

This Class - Organization

• All the information and material is on my webpage: https://www.bauer.uh.edu/rsusmel/4397/4397.htm

• Textbook: Introductory Econometrics for Finance, 4th edition or older, by Chris Brooks.

• Exams

- Midterms: September 26 and October 24 (tentative)
- Final: According to UH Schedule (December 12, 5PM-8PM)
- Research Project (Paper): November 5
- Case Presentations (presentation and discussion of project): After project, to be scheduled during Office Hours.

- Homework: Aug 29, Sep 12, Sep 24 & Nov 26 (or Dec 3, depending on progress of class) 2

This Class – Organization

• Final Grade

A weighted average: Midterms (2): **40%**

Final (substitution with a short paper encouraged): 30%

Homework: 10%

Class Project: 15%

Presentation: 5%

• R Program

After this class, **install R** in your machine. Previous students had a strong preference for **R Studio**.

Next class, we will introduce R and run some simple R programs. More advance programs will be run throughout the semester.

This Class – Comments from Previous Classes

• Class is very technical

We will go over many stats concepts and definitions, some derivations and lots of formulas. But, we will apply each topic to a finance setting.

• Class covers a lot of material

We will cover as much as possible of Chris Brooks's textbook. Last year we were able to cover 7 chapters (& previous years, 8 chapters).

• Instructor (me) goes fast

Questions are a great way to slow me down. Ask questions, please. All questions and interruptions are appreciated

This Class - Comments from Previous Classes

• More comments from previous classes

"It was difficult to keep awake in the class."

"Very technical course."

"Very organized lectures and course."

"I had problems with R almost all semester. Only at the end, I was able to understand what was going on."

"Learned a lot. Good course. Enjoyed the exams. One of the good courses (in program)."

"This course is much too quantitative."

"We covered too much info too fast."

"This is one of the few courses that I feel I've truly earned what I'm paying the university."

"He fried my brain."

This class – Overview

• This is an applied technical class, with some econometric theory and many stats concepts, followed by related financial applications.

• We will review many math and statistical topics.

• Some technical material may be new to you, for example Linear Algebra. The new material is introduced to simplify exposition. You will not be required to have a deep understanding of the new material, but you should be able to follow the intuition.

• This is not a programming class, but we will use R to estimate models. I will cover some of the basics in class and I will run in class all the programs you need to run.

• For some students, the class will be dry ("*He fried my brain*," a student said in 2020.)

This class – Main Applied Topics

- How do we measure returns and risks of financial assets?
- Is the equity premium (excess returns of stocks over bonds) really that high?
- Can we measure left-tail (unusual/extreme negative) risk?
- How do we determine a good model for financial assets? Is the CAPM a good model for stock returns? What about Fama-French?
- Can we explain asset returns?
- How can one explain variations in stock returns across various stocks?
- Are asset returns predictable? In the short run? In the long run?
- Are markets efficient?
- Does the risk of an asset vary with time? What are the implications? How can one model time-varying risk?

7

This class – Main Technical TopicsUnderstanding Distributions and Moments Testing and Confidence Intervals Bootstrap Linear Regression Testing Hypothesis in the Classical Linear (Regression) Model Finding a Good Statistical and Financial Model Forecasting Time Series Models Efficiency & EMH. (Application of Many Concepts) Time-varying Volatility (if time allows). Integrated Time Series and Co-integration (if time allows).

This class - Goals

• <u>Goal of the class</u>: Students should be comfortable with applied regressions, testing financial hypothesis and forecasting.

• Secondary goal: Get students familiar with R & running R programs.



What is Econometrics?

Example: We want to estimate expected annual excess returns for Exxon, $E[r_{XOM} - r_f]$.

- Simple approach: Compute the average excess return of XOM in the past 50 years to estimate the annual expected return. We get an annualized **2.81%** estimate. Good, we use data & statistics.

- More sophisticated approach: Add economic theory. That is, use econometrics. For example, we use the Capital Asset Pricing Model (**CAPM**) that states a linear relation, in equilibrium, between excess market returns, $r_M - r_f$, & excess returns, $r_i - r_f$, for any asset *i*:

$$\mathbf{E}[r_i - r_f] = \beta_i \mathbf{E}[(r_M - r_f)]$$

We get data on r_i , r_f , and r_M . Then, we use a linear regression to estimate β_i .

