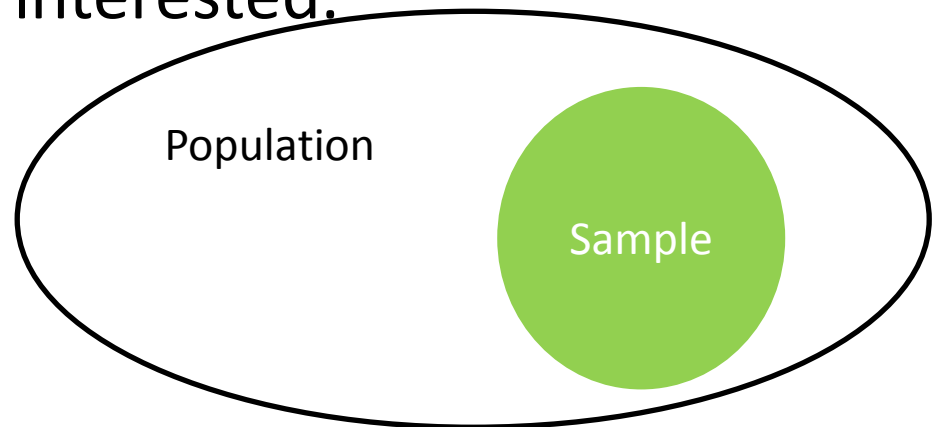# Statistical concepts and market returns

# Populations and samples

- The subset of data used in statistical inference is known as a **sample** and the larger body of data is known as the **population.**
  - The **population** is defined as all members of the group in which we are interested.

# Parameters and Sample Statistics

**A population has parameters, and a sample has statistics.**

- Descriptive statistics that characterize population values are called **parameters**.
  - Examples: mean, median, mode, variance, skewness, kurtosis
- Descriptive statistics that characterize samples are known as **sample statistics**.
  - Examples: sample mean, sample median, sample variance
- By convention, we often omit the term "sample" in front of sample statistics, a practice that can lead to confusion when discussing both the sample and the population.

# Measurement Scales

**Statistical inference is affected by the type of data we are trying to analyze.**

Weak Scales

- **Nominal scales** categorize data but do not rank them.
  - Examples: fund style, country of origin, manager gender
- **Ordinal scales** sort data into categories that are ordered with respect to the characteristic along which the scale is measured.
  - Examples: "star" rankings, class rank, credit rating
- **Interval scales** provide both the relative position (rank) and assurance that the differences between scale values are equal.
  - Example: temperature
- **Ratio scales** have all the characteristics of interval scales and a zero point at the origin.
  - Examples: rates of return, corporate profits, bond maturity

Strong Scales

# Identifying Scales of Measurement

- **State the scale of measurement for each of the following:**

- Credit ratings for bond issues

- Cash dividends per share

- Hedge fund classification types

- Bond maturity in years

# Holding period returns

**Holding period returns are a fundamental building block of the statistical analysis of investments.**

- Holding period returns (HPR) are calculated as the price at the end of the period plus any cash distribution during the period minus the beginning of period price, all divided by the beginning period price.

$$R_t = \frac{P_t - P_{t-1} + D_t}{P_{t-1}}$$

- For this stock, which is nondividend paying, the HPRs are:

| Time | Price | HPR | Time | Price | HPR |
|------|-------|------|------|-------|--------|
| 0 | 27.00 | — | 7 | 25.90 | 2.38% |
| 1 | 25.77 | −4.57% | 8 | 27.01 | 4.28% |
| 2 | 24.73 | −4.04% | 9 | 28.20 | 4.42% |
| 3 | 24.32 | −1.64% | 10 | 29.52 | 4.68% |
| 4 | 24.39 | 0.28% | 11 | 31.63 | 7.16% |
| 5 | 24.71 | 1.34% | 12 | 35.25 | 11.43% |
| 6 | 25.30 | 2.35% | | | |

# Frequency distributions

**A tabular display of data summarized into intervals is known as a frequency distribution.**

Constructing a frequency distribution:
1. Sort the data in ascending order.
2. Calculate the range of the data, defined as
   Range = Maximum value – Minimum value.
3. Decide on the number of intervals in the frequency distribution, $k$.
4. Determine interval width as Range/$k$.
5. Determine the intervals by successively adding the interval width to the minimum value to determine the ending points of intervals, stopping after reaching an interval that includes the maximum value.
6. Count the number of observations falling in each interval.
7. Construct a table of the intervals listed from smallest to largest that shows the number of observations falling in each interval.

# Frequency Distributions

## Focus on: Holding Period Returns

- Suppose we have 12 holding period return observations from a non-dividend-paying stock, sorted in ascending order:

  −4.57, −4.04, −1.64, 0.28, 1.34, 2.35, 2.38, 4.28, 4.42, 4.68, 7.16, and 11.43.

- Using $k$ = 4, we have intervals with width of 4.

