

# Introduction to Sampling & Statistical Inference

BU609-4

## Objectives

- Understand sampling distributions
- Understand estimation concept
- Apply sampling and estimation concepts when working with means, proportions, variances (one or two populations)

1

2

## Parameters vs Statistics

- |            |                     |
|------------|---------------------|
| ● $\mu$    | ● $\bar{x}$ (x-bar) |
| ● $\sigma$ | ● $s$               |
| ● $p$      | ● $\hat{p}$ (p-hat) |

Also need to consider number of distinct populations; can work with one, two or more populations

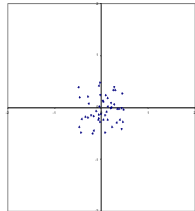
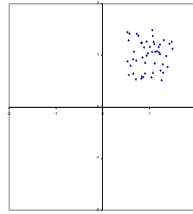
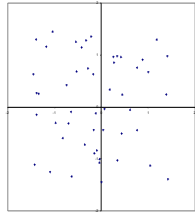
3

## Selecting Statistic

- Selecting sample statistic to estimate parameter value depends on characteristics of statistic
- Estimator's desirable characteristics
  - **Unbiasedness:** unbiased estimator is one whose expected value is equal to parameter it estimates
  - **Consistency:** unbiased estimator is said to be consistent if difference between estimator and parameter grows smaller as sample size increases
  - **Relative efficiency:** For two unbiased estimators, one with smaller variance is said to be relatively efficient
- Fortunately, statisticians do this for us

4

## Bias vs Efficiency



Upper Left: Unbiased, Inefficient

Upper Right: Biased, Efficient

Lower Left: Unbiased, Efficient

5

## Sampling Distributions: Introduction

- In real life calculating **parameters** of populations is usually impossible
- Rather than investigating entire population, take a sample, calculate a **statistic** related to **parameter** of interest, and make an inference
- **Sampling distribution** of **statistic** is prob. dist. of statistic (mean, std dev, proportion)
  - can be Normal, or another known distribution
  - can use this to estimate how close statistic is to parameter

6

## Opinion Polls – Example

- “Rich get Richer, Poor get Poorer, Poll Shows” (*The Record*, 24-10-01)
  - “About 40 per cent of Canadians believe they got poorer during the last 10 years, an opinion poll suggests. The Leger Marketing survey indicated 39.6 per cent of respondents thought they were not as rich as 10 years ago, compared with 32.6 per cent who believed they were richer. Another 26.2 per cent reported no change. The survey of 1,508 Canadians has a margin of error of 2.6 percentage points, 19 times out of 20.”

7

## Sampling Distribution of Mean

- An example
  - A die is thrown infinitely many times. Let  $x$  represent the number of spots showing on any throw.
  - Probability distribution of  $x$  is shown below
    - » What type of distribution is this?

$x$	1	2	3	4	5	6
$p(x)$	1/6	1/6	1/6	1/6	1/6	1/6

$$E(X) = 1(1/6) + 2(1/6) + 3(1/6) + \dots = 3.5$$

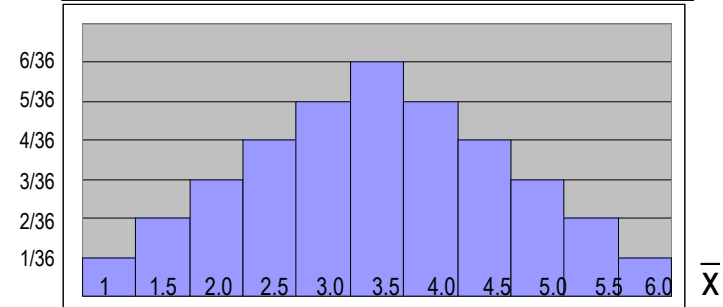
$$V(X) = (1-3.5)^2(1/6) + (2-3.5)^2(1/6) + \dots = 2.92$$

8

- Suppose we want to estimate  $\mu$  from the mean  $\bar{x}$  of a sample of size  $n = 2$
- What is distribution that  $\bar{x}$  can follow?

