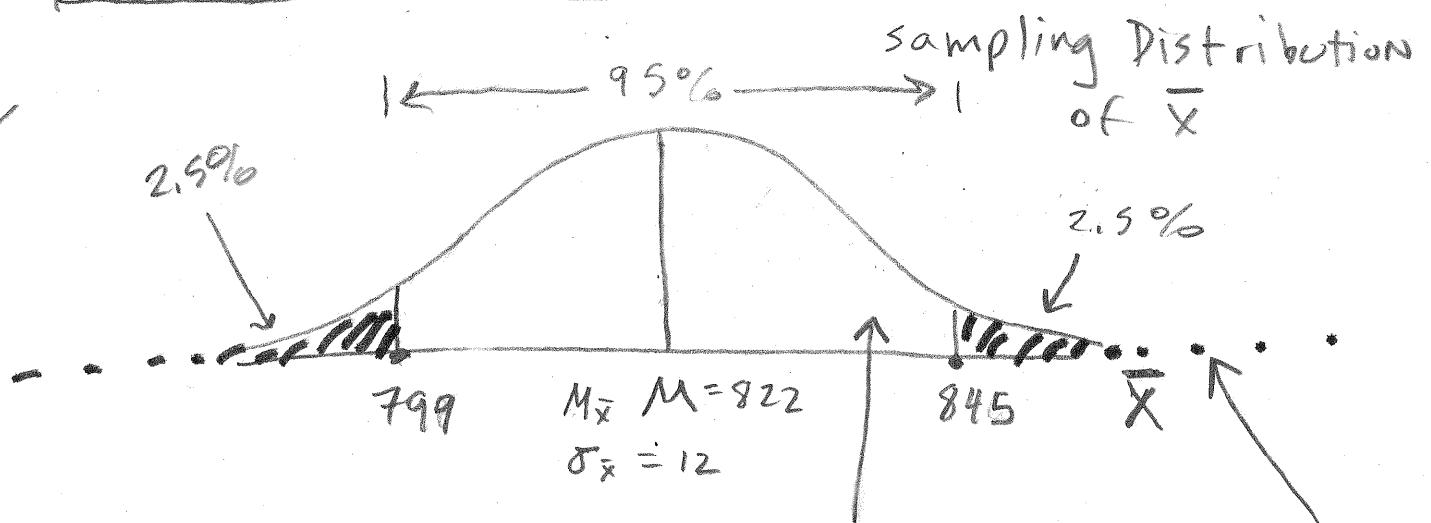


Chapter 7: Know Pop. Mean, M

Ch. 7

S DoX



took sample & compare
to Sampling Distribution
of \bar{X}

$$\bar{X}_1 = 837$$

\bar{X}_1 inside
so we say
original claim
of 822 seems
reasonable

$$\bar{X}_2 = 850$$

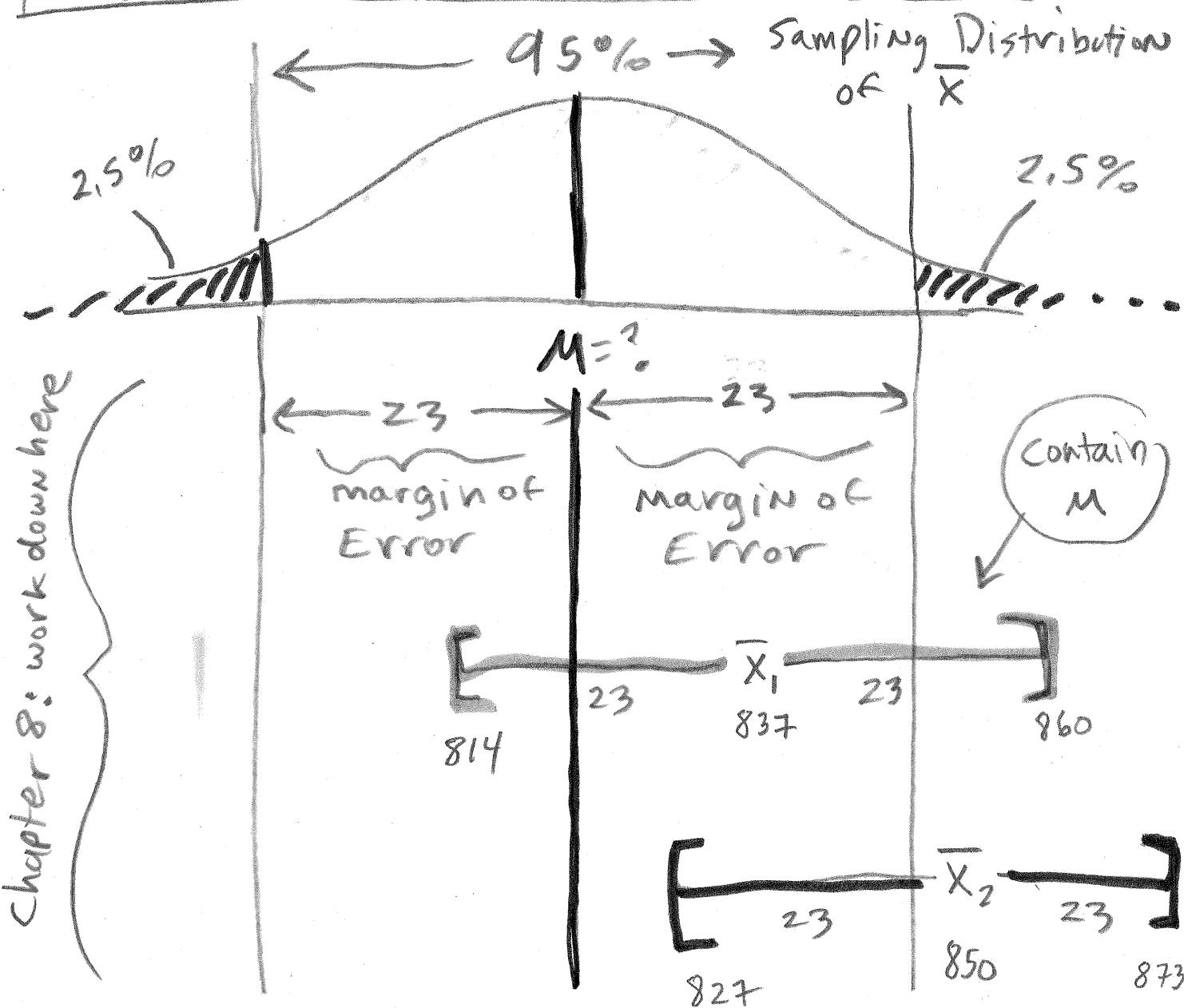
\bar{X}_2 outside
so we say
original claim
of 822 seems
unreasonable