11

What is Econometrics? • Steps: (1) Economic Theory: The CAPM: $E[r_{i=XOM} - r_{f}] = \beta_{i} E[(r_{M} - r_{f})]$ (2) Data: Collect data, 1973-now for r_{XOM} , r_{f} , & r_{M} . (3) Mathematical Statistics: Use a linear regression to estimate β_{i} : $r_{XOM} - r_{f} = \alpha_{XOM} + \beta_{XOM} (r_{M} - r_{f}) + \varepsilon_{XOM}$ \Rightarrow Compute b_{XOM} (the regression estimator of β_{XOM}), say 0.665. • Now, we are ready to compute the expected excess return for XOM: Expected excess XOM return: $b_{XOM} * \text{Average}(r_{M} - r_{f})$. : 0.665 * 0.0727 = 0.0483 (= 4.83%)

What is Econometrics?

Issues

Of course, there are many potential problems (& assumptions) behind our estimation of expected excess returns for XOM.

We can raise many issues regarding what we have done:

- Economic Theory question: Is the CAPM a good model?
- Data question: Is 50 years enough data?
- Stats question: Do the assumptions behind the linear CAPM hold?

What is Financial Econometrics?

• Financial Econometrics is applied econometrics to financial data. That is, we study the statistical tools that are needed to analyze and address the specific types of questions and modeling challenges that appear in analyzing financial data.

• Always keep in mind that almost in all cases, financial data is not *"experimental data.*" We have no control over the data. We have to learn how to deal with the usual problems in financial data.

- Typical applications of econometric tools to finance:
 - Describe data. For example, expected returns & volatility.
 - Test hypothesis. For example, are stocks riskier than bonds?
 - Build and test models. For example, the different Fama-French factor models used to estimate expected returns.



• In general, in finance we deal with **trade-offs**. The usual trade-off: Risk & Return.

- How do we measure risk and return?
- Can we predict them?
- How do we measure the trade-off?
- How much should I be compensated for taking a given risk?

• Thus, we will be concerned with quantifying rewards and risks associated with uncertain outcomes.

What is Financial Econometrics?

This Lecture

We will review some basic concept of Probability and Statistics:

- Random Variable
- Distribution Functions
- Descriptive Statistics: Moments
- Population & Sample
- Sample Statistics & Estimators
- Law of Large Numbers (LLN)
- Central Limit Theorem (CLT)
- Sampling Distributions

Review – Random Variable

• In probability, a *random variable* (RV), or *stochastic variable*, is described informally as a variable whose values depend on outcomes of an *experiment*. (Experiment: Act/process with an unknown outcome).

Examples:

1. We throw two coins and count the number of heads.

2. We define X = 1 if the economy grows two consecutive quarters and

X = 0, otherwise. (This is an example of a *Bernouille* (or *indicator*) RV.)

3. We read comments from IBM's CEO and compute IBM's return.

4. We count the days in a week that XOM has a positive return.

5. We look at a CEO and write his/her highest education degree.

6. We compute the weekly sign of stock returns of two unrelated firms: Positive (U: up) or negative (D: down). We count the times at least one stock is up: {D,U}, {U,D}, {U,U}.

17

Review – Random Variable

• For some RVs, it is easy to enumerate all possible outcomes. For instance, for the fourth (XOM) above: {0, 1, 2, 3, 4, 5}. But, for some RVs, it can be complicated. For example, for the third (IBM) example: $\{-100\%, K\}$, where *K* is a large positive number.

• The set of all possible outcomes is called *sample space*, denoted by Ω .

• An event A is a set containing outcomes from the sample space. For example, for the IBM return example, the return is between -1.64% and 1.64% is an event.

• The collection of all possible events is Σ . For example, for the IBM return example, {(1.1%, 1.2%), (-0.02%, -0.001%), (2.00%, 12.657%), (-5%, 5%), (-100%, -13.95%),(-1.64%, 1.64%), (0%, 350%), ... }

Review – Random Variable

• In general, a RV is a *function* whose domain is the sample space, Ω . It produces numbers. For instance, in example 6 above, instead of using $\{U, U\}$ when both stocks go up, we use 2.