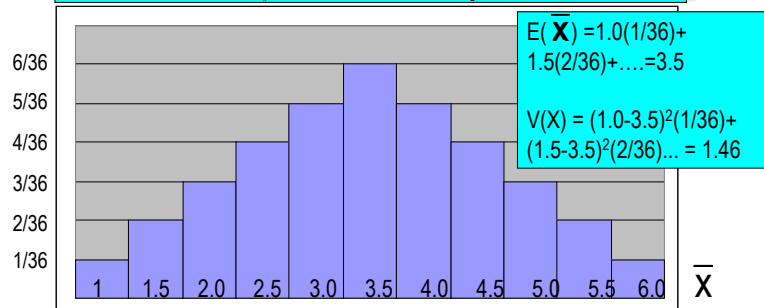
Sample	Mean	Sample	Mean	Sample	Mean
1	1,1	1	13	3,1	2
2	1,2	1,5	14	3,2	2,5
3	1,3	2	15	3,3	3
4	1,4	2,5	16	3,4	3,5
5	1,5	3	17	3,5	4
6	1,6	3,5	18	3,6	4,5
7	2,1	1,5	19	4,1	2,5
8	2,2	2	20	4,2	3
9	2,3	2,5	21	4,3	3,5
10	2,4	3	22	4,4	4
11	2,5	3,5	23	4,5	4,5
12	2,6	4	24	4,6	5

Sample	Mean	Sample	Mean	Sample	Mean
1	1,1	1	13	3,1	2
2	1,2	1,5	14	3,2	2,5
3	1,3	2	15	3,3	3
4	1,4	2,5	16	3,4	3,5
5	1,5	3	17	3,5	4
6	1,6	3,5	18	3,6	4,5
7	2,1	1,5	19	4,1	2,5
8	2,2	2	20	4,2	3
9	2,3	2,5	21	4,3	3,5
10	2,4	3	22	4,4	4
11	2,5	3,5	23	4,5	4,5
12	2,6	4	24	4,6	5



Sample	Mean	Sample	Mean	Sample	Mean
1	1,1	1	13	3,1	2
2	1,2	1,5	14	3,2	2,5
3	1,3	2	15	3,3	3
4	1,4	2,5	16	3,4	3,5
5	1,5	3	17	3,5	4
6	1,6	3,5	18	3,6	4,5
7	2,1	1,5	19	4,1	2,5
8	2,2	2	20	4,2	3
9	2,3	2,5	21	4,3	3,5
10	2,4	3	22	4,4	4
11	2,5	3,5	23	4,5	4,5
12	2,6	4	24	4,6	5

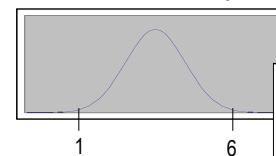
Note :  $\mu_{\bar{x}} = \mu_x$  and  $\sigma_{\bar{x}}^2 = \frac{\sigma_x^2}{2}$



$n = 5$

$\mu_{\bar{x}} = 3.5$

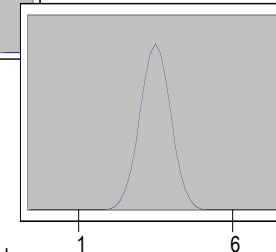
$\sigma_{\bar{x}}^2 = .5833 (= \frac{\sigma_x^2}{5})$



$n = 10$

$\mu_{\bar{x}} = 3.5$

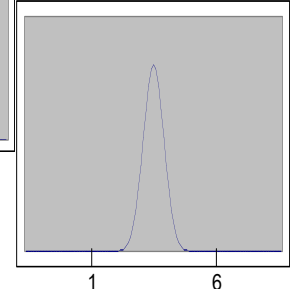
$\sigma_{\bar{x}}^2 = .2917 (= \frac{\sigma_x^2}{10})$



$n = 25$

$\mu_{\bar{x}} = 3.5$

$\sigma_{\bar{x}}^2 = .1167 (= \frac{\sigma_x^2}{25})$



Notice that  $\sigma_{\bar{x}}^2$  is smaller than  $\sigma_x^2$ . The larger the sample size the smaller  $\sigma_{\bar{x}}^2$ . Therefore,  $\bar{X}$  tends to fall closer to  $\mu$ , as the sample size increases.