- ① Even though 2.5% of samples would occur out here, because getting a sample of 850 is so unlikely, we say that from our sample evidence it seems unlikely, and thus original claim seems unreasonable.
- ② Because samples will almost always have Sampling Error, the sampling Error with $\bar{X} = 850$ is Not acceptable, whereas with $\bar{X} = 837$ the sampling Error is acceptable.
- ③ we create point estimates & compare them directly to S DoX.
- ④ we make statements like:

We are 95% sure that \bar{X} will lie between 799 & 845

Chapter 8: Don't Know Pop Mean, M

$SDo\bar{X}$ ch.8



- ① we run the risk that 2.5% of our intervals will not contain M & will be on upper side.
- ② we develop confidence intervals to estimate the pop. parameter (we will estimate $\bar{X} \pm \bar{E}$)
- ③ we make statements like: we are 95% sure that M will be in our interval

Chapter 8

Confidence Intervals (Interval Estimation)

(1)

Chapter 7

Know what pop. Mean, M , is or what it should be

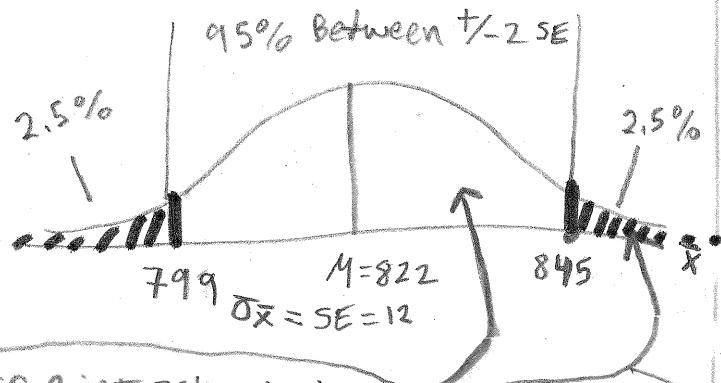
when insurance company says "average cost for policy is \$822, we know what M should be



we develop point estimates

& compare them directly

to Sampling Distribution of \bar{X}



If point estimate is $\bar{X}_1 = \$831$
original claim seems reasonable

If point est. is $\bar{X}_2 = \$850$
original claim seems unreasonable

Make statements about \bar{X} :

We are 95% sure that \bar{X} will lie between 799 and 845

Chapter 8

Don't know what pop. Mean, M , is or what it should be

when printer cartridge manufacturer wants an "average number of pages per cartridge", they can't calculate pop. mean, M , because they can't print all pages from all cartridges.

we develop confidence intervals

if our point estimate is $\bar{X} = 2409$ pages

$$\bar{X} = 2409 \text{ pages}$$

we can't compare \bar{X} directly to Sampling Distribution of \bar{X}
but we can:

Add a "Margin of Error" to our point estimate:

$$+/- 217 \text{ pages}$$

to create a 95% confidence interval

$$[\underline{\bar{X}} \overline{\bar{X}}]$$

2192 2409 2626

And say:

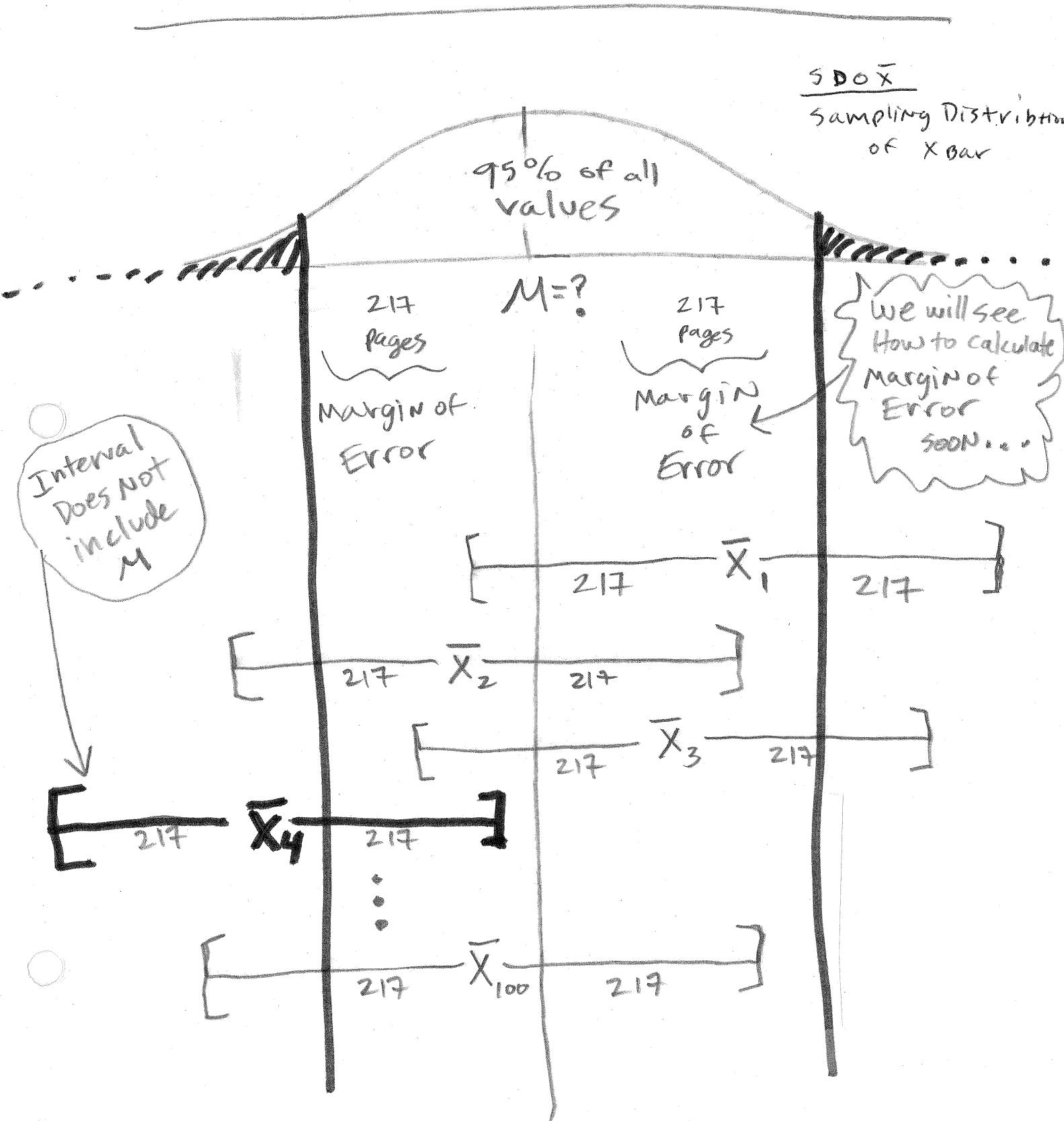
We are 95% sure that M will lie between 2192 & 2626 pages
or
Avg pages will be 2409 w/ Margin of Error 217

