

Raoul LePage

Professor

STATISTICS AND PROBABILITY

www.stt.msu.edu/~lepage

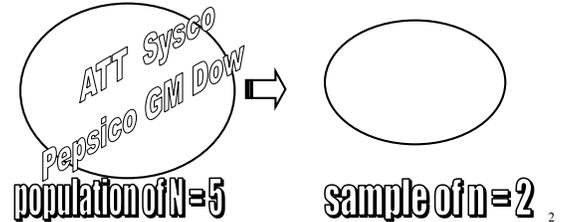
click on STT461_Sp05

Slides 38-41 revised 1-7-05.

1

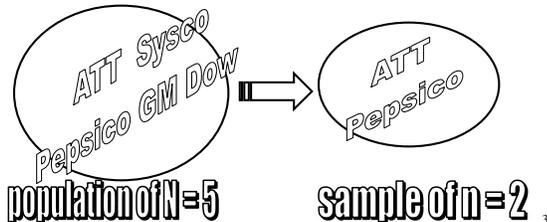
THE GREAT TRICK OF STATISTICS

The overwhelming majority of samples of n from a population of N can stand-in for the population.



THE GREAT TRICK OF STATISTICS

The overwhelming majority of samples of n from a population of N can stand-in for the population.



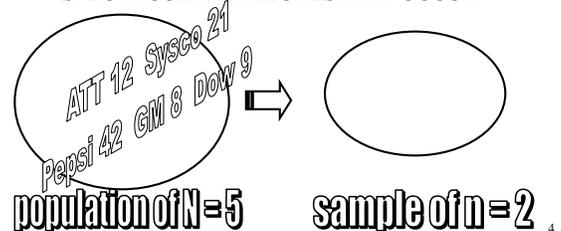
3

GREAT TRICK: SOME CAVEATS

For a few characteristics at a time, such as profit, sales, dividend.

Sample size n must be "large."

SPECTACULAR FAILURES MAY OCCUR!



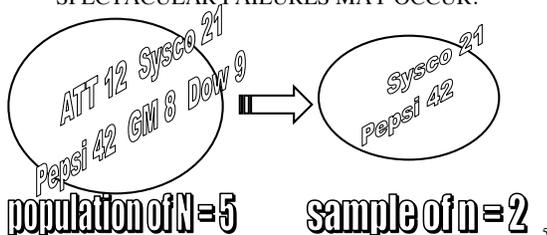
4

GREAT TRICK: SOME CAVEATS

For a few characteristics at a time, such as profit, sales, dividend.

Sample size n must be "large."

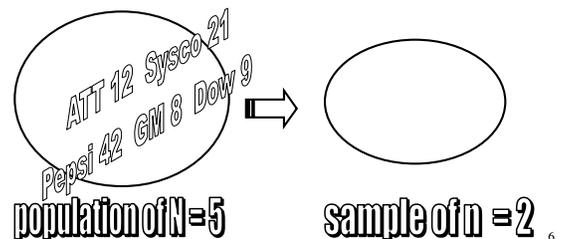
SPECTACULAR FAILURES MAY OCCUR!



5

HOW ARE WE SAMPLING?

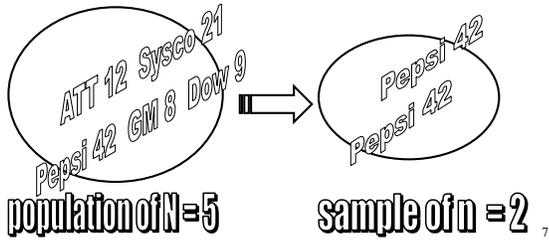
With-replacement vs without replacement.



6

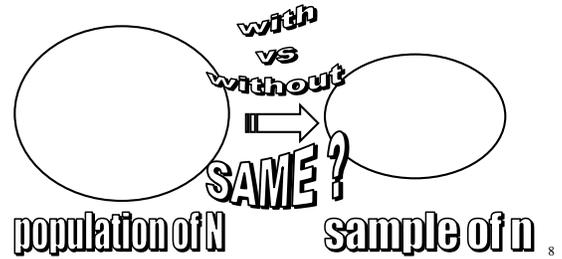
HOW ARE WE SAMPLING?

With-replacement vs without replacement.



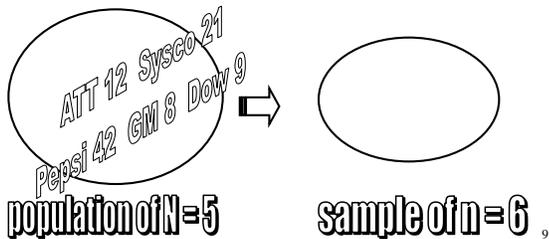
DOES IT MAKE A DIFFERENCE?

