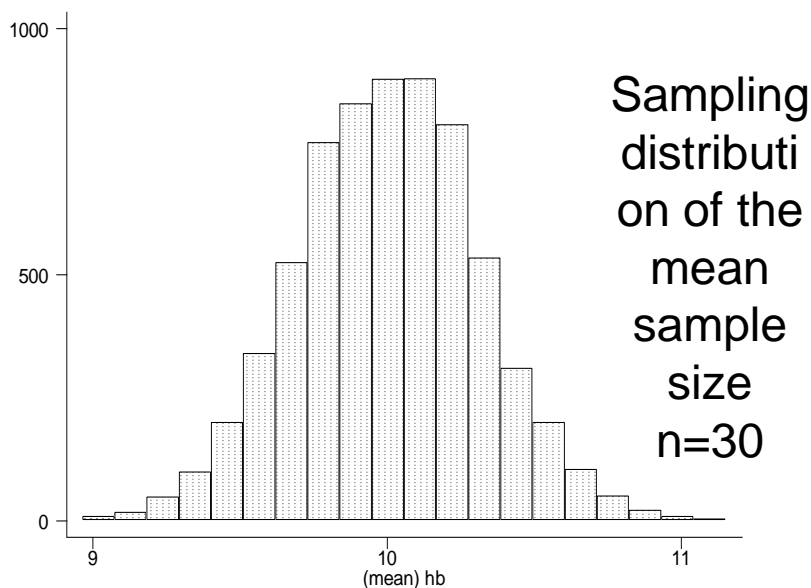
Frequency



Frequency



Frequency

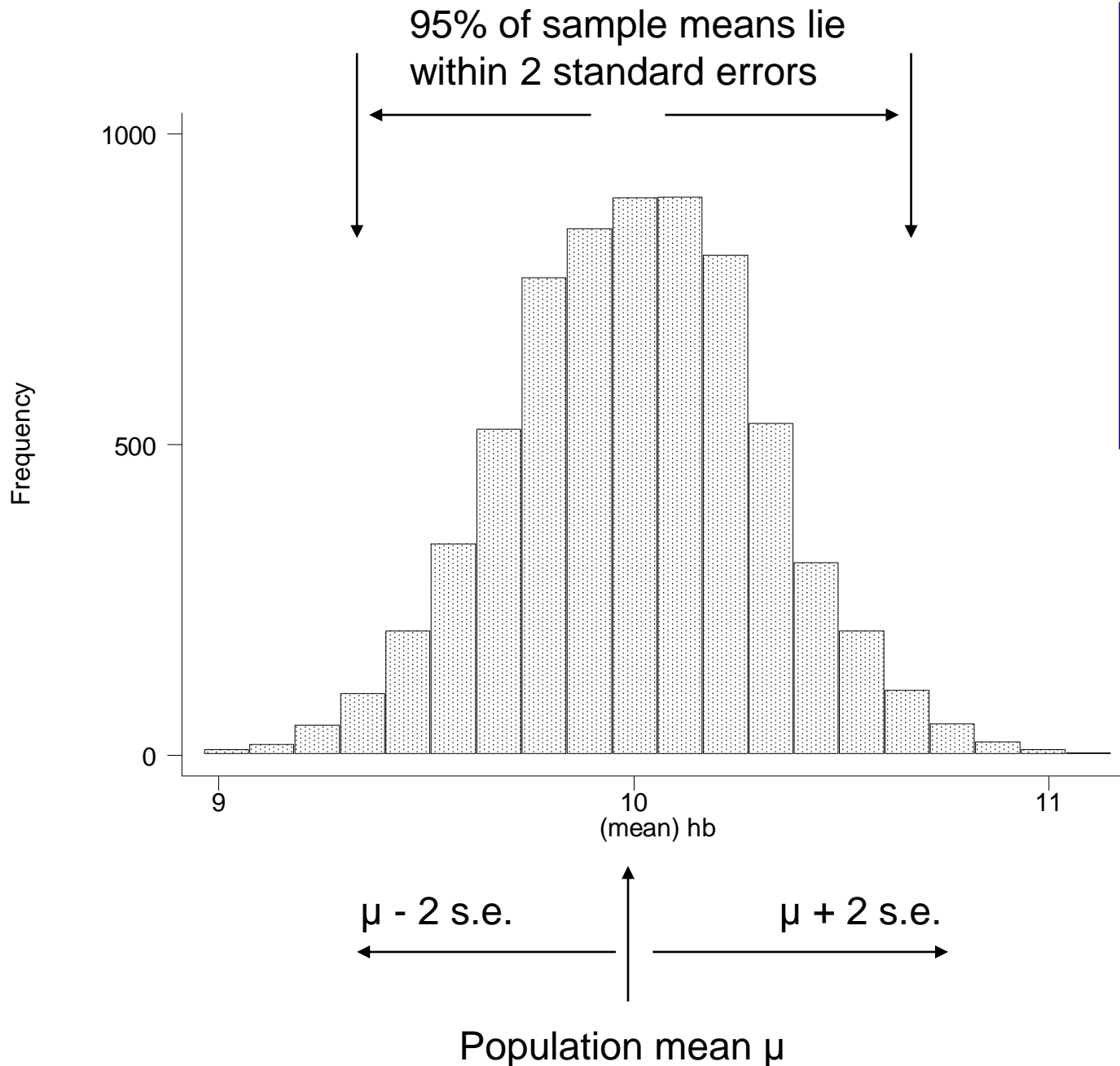


□ Bigger sample size, less variability in sample means: $s.e. = SD/\sqrt{n}$

□ Average value of sample means is equal to the population mean

□ Sample means follow normal distribution if the distribution in the population is normal