For 95% Confidence Intervals

(2)

Because we don't know M

technically it would look like this:



→ Notice for printer cartridge example: ③

we CAN NOT calculate population mean, M .

why?

Because we would have to test every cartridge made or count every page from every cartridge in every house...

other Examples:

① we do not know how many meals married couples eat out each week
- we can take a sample and our \bar{x} (mean) will be the point estimate for the unknown pop. Mean.

② we do not know what the mean amount of profit is per auto sold in USA.
- we can take a sample and our \bar{x} will be the point estimate for the unknown pop. mean.

Since \bar{x} rarely (if ever) equals M we can't just take sample & use \bar{x} , we must add Margin of Error to each side of \bar{x} & create Confidence Interval for our estimate of pop. mean, M .

Confidence Intervals (Interval Estimation)

- ① An estimate of the population parameter (4)
that provides an interval believed to contain
the value of the parameter.
- ② estimate of population parameter = \bar{x}
point estimate \pm Margin of Error

Confidence Level or Confidence Coefficient

probability that population parameter will
occur within our interval.

Examples: 90%

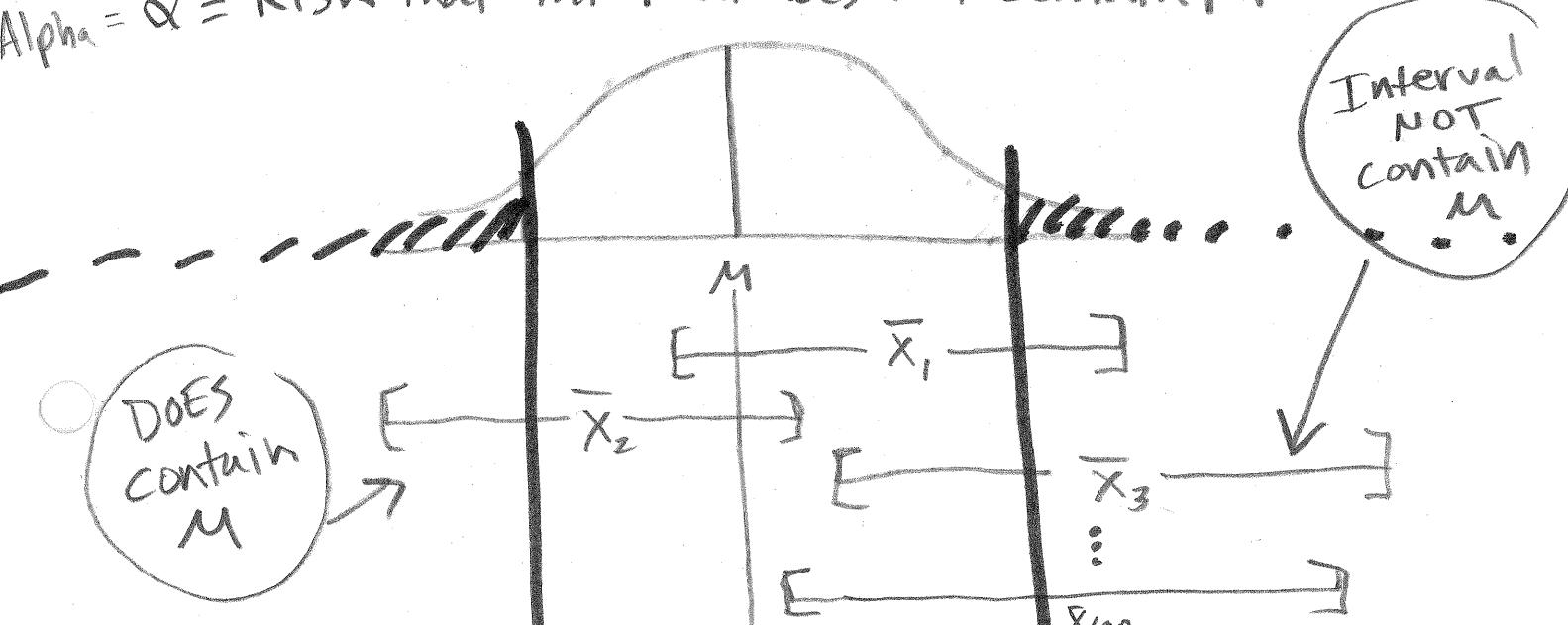
95%

99%

Level of Significance = α = Alpha

$\alpha = \text{Level of Significance} = 1 - \text{"Confidence Interval"}$

$\alpha = \text{Risk that interval does not contain pop. parameter.}$

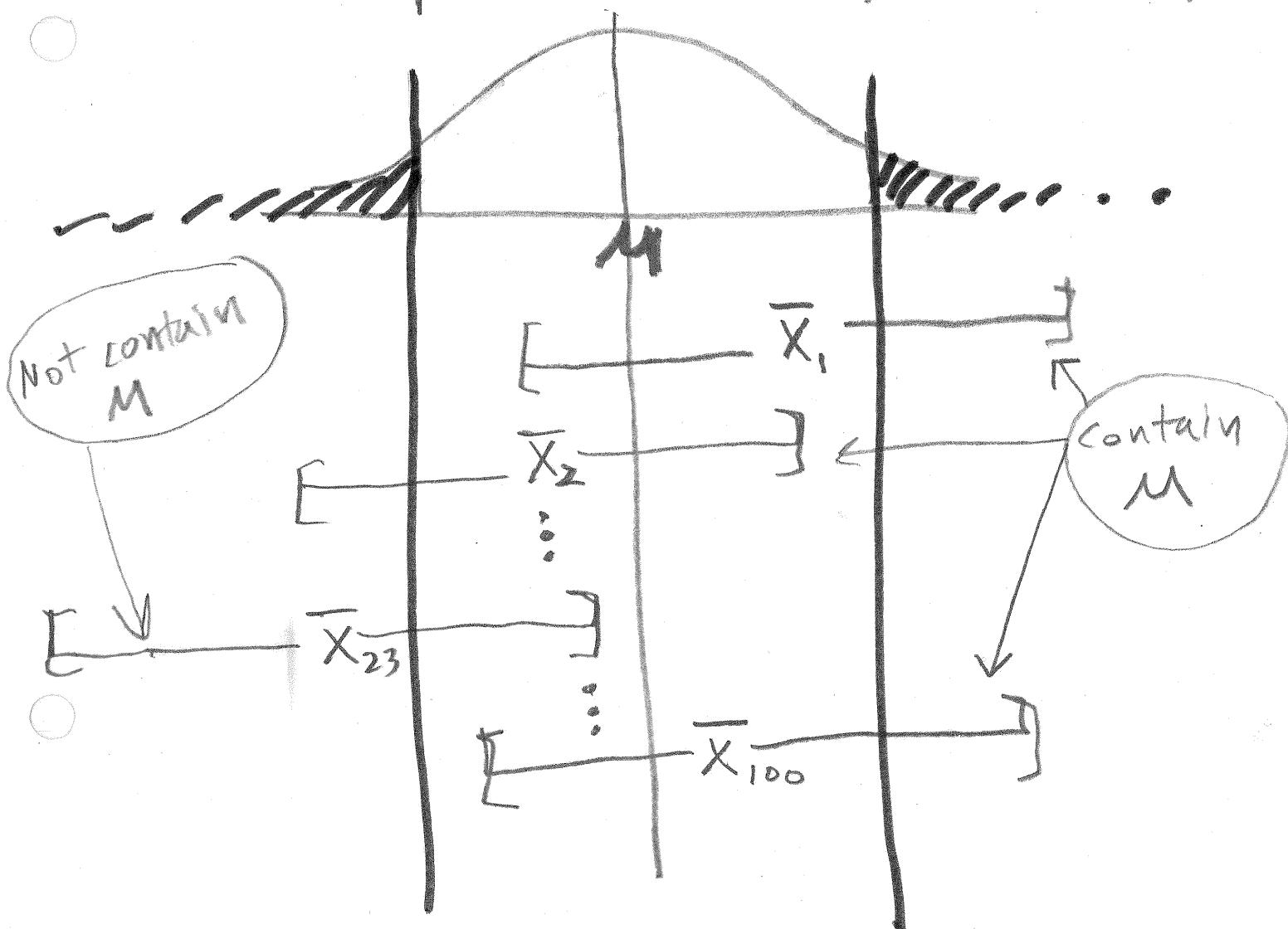


Note:

95% confidence

(5)

Interval



★ If population Distribution follows a Normal Distribution, and you construct a 95% confidence Interval, 95 would contain M , & 5 would not contain M .