- The resulting frequency distribution is

| Interval | Absolute Frequency |
|---|---|
| −4.57 ≤ observation < −0.57 | 3 |
| −0.57 ≤ observation < 3.43 | 4 |
| 3.43 ≤ observation < 7.43 | 4 |
| 7.43 ≤ observation ≤ 11.43 | 1 |

# Relative and cumulative frequency
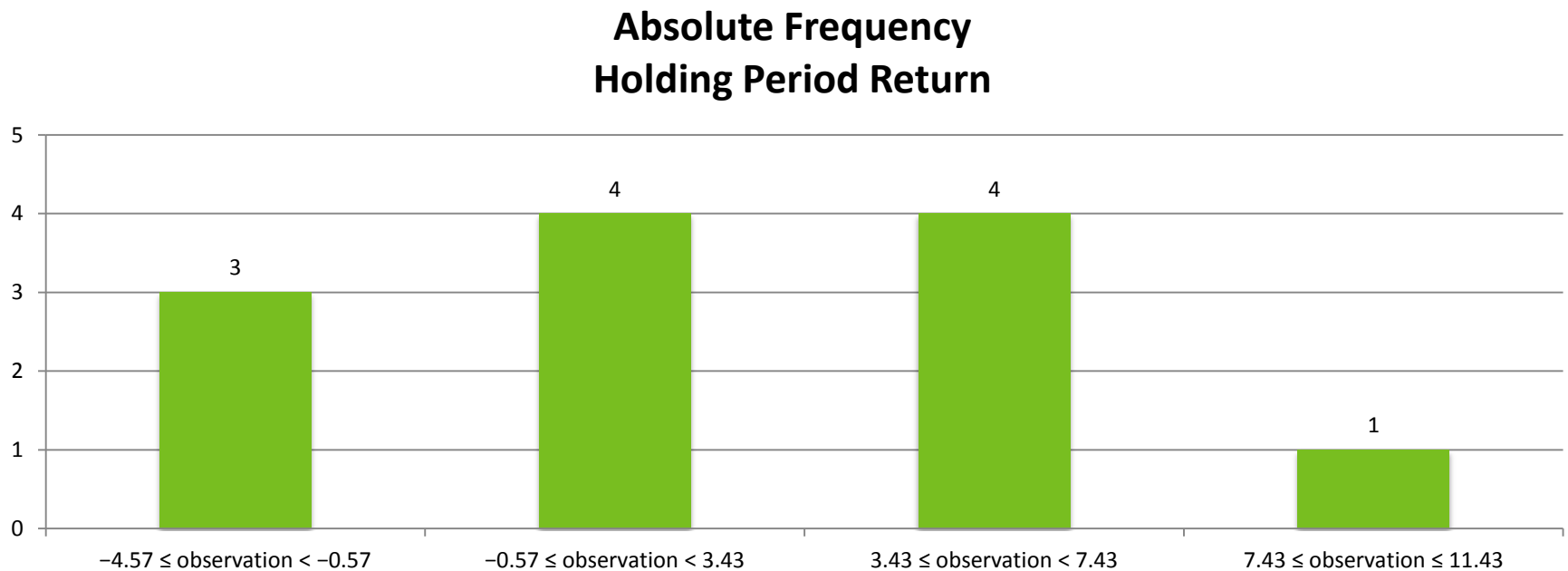
## Focus on: Holding Period Returns

- **Relative frequency** is the absolute frequency divided by the total number of observations.
- **Cumulative (relative) frequency** is the relative frequency of all observations occurring before a given interval.

| Interval | Absolute Frequency | Relative Frequency | Cumulative Frequency |
|---|---|---|---|
| −4.57 ≤ observation < −0.57 | 3 ÷ 12 | 0.250 | 0.250 |
| −0.57 ≤ observation < 3.43 | 4 | 0.333 | 0.583 |
| 3.43 ≤ observation < 7.43 | 4 | 0.333 | 0.917 |
| 7.43 ≤ observation ≤ 11.43 | 1 | 0.083 | 1.000 |

# Histograms

## Focus on: Holding Period Returns

- Histograms are the graphical representation of a frequency distribution.

**Absolute Frequency**
**Holding Period Return**



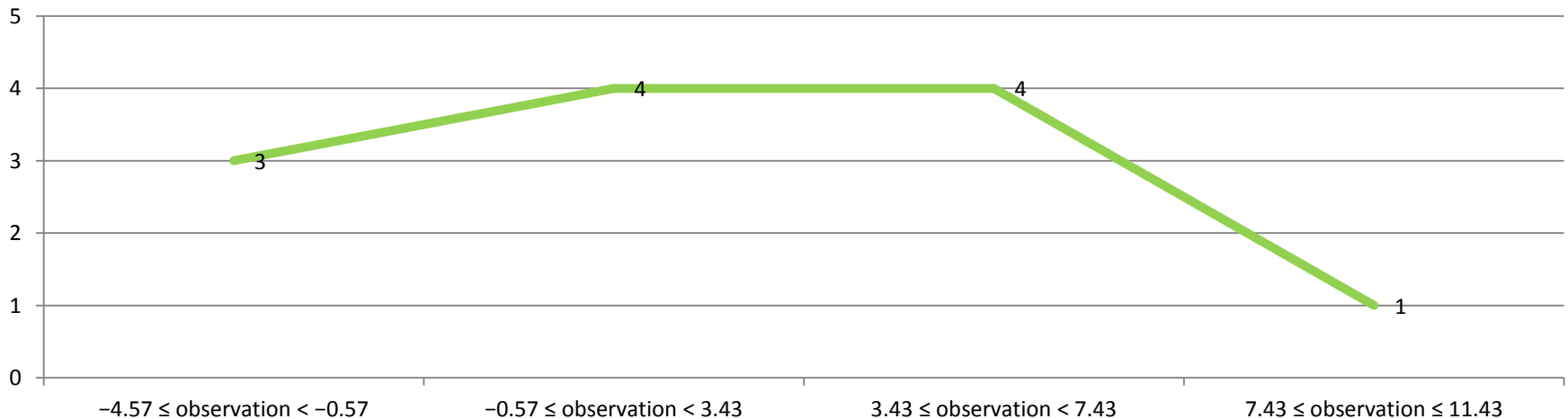| | | | |
|---|---|---|---|
| 3 | 4 | 4 | 1 |
| −4.57 ≤ observation < −0.57 | −0.57 ≤ observation < 3.43 | 3.43 ≤ observation < 7.43 | 7.43 ≤ observation ≤ 11.43 |

