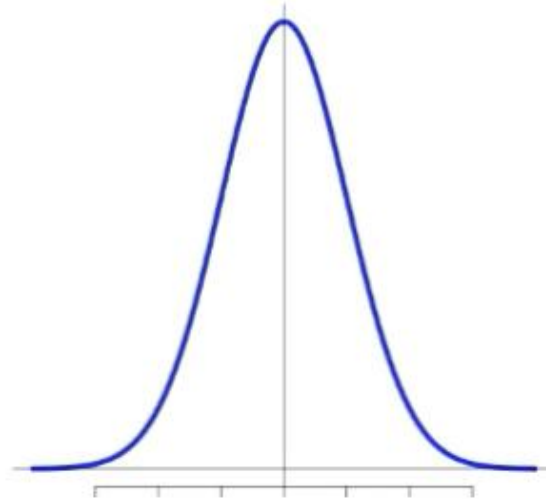# Learning Objectives

- ► Upon completion of this lecture, you will be able to:
    - ► Describe the basic properties of the normal curve
    - ► Describe how the normal distribution is completely defined by its mean and standard deviation
    - ► Recite the 68–95–99.7% rule for the normal distribution with regards to standard deviations

# The Normal Distribution—1

► The normal distribution is a theoretical probability distribution that is perfectly symmetric about its mean (and median and mode)
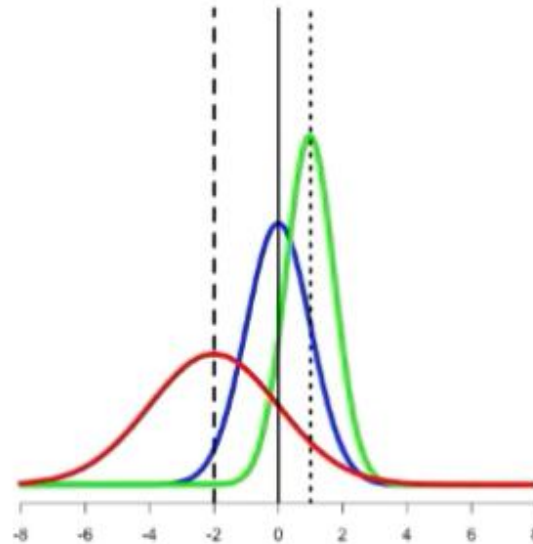  ► A "bell"-like shape

# The Normal Distribution—2

► The normal distribution is also called the "Gaussian distribution" in honor of its inventor Carl Friedrich Gauss

# Defining Quantities for any Normal Distribution

▶ Normal distributions are uniquely defined by two quantities: a mean ($\mu$) and standard deviation ($\sigma$)

▶ There are literally an infinite number of possible normal curves for every possible combination of ($\mu$) and ($\sigma$)
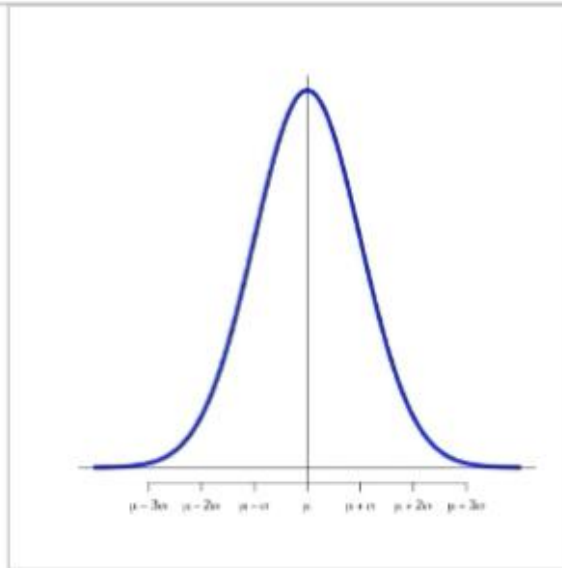
# Underlying Formula for Normal Distribution

► This function defines the normal curve for any given (μ) and (σ)

► The proportion of values falling between a and b under a normal curve is given by:

$$\int_a^b \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(x-\mu)^2}{2\sigma^2}} \, dx$$

# Structural Properties of the Normal Distribution (Curve)—1

▶ All normal distributions, regardless of mean and standard deviation values, have the same structural properties:

    ▶ Mean = median (= mode)