★ If pop. distribution Not Normal, the Sample size must be sufficiently big (central Limit Theorem) so that approximately 95 have M & 5 ^{not have} M .

Confidence Interval for μ σ Known*

* Known or have good estimate

$$\bar{X} \pm Z_{\alpha/2} * \frac{\sigma}{\sqrt{n}} \quad \left. \begin{array}{l} \text{Standard} \\ \text{Error} \end{array} \right\}$$

(6)

σ = sigma = Margin of Error

\bar{X} = sample mean = Population Standard Deviation

Z = Standardize value = Number of Standard Error

$\sigma_{\bar{X}}$ = Standard Error = Standard Deviation of Sampling Distribution

$1 - \alpha$ = Confidence level (confidence coefficient) of \bar{X}

α = alpha = Risk of Not in Interval

Margin of Error = Amount on either side of \bar{X} to create Interval

$Z_{\alpha/2}$ = Z on upper end

n = sample size

Excel

Calculate Z

Method 1:

$$\text{Upper Limit} = \text{NORM.INV}(1 - \frac{\alpha}{2}, \bar{X}, \sigma_{\bar{X}})$$

$$\text{Lower Limit} = \text{NORM.INV}(\frac{\alpha}{2}, \bar{X}, \sigma_{\bar{X}})$$

Method 2: $Z_{\text{upper}} = \text{NORM.S.INV}(1 - \frac{\alpha}{2})$

Calculate Margin of Error

{Margin of Error}

Method 3:

$$= \text{CONFIDENCE.NORM}(\alpha, \sigma, n)$$

σ known or at least reliably estimated

- ① when there is large amounts of historic data, sigma (σ) may be known.
- ② Quality control applications where process is assumed to be operating correctly, it is appropriate to treat sigma as known.
- ③ Most situations σ is not known, so we don't use Normal (z) distribution, we use t-distribution. In this case we will have to use the same sample to estimate M & σ . We will use \bar{x}_{bar} & s .
(more later...)

