

# Summarizing Measured Data

Raj Jain

Washington University in Saint Louis

Saint Louis, MO 63130

Jain@cse.wustl.edu

These slides are available on-line at:

<http://www.cse.wustl.edu/~jain/cse567-11/>



- ❑ Basic Probability and Statistics Concepts: CDF, PDF, PMF, Mean, Variance, CoV, Normal Distribution
- ❑ Summarizing Data by a Single Number: Mean, Median, and Mode, Arithmetic, Geometric, Harmonic Means
- ❑ Mean of A Ratio
- ❑ Summarizing Variability: Range, Variance, percentiles, Quartiles
- ❑ Determining Distribution of Data: Quantile-Quantile plots

## Part III: Probability Theory and Statistics

1. How to report the performance as a single number? Is specifying the mean the correct way?
2. How to report the variability of measured quantities? What are the alternatives to variance and when are they appropriate?
3. How to interpret the variability? How much confidence can you put on data with a large variability?
4. How many measurements are required to get a desired level of statistical confidence?
5. How to summarize the results of several different workloads on a single computer system?
6. How to compare two or more computer systems using several different workloads? Is comparing the mean sufficient?
7. What model best describes the relationship between two variables? Also, how good is the model?

# Basic Probability and Statistics Concepts

- **Independent Events:** Two events are called independent if the occurrence of one event does not in any way affect the probability of the other event.
  - Examples: Successive throws of a coin
  - Price of a stock on successive days
- **Random Variable:** A variable is called a random variable if it takes one of a specified set of values with a specified probability.
  - Examples: Weather, temperature, execution time

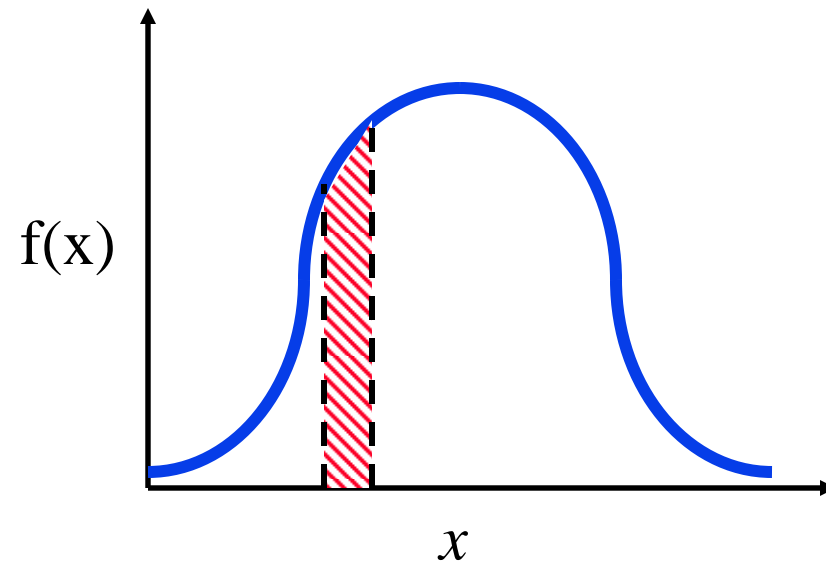
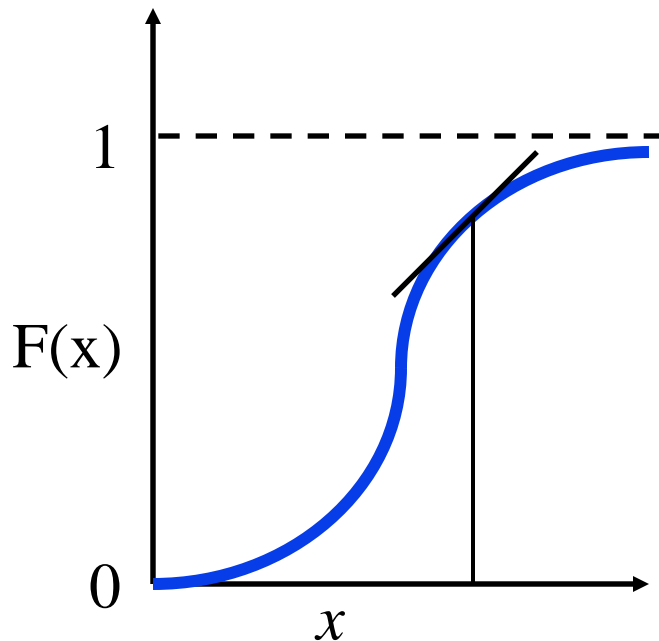
# CDF, PDF, and PMF