□ and if n is large, is roughly normal even if the population is not normal



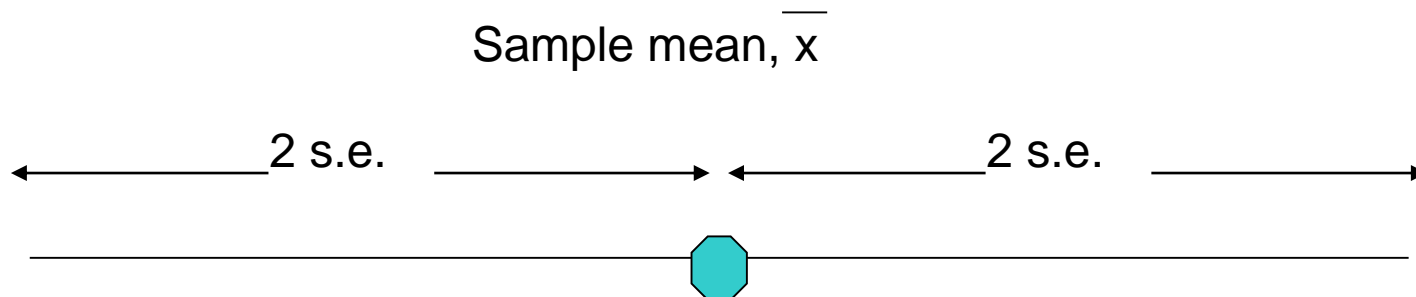
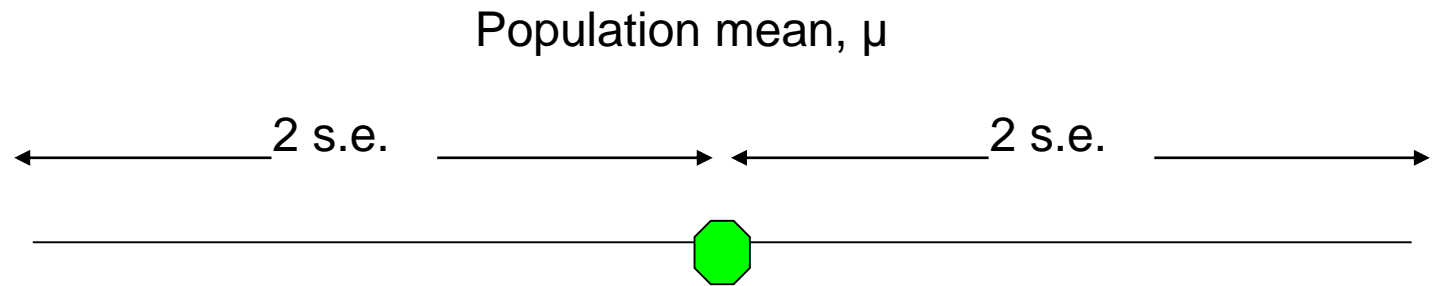
If we repeated our sampling many times, 95% of sample means would be within 2 standard errors of the population mean.