    ▶ Values are symmetrically distributed around the mean

    ▶ Values "closer" to the mean are more frequent than values "farther" from the mean
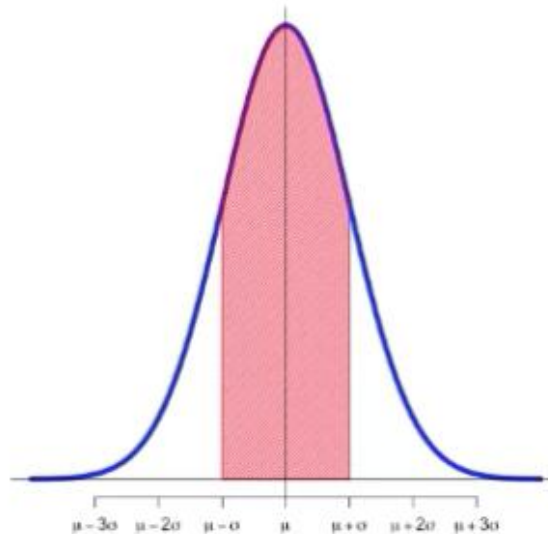
# Structural Properties of the Normal Distribution (Curve)—2

- ► All normal distributions, regardless of mean and standard deviation values, have the same structural properties:
  - ► The entire distribution of values described by a normal distribution can be completely specified by knowing just the mean and standard deviation
  - ► Since all normal distributions have the same structural properties, we can use a reference distribution, called the *standard normal distribution*, to elaborate on some of these properties
  - ► In the next section, we'll show that any normal distribution can be easily rescaled to this standard normal distribution

# The 68–95–99.7 Rule for the Normal Distribution—1

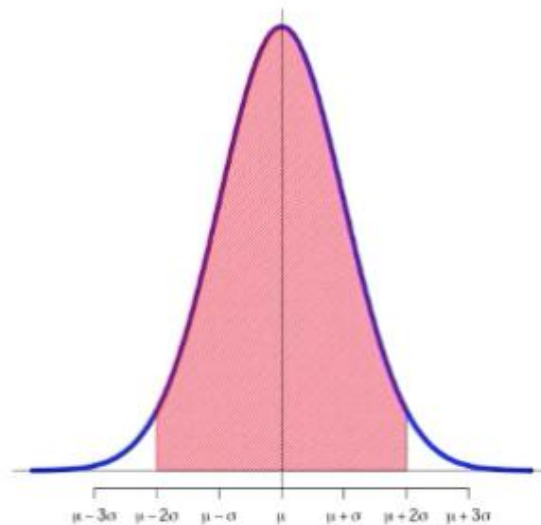▶ 68% of the observations in a normal distribution fall within one standard deviation of the mean

# The 68–95–99.7 Rule for the Normal Distribution—2

► There are several ways to state this. For data whose distribution is approximately normal:
  ► 68% of the observations fall within one standard deviation of the mean
  ► The probability that any randomly selected value is within one standard deviation of the mean is 0.68 or 68%

# The 68–95–99.7 Rule for the Normal Distribution—3

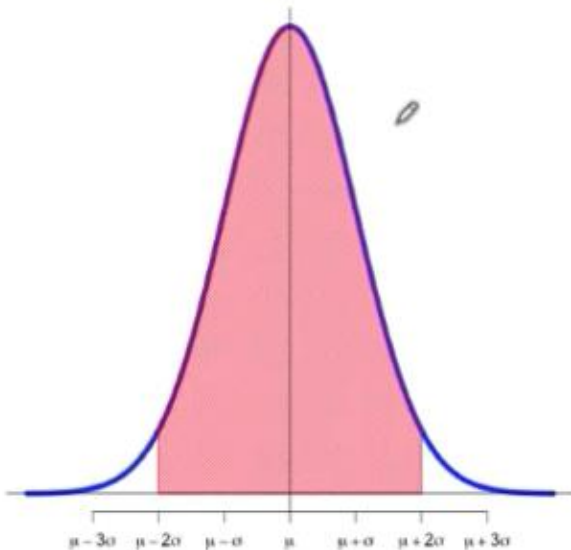▶ 95% of the observations fall within two standard deviations of the mean (truthfully, within 1.96)

# The 68–95–99.7 Rule for the Normal Distribution—4

▶ 99.7% of the observations fall within three standard deviations of the mean

# 2.5th and 97.5th Percentiles of a Normal Distribution

▶ 95% of the observations fall within two standard deviations of the mean (truthfully, within 1.96)