- **Cumulative Distribution Function:** Starts at 0. Ends at 1

$$F_x(a) = P(x \leq a)$$

- **Probability Density Function:** Starts at 0 and ends at 0

$$f(x) = \frac{dF(x)}{dx}$$



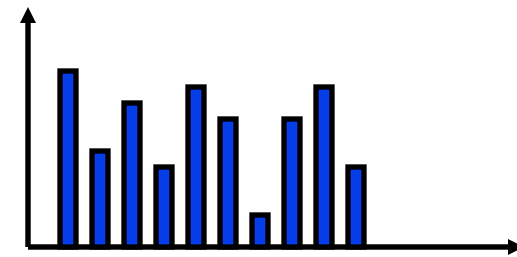
## CDF, PDF, and PMF (Cont)

- Given a pdf  $f(x)$ :

$$P(x_1 < x \leq x_2) = F(x_2) - F(x_1) = \int_{x_1}^{x_2} f(x) dx$$

- Probability Mass Function:** For discrete random variables:

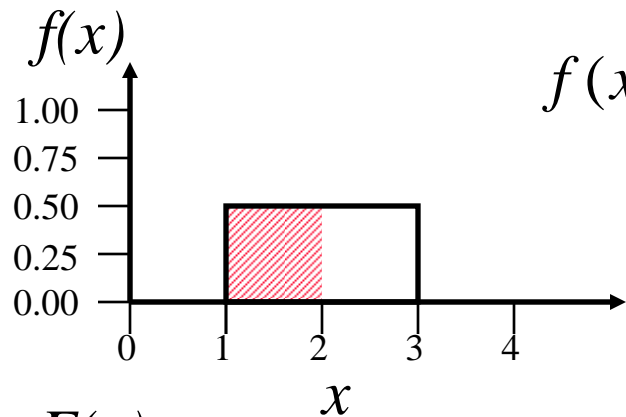
$$f(x_i) = p_i$$



$$\begin{aligned} P(x_1 < x \leq x_2) &= F(x_2) - F(x_1) \\ &= \sum_{x_1 < x_i \leq x_2} p_i \end{aligned}$$

# CDF, PDF, and PMF (Cont)

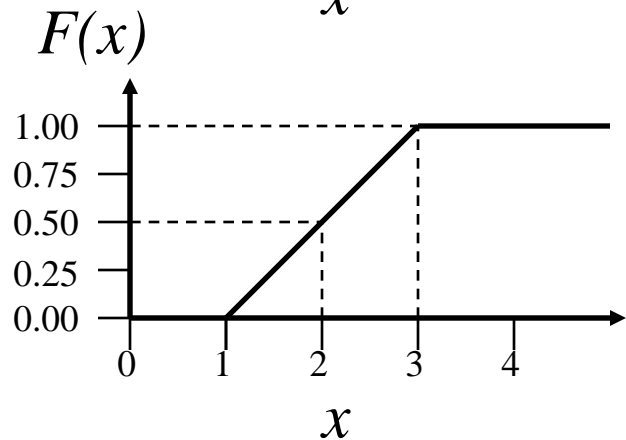
- Example: Response time – uniformly distributed between 1 and 3 seconds.



$$f(x) = \frac{1}{2} \quad 1 \leq x \leq 3$$

$$P(1 \leq x \leq 3) = 0.5$$

$$P(x \leq 3) = 0.5$$



$$F(x) = \int_{-\infty}^x f(x) dx$$

# Mean, Variance, CoV

## □ Mean or Expected Value:

$$\text{Mean } \mu = E(x) = \sum_{i=1}^n p_i x_i = \int_{-\infty}^{+\infty} x f(x) dx$$

## □ Variance: The expected value of the square of distance between x and its mean

$$\begin{aligned} \text{Var}(x) &= \sigma^2 = E[(x - \mu)^2] = \sum_{i=1}^n p_i (x_i - \mu)^2 \\ &= \int_{-\infty}^{+\infty} (x_i - \mu)^2 f(x) dx \end{aligned}$$

## □ Coefficient of Variation:

$$\text{C.O.V.} = \frac{\text{Standard Deviation}}{\text{Mean}} = \frac{\sigma}{\mu}$$