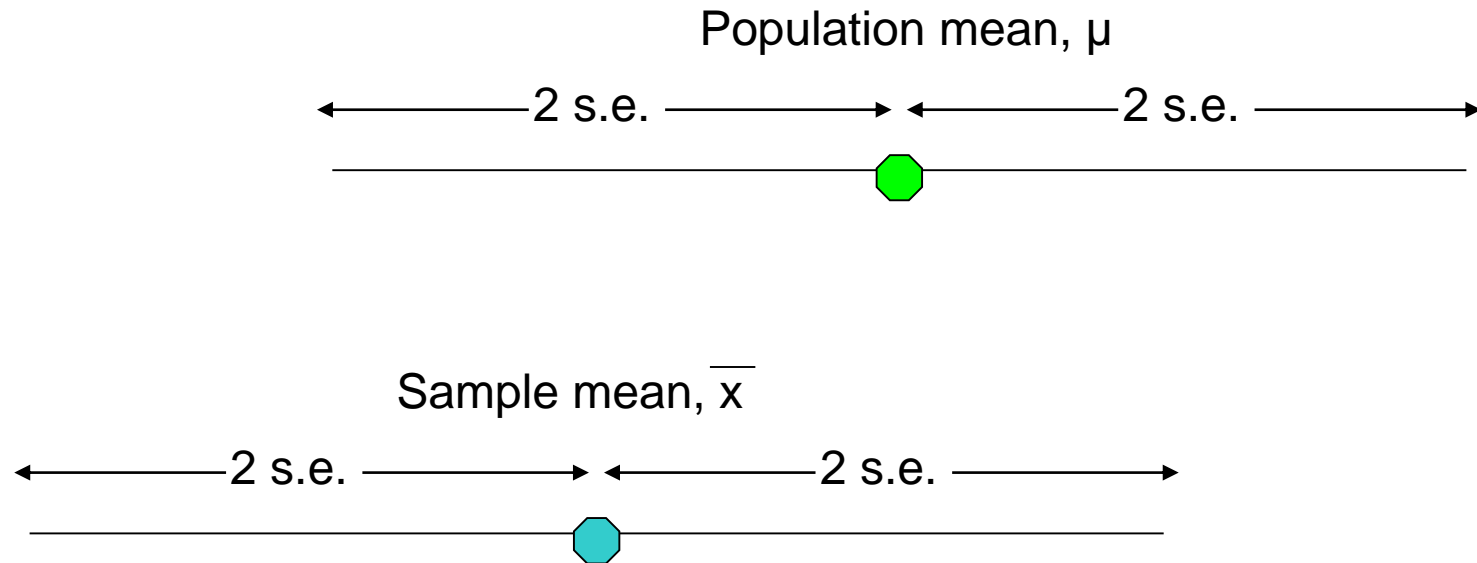
In practice we usually have only one sample but we can use the sample data to quantify the uncertainty in our single estimate

95% of sample means lie within 2 s.e. of the population mean

so we can also say:

95% of the time the population mean will lie within 2 s.e. of
the our sample mean