Rule of thumb: With and without replacement are about the same if $\sqrt{(N-n)/(N-1)} \sim 1$.



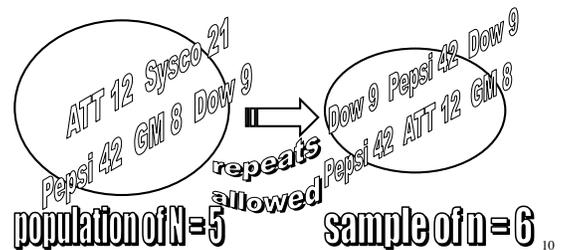
UNLIMITED SAMPLING

WITH-replacement samples have no limit to the sample size n .

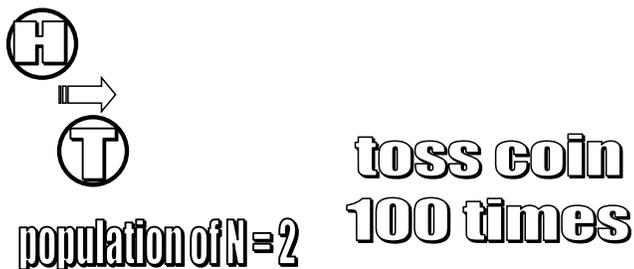


UNLIMITED SAMPLING

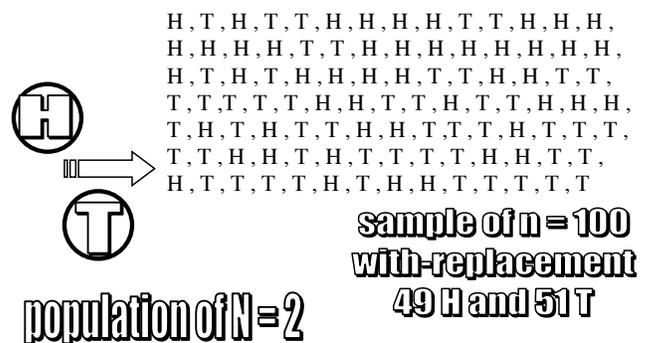
WITH-replacement samples have no limit to the sample size n .



TOSS COIN 100 TIMES

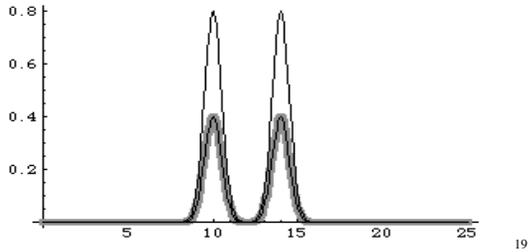


TOSS COIN 100 TIMES



NARROWER TENTS = MORE DETAIL

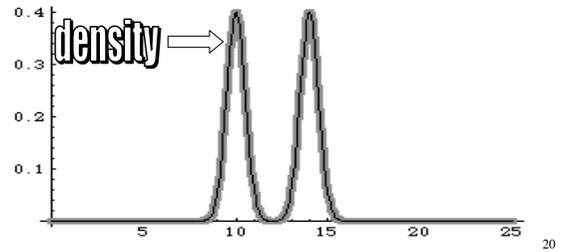
Making the tents narrower isolates different parts of the data and reveals more detail.



19

THE DENSITY BY ITSELF

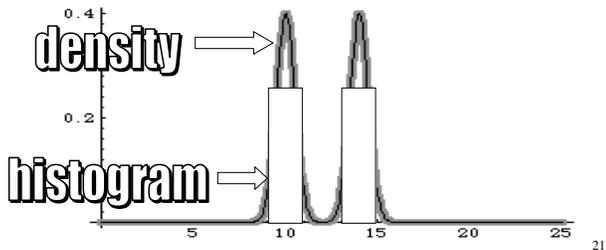
With narrow tents.



20

DENSITY OR HISTOGRAM ?

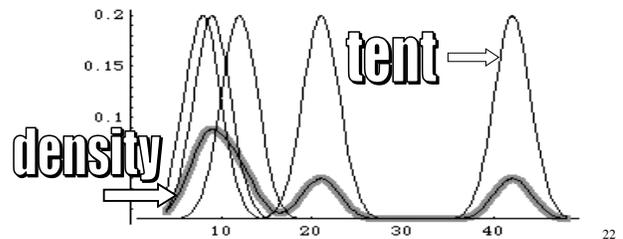
Histograms lump data into categories (the black boxes), not as good for continuous data.