# Frequency Polygon

## Focus on: Holding Period Returns

- Frequency polygons are often used to provide higher visual continuity than histograms.

**Absolute Frequency**
**Holding Period Return**



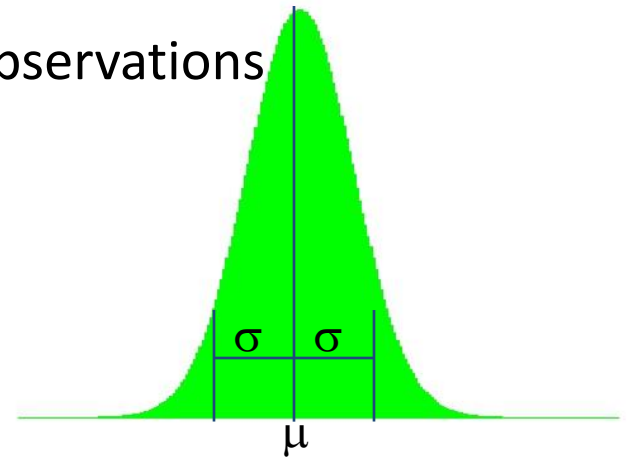| | | | |
|---|---|---|---|
| −4.57 ≤ observation < −0.57 | −0.57 ≤ observation < 3.43 | 3.43 ≤ observation < 7.43 | 7.43 ≤ observation ≤ 11.43 |

# Measures of central tendency

**These measures describe where the data are centered.**

- Arithmetic Mean
  - The **arithmetic mean** is the sum of the observations divided by the number of observations.

    - **Population mean** → $\mu = \dfrac{\sum_{i=1}^{N} X_i}{N}$

    - **Sample mean** → $\bar{X} = \dfrac{\sum_{i=1}^{N} X_i}{N}$

    - The sample mean is often interpreted as center of gravity, for a given set of data.
    - **Cross-sectional data** occur across different observation types at one point in time, and **time-series data** occur for the same unit of observation across time.