## Central Limit Theorem (CLT)

- If a truly random sample is independently drawn from **any** population,
- **Sampling distribution** of **sample mean** is approximately Normal (for sufficiently large sample size)
- Larger the sample size, more closely sampling distribution of sample mean will resemble Normal distribution

13

## Guidelines for C.L.T.

- If sampling from a relatively symmetric dist. (Uniform, binomial with  $p=.5$ ),  $n$  can be small (5 - 15)
- If sampling from a relatively skewed dist (binomial with  $p=.1$ , or bimodal),  $n$  should be  $\geq 30$
- When  $n \geq 30$ , variability of sample mean is typically reasonable

14

## Sampling Dist. of Sample Mean

1.  $\mu_{\bar{x}} = E(\bar{x}) = \mu_x$
2.  $\sigma_{\bar{x}}^2 = s^2 = \frac{\sigma_x^2}{n} \quad s = \sigma / \sqrt{n}$
3. If  $x$  is Normal,  $\bar{x}$  is Normal. If  $x$  is nonnormal  $\bar{x}$  is approx Normal for sufficiently large sample size ( $n \geq 30$ )

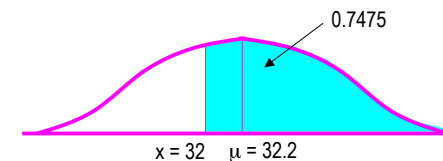
Question : How do we know  $\sigma$  ?

Note : SAMPLE std dev called "std error of mean"

15

## Bottling Example

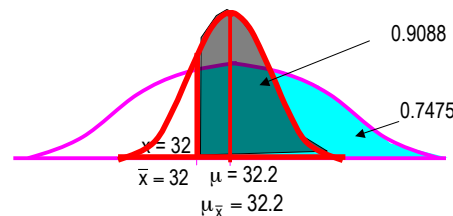
- Amount of beverage in each bottle is normally distributed with mean of 32.2 ounces and standard deviation of .3 ounces.
  - Find probability that bottle bought by customer will contain more than 32 ounces.



16

## Exercise (cont.)

- Find probability that carton of four bottles will have mean of more than 32 ounces of pop per bottle
- What about case of 8 bottles?



17

## Finance Application of CLT

- An initial investment  $I$  grows over time
  - In first month grew by factor of 1.1, then shrank by factor of .97 in second month, and grew by factor of 1.2 in third month
- Value at end of 3 months =  $I * 1.1 * .97 * 1.2$ 
  - If you studied logarithms, may remember that they convert multiplications to additions
  - Logarithm of value after 3 months is  $\log(I) + \log(1.1) + \log(.97) + \log(1.2)$
- Therefore log of value of investment after large number of periods is sum of bunch of random numbers, and is hence normally distributed
  - An uncertain number whose logarithm is normally distributed is known as log normal
  - This is why investment prices are often modeled as log normal random variables

18

## Sampling Distributions Recap

- Basis of **Statistical Inference**
- Particular sampling distribution depends on parameter (mean, variance, proportion, etc.)
  - Can take on many different forms
- Always based on **random sampling!**
- Knowing sampling distribution, and our sample statistics, we can make inferences about population parameter

19

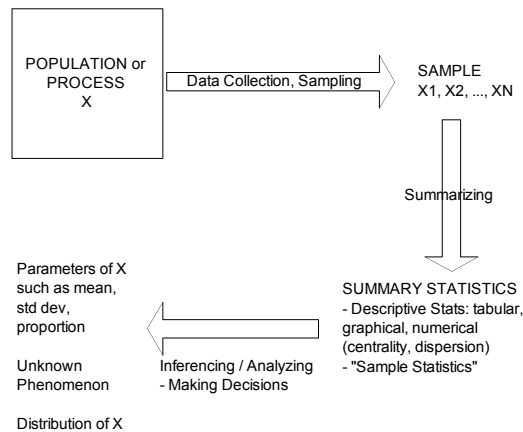
## Statistical Inference: Introduction

- Statistical inference is process by which we acquire information about **populations** from **samples**
- Two procedures for making inferences:
  - Estimation
  - Hypotheses testing
- Critical that sample is a **random sample**