21

DENSITY FOR {12, 21, 42, 8, 9}

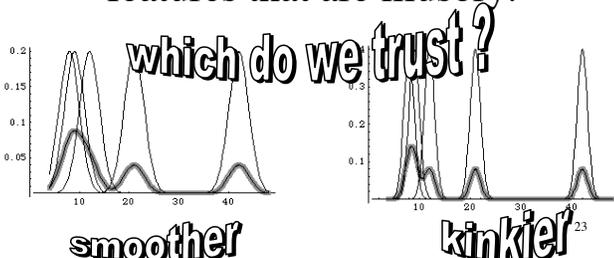
Plot of average heights of 5 tents placed at data {12, 21, 42, 8, 9}.



22

IS DETAIL ILLUSORY ?

Narrower tents operate at higher resolution but they may bring out features that are illusory.



23

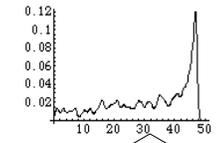
THE MEAN OF A DENSITY IS THE SAME AS THE MEAN OF THE DATA FROM WHICH IT IS MADE.

24

BEWARE OVER-FINE RESOLUTION

Population of $N = 500$ compared with two samples of $n = 30$ each.

POP mean = 32.02



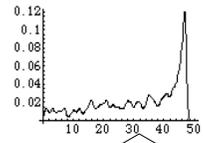
population of $N = 500$

with 2 samples of $n = 30$ ²⁵

BEWARE OVER-FINE RESOLUTION

Population of $N = 500$ compared with two samples of $n = 30$ each.

POP mean = 32.02



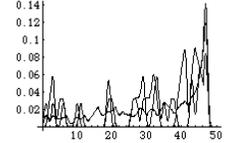
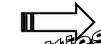
population of $N = 500$

sample means

SAM1 mean = 33.03
SAM2 mean = 30.60

are close

densities not good at fine resolution



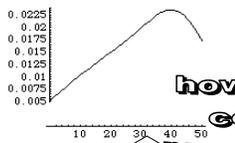
with 2 samples of $n = 30$ ²⁶

WE DO BETTER AT COARSE RESOLUTION

The same two samples of $n = 30$ each from the population of 500.

POP mean = 32.02

SAM1 mean = 33.03
SAM2 mean = 30.60



how about coarse resolution?



population of $N = 500$

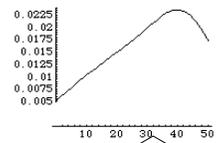
with 2 samples of $n = 30$ ²⁷

WE DO BETTER AT COARSE RESOLUTION

The same two samples of $n = 30$ each from the population of 500.

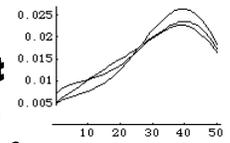
POP mean = 32.02

SAM1 mean = 33.03
SAM2 mean = 30.60



population of $N = 500$

good at coarse resolution



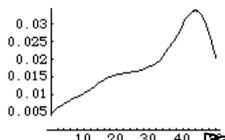
with 2 samples of $n = 30$ ²⁸

HOW ABOUT MEDIUM RESOLUTION?

The same two samples of $n = 30$ each from the population of 500.

POP mean = 32.02

SAM1 mean = 33.03
SAM2 mean = 30.60



medium resolution?



population of $N = 500$

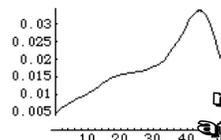
with 2 samples of $n = 30$ ²⁹

HOW ABOUT MEDIUM RESOLUTION?

The same two samples of $n = 30$ each from the population of 500.

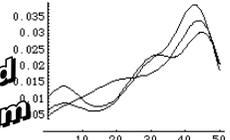
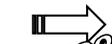
POP mean = 32.02

SAM1 mean = 33.03
SAM2 mean = 30.60



population of $N = 500$

not good at medium resolution

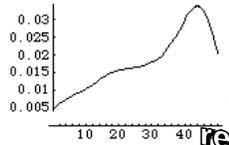


with 2 samples of $n = 30$ ³⁰

SAMPLING ONLY 600 FROM 500 MILLION?

A sample of only $n = 600$ from a population of $N = 500$ million.
(medium resolution)

POP mean = 32.02



medium resolution?

large sample of $n = 600$?

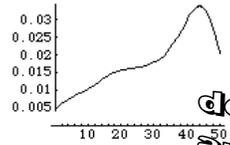
population of $N = 500,000$

with a sample of $n = 600$

SAMPLING ONLY 600 FROM 500 MILLION?