# Mean, Variance, CoV: Examples

## □ Mean or Expected Value:

$$\begin{aligned}\text{Mean } \mu &= E(x) = \int_{-\infty}^{+\infty} x f(x) dx \\ &= \int_1^3 x \frac{1}{2} dx = \frac{x^2}{4} = \frac{9}{4} - \frac{1}{4} = 2\end{aligned}$$

## □ Variance:

$$\begin{aligned}\text{Var}(x) &= \int_{-\infty}^{+\infty} (x_i - \mu)^2 f(x) dx \\ &= \int_1^3 (x - 2)^2 \frac{1}{2} dx = \frac{(x-2)^3}{6} \\ &= \frac{1}{6} + \frac{1}{6} = \frac{1}{3}\end{aligned}$$

## □ Standard Deviation: $\sigma = \frac{1}{\sqrt{3}}$

## □ Coefficient of Variation:

$$\text{C.O.V.} = \frac{\text{Standard Deviation}}{\text{Mean}} = \frac{\sigma}{\mu} = \frac{1}{2\sqrt{3}}$$

# Covariance and Correlation

## □ Covariance:

$$\begin{aligned} Cov(x, y) &= \sigma_{xy} = E[(x - \mu_x)(y - \mu_y)] \\ &= E(xy) - E(x)E(y) \end{aligned}$$

- For independent variables, the covariance is zero:

$$E(xy) = E(x)E(y)$$

- Although independence always implies zero covariance, the reverse is not true.

- **Correlation Coefficient:** normalized value of covariance

$$\text{Correlation}(x, y) = \rho_{xy} = \frac{\sigma_{xy}^2}{\sigma_x \sigma_y}$$

The correlation always lies between -1 and +1.

# Mean and Variance of Sums

□ If  $x_1, x_2, \dots, x_k$

are  $k$  random variables and if  $a_1, a_2, \dots, a_k$  are  $k$  arbitrary constants (called weights), then:

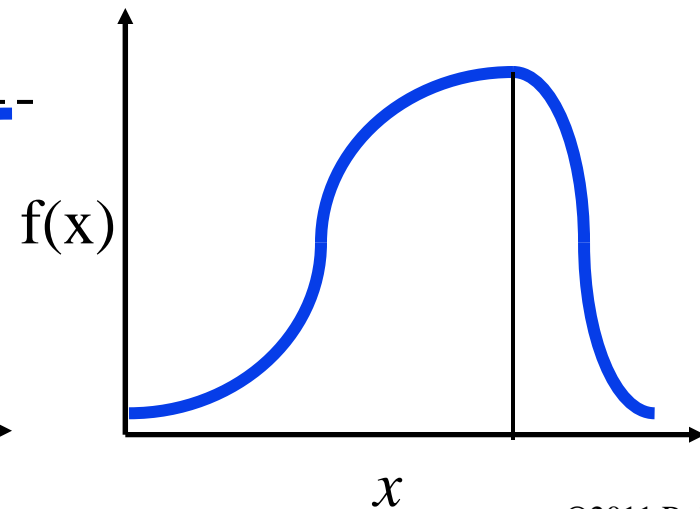
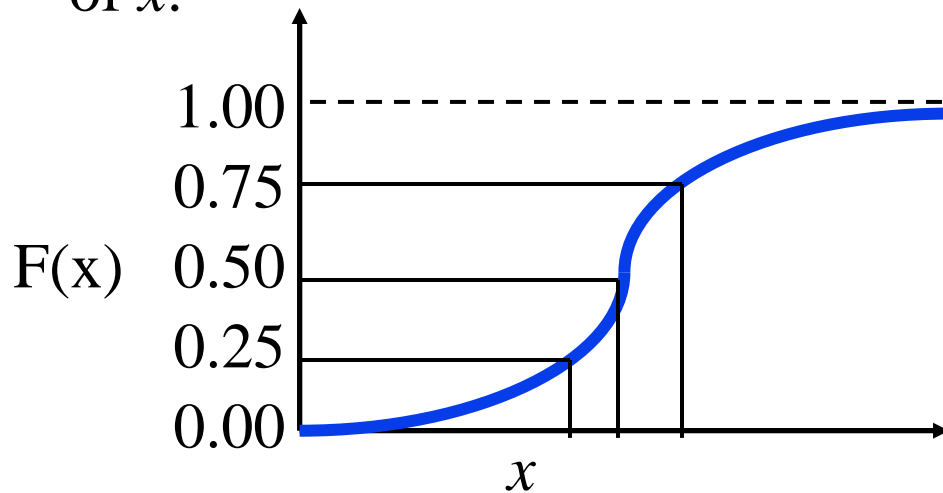
$$\begin{aligned} E(a_1x_1 + a_2x_2 + \dots + a_kx_k) \\ = a_1E(x_1) + a_2E(x_2) + \dots + a_kE(x_k) \end{aligned}$$