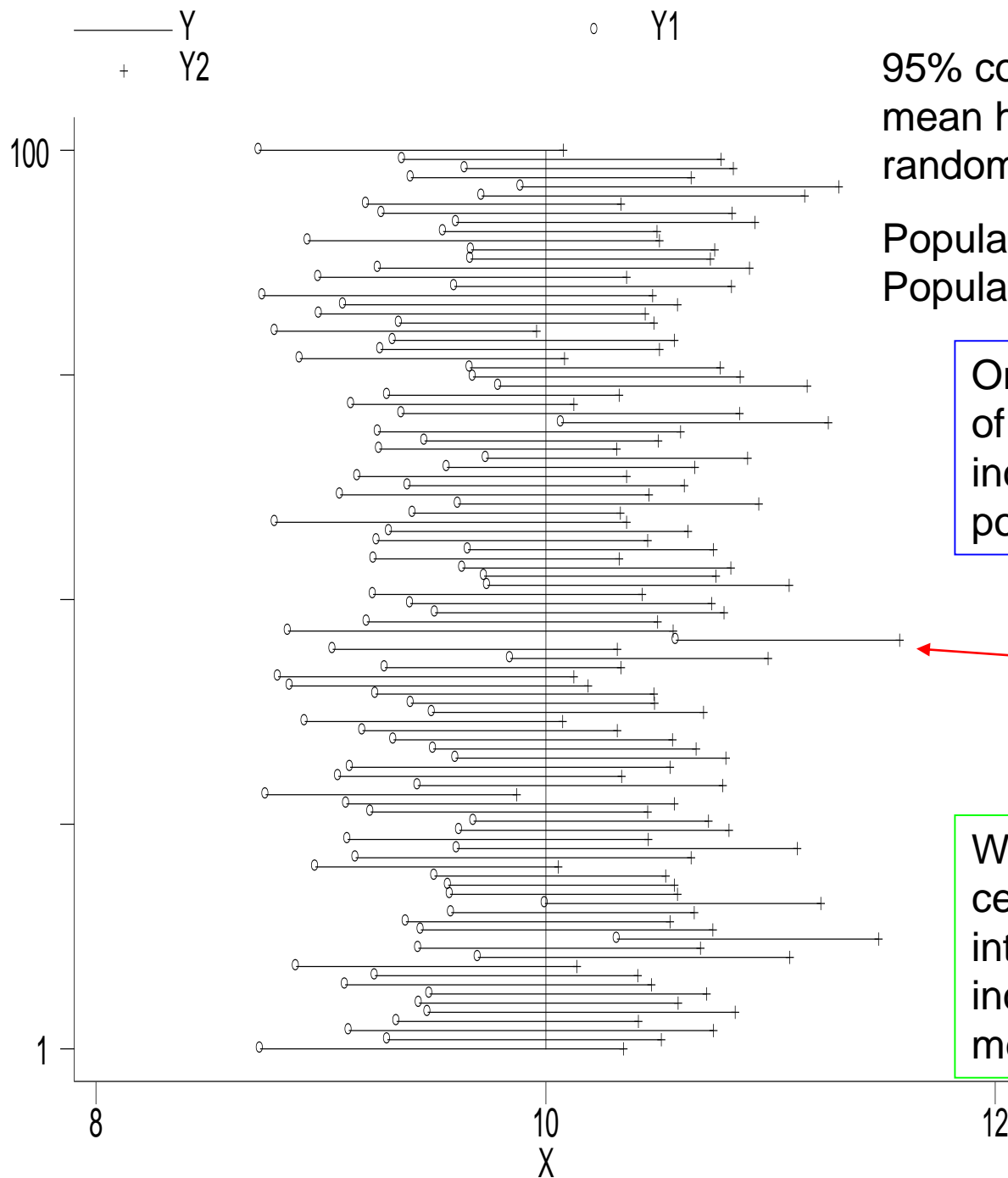
Sample estimate of the SD: 1.5 g/dl

$$\text{Standard error s.e.} = \text{SD}/\sqrt{n}$$

$$\text{s.e.} = 1.5/\sqrt{30} = 1.5/5.477 = 0.274$$

95% confidence interval $9.6 - 2 \times 0.274$ to $9.6 + 2 \times 0.274$
9.1 g/dl to 10.1 g/dl