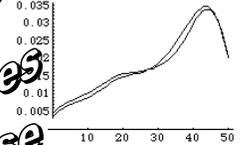
A sample of only $n = 600$ from a population of $N = 500$ million.
(MEDIUM resolution)

POP mean = 32.02



population of $N = 500,000$

sample of $n = 600$
sample mean = 32.84



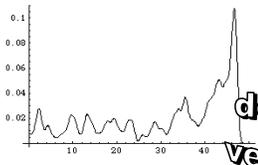
mean very close densities are close

with a sample of $n = 600$

SAMPLING ONLY 600 FROM 500 MILLION?

A sample of only $n = 600$ from a population of $N = 500$ million.
(FINE resolution)

POP mean = 32.02



population of $N = 500,000$

sample of $n = 600$
sample mean = 32.84



densities very close

FINE resolution

with a sample of $n = 600$

THE ROLE OF RANDOM SAMPLING

IF THE OVERWHELMING MAJORITY OF SAMPLES ARE "GOOD SAMPLES" THEN WE CAN OBTAIN A "GOOD" SAMPLE BY RANDOM SELECTION.

34

HOW TO SAMPLE RANDOMLY? SELECTING A LETTER AT RANDOM

With-replacement:

a = 00-02 b = 03-05 z = 75-77

From Table 14 pg. 869:

1559 9068 9290 8303 etc...

15 59 90 68 etc... (split into pairs)

we have 15 = f, 59 = t, 90 = none, etc...

(for samples without replacement just pass over any duplicates).

35

TALKING POINTS

- The Great Trick of Statistics.
 - The overwhelming majority of all samples of n can "stand-in" for the population to a remarkable degree.
 - Large n helps.
 - Do not expect a given sample to accurately reflect the population in many respects, it asks too much of a sample.
- The Law of Averages is one aspect of The Great Trick.
 - Samples typically have a mean that is close to the mean of the population.
 - Random samples are nearly certain to have this property since the overwhelming majority of samples do.
- A density is controlled by the width of the tents used.
 - Small samples zero-in on coarse densities fairly well.
 - Samples in hundreds can perform remarkably well.
 - Histograms are notoriously unstable but remain popular.
- Making a density from two to four values; issue of resolution.
- With-replacement vs without; unlimited samples.
- Using Table 14 to obtain a random sample.

36

The Great Trick is far more powerful than we have seen.

A typical sample closely estimates such things as a population mean or the shape of a population density.

But it goes beyond this to reveal how much variation there is among sample means and sample densities.

A typical sample not only estimates population quantities.

It estimates the sample-to-sample variations of its own estimates. ³⁷

EXAMPLE: ESTIMATING A MEAN

The average account balance is \$421.34 for a random with-replacement sample of 50 accounts.

We estimate from this sample that the average balance is \$421.34 for all accounts.

From this sample we also estimate and display a “margin of error”

$$\$421.34 \pm \$65.22 = \bar{x} \pm 1.96 \frac{s}{\sqrt{n}}$$

s denotes "sample standard deviation"

SAMPLE STANDARD DEVIATION

$$s = \sqrt{\frac{\sum_{i=1}^{i=n} (x_i - \bar{x})^2}{n - 1}}$$

NOTE: Sample standard deviation *s* may be calculated in several equivalent ways, some sensitive to rounding errors, even for $n = 2$. ³⁹

EXAMPLE: MARGIN OF ERROR CALCULATION

The following margin of error calculation for $n = 4$ is only an illustration. A sample of four would not be regarded as large enough.

Profits per sale = {12.2, 15.3, 16.2, 12.8}.

Mean = 14.125, $s = 1.92765$, $\text{root}(4) = 2$.

Margin of error = $\pm 1.96 (1.92765 / 2)$

Report: 14.125 \pm 1.8891.

A precise interpretation of margin of error will be given later in the course, including the role of 1.96. The interval 14.125 \pm 1.8891 is called a “95% confidence interval for the population mean.”

We used: $(12.2-14.125)^2 + (15.3-14.125)^2 + (16.2-14.125)^2 + (12.8-14.125)^2 = 11.1475$. ⁴⁰

EXAMPLE: ESTIMATING A PERCENTAGE

A random with-replacement sample of 50 stores participated in a test marketing. In 39 of these 50 stores (i.e. 78%) the new package design outsold the old package design.

We estimate from this sample that 78% of all stores will sell more of new vs old.

We also estimate a “margin of error”

$\pm 11.6\%$

percentage is the mean for 100 = new sale, 0 = old sale ⁴¹