□ For independent variables:

$$\begin{aligned} Var(a_1x_1 + a_2x_2 + \dots + a_kx_k) \\ = a_1^2Var(x_1) + a_2^2Var(x_2) + \dots + a_k^2Var(x_k) \end{aligned}$$

# Quantiles, Median, and Mode

- **Quantile:** The  $x$  value at which the CDF takes a value  $\alpha$  is called the  $\alpha$ -quantile or  $100\alpha$ -percentile. It is denoted by  $x_\alpha$ :  
$$P(x \leq x_\alpha) = F(x_\alpha) = \alpha$$
- **Median:** The 50-percentile or (0.5-quantile) of a random variable is called its median.
- **Mode:** The most likely value, that is,  $x_i$  that has the highest probability  $p_i$ , or the  $x$  at which pdf is maximum, is called mode of  $x$ .



# Normal Distribution

- **Normal Distribution:** The sum of a large number of independent observations from any distribution has a normal distribution.

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad -\infty \leq x \leq +\infty$$

- A normal variate is denoted as  $N(\mu, \sigma)$ .
- **Unit Normal:** A normal distribution with zero mean and unit variance. Also called **standard normal distribution** and is denoted as  $N(0, 1)$ .

# Normal Quantiles

- An  $\alpha$ -quantile of a unit normal variate  $z \sim N(0,1)$  is denoted by  $z_\alpha$ . If a random variable  $x$  has a  $N(\mu, \sigma)$  distribution, then  $(x-\mu)/\sigma$  has a  $N(0,1)$  distribution.

$$P\left(\frac{x-\mu}{\sigma} \leq z_\alpha\right) = \alpha$$

or

$$P(x \leq \mu + z_\alpha \sigma) = \alpha$$

# Why Normal?

- There are two main reasons for the popularity of the normal distribution:
  1. *The sum of  $n$  independent normal variates is a normal variate.* If,  $(x_i \sim N(\mu_i, \sigma_i))$   
then  $x = \sum_{i=1}^n a_i x_i$  has a normal distribution with  
mean  $\mu = \sum_{i=1}^n a_i \mu_i$  and variance  $\sigma^2 = \sum_{i=1}^n a_i^2 \sigma_i^2$ .
  2. *The sum of a large number of independent observations from any distribution tends to have a normal distribution.*  
This result, which is called **central limit theorem**, is true for observations from all distributions  
 $\Rightarrow$  Experimental errors caused by many factors are normal.