Note: $\bar{x} \pm Z_{\alpha/2} * \frac{\sigma}{\sqrt{n}}$

If particular sample size provides an interval that is too wide for any practical use, increase sample size, n .

$$SE = \frac{\sigma}{\sqrt{n}} \leftarrow \text{As this increases (sample size)}$$

$\frac{\sigma}{\sqrt{n}}$ decreases (Standard Error)

$Z_{\alpha/2} * \frac{\sigma}{\sqrt{n}}$ decreases (Margin of Error)

Narrower interval, & thus greater precision of estimate of M .

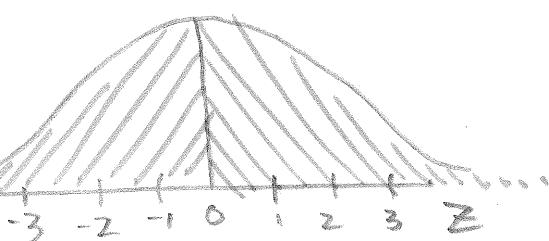
Q: But what if we do not know the population standard deviation, sigma (σ)? Can we still construct confidence intervals?

A: Yes we can! But will have to use the point estimate sample standard deviation, s , in place of σ .

When we use the sample standard deviation, s , instead of the population standard deviation, σ , we cannot use the Standard Normal curve & the Z-score. When we use the sample standard deviation, s , we must use the t-distribution.

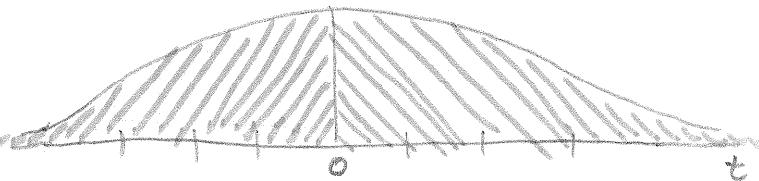
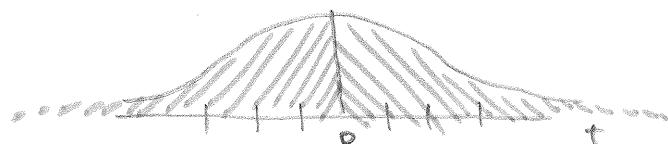
σ known

σ not known



1 standard Normal

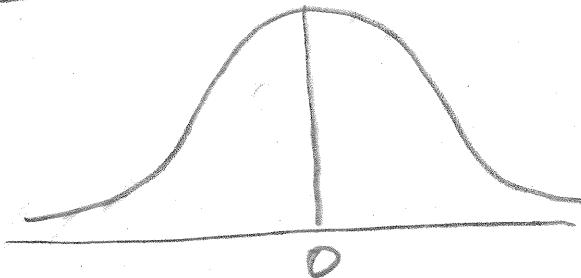
Curve using Z



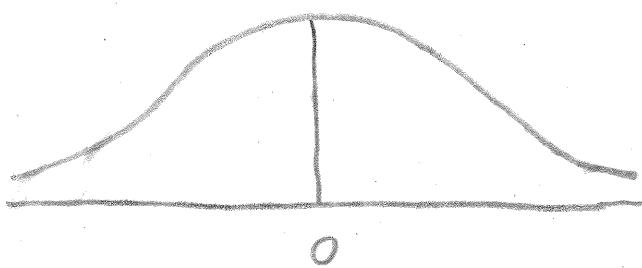
many t-distribution curves

Many t -Distributions

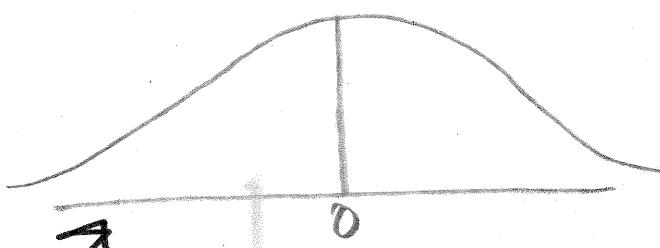
9



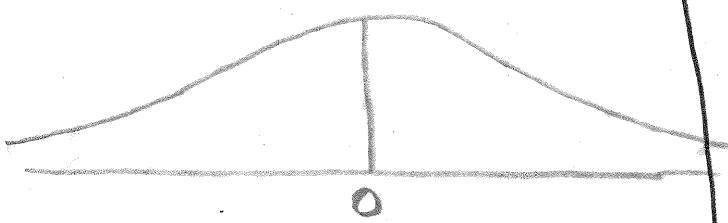
$n = 100$
 $df = 99$



$n = 50$
 $df = 49$



$n = 30$
 $df = 29$



$n = 15$
 $df = 14$

One for each sample size, n , or
Degrees of Freedom, $n - 1$.

This means that when creating our Confidence Interval, we will have to choose the correct t -Distribution

• IF we have $n = 50$, we choose

• IF we have $n = 30$, we choose

Degrees of freedom (statistics)

From Wikipedia, the free encyclopedia

In statistics, the number of **degrees of freedom** is the number of values in the final calculation of a statistic that are free to vary.^[1]

Estimates of statistical parameters can be based upon different amounts of information or data. The number of independent pieces of information that go into the estimate of a parameter is called the degrees of freedom (df). In general, the degrees of freedom of an estimate of a parameter is equal to the number of independent scores that go into the estimate minus the number of parameters used as intermediate steps in the estimation of the parameter itself (which, in sample variance, is one, since the sample mean is the only intermediate step).^[2]

Mathematically, degrees of freedom is the dimension of the domain of a random vector, or essentially the number of 'free' components: how many components need to be known before the vector is fully determined.

The term is most often used in the context of linear models (linear regression, analysis of variance), where certain random vectors are constrained to lie in linear subspaces, and the number of degrees of freedom is the dimension of the subspace. The degrees-of-freedom are also commonly associated with the squared lengths (or "Sum of Squares") of such vectors, and the parameters of chi-squared and other distributions that arise in associated statistical testing problems.