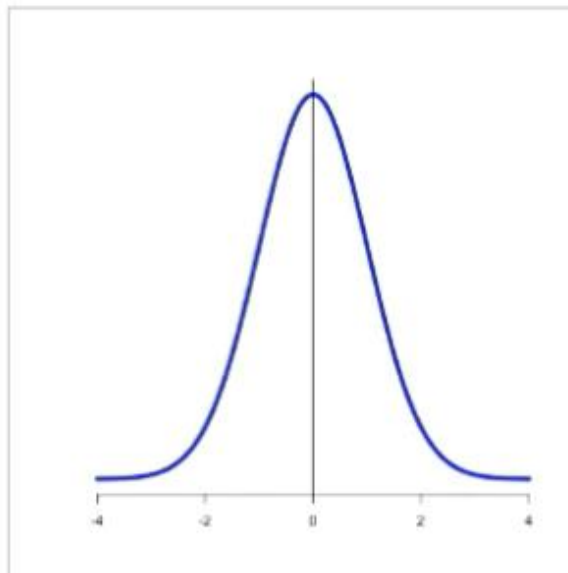▶ The middle 95% of values fall between $\mu -2\sigma$ and $\mu+2\sigma$

▶ 2.5% of the values are smaller than (and hence 97.5% are greater than) $\mu -2\sigma$

▶ 97.5% of the values are smaller than (and hence 2.5% are greater than) $\mu+2\sigma$

# Percentage of Observations Under the Normal Distribution

► Where did this rule come from: in other words, how do I know these relationships?

► What about the percentages under the curve for other standard deviation distances from the mean?

► All of the information I quoted, and much more, can be found in a "standard normal table"

# The Standard Normal Distribution

▶ The standard normal distribution is a normal distribution with mean $\mu = 0$, and standard deviation $\sigma = 1$

▶ Any normal distribution with mean $\mu$ and standard deviation $\sigma$ can be rescaled to a standard normal distribution.

# Percentage of Observations Under the Normal Distribution: Exhibit A

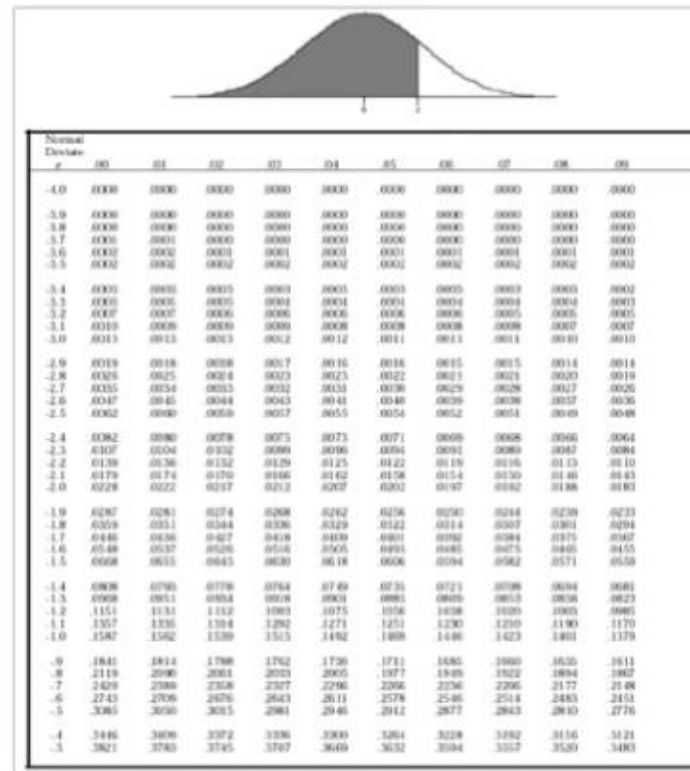# Percentage of Observations Under the Normal Distribution: Exhibit B

# Percentage of Observations Under the Normal Distribution—2

▶ In this class, I will only have you find relevant percentages under a normal curve for some early activities, and this will be done easily using R

▶ Generally speaking, I only want you to be familiar with the 68–95–99.7 rule

▶ Such computations will be wrapped into other analyses later in the course and completely handled by a computer

# Using R to Compute Normal Curve Percentages—1

► We can use R as a calculator, i.e., an automatic standard normal table

► The relevant command that "looks up" values in a standard normal table is:

► For converting any standard deviation value (above or below the mean), $z$, to a corresponding proportion under a normal curve, the syntax is:

```
pnorm(z)
```