## Homework 12A: Exercise 12.7

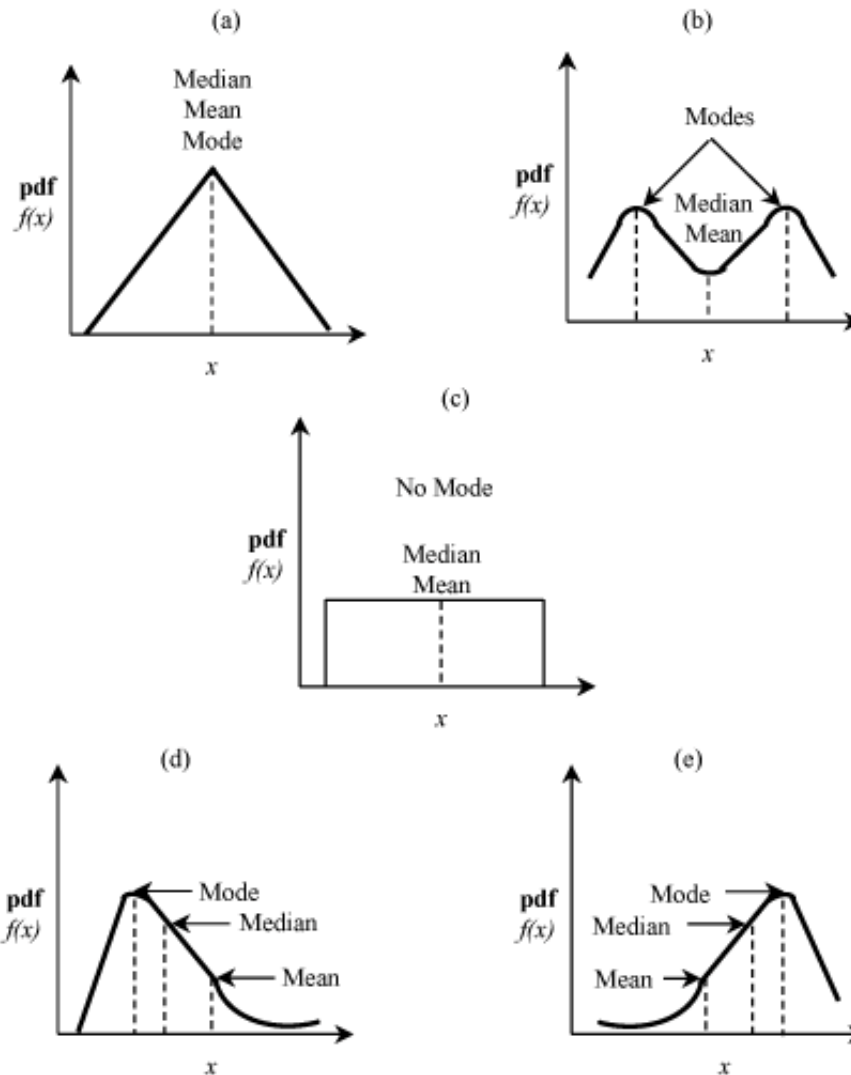
- The execution times of queries on a database is normally distributed with a mean of 5 seconds and a standard deviation of 1 second. Determine the following:
  - a. What is the probability of the execution time being more than 8 seconds.
  - b. What is the probability of the execution time being less than 6 seconds.
  - c. What percent of responses will take between 4 and 7 seconds?
  - d. What is the 95-percentile execution time?