20

## Statistical Model

### KEY STATISTICAL CONCEPTS



21

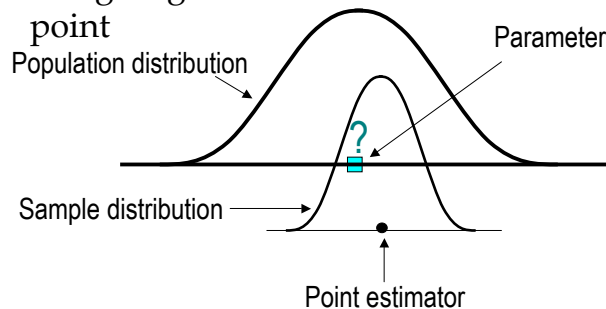
## Concept of Estimation

- Objective of estimation is to determine value of a **population parameter** on the basis of a **sample statistic**
- Two types of estimators
  - Point Estimator
    - » Single value
  - Interval estimator
    - » Interval likely to contain true parameter value

22

## Point Estimator Example (Mean)

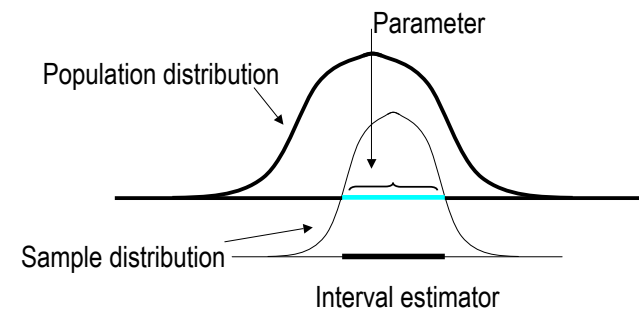
- Draws inference about parameter by using single value or point



23

## Interval Estimator Example

- Interval estimator draws inferences about parameter by using an interval
- Interval estimator is affected by sample size



24

## Confidence Statements

- In making estimations, **want to state how sure/unsure we are**
  - **Level of significance** =  $\alpha$ 
    - » used for hypothesis tests; .01, .05, .1 common
  - **Confidence level** =  $(1 - \alpha) \times 100\%$ 
    - » used for confidence intervals; 90%, 95%, 99% common
- Random sample from pop. has  $(1 - \alpha) \times 100\%$  probability (or confidence) that confidence interval contains unknown pop. Parameter
  - 95% most commonly used
- In repeated samplings,  $(1 - \alpha) \times 100\%$  of conf. intervals will contain unknown pop. parameter

25

## Commonly used Values

- Confidence level gives  $\alpha$  (prob. of being wrong; tail area)
- Corresponds to a particular **# of standard deviations** around point estimate
  - commonly called “z-value”; obtain from Normal table or use Excel

Confidence level	$\alpha$	$\alpha/2$	(z) # std dev
0.90	0.10	0.05	1.645
0.95	0.05	0.025	1.96
0.98	0.02	0.01	2.33
0.99	0.01	0.005	2.575

27

## Est. Population Mean ( $\sigma$ Known)

- How is interval estimator produced from sampling distribution?
  - To estimate  $\mu$ , sample of size  $n$  is drawn from population, and its mean ( $\bar{x}$ ) is calculated
  - By CLT,  $\bar{x}$  is  $N(\mu, \sigma/\sqrt{n})$
  - Point estimate is sample mean ( $\bar{x}$ )
  - Interval is  $\bar{x} \pm \text{so many std deviations}$ 
    - » # of std dev depends on desired confidence level (how confident we want to be that interval contains population mean  $\mu$ )

26

## Excel & z-values

- Based on standard Normal distribution ( $\mu = 0, \sigma = 1$ )
- **NORMSDIST(z)**
  - Returns LHS probability for given z-value
- **NORMSINV( $\alpha$ )**
  - Returns z value for given LHS area
- **Drawing a diagram** helps, since tables are symmetrical and these Excel functions are not; **watch this!**

28

## Confidence Interval for Mean ( $\sigma$ known)

$$\left[ \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

- How do we know  $\bar{x}$ ,  $\sigma$ ,  $n$ ,  $z_{\alpha/2}$ ?
- If we want to be 95% confident that this interval encapsulates population mean, what value of  $z$  would we use?
- What proportion of time will mean be outside this interval?