# Using R to Compute Normal Curve Percentages—2

► As with any  print version of a standard normal table, it is important to know what information pnorm(z) returns

- ▶ The normal distribution is a theoretical probability distribution that is symmetric and "bell-shaped"

- ▶ There are literally an infinite number of normal distributions, and each can be completely specified by only two quantities: the mean and standard deviation

- ▶ For all normal distributions
  - ▶ 68% of observations described by a normal distribution fall within 1 sd of the mean
  - ▶ 95% of observations described by a normal distribution fall within 2 sds of the mean
  - ▶ 99.7% of observations described by a normal distribution fall within 3 sds of the mean

- ▶ Other such percentages can be found using a standard normal table (available via R)

# Learning Objectives

- ► Upon completion of this lecture, you will be able to:
  - ► Create ranges containing a certain percentage of observations in an (approximately normal) distribution using only an estimate of the mean and standard deviation
  - ► Figure out how far any individual data point is from the mean of its distribution in standardized units (compute a z-score)
  - ► Convert z-scores to statements about relative proportions/probabilities for values that have an (approximately) normal distribution

# The Normal Distribution, Generally Speaking—1

- The normal distribution is a theoretical probability distribution
  - No real data is perfectly described by this distribution

- For example, in a true normal distribution, the tails go on to negative and positive infinity, respectively

# The Normal Distribution, Generally Speaking—2

► However, the distributions of some data will be well approximated by a normal distribution
  ► In such situations we can use the properties of the normal curve to characterize aspects of the data distribution

# Using Sample Estimates and Properties of Normal Distribution to Estimate Percentiles

▶ Using only the sample mean and standard deviation, and assuming normality, let's estimate the 2.5th and 97.5th percentiles SBP in this population

$$2.5^{\text{th}} \%\text{ile: } \bar{x} - 2s = 123.6 - (2 \times 12.9) = 97.8 \text{ mmHg}$$

$$97.5^{\text{th}} \%\text{ile: } \bar{x} + 2s = 123.6 + (2 \times 12.9) = 149.4 \text{ mmHg}$$

▶ Based on this sample data, we estimate that *most (95%) of the men* in this clinical population have systolic blood pressures *between 97.8 and 149.4 mmHg*

▶ Note: the observed 2.5th and 97.5th percentiles of the 113 sample value are 100.7 mmHg and 151.2 mmHg, respectively

# Example: The "z-score"—1

► Suppose you want to use the results from this sample of 113 men from a clinic to evaluate individual male patients relative to the population of all such patients

► For example, suppose a patient in your clinic has a SBP measurement of 130 mmHg. What proportion of men at the clinic have SBP measurements greater than this patient?

# Example: The "z-score"—1

▶ Suppose you want to use the results from this sample of 113 men from a clinic to evaluate individual male patients relative to the population of all such patients

▶ For example, suppose a patient in your clinic has a SBP measurement of 130 mmHg. What proportion of men at the clinic have SBP measurements greater than this patient?

# Example: The "z-score"—3

- ▶ If we translate this measurement of 130 mmHg to units of standard deviation, we can find out how many sample standard deviations this person's SBP is above the sample mean. To do this:

$$\text{Take } \frac{individual\ observation - \bar{x}}{s} = \frac{130 - 123.6}{12.9} = \frac{6.4 \text{ mmHg}}{12.9 \text{ mmHg}/s}$$

$$\approx 0.5 \ standard\ deviations$$

- ▶ Now, the same question can be rephrased as, "What percentage of observations in a normal curve are more than 0.5 SD above its mean?"

# Example: The "z-score"—2

► We want to figure out:

$$\overline{x} = 123.6 \text{ mmHg}$$
$$s = 12.9 \text{ mmHg}$$
$$\widehat{m} = 123.0 \text{ mmHg}$$

**Systolic Blood Pressure (SBP) Measurements**



SPB (mmHg)
Random Sample of 113 Men

# Example: The "z-score"—4

▶ "What percentage of observations in a normal curve are greater than 0.5 SD above its mean?"