Mathematically, **X**: $\Omega \rightarrow R$.

Remark: The name "random variable" is confusing; it is just a function!

• We put some mathematical structure (pdf, pmf, CDF) to the concept of RV to describe what is more/less likely to happen to the (randomly determined) events.

For example, we would like to know which event is more/less likely for the IBM example: Is (1.1%-1.2%) more likely than (-0.02%, -0.001%)?₁₉

Review – PMF for a Discrete RV

• **Definition**: Let *X* be a discrete RV. Let p(x) be a function with the following properties:

- $1.0 \le p(x) \le 1$
- 2. $\sum_{i=1}^{\infty} p(x_i) = 1$
- 3. $P[a \le X \le b] = \sum_{a \le x \le b} p(x)$

Then, p(x) is called the *probability function* or *probability mass function* (pmf) of *X*. We use p(x) to describe the behavior of a discrete RV.

Example: Suppose the discrete RV X is the number of days in a week that XOM has a positive return. Using Property 3, we can compute the probability that XOM's has a positive return in 3 or more days in a week:

 $P[a = 3 \le X \le b = 5] = p(x = 3) + p(x = 4) + p(x = 5)$ ²⁰

Review – PDF for a Continuous RV

Definition: Let *X* be a continuous RV, like stock returns. Let f(x) be a function defined for $-\infty < x < \infty$ with the following properties:

1. $f(x) \ge 0$. 2. $\int_{-\infty}^{\infty} f(x) dx = 1$. 3. $P[a \le X \le b] = \int_{a}^{b} f(x) dx$

Then, f(x) is called the *probability density function* (pdf) of X. We use the pdf to describe the behavior of a continuous RV.

Example: Suppose the continuous RV *X* is IBM's daily stock returns and we know the pdf. Then, using Property 3, we can compute the probability that IBM's daily return is between a = -1.64% and b = 1.64%:

$$P[-1.64\% \le X \le 1.64\%] = \int_{a=-1.64}^{b=1.04} f(x)dx$$





Review – Popular PDFs: Normal Distribution

• When $\mu = 0$ and $\sigma^2 = 1$, we call the distribution *standard normal*. We write $X \sim N(0, 1)$. This is the distribution that is tabulated.

The normal distribution is often used to describe or approximate any variable that tends to cluster around the mean. It is the most assumed distribution in economics and finance: rates of return, growth rates, IQ scores, observational errors, etc.

• The central limit theorem (CLT) provides a justification for the normality assumption when the sample size, *N*, is large.

Notation: PDF: $X \sim N(\mu, \sigma^2)$ CDF: $\Phi(x)$



• Let the continuous RV *X* have density function):

$$f(x) = \begin{cases} \frac{\lambda^{\alpha}}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} & x \ge 0\\ 0 & x < 0 \end{cases}$$

where α , $\lambda > 0$ and $\Gamma(\alpha)$ is the gamma function evaluated at α .

Then, *X* is said to have a *Gamma distribution* with parameters α and λ , denoted as $X \sim \text{Gamma}(\alpha, \lambda)$ or $\Gamma(\alpha, \lambda)$.

It is a family of distributions, with special cases:

- Exponential Distribution, or $\text{Exp}(\lambda)$: $\alpha = 1$.
- Chi-square Distribution, or χ^2_{ν} : $\alpha = \nu/2$ and $\lambda = \frac{1}{2}$.



Review – Popular PDFs: Other Distributions

• Other distributions we will use in class:

(1) t-distribution: A ratio of a standard normal and the square root of a χ^2_{ν} divided by ν . That is, let Y ~ N(0, 1) and W ~ χ^2_{ν} , then

$$t=\frac{Y}{\sqrt{W/\nu}}\sim t_{\nu}.$$

Below, we plot a simulated t-distribution with $\nu = 5$ (in red), along a normal distribution (in blue). It has thicker tails. As ν increases, t_{ν} converges to a N(0, 1) distribution.