# Measures of central tendency
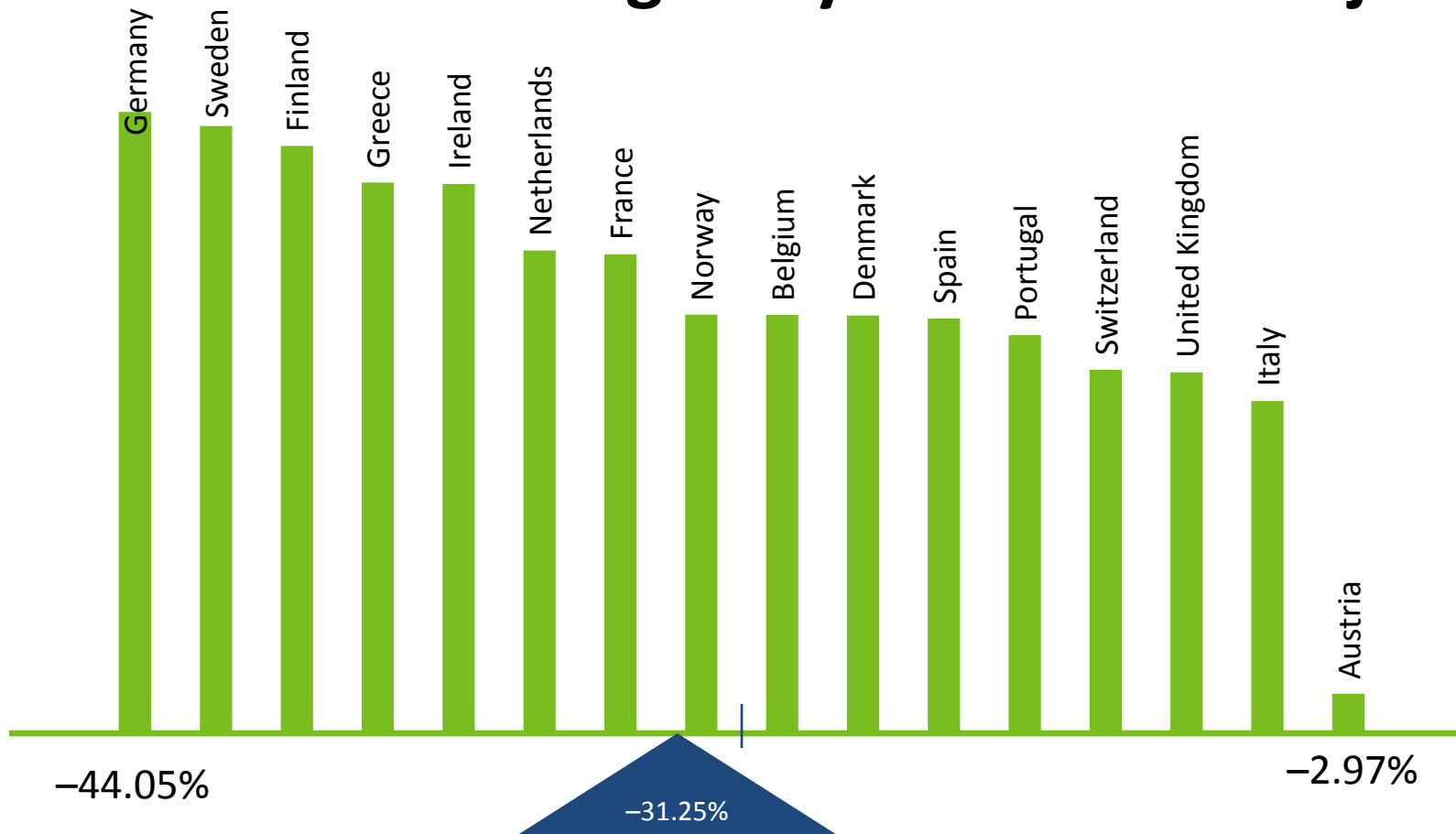
## Focus on: Cross-Sectional Sample Mean Return

| Country | Return | Country | Return |
|---------|--------|---------|--------|
| Austria | −2.97% | Italy | −23.64% |
| Belgium | −29.71% | Netherlands | −34.27% |
| Denmark | −29.67% | Norway | −29.73% |
| Finland | −41.65% | Portugal | −28.29% |
| France | −33.99% | Spain | −29.47% |
| Germany | −44.05% | Sweden | −43.07% |
| Greece | −39.06% | Switzerland | −25.84% |
| Ireland | −38.97% | United Kingdom | −25.66% |

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

$$\bar{X} = \frac{-500.04}{16} = -31.25\%$$

*Source*: www.msci.com.

# Measures of central tendency

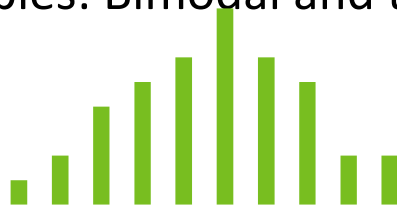**Mean as a center of gravity for the data object**

# Measures of central tendency

**These measures also describe where the data are centered.**

- Weighted Mean → $$\bar{X}_W = \sum_{i=1}^{n} w_i X_i$$

  - The sum of the observations times each observation's weight (proportional representation in the sample), where the weight is chosen to meet a statistical or financial goal. Example: Portfolio return

- Geometric Mean → $$G = \sqrt[n]{\prod_{i=1}^{n} X_i}$$

  - Represents the growth rate or compounded return on an investment when $X$ is $1 + R$

- Harmonic Mean → $$\bar{X}_H = {}^{n} \Big/ {\sum_{i=1}^{n} {}^{1}/{X_i}}$$

  - A weighted mean in which each observation's weight is inversely proportional to its magnitude. Example: Cost averaging
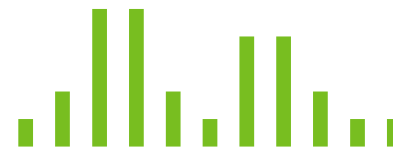
# Measures of central tendency

**These measures also describe where the data are centered.**

- The **median** is the middle observation by rank.
  - When we have an odd number of observations, the median will be the closest to the middle. When we have an even number, the median will be the average of the two middle values.
- The **mode** is the most frequently occurring value in a distribution.
  - Distributions are unimodal when there is a single most frequently occurring value and multimodal if there is more than one frequently occurring value.
  - Examples: Bimodal and trimodal

Unimodal

Bimodal

# Measures of central tendency

**Focus on: Calculating a Median or Mode**

$$\text{Median} = \frac{-29.73\% + (-29.71\%)}{2} = -29.72\%$$

| Rank | Country | Return | Rank | Country | Return |
|------|---------|--------|------|---------|--------|
| 1 | Germany | −44.05% | **9** | **Belgium** | **−29.71%** |
| 2 | Sweden | −43.07% | 10 | Denmark | −29.67% |
| 3 | Finland | −41.65% | 11 | Spain | −29.47% |
| 4 | Greece | −39.06% | 12 | Portugal | −28.29% |
| 5 | Ireland | −38.97% | 13 | Switzerland | −25.84% |
| 6 | Netherlands | −34.27% | 14 | United Kingdom | −25.66% |
| 7 | France | −33.99% | 15 | Italy | −23.64% |
| **8** | **Norway** | **−29.73%** | 16 | Austria | −2.97% |

# Median: The Case of the Price–Earnings Ratio

**ble 10.    P/Es for a Client Portfolio**

| Stock | Consensus Current EPS | Consensus Current P/E |
|---|---|---|
| Caterpillar, Inc. | 6.34 | 13.15 |
| Ford Motor Company | 1.55 | 10.97 |
| General Dynamics | 6.96 | 12.15 |
| Green Mountain Coffee Roasters | 3.25 | 25.27 |
| McDonald's Corporation | 5.61 | 17.16 |
| Qlik Technologies | 0.17 | 204.82 |
| Questcor Pharmaceuticals | 4.79 | 13.94 |

e: Consensus current P/E was calculated as price as of 9 September 2013 divided by consensus EPS as of the same date.