▶ Using pnorm in R:



```
> pnorm(.5)
[1] 0.6914625
> |
```

# Example: The "z-score" —5

- So 69% of the observations described by a standard normal curve are less than or equal to 0.5 standard deviations above the mean of 0

- Hence, the remaining 31% are more than 0.5 standard deviations above 0

- In terms of the original question posed, this means that an estimated 31% of the males in this population have blood pressures greater than 130 mmHg (i.e., using only the mean and sd, we have estimated the 69th percentile to be 130 mmHg)

- Just for context/comparison: the 70[th] percentile of the observed 113 values is 130 mmHg

# Example: The "z-score" —6

► Another way to interpret this is as an (estimated) probability: the probability that any males in the population has a blood pressure measurement more than .5 standard deviations above the mean is .31 or 31%

# The "z-score," Generally Speaking—1

▶ The type of computation we did to convert the SBP value of 130 to the number of SDs above (or below) the sample mean is sometimes called a *z-score*

▶ There is nothing special about a z-score; it is simply a measure of the relative distance (and direction) of a single observation in a data distribution relative to the mean of the distribution
   ▶ This distance is converted to units of standard deviation

▶ This is akin to converting kilometers to miles, or dollars to rupees

# The "z-score": A Parallel Example—1

▶  You are an American who is apartment hunting in an unnamed European city. You wish to find an apartment within walking distance (+/- 1.5 miles) of the large organic supermarket, which is on Main Boulevard (E/W). You are only considering apartments on Main Blvd.

▶ The supermarket is 2 km west of the main city square. You are interested in 3 apartments:

    Apt 1 is 6 km west of the city square

    Apt 2 is .75 km west of the city square

    Apt 3 is 1 km *east* of the city square

# Summary—1

▶ The normal distribution is a theoretical probability distribution, which can be completely defined by two characteristics: the mean and standard deviation

▶ No real-world data has a perfect normal distribution; however, some continuous measures are reasonably approximated by a normal distribution

## Summary—2

- When dealing with samples from populations of (approximately) normally distributed data, the distribution of sample values will also be approximately normal. We can use the sample mean and standard deviation estimates, $\bar{x}$ and $s$, to:
  - Create ranges containing a certain percentage of observations, or in other words: estimate the probability that an observed data point falls within a certain range of values
  - Figure out how far any individual data point is from the mean of its distribution in standardized units (compute a z-score)
  - Convert z-scores to statements about relative proportions/probabilities (and, hence, percentiles) for values that have an (approximately) normal distribution

# Learning Objectives

- ▶ Upon completion of this lecture, you will be able to:
    - ▶ Describe situations in which using only the mean and standard deviation of a distribution of values to characterize the entire distribution will not work well
    - ▶ Realize that z-scores are nothing "special"; z-scores are just a (standardized) measure of distance
    - ▶ Understand that z-scores do not necessarily align with the corresponding percentiles for a normal distribution for data that do not follow a normal distribution
    - ▶ Choose the right approach to estimating ranges for individual values, and to computing percentage greater (or less) than a specific value using non-normal data distributions

# The Theoretical Normal Distribution

- ► The normal distribution is a theoretical probability distribution
  - ► No real data is perfectly described by this distribution

- ► For example, in a true normal distribution, the tails go on to negative and positive infinity, respectively

# Applying the Normal Distribution Properties to Sample Data

► The distributions of some data will be well approximated by a normal distribution
  ► In such situations, we can use the properties of the normal curve to characterize aspects of the data distribution

► But the distributions of much data _will not_ be well approximated by a normal distribution
  ► In such situations using the properties of the normal curve to characterize aspects of the data distribution will yield invalid results

# Estimating Percentiles Based on Normality Assumption with Skewed Data (Length of Stay)—2

▶ In this example, using the properties of the normal curve to estimate an interval containing the "middle 95%" of length of stay values for the claims population yields useless results

▶ Better to take the observed 2.5th and 97.5th percentiles of the sample data and report these as an estimate of the "middle 95%"

▶ "Based on this sample data, we estimate that *most (95%)* of the persons making claims in this health care population had length of stays between *1 and 21 days* in 2011."

# Relative Proportions Based on Normality Assumption, Skewed Data (Length of Stay)—3

► Based on these analyses, we estimate that approximately 25% of the claims had total length of stay greater than 5 days

► The above percentage (25%) is a lot smaller than the estimate of 45% we got using the mean and standard deviation to compute a z-score

# Summary

- While sample means and sample standard deviations are useful summary measures regardless of the data for which they are computed, these two quantities do not always help to characterize the data distribution (so far, this has worked only when the data distribution is approximately normal)

- For skewed distributions and others that are not approximately normal, using only the mean and standard deviation to characterize the entire underlying distribution can result in, at best, incorrect results and at worst, nonsensical results