Review – Popular PDFs: Other Distributions (2) F-distribution: A ratio of two independent χ^2 distributions, divided by their degrees of freedom. That is, let $Z_1 \sim \chi^2_{\nu_1}$ and $Z_2 \sim$ $\chi^2_{\nu_2}$, then $F = \frac{Z_1 / \nu_1}{Z_2 / \nu_2} \sim F_{\nu, \nu_2}$ Simulated F-distribution with df=5 & 10 80 99 4 Density 8 00 5 8 5 10 15 0 N = 1000 Bandwidth = 0.1719 28





Review – The Empirical Distribution

• The empirical distribution (ED) of a dataset is simply the distribution that we observe in the data.

The ED is a discrete distribution that gives equal weight to each data point, assigning a 1/N probability to each of the original N observations.

We form a cumulative distribution function, F^* , that is a step function that jumps up by 1/N at each of the N data point:

$$F^*(x) = 1/N \sum_{i=1}^N I(x_i \le x),$$

where I(.) is the indicator function:

$$I(x_i \le x) = 1, \qquad \text{if } x_i \le x$$

$$I(x_i \le x) = 0, \qquad \text{if } x_i > x$$





Example: We use a histogram to estimate the distribution (pdf) of a RV. Let X = Percentage changes in the CHF/USD exchange rate = e_f Data: Monthly - January 1973 to March 2024 (N = 615 observations).



Review - Moments of Random Variables

• The moments of a random variable X are used to describe the behavior of the RV (discrete or continuous).

Definition: k^{tb} Moment

Let X be a RV (discrete or continuous), then the k^{tb} moment of X is:

$$\mu_{k} = E(X^{k}) \qquad = \begin{cases} \sum_{x} x^{k} p(x) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{x} x^{k} f(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

• The first moment of X, $\mu = \mu_1 = E(X)$ is the center of gravity of the distribution of X.

• The higher moments give different information regarding the shape of the distribution of X.

Review – Moments of Random Variables

Example: Suppose X is the number of days in a week that XOM has a positive return. We want to know the first moment, the mean, of the distribution. That is,

$$\mu_1 = \sum_x x \, p(x) = 0 * p(x = 0) + 1 * p(x = 1) + 2 * p(x = 2) + + 3 * p(x = 3) + 4 * p(x = 4) + 5 * p(x = 5)$$

Suppose we can describe *X* with a Binomial distribution, with p=0.52. That is, XOM has a 52% probability of having a positive return. Then,

<u>Interpretation</u>: The expected number of days in week with positive returns for XOM is 2.6 days.

Note: For a continuous RV, we need to integrate to get moments.

35

Review – Moments of a RV

Definition: Central Moments

Let X be a RV (discrete or continuous). Then, the k^{tb} central moment of X is defined to be:

$$\mu_k^0 = E[(X - \mu)^k] = \begin{cases} \sum_{\substack{x \\ \infty \\ -\infty}} (x - \mu)^k p(x) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} (x - \mu)^k f(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

where $\mu = \mu_1 = E(X)$ = the first moment of *X*.

• The central moments describe how the probability distribution is distributed about the center of gravity, μ .

Review - Moments of a RV

• The first central moments is given by:

 $\mu_1^0 = E[X - \mu] = 0$

The second central moment depends on the *spread* of the probability distribution of X about μ . It is called the variance of X and is denoted by the symbol $\sigma^2 = \operatorname{var}(X)$:

$$\mu_2^0 = E[(X - \mu)^2] = \operatorname{Var}[X] = \sigma^2$$

The square root of var(X) is called the *standard deviation* of X and is denoted by the symbol $\sigma = SD(X)$. We also refer to it as *volatility*:

$$\sqrt{\mu_2^0} = \sqrt{E[(X-\mu)^2]} = \sigma$$

37

Review – Moments of a RV

Example: Suppose *X* is the number of days in a week that XOM has a positive return. We want to know the second central moment, $\mu_2^0 = \sigma^2$ (& volatility, σ). (Recall that $\mu_1 = \mu = 2.6$ days). Then, $\sigma^2 = \sum_x (x - \mu)^2 p(x) = (0 - 2.6)^{2*} p(x = 0) + (1 - 2.6)^{2*} p(x = 1) + (2 - 2.6)^2 * p(x = 2) + (3 - 2.6)^{2*} p(x = 3) + (4 - 2.6)^2 * p(x = 4) + (5 - 2.6)^{2*} p(x = 5)$ Again, assume *X* follows a Binomial distribution, with p=0.52. Then, $\sigma^2 = (0 - 2.6)^2 * 0.0255 + (1 - 2.6)^2 * 0.1380 + (2 - 2.6)^2 * 0.2990 + (3 - 2.6)^2 * 0.3240 + (4 - 2.6)^2 * 0.1755 + (5 - 2.6)^2 * 0.0380 = 1.24802 \implies \sigma = \operatorname{sqrt}(1.24802) = 1.117148$ Interpretation: The volatility of *X* is 1.12 days. Note: Again, for a continuous RV, we need to integrate to get moments³⁸





Review - Moments of a RV

Example (continuation): We want to know $\mu_3^0 \& \gamma_1$ for X = the number of days in a week that XOM has a positive return. Then, $\mu_3^0 = \sum_x (x - \mu)^3 p(x) = (0 - 2.6)^{3*} p(x = 0) + (1 - 2.6)^{3*} p(x = 1)$ $+ (2 - 2.6)^3 * p(x = 2) + (3 - 2.6)^{3*} p(x = 3)$ $+ (4 - 2.6)^3 * p(x = 4) + (5 - 2.6)^{3*} p(x = 5)$ Again, assume X follows a Binomial distribution, with p=0.52. Then, $\mu_3^0 = (0 - 2.6)^3 * 0.0255 + (1 - 2.6)^3 * 0.1380 + (2 - 2.6)^3 * 0.2990$ $+ (3 - 2.6)^3 * 0.3240 + (4 - 2.6)^3 * 0.1755 + (5 - 2.6)^3 * 0.0380$ $= -0.04989 \implies \gamma_1 = \frac{\mu_3^0}{\sigma^3} = \frac{-0.04989}{(1.11715)^3} = -0.03578522$ Interpretation: X has a small, but negative skewness. The left tail is a little bit longer.

41

Review – Moments of a RV: Skewness

Skewness and Economics

For returns:

- Zero skew means symmetrical gains and losses.
- Positive skew suggests many small losses and few rich returns.
- Negative skew indicates lots of minor wins offset by rare major losses.

• In financial markets, stock returns at the firm level show positive skewness, but at the aggregate (index) level show negative skewness.

• From horse race betting and from U.S. state lotteries there is evidence supporting the contention that gamblers are not necessarily risk-lovers but **skewness-lovers**: Long shots are overbet (positive skewness loved!).





Review - Moments of a RV: Kurtosis

• Typical financial returns series has $\gamma_2 > 0$. Below, I simulate a series with $\mu=0, \sigma=1, \gamma_1=0$ & kurtosis = 6 ($\gamma_2=3$), overlaid with a standard normal distribution. Fat tails are seen on both sides of the distribution.



Review – Moments and Expected Values

• Note that moments are defined by expected values. We define the expected value of a function of a continuous RV X, g(X), as

$$E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x)dx$$

• If X is *discrete* with probability function p(x)

$$E[g(X)] = \sum_{x} g(x)p(x) = \sum_{i} g(x_i)p(x_i)$$

Examples:

 $g(x) = (x - \mu)^2 \qquad \Rightarrow \operatorname{E}[g(x)] = \operatorname{E}[(x - \mu)^2]$ $g(x) = (x - \mu)^k \qquad \Rightarrow \operatorname{E}[g(x)] = \operatorname{E}(x - \mu)^k]$

• We estimate expected values with sample averages. The Law of Large Numbers (LLN) tells us they are *consistent* estimators of expected values. 46

Review – Population and Sample

Definition: Population

A population is the totality of the elements under study. We are interested in learning something about this population.

Examples: Number of alligators in Texas, percentage of unemployed workers in cities in the U.S., total return of all stocks in the U.S., 10-year Japanese government bond yields from 1960-2024.

Usually, a complete enumeration of all the values in the population is impractical or impossible. Thus, the descriptive statistics describing/generating the population –i.e., the *population parameters*– will be considered unknown.

<u>Note</u>: A Random Variable (RV) X defined over a population is called the population RV. The population RV generates the data. We call the population RV the "*Data Generating Process*," or DGP.