Interpretation: we can say with 95% confidence that the mean Haemoglobin concentration in children in the population could be as small as 9.1 or as big as 10.1 g/dl



95% confidence intervals for
mean haemoglobin from 100
random samples of size 30

Population mean 10g/dl
Population SD 1.7g/dl

On average 95%
of the intervals
include the
population mean

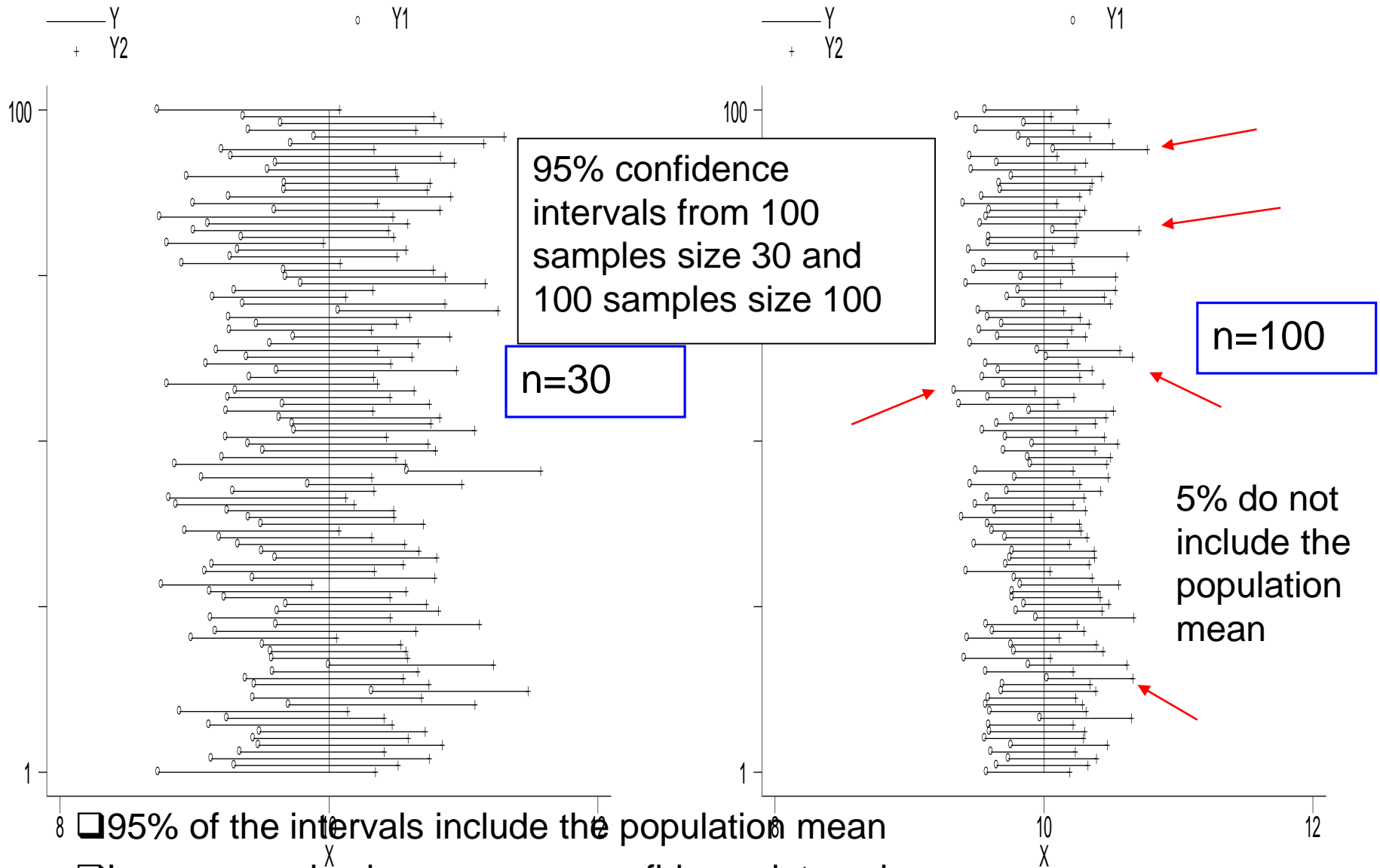
Small chance the
interval will not
include the
population mean