# Weighted average

**Also known as a weighted mean, the most common application of this measure in investments is the weighted mean return to a portfolio.**

- Consider again the country-level data. You have constructed a portfolio that has 50% of its weight in Portugal, Ireland, Greece, and Spain and 50% of its weight in Germany and the UK. Each of the first four countries is equally weighted within the 50%, as are Germany and the UK within their 50%. What is the weighted average return to the portfolio?

$$\bar{X}_W = \sum_{i=1}^{n} w_i X_i$$

| Country | Weight | Return | Component Return |
|---|---|---|---|
| Portugal | 12.50% | −28.29% | −3.54% |
| Ireland | 12.50% | −23.64% | −2.96% |
| Greece | 12.50% | −39.06% | −4.88% |
| Spain | 12.50% | −29.47% | −3.68% |
| Germany | 25.00% | −44.05% | −1.01% |
| UK | 25.00% | −25.66% | −6.42% |
| Sum | 100% | Weighted Mean = | −32.49% |

# Measures of dispersion

**Dispersion measures variability around a measure of central tendency. If mean return represents reward, then dispersion represents risk.**

- # Mean Absolute Deviation (MAD) $\rightarrow$ $\text{MAD} = \frac{\sum_{i=1}^{n} |X_i - \bar{X}|}{n}$
  - The arithmetic average of the absolute value of deviations from the mean.

# Measures of dispersion

**Dispersion measures variability around a measure of central tendency. If mean return represents reward, then dispersion represents risk.**

- Variance is the average squared deviation from the mean.

  – Population variance $\rightarrow \sigma^2 = \frac{\sum_{i=1}^{n}(X_i - \mu)^2}{n}$

  – Sample variance $\rightarrow s^2 = \frac{\sum_{i=1}^{n}(X_i - \bar{X})^2}{n-1}$

- Sample variance is "penalized" by dividing by $n-1$ instead of $n$ to account for the fact that the measure of central tendency used, $\bar{X}$, is an estimate of the true population parameter, $\mu$, and so has some uncertainty associated with it.

- Standard deviation is the square root of variance.

# Measures of dispersion

## Focus on: Sample Standard Deviation

| Country | Return | Squared Deviation from Mean |
|---|---|---|
| Germany | −44.05% | 0.016384 |
| Sweden | −43.07% | 0.013971 |
| Finland | −41.65% | 0.010816 |
| Greece | −39.06% | 0.00610 |
| … | … | … |
| Austria | −2.97% | 0.0780 |
| | Sum= | 0.1486 |
| | $s^2 =$ | 0.0099 |
| | $s =$ | 9.95% |

# Semivariance

**We are often concerned with measures of risk that focus on the "downside" of the possible outcomes—in other words, the losses.**

- Semivariance is the average squared deviation below the mean.
  - Semideviation is the square root of semivariance.
  - Both are a measure of dispersion focusing only on those observations below the mean.

$$\sum_{for\ all\ Xi<\bar{X}} \frac{(X_i - \bar{X})^2}{n^* - 1}$$

  - Target semivariance, by analogy, is the average squared deviation below some specified target rate, *B*, and represents the "downside" risk of being below the target, *B*.