29

## Marketing Example

- Number & types of television programs and commercials targeted at children are affected by amount of time children watch TV
  - Survey conducted among 100 North American children, in which they were asked to record # of hours they watched TV/week
  - Sample mean is 27.191 hrs
  - Population standard deviation of TV watch was known to be  $\sigma = 8.0$  hrs
  - Estimate watch time with 95% confidence

30

## Solution

- Parameter to be estimated is  $\mu$ , mean time of TV watched per week per child
  - Need to compute interval estimator for  $\mu$

$$\begin{aligned} \bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} &= 27.191 \pm z_{.025} \frac{8.0}{\sqrt{100}} \\ &= 27.191 \pm 1.96 \frac{8.0}{\sqrt{100}} = 27.191 \pm 1.57 = [25.621, 28.761] \\ \text{LCL} &= 25.621, \text{UCL} = 28.761 \end{aligned}$$

Since  $1 - \alpha = .95$ ,  $\alpha = .05$ .  
Thus  $\alpha/2 = .025$ .  $Z_{.025} = 1.96$

31

## Width of Interval Estimate

- Width of interval estimate is function of:
  - population standard deviation ( $\sigma$ )
  - confidence level ( $1 - \alpha$ ); z-value ( $z_{\alpha/2}$ )
  - sample size ( $n$ )
  - has nothing to do with  $N$  (pop. size)!
- Smaller is better (better understanding of where  $\mu$  is)
  - *How can we influence this?*

32



## Inference About Pop. Mean ( $\sigma$ Unknown)

- How can we know  $\sigma$  when we are only working with a sample?
- When  $\sigma$  is unknown, use its point estimator  $s$ , and **sampling distribution** is **Student t** distribution
- Requires Normal (or approx Normal) population
- As sample size  $n$  increases, t-dist approaches N-dist ( $n > 200$ )

33

## t-Statistic

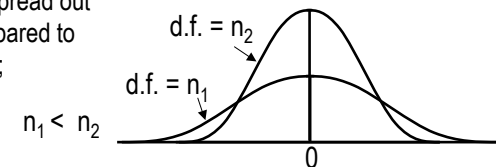
$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \quad t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \quad df = n - 1$$

When sampled population is **normally distributed**, statistic **t** is **Student t** distributed.

### "Degrees of freedom"

- function of sample size
- determines how spread out distribution is (compared to normal distribution);
- $= n - 1$

**t distribution** is mound-shaped, and symmetrical around zero.



34

## Confidence Interval for Mean ( $\sigma$ unknown, $s$ known)

$$\left[ \bar{x} - t_{\alpha/2, df} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2, df} \frac{s}{\sqrt{n}} \right]$$

- This application of t-distribution is **robust**. It is an adequate approximate sampling distribution for moderately non-normal populations.
- Use Excel (or tables) for t-values

35

## Excel & Student t-Distribution

- TDIST(x, df, tails)
  - returns "tail" probability (area under curve)
  - $x$  = specific t-value
  - tails must be 1 or 2
- TINV(prob, df)
  - returns t-value for specified probability & degrees of freedom (**based on 2-tails**)
  - df** = **degrees of freedom** (sample size minus one,  $n-1$ )
  - prob** =  $\alpha$  (tails area)

36

## *Practical vs Statistical Significance*

- Take a large enough sample, and you can usually show statistically significant difference between two sample means (or other parameters such as variance or proportion)
- But is this difference of practical/economic significance?
  - Managerial decision re what is of practical significance

37

## *Paired Difference Test*

- Use for difference of means, when samples are **not independent** (matched pairs)
  - More efficient!
  - Calculate difference for each sample pair
  - Then work with this new series of “paired differences”; same as working with a single sample
  - Symbols:  $\mu_D$  and  $\sigma_D$ ,  $\bar{x}_D$  and  $s_D$ ,  $n_D$

38

## *Recap re Confidence Intervals*

- Formula depends on **Parameter** of interest & its **Sampling Distribution**
  - Parameter defines sample statistic & sampling distribution (sometimes requirements re parent distribution)
  - We focused on  $\mu$ ; text also covers  $\sigma^2$ ,  $p$  for 1 or 2 populations
- Point Estimate  $\pm$  **so many std deviations** of sampling distribution
  - # of std dev depends on desired confidence level  $(1-\alpha)*100\%$ ; 2-tailed, hence  $\alpha/2$

39

## *Additional Topics (Level 1)*

- Sampling & Estimation
  - Sampling in investment analysis
  - Data-snooping bias
  - Sample selection bias: survivorship, delisting
  - Using t-value tables
  - Conf. intervals when sampling from non-N population

40