While introductory texts may introduce degrees of freedom as distribution parameters or through hypothesis testing, it is the underlying geometry that defines degrees of freedom, and is critical to a proper understanding of the concept. Walker (1940)^[3] has stated this succinctly:

For the person who is unfamiliar with N -dimensional geometry or who knows the contributions to modern sampling theory only from secondhand sources such as textbooks, this concept often seems almost mystical, with no practical meaning.

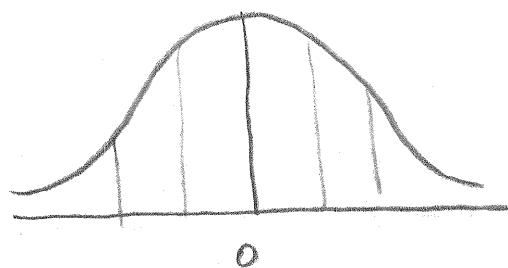
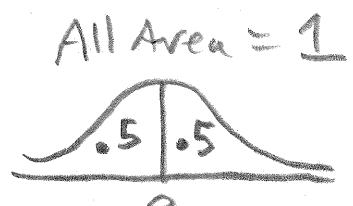
t-Distribution (Student's t-Distribution)

When are we allowed to use t-Distribution?

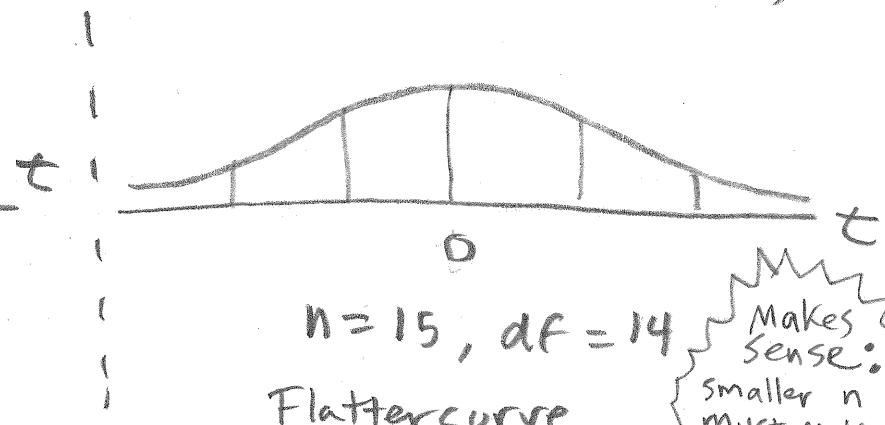
Ans: When population Distribution is Normal Shaped or near Normal, or n is sufficiently large (Central Limit Theorem)

Characteristics of t-Distribution

- ① Continuous Distribution
- ② Normal (Bell) shape & Symmetrical
- ③ The smaller n or df , the larger the Standard Deviation, the flatter the curve,



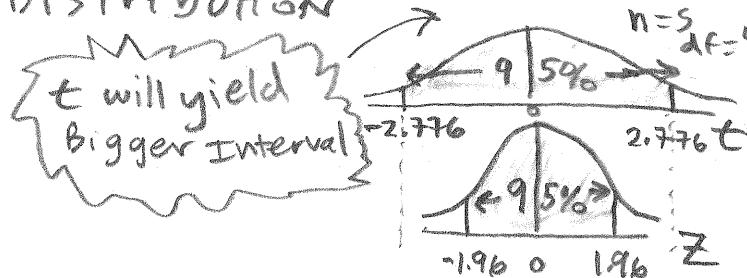
$n = 30, df = 29$
taller curve
smaller standard deviation



$n = 15, df = 14$
Flatter curve
Bigger standard deviation

Makes Sense:
Smaller n must make Interval Bigger

- ④ At center, t-Distribution is Flatter & more spread out than Normal Z Distribution



- ⑤ As n or df increase, t-Distribution approaches Z-Distribution.

Confidence Interval for μ & σ Not Known

$$\bar{X} \pm t_{\alpha/2} * \frac{s}{\sqrt{n}}$$

Margin of Error

Standard Error

\bar{X} = Sample mean

s = Sample Standard Deviation

$1 - \alpha$ = confidence level (confidence coefficient)

α = alpha = risk μ not in interval

$t_{\alpha/2}$ = t value on upper end

t = standardize t value = # of standard errors

n = sample size

Margin of Error = Amount on either side of \bar{X} to create Interval

s/\sqrt{n} = Standard Error = SD of Sampling Distribution of \bar{X}

df = Degrees of Freedom = $n - \#$ of samples.

Excel

calculate t

Method 1:

$$t_{\text{upper}} = \text{T.INV}(1 - \frac{\alpha}{2}, df)$$

$$t_{\text{lower}} = \text{T.INV}(\frac{\alpha}{2}, df)$$

calculate MoE

Method 2:

$$\left\{ \begin{array}{l} \text{Margin} \\ \text{of} \\ \text{Error} \end{array} \right\} = \text{CONFIDENCE.T}(\alpha, s, n)$$

calculate
 \bar{X}, s, MoE

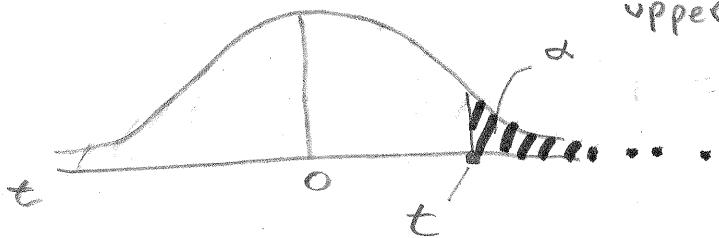
Method 3:

$$\left\{ \begin{array}{l} \bar{X} \pm s \\ \text{and} \\ \text{Margin of Error} \end{array} \right\}$$

Data Analysis, Descriptive Statistics
(Data Ribbon Tab)