We will not know for
certain whether the
interval we calculate
includes the population
mean

Now consider what will change if we increase the sample size to $n=100$:

What happens to the width of the interval?

What proportion of intervals will include the population mean?



95% of the intervals include the population mean

Larger sample size \rightarrow narrow confidence interval

We can improve the probability of including the population mean by using 99% confidence intervals, at the cost of having wider intervals and so more uncertainty

- For large samples: $\bar{x} \pm 1.96 \text{ s.e.}$
- With continuous data, if samples are small the estimated standard error tends to be an underestimate. To correct for this we use t-tables to find the multiplier:

degrees of freedom	Two-tailed probability		
	0.1	0.05	0.01
1	6.314	12.706	63.656
2	2.920	4.303	9.925
3	2.353	3.182	5.841
4	2.132	2.776	4.604
5	2.015	2.571	4.032
10	1.812	2.228	3.169
20	1.725	2.086	2.845
29	1.699	2.045	2.756
30	1.697	2.042	2.750
100	1.660	1.984	2.626
150	1.655	1.976	2.609

Sample mean 9.6 Sample SD 1.5

Standard error $1.5/\sqrt{30}=0.274$

95% confidence interval:

$9.6 - \mathbf{2.045} \times 0.274$ **to** $9.6 + \mathbf{2.045} \times 0.274$

9.04 to 10.16 g/dl

99% confidence interval:

$9.6 - \mathbf{2.756} \times 0.274$ **to** $9.6 + \mathbf{2.756} \times 0.274$

8.85 to 10.35 g/dl

Confidence interval on μ (σ known)

$$\Pr\left\{-1.96 \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq 1.96\right\} = 0.95$$

$$\Pr\left\{\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right\} = 0.95$$

$$\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right)$$

is a 95% confidence interval for μ .

A survey of haemoglobin status in children <5yrs in Kilifi district.
30 children gave a finger prick blood sample and the mean Hb was
9.6g/dl with a standard deviation of 1.5 g/dl.

Sample estimate of the SD: 1.5 g/dl

Standard error s.e. = SD/ \sqrt{n}

$$\text{s.e.} = 1.5 / \sqrt{30} = 1.5 / 5.477 = 0.274$$

95% confidence interval $9.6 - 2 \times 0.274$ to $9.6 + 2 \times 0.274$
9.1 g/dl to 10.1 g/dl

Interpretation: we can say with 95% confidence that the mean Haemoglobin concentration in children in the population could be as small as 9.1 or as big as 10.1 g/dl

Confidence Intervals: Factors affecting width

1. Data dispersion

- Measured by σ

Intervals extend from

$$\bar{X} - Z\sigma_{\bar{X}} \text{ to } \bar{X} + Z\sigma_{\bar{X}}$$

2. Sample size

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

3. Level of confidence

$(1 - \alpha)$

- Affects Z

Confidence intervals: summary

- ❑ The confidence interval expresses the **uncertainty** in sample estimates of means, proportions, treatment efficacy, dose response,,,
- ❑ Larger the sample, the narrower the CI
- ❑ Can improve the probability of including the population mean by calculating 99% interval but at a cost of having a wider interval and thus greater uncertainty
- ❑ Should design studies to yield CI's that are narrow enough to draw replicative conclusions
- ❑ Interpretation depends on the assumption that the sample was **representative** of the population