Review – Population and Sample

Example: The total returns of the stocks on the S&P 500 index is *not* a random sample.

• In general, in finance and economics, we do not deal with random samples. The collected observations will have issues that make the sample not a true random sample.

<u>Remark</u>: In mathematical terms, given a random variable X with distribution F, a *random sample* of length N is a set of N **independent, identically distributed** (*i.i.d.*) random variables with distribution F.

• We will estimate population parameters using sample analogues: mean, sample mean; variance, sample variance; β , **b**; etc.

49

Review - Samples and Types of Data

• The samples we collect are classified in three groups:

• **Time Series Data**: Collected over time on one or more variables, with a particular *frequency* of observation. Example: we record for 10 years the monthly S&P 500 returns, or 10' IBM returns.

<u>Usual notation</u>: x_t , t = 1, 2, ..., T.

• **Cross-sectional Data**: Collected on one or more variables collected at a single point in time. Example: today we record all closing returns for the members of the S&P 500 index.

<u>Usual notation</u>: x_i , i = 1, 2, ..., N.

• **Panel Data**: Cross-sectional data collected over time. Example: the CRSP database collects daily prices of all U.S. traded stocks since 1962. Usual notation: $x_{i,t}$, i = 1, 2, ..., N & t = 1, 2, ..., T. 50







Review – Sample Statistic: Estimators

• The definition of a sample statistic is very general. For example, by definition $(x_1 + x_N)/2$ is a statistic; we could claim that it estimates the population mean of the variable X. However, this is probably not a good estimate.

• We would like our estimators, $\hat{\theta}$, to have certain desirable properties, for example, low bias and low variance, where bias and variance are: - Bias[$\hat{\theta}$] = E[$\hat{\theta}$] - θ

$$-\operatorname{Var}[\widehat{\theta}] = \operatorname{E}[(\widehat{\theta} - \operatorname{E}[\theta])^2]$$

Ideally, we would like to have $\hat{\theta}$ with both low bias and low variance, but as we would see later, in general, we have a trade-off between these two properties.



Review – Estimators: Properties

• The first two properties for estimators hold for samples of any size, not just large samples –i.e., when $N \rightarrow \infty$.

We associate bias with lack of accuracy and efficiency/variance with uncertainty.

• It is common to evaluate an estimator using the Mean Squared Error (MSE), which combines bias and variance:

$$MSE[\widehat{\theta}] = E[(\widehat{\theta} - \theta)^2] = Bias[\widehat{\theta}]^2 + Var[\widehat{\theta}].$$

RS 2024 copyright. Not to be posted online/shared without written consent from author



• Long history: Gerolamo Cardano (1501-1576) stated it without proof. Jacob Bernoulli published a rigorous proof in 1713.

Theorem (Weak LLN)

Let X_1, \ldots, X_N be N mutually independent & identically distributed RVs, each having mean, μ , and SD, σ , finite. We say $\{X_N\}$ is *i.i.d.*

Let $\overline{X} = \frac{\sum_{i=1}^{N} X_i}{N}$. Then for any $\delta > 0$ (no matter how small) $P[|\overline{X} - \mu| < \delta] = P[|\mu - \delta < \overline{X} < \mu + \delta] \rightarrow 1$, as $N \rightarrow \infty$ • There are many versions of the LLN. It is a general result: A sample average as the sample size goes to infinite tends to its expected value. Also written as $\overline{X}_N \xrightarrow{p} \mu$. (convergence in probability) ⁵⁷



Review – Central Limit Theorem (CLT)

• Let X_1, \ldots, X_N be a sequence of *i.i.d.* RVs with finite mean μ , and finite variance σ^2 . Then, as N increases, the distribution of the (normalized) sample mean, \overline{X}_N , approaches the normal distribution with mean μ and variance σ^2/N .

This theorem is sometimes stated as $\frac{\sqrt{N}(\bar{X}-\mu)}{\sigma} \xrightarrow{d} N(0,1)$

where \xrightarrow{d} means "the limiting distribution (asymptotic distribution) is" (or *convergence in distribution*).

• Many versions of the CLT. This one is the Lindeberg-Lévy CLT.