Excel t functions & NOT KNOWN

1 tail to right



$$\text{upper } t = \text{T.INV}(1-\alpha, df)$$

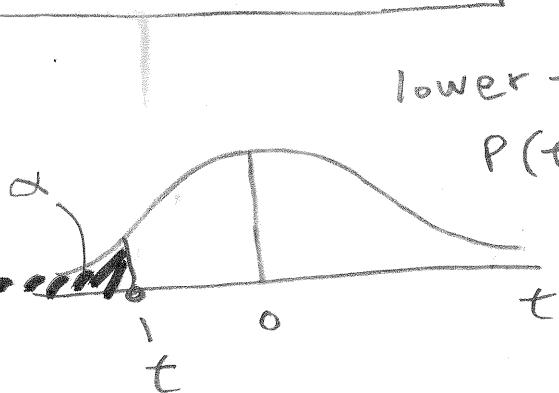
$$P(t \text{ or greater}) = \text{T.DIST.RT}(t, df)$$

OR
↓

$$1 - \text{T.DIST}(t, df, 1)$$

Neg. or Pos.
& will give
Prob. to Right

1 tail to left



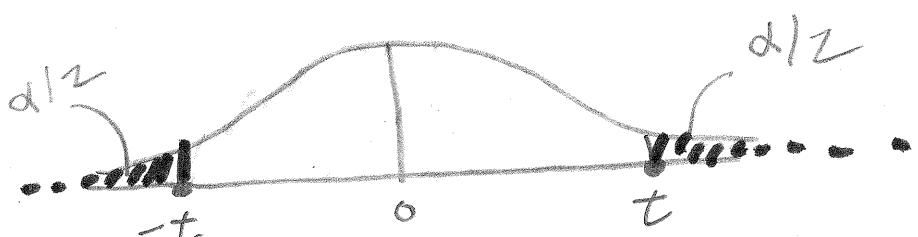
$$\text{lower } t = \text{T.INV}(\alpha, df)$$

$$P(t \text{ or less}) = \text{T.DIST}(t, df, 1)$$

Neg. or Pos.
will give
Prob. to Left

1 = Probability
0 = height of
curve
(used for
charting)

2 tail



$$\text{upper } t = \text{T.INV.2T}(\alpha, df)$$

$$P(t \text{ or greater} \text{ OR } -t \text{ or less}) = \text{T.DIST.2T}(t, df)$$

$$P(\text{Between } -t \text{ & } t) = \text{T.DIST}(t, df, 1) - \text{T.DIST}(-t, df, 1)$$

Charting t Distributions

$$\left\{ \begin{array}{l} \text{Height of} \\ \text{curve} \end{array} \right\} = T.DIST(t, df, 0)$$

t = cumulative
prob. from
 $-\infty$ to t

0 = Height of
curve

Histograms from Sample Data

- ① For a clue about shape of population distribution
- ② Not conclusive, but sometimes it is best into we have.
- ③ Helps to determine sample size & whether you can use t distribution.

Note: Create of t-distribution

- created by William Sealy Gosset
- Gosset wrote under name "Student"
- Gosset was Oxford graduate in Mathematics
- Gosset worked for Guinness Brewery (Dublin, Ireland)
- Developed t Distribution while working on small scale materials & temperature experiments
- Gosset's mathematical development of t distribution assumed that populations had normal distributions.
- Research shows that t distribution can be applied in some situations where pop. Distribution deviates from normal distribution

Guidelines for using t distributions when pop not

Normal

- ① If pop. is normal, if you create 100 similar 95% confidence intervals, 95 should contain pop. M & S not.
- ② If pop Not normal, if you create 100 similar 95% confidence intervals, about 95 should contain M & about S will not if:

- Art form
- Professional
Judgement.

- Distribution almost symmetrical $n \geq 15$
- Distribution near normal $n \geq 30$
- Highly skewed or has outliers $n \geq 50$

Example 2 σ NOT KNOWN

P. 16

Printers Manufacturer

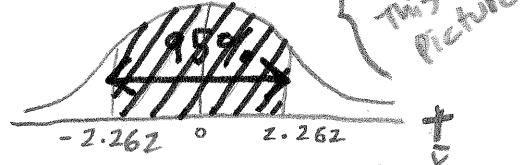
Printers Manufacturer is a manufacturer of ink jet printers. They would like to include as part of their advertising the number of pages a user can expect from an ink cartridge. A sample of 10 cartridges was taken with a mean of 2409 pages and a standard deviation of 304 pages.

8. List variables: $n = 10$
 $df = 10 - 1 = 9$

$$\bar{X} = 2409 \text{ pages}$$

$$s = 304 \text{ pages}$$

$$\frac{s}{\sqrt{n}} = \frac{304}{\sqrt{10}} = 96.13$$



9. What is the best estimation for the population mean?

$\bar{X} = 2409 \text{ pages}$ is the best estimate for our pop. mean.
 'Point estimate'

10. Determine a 95% Confidence Interval

a. State the level of confidence, state the df, then look up t in back of book.

$$\text{level of confidence} = .95 \quad df = 10 - 1 = 9 \quad \Rightarrow \quad t = 2.262$$

$$\text{Excel} = TINV(1 - .95, 9)$$

$$2409 \pm 217.45$$

Margin of error

b. Using the correct confidence interval formula, calculate the confidence limits

$$\bar{X} \pm t \frac{s}{\sqrt{n}}$$

$$2409 \pm 2.262 * 96.13$$

$$2192 \pm 2626$$

"The limits for the 95% CI are
 2192 pages and 2626 pages.

11. Conclusions:

- The customer can expect an average of 2409 pages/cartridge.
- The typical usage ranges from 2192 pages to 2626 pages (given a 95% CI). The margin of error is .

The Margin of error is 217.45 pages.

We are 95% sure that the population mean lies between 2192 & 2626 pages.

- If we constructed 100 similar intervals, 95 of them would contain pop. mean.

Proportions

17

Proportions: The fraction, ratio, or percent indicating the part of the sample or the population having a particular trait of interest.