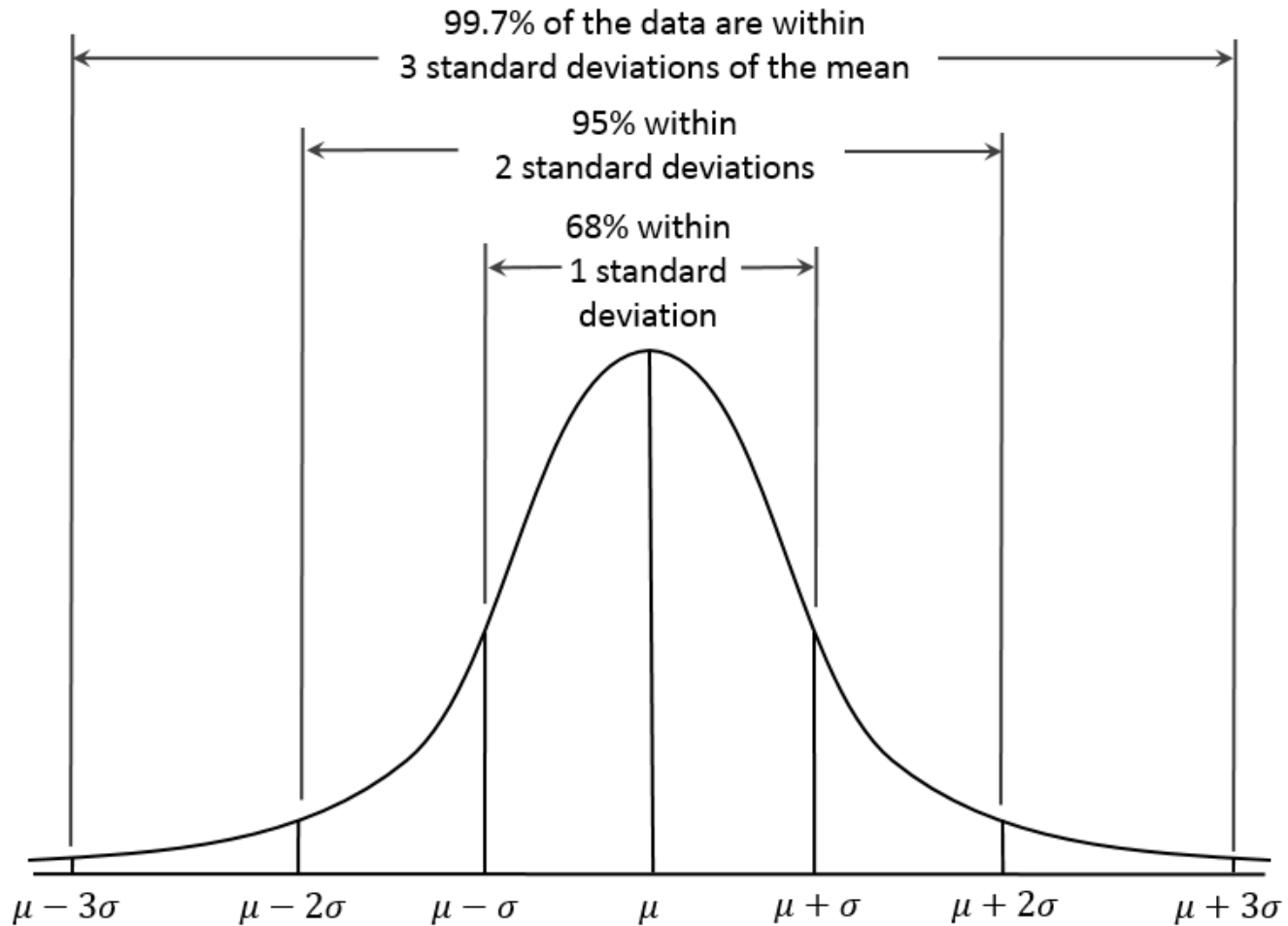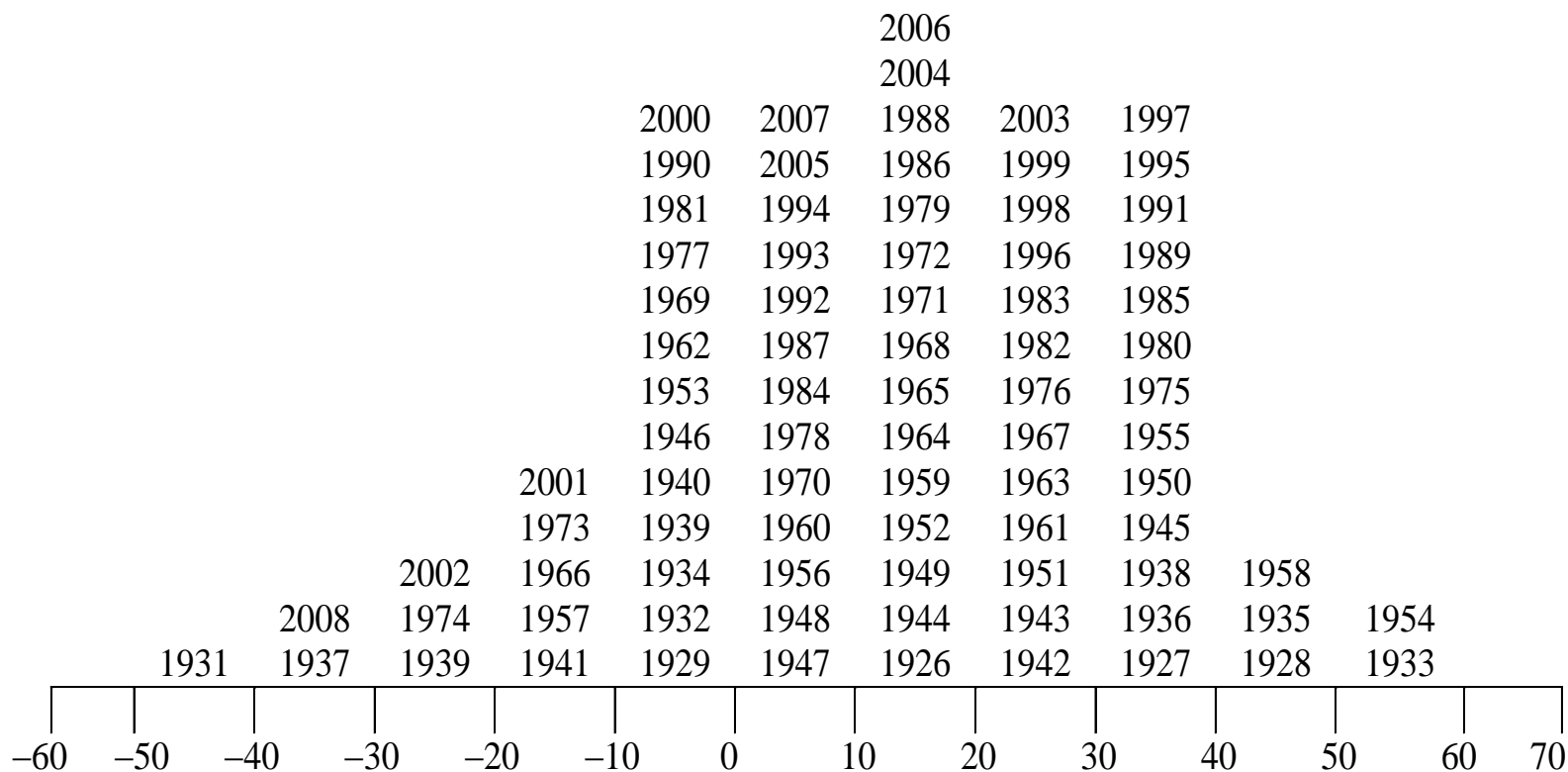$$\sum_{for\ all\ Xi<B} \frac{(X_i - B)^2}{n^* - 1}$$