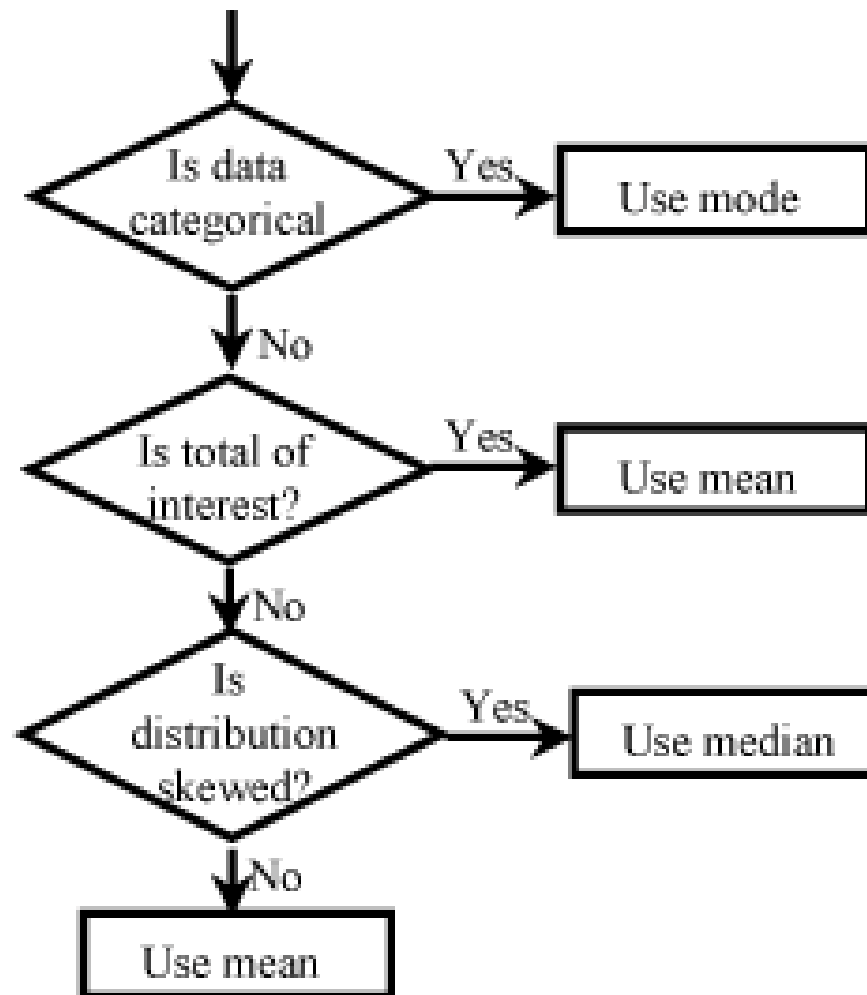
# Summarizing Data by a Single Number

- ❑ **Indices of central tendencies:** Mean, Median, Mode
- ❑ **Sample Mean** is obtained by taking the sum of all observations and dividing this sum by the number of observations in the sample.
- ❑ **Sample Median** is obtained by sorting the observations in an increasing order and taking the observation that is in the middle of the series. If the number of observations is even, the mean of the middle two values is used as a median.
- ❑ **Sample Mode** is obtained by plotting a histogram and specifying the midpoint of the bucket where the histogram peaks. For categorical variables, mode is given by the category that occurs most frequently.
- ❑ Mean and median always exist and are unique. Mode, on the other hand, may not exist.

# Mean, Median, and Mode: Relationships



# Selecting Mean, Median, and Mode



# Indices of Central Tendencies: Examples

- ❑ Most used resource in a system: Resources are categorical and hence mode must be used.
- ❑ Interarrival time: Total time is of interest and so mean is the proper choice.
- ❑ Load on a Computer: Median is preferable due to a highly skewed distribution.
- ❑ Average Configuration: Medians of number devices, memory sizes, number of processors are generally used to specify the configuration due to the skewness of the distribution.

# Common Misuses of Means

- ❑ *Using mean of significantly different values:*  
 $(10+1000)/2 = 505$
- ❑ *Using mean without regard to the skewness of distribution.*

System A	System B
10	5
9	5
11	5
10	4
10	31
Sum=50	Sum=50
Mean=10	Mean=10
Typical=10	Typical=5

## Misuses of Means (cont)

- ❑ *Multiplying means to get the mean of a product*

$$E(xy) \neq E(x)E(y)$$

- ❑ Example: On a timesharing system,

Average number of users is 23

Average number of sub-processes per user is 2

What is the average number of sub-processes?

Is it 46? No!

The number of sub-processes a user spawns depends upon how much load there is on the system.

- ❑ *Taking a mean of a ratio with different bases.*

Already discussed in Chapter 11 on ratio games and is discussed further later

# Geometric Mean

$$\hat{x} = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

- Geometric mean is used if the product of the observations is a quantity of interest.

# Geometric Mean: Example

- The performance improvements in 7 layers:

Protocol Layer	Performance Improvement
7	18%
6	13%
5	11%
4	8%
3	10%
2	28%
1	5%

Average improvement per layer

$$= \{(1.18)(1.13)(1.11)(1.08)(1.10)(1.28)(1.05)\}^{\frac{1}{7}} - 1$$
$$= 0.13$$



# Examples of Multiplicative Metrics

- ❑ Cache hit ratios over several levels of caches
- ❑ Cache miss ratios
- ❑ Percentage performance improvement between successive versions
- ❑ Average error rate per hop on a multi-hop path in a network.

# Geometric Mean of Ratios

$$\begin{aligned} gm\left(\frac{x_1}{y_1}, \frac{x_2}{y_2}, \dots, \frac{x_n}{y_n}\right) &= \frac{gm(x_1, x_2, \dots, x_n)}{gm(y_1, y_2, \dots, y_n)} \\ &= \frac{1}{gm\left(\frac{y_1}{x_1}, \frac{y_2}{x_2}, \dots, \frac{y_n}{x_n}\right)} \end{aligned}$$

- The geometric mean of a ratio is the ratio of the geometric means of the numerator and denominator  
=> the choice of the base does not change the conclusion.
- It is because of this property that sometimes geometric mean is recommended for ratios.
- However, if the geometric mean of the numerator or denominator do not have any physical meaning, the geometric mean of their ratio is meaningless as well.