<u>Remark:</u> The CLT gives only an asymptotic distribution. We usually take it as an approximation, since N is finite. In these cases, the notation goes from $\stackrel{d}{\rightarrow}$ to $\stackrel{a}{\rightarrow}$.

Review – Expected Values & Sample Averages

• We estimate expected values with sample averages. For example, the first moment, μ , & the second central moment, σ^2 , are estimated by:

$$\bar{X} = \frac{\sum_{i=1}^{N} X_i}{N}$$

$$s^2 = \frac{\sum_{i=1}^{N} (X_i - \bar{X})^2}{N-1} \quad (N - 1 \text{ adjustment needed for } \mathbb{E}[s^2] = \sigma^2)$$

• They are both *unbiased* estimators of their respective population moments That is,

 $E[\overline{X}] = \mu$ $E[s^{2}] = \sigma^{2} \qquad ``\mu \& \sigma^{2} \text{ are population parameter''}$

<u>Note</u>: Unbiased estimator = "On average, we get the population parameter."

RS 2024 copyright. Not to be posted online/shared without written consent from author

Review – Sampling Distributions: \bar{X}

• All statistics, T(X), are functions of RVs and, thus, they have a distribution. Depending on the sample, we can observe different values for T(X), thus, the finite sample distribution of T(X) is called the *sampling distribution*.

For the sample mean \overline{X} , if the X_i 's are normally distributed, then the sampling distribution is normal with mean μ and variance σ^2/N . Or $\overline{X} \sim N(\mu, \sigma^2/N)$.

<u>Note</u>: If the data is not normal, the CLT is used to approximate the sampling distribution by the asymptotic one, usually after some

manipulations. Again, in those cases, the notation goes from \xrightarrow{d} to \xrightarrow{a} .

• The SD of the sampling distribution is called the *standard error* (SE). Then, $SE(\overline{X}) = \sigma/sqrt(N)$.

61

62

Review – Sampling Distributions: \overline{X}

• Summary for \overline{X} :

Sampling distribution: $\overline{X} \sim N(\mu, \sigma^2/N)$ (if data normal)

Mean: $E[\overline{X}] = \mu$

Variance: $\operatorname{Var}[\overline{X}] = \sigma^2 / N$.

<u>Note</u>: If the data is not normal (& N is large), the CLT can be used to approximate the sampling distribution by the asymptotic one:

$$\bar{X} \xrightarrow{a} N(\mu, \sigma^2/N)$$





Review – Sampling Distributions: s^2

• For the sample variance s^2 , if the X_i 's are normally distributed, the sampling distribution is derived from this result:

 $(N-1) s^2/\sigma^2 \sim \chi^2_{N-1}.$

We use the properties of a χ^2_{ν} to derive the mean & variance of s^2 :

Property 1. Let $Z \sim \chi_{\nu}^2$. Then, $E[Z] = \nu$. **Property 2.** Let $Z \sim \chi_{\nu}^2$. Then, $Var[Z] = 2 * \nu$.

Application: Let $Z = (N - 1) s^2 / \sigma^2 \sim \chi^2_{N-1}$ From Property 1: $E[(N - 1) s^2 / \sigma^2] = N - 1$ $\Rightarrow E[s^2] = \sigma^2$

From Property 2:
$$\operatorname{Var}[(N-1) s^2/\sigma^2] = 2 * (N-1)$$

 $\Rightarrow \operatorname{Var}[s^2] = 2 * \sigma^4/(N-1)$
 $\Rightarrow \operatorname{SE}(s^2) = \operatorname{SD}(s^2) = \sigma^2 * \operatorname{sqrt}[2/(N-1)]$ 65

Review – Sampling Distributions: s^2 • Summary for s^2 of normal variates: Sampling distribution: $(N-1) s^2/\sigma^2 \sim \chi^2_{N-1}$. Mean: $E[s^2] = \sigma^2$ Variance: $Var[s^2] = 2 * \sigma^4/(N-1)$. Note: If the data is not normal (& N is large), the CLT can be used to approximate the sampling distribution by the asymptotic one: $s^2 \xrightarrow{a} N(\sigma^2, \sigma^4 * (\kappa - 1)/N)$ where $\kappa = \frac{\mu_q^4}{\sigma^4}$ (recall when data is normal, $\kappa = 3$).