# Normal Distribution Function

# EXHIBIT 5-9 Histogram of U.S. Large Company Stock Returns, 1926-2008

_____

```
                                                2006
                                                2004
                            2000   2007   1988   2003   1997
                            1990   2005   1986   1999   1995
                            1981   1994   1979   1998   1991
                            1977   1993   1972   1996   1989
                            1969   1992   1971   1983   1985
                            1962   1987   1968   1982   1980
                            1953   1984   1965   1976   1975
                            1946   1978   1964   1967   1955
                     2001   1940   1970   1959   1963   1950
                     1973   1939   1960   1952   1961   1945
              2002   1966   1934   1956   1949   1951   1938   1958
       2008   1974   1957   1932   1948   1944   1943   1936   1935   1954
  1931 1937   1939   1941   1929   1947   1926   1942   1927   1928   1933

 ├────┼────┼────┼────┼────┼────┼────┼────┼────┼────┼────┼────┼────┤
-60  -50  -40  -30  -20  -10   0    10   20   30   40   50   60   70
```

# Chebyshev's inequality

**This expression gives the minimum proportion of values, *p*, within *k* standard deviations of the mean for any distribution whenever *k* > 1.**

$$p \geq 1 - \frac{1}{k^2}$$

| *k* | Interval around the Mean | *p* |
|---|---|---|
| 1.25 | $\bar{X} \pm 1.25s$ | 0.36 |
| 1.50 | $\bar{X} \pm 1.50s$ | 0.56 |
| 2.00 | $\bar{X} \pm 2.00s$ | 0.75 |
| 2.50 | $\bar{X} \pm 2.50s$ | 0.84 |
| 3.00 | $\bar{X} \pm 3.00s$ | 0.89 |
| 4.00 | $\bar{X} \pm 4.00s$ | 0.94 |

# Chebyshev's inequality

**Focus on: Calculating Proportions Using Chebyshev's Inequality**

- For our country data, the mean is −31.25% and the sample standard deviation is 9.95%.

- Lower cutoff at 1.25 standard deviations:

$$-31.25\% - 1.25\,(9.95\%) = -43.6875\%$$

- Upper cutoff at 1.25 standard deviations:

$$-31.25\% + 1.25\,(9.95\%) = -18.8125\%$$

| $k$ | Lower Cutoff | Upper Cutoff | Actual $p$ | Chebyshev's $p$ |
|------|-------------|-------------|-----------|-----------------|
| 1.25 | −43.69% | −18.81% | 0.875 | 0.36 |
| 1.50 | −46.18% | −16.32% | 0.938 | 0.56 |
| 2.00 | −51.16% | −11.34% | 0.95 | 0.75 |
| 2.50 | −56.13% | −6.37% | 0.97 | 0.84 |
| 3.00 | −61.11% | −1.39% | 0.997 | 0.89 |
| 4.00 | −71.07% | 8.57% | 1.000 | 0.94 |

# Chebyshev's inequality

## Applying Chebyshev's Inequality

According to Table 22, the arithmetic mean monthly return and standard deviation of monthly returns on the S&P 500 were 0.94 percent and 5.50 percent, respectively, during the 1926–2012 period, totaling 1,044 monthly observations. Using this information, address the following:

1. Calculate the endpoints of the interval that must contain at least 75 percent of monthly returns according to Chebyshev's inequality.

# Combining risk and return

**Measures of relative dispersion are used to compare risk and return across differing sets of observations.**

- The **coefficient of variation** is the ratio of the standard deviation of a set of observations to their mean value.
  – This ratio can be thought of as the units of risk per unit of mean return.

- The **Sharpe Ratio** is the ratio of the mean excess return (mean return minus the mean risk-free rate) per unit of standard deviation.
  – This ratio can be thought of as units of risky return (excess return) per unit of risk.
  – This will also be the slope of a line in expected return/standard deviation space.

$E(r)$

$S_p$

$r_f$

$\sigma$

# Combining risk and return

**Focus on: Coefficient of Variation and the Sharpe Ratio**

- Consider a portfolio with a mean return of 25.26% and a standard deviation of returns of 9.95%.