Example: A recent survey of Highline students indicated that 98 out of 100 surveyed thought textbooks were too expensive.

$$x = \# \text{ of successes} = \text{particular trait of interest} = 98$$

$$n = \text{sample size} = 100$$

$$\bar{p} = \frac{x}{n} = \text{sample proportion}$$

Notes about proportions:

- ① The sample proportion is our best estimate of our population proportion, π .
- ② Proportions are Nominal Level data.
- ③ Success or failure sounds Binomial, right? so...

In order to build a confidence interval for proportions:

must verify:

- ① Are there a fixed # of trials?
- ② Are results independent?
- ③ Does each trial result in success or failure?
- ④ π stay same each trial?
- ⑤ $n * \pi > 5$
- $n * (1 - \pi) > 5$

construct of Confidence Interval for Proportion

alternative notation

$X = \# \text{ of successes}$

$n = \text{sample size}$

$\frac{X}{n} = \bar{p} = \text{sample proportion} = \text{best estimate for } \pi \text{ or } p$

{Confidence
Limits}

$$\text{Confidence Limits} = \bar{p} \pm Z \sqrt{\frac{\bar{p} * (1 - \bar{p})}{n}}$$

Margin
of
error

$p = \pi = \text{population proportion}$

Sampling error = $\bar{p} - \pi$ or $\bar{p} - p$

\bar{p} = point estimate

or books Notations:

$$\bar{p} \pm Z_{\alpha/2} \sqrt{\frac{\bar{p} * (1 - \bar{p})}{n}}$$

area on upper side

Excel confidence Interval for Proportions

(19)

- $\{ \text{sample size } n \} = \text{COUNTA}(\text{Range})$
- $\{ \begin{matrix} \text{Count} \\ \text{Successes} \end{matrix} \} = \text{COUNTIF}(\text{Range}, \text{criteria})$
- $\{ Z \text{ value} \} = \text{NORM.S.INV}(1 - \alpha/2)$
- $\{ \begin{matrix} \text{Standard} \\ \text{Error} \end{matrix} \} = \text{SQRT}(\bar{p} * (1 - \bar{p}) / n)$

Furniture Land South

Furniture Land South surveyed their customers ($n = 600$) to see if they liked the new line of durable foam furniture decorated in bright colors. 414 said they were excited about the new line. All the binomial tests are met.

1. List variables:

best estimate for π :

$$n = 600$$

$$x = 414$$

$$P = \frac{x}{n} = \frac{414}{600} = .69$$

$$\rightarrow \text{for } \pi$$

success = excited

Failure = not excited

1 Fixed # trials = yes $n = 600$

2 Independent ✓

3 S/I/F ✓

4 $\pi = .69$ same each time ✓

5 $\pi * n > 5 \Rightarrow .69 * 600 = 414 > 5$

6 $n * (1 - \pi) > 5 \Rightarrow 600 * .31 = 186 > 5$

2. What is the best estimation for the population proportion?

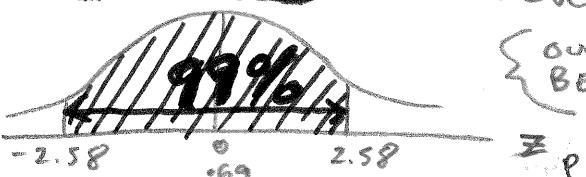
The point estimate $P = .69$ is the best estimate for our population proportion

3. Determine a 99% Confidence Interval

- a. State the level of confidence,

then find z

$$\text{Level of confidence} = .99 \quad .99/2 =$$



$$\left\{ \begin{array}{l} \text{our interval} \\ \text{below this} \\ \text{picture} \end{array} \right. \quad Z = .258 \quad .495$$

- b. Using the correct confidence interval formula, calculate the confidence limits

$$P \pm Z \sqrt{\frac{P * (1-P)}{n}} = .69 \pm 2.58 \sqrt{\frac{.69 * .31}{600}} \Rightarrow .69 \pm .0487$$

.6413 and .7387 are the confidence limits.

4. Conclusions:

- The owner can be 99% sure that the population proportion (% of customers excited about new product) is between .6413 and .7387. The new product will probably be popular.
- .0487 is Margin of error

- 99 of 100 similarly constructed intervals would contain pop. proportion.

Confidence Interval:

Point estimator \pm Margin of Error



If we know this
we can solve
for n , sample
size.

Point Estimator \pm Error (E)

$$E = Z_{\alpha/2} * \frac{\sigma}{\sqrt{n}}$$

$$\sqrt{n} = \frac{Z_{\alpha/2} * \sigma}{E}$$

Sample size
for Interval
Estimation for
Pop. Mean

$$n = \left(\frac{Z_{\alpha/2} * \sigma}{E} \right)^2$$

These are
both
estimates
Judgment
Must be
Used

Sample size
for Interval
Estimation for
Pop. Proportion

$$n = p * (1-p) * \left(\frac{Z_{\alpha/2}}{E} \right)^2$$

sample
size
for
confidence
interval
pop mean

$$N = \frac{(Z_{\alpha/2} * \sigma)^2}{E^2}$$

Requires estimate: ways to get estimate

- ① use estimate for sigma from previous studies (or gov. studies)
- ② pilot study before you run experiment
- ③ estimate large & small values:

$$\frac{\text{large} - \text{small}}{4} = \text{approximate } \sigma$$

* Always Round up when calculating $n^{!!}$

Excel: = ROUNDUP(n estimate, 0)

sample size
for
confidence
Interval
pop
proportion

$$N = \left(\frac{Z_{\alpha/2}}{E} \right)^2 * P * (1 - P)$$

P.23

Requires estimate

How to estimate:

- ① use estimate from previous studies or available data or gov. data
- ② Pilot study before you run experiment
- ③ Best guess (requires judgment)
- ④ None of these apply use $P = .50$
 - * $.5 * (1 - .5) = .25$ and as P increases it will never get bigger than $.25$). This will always yield largest sample size.

Note: Margin of Error for proportion

* margin of error for estimating pop. proportion is almost always less than 0.10. Gallup & Harris commonly use $E = .03$ or $E = .04$