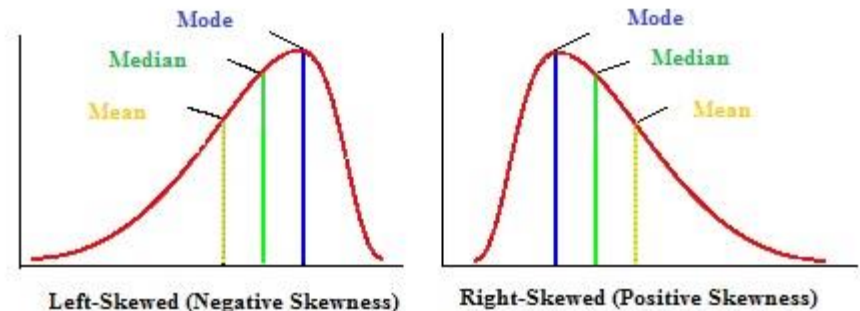$$CV = \frac{s}{\bar{\bar{X}}} = \frac{9.95\%}{25.26\%} = 0.3939$$

  – The coefficient of variation is

$$S_p = \frac{\bar{R}_p - \bar{R}_f}{s_p} = \frac{25.26\% - 3\%}{9.95\%} = 2.2372$$

  – If the risk-free rate is 3%, then the Sharpe Ratio is
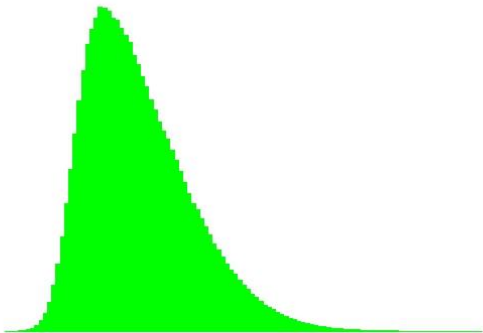
# Combining centrality, dispersion, and symmetry

- For a symmetrical distribution, the mean, median, and mode (if it exists) will all be at the same location.

- If the distribution is positively skewed, then the mean will be greater than the median, which will be greater than the mode (if it exists).

- If the distribution is negatively skewed, then the mean will be less than the median, which will be less than the mode (if it exists).
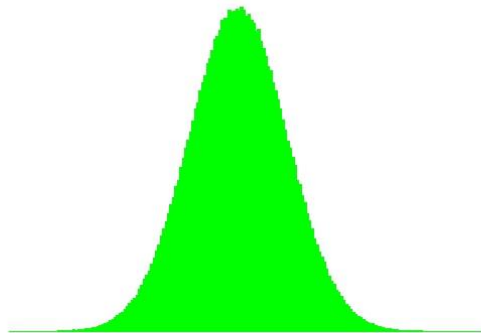


Left-Skewed (Negative Skewness)    Right-Skewed (Positive Skewness)

31

# Skewness

**The degree of symmetry in the dispersion of values around the mean is known as skewness.**
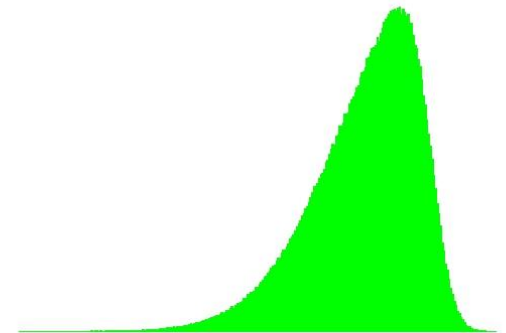
- If observations are equally dispersed around the mean, the distribution is said to be symmetrical.

- If the distribution has a long tail on one side and a "fatter" distribution on the other side, it is said to be skewed in the direction of the long tail.
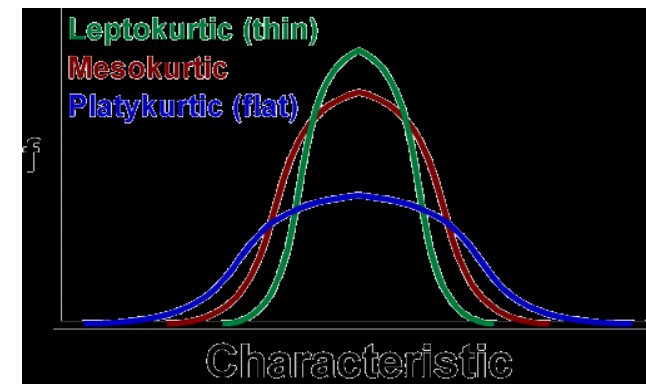
| Skew Right | No Skew | Skew Left |

# Kurtosis

- Kurtosis measures the relative amount of "peakedness" as compared with the normal distribution, which has a kurtosis of 3.
    - We typically express this measure in terms of excess kurtosis being the observed kurtosis minus 3.
    - Distributions are referred to as being
  1. Leptokurtic (more peaked than the normal; fatter tails)
  2. Platykurtic (less peaked than the normal; thinner tails) or
  3. Mesokurtic (equivalent to the normal).

# Summary

- The underlying foundation of statistically based quantitative analysis lies with the concepts of a sample versus a population.
  - We use sample statistics to describe the sample and to infer information about its associated population.
  - Descriptive statistics for samples and populations include measures of centrality, location, and dispersion, such as mean, range, and variance, respectively.
  - We can combine traditional measures of return (such as mean) and risk (such as standard deviation) to measure the combined effects of risk and return using the coefficient of variation and the Sharpe Ratio.
- The normal distribution is of central importance in investments, and as a result, we often compare statistical properties, such as skewness and kurtosis, with those of the normal distribution.