# Harmonic Mean

$$\ddot{x} = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

- Used whenever an arithmetic mean can be justified for  $1/x_i$   
E.g., Elapsed time of a benchmark on a processor
- In the  $i^{\text{th}}$  repetition, the benchmark takes  $t_i$  seconds. Now suppose the benchmark has  $m$  million instructions, MIPS  $x_i$  computed from the  $i^{\text{th}}$  repetition is:  $x_i = \frac{m}{t_i}$
- $t_i$ 's should be summarized using arithmetic mean since the sum of  $t_i$  has a physical meaning  
 $\Rightarrow x_i$ 's should be summarized using harmonic mean since the sum of  $1/x_i$ 's has a physical meaning.

## Harmonic Mean (Cont)

- The average MIPS rate for the processor is:

$$\begin{aligned}\ddot{x} &= \frac{n}{\frac{1}{m/t_1} + \frac{1}{m/t_2} + \dots + \frac{1}{m/t_n}} \\ &= \frac{m}{\frac{1}{n}(t_1 + t_2 + \dots + t_n)}\end{aligned}$$

- However, if  $x_i$ 's represent the MIPS rate for  $n$  different benchmarks so that  $i^{\text{th}}$  benchmark has  $m_i$  million instructions, then harmonic mean of  $n$  ratios  $m_i/t_i$  cannot be used since the sum of the  $t_i/m_i$  does not have any physical meaning.
- Instead, as shown later, the quantity  $\sum m_i/\sum t_i$  is a preferred average MIPS rate.

# Weighted Harmonic Mean

- The weighted harmonic mean is defined as follows:

$$\bar{x} = \frac{1}{\frac{w_1}{x_1} + \frac{w_2}{x_2} + \dots + \frac{w_n}{x_n}}$$

where,  $w_i$ 's are weights which add up to one:

$$w_1 + w_2 + \dots + w_n = 1$$

- All weights equal  $\Rightarrow$  Harmonic, I.e.,  $w_i = 1/n$ .
- In case of MIPS rate, if the weights are proportional to the size of the benchmark:

$$w_i = \frac{m_i}{m_1 + m_2 + \dots + m_n}$$

- Weighted harmonic mean would be:

$$\bar{x} = \frac{m_1 + m_2 + \dots + m_n}{t_1 + t_2 + \dots + t_n}$$

# Mean of A Ratio

- 1. If the sum of numerators and the sum of denominators, both have a physical meaning, the average of the ratio is the ratio of the averages.*

For example, if  $x_i = a_i/b_i$ , the average ratio is given by:

$$\begin{aligned} \text{Average}\left(\frac{a_1}{b_1}, \frac{a_2}{b_2}, \dots, \frac{a_n}{b_n}\right) &= \frac{a_1 + a_2 + \dots + a_n}{b_1 + b_2 + \dots + b_n} \\ &= \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n a_i}{\frac{1}{n} \sum_{i=1}^n b_i} = \frac{\bar{a}}{\bar{b}} \end{aligned}$$

# Mean of a Ratio: Example

## □ CPU utilization

Measurement Duration	CPU Busy
1	45%
1	45%
1	45%
1	45%
100	20%
Sum	200
Mean	$\neq 200/5$ or 40%

## Example (Cont)

$$\begin{aligned} \text{Mean CPU utilization} &= \frac{\text{Sum of CPU busy times}}{\text{Sum of measurement durations}} \\ &= \frac{0.45 + 0.45 + 0.45 + 0.45 + 20}{1 + 1 + 1 + 1 + 100} \\ &= 21\% \end{aligned}$$

- ❑ Ratios cannot always be summarized by a geometric mean.
- ❑ A geometric mean of utilizations is useless.



## Mean of a Ratio: Special Cases

- a. *If the denominator is a constant and the sum of numerator has a physical meaning, the arithmetic mean of the ratios can be used.*

That is, if  $b_i = b$  for all  $i$ 's, then:

$$\begin{aligned} & \text{Average}\left(\frac{a_1}{b}, \frac{a_2}{b}, \dots, \frac{a_n}{b}\right) \\ &= \frac{1}{n} \left( \frac{a_1}{b} + \frac{a_2}{b} + \dots + \frac{a_n}{b} \right) \\ &= \frac{\sum_{i=1}^n a_i}{nb} \end{aligned}$$

- Example: mean resource utilization.

## Mean of Ratio (Cont)

- b. If the sum of the denominators has a physical meaning and the numerators are constant then a harmonic mean of the ratio should be used to summarize them.*

That is, if  $a_i = a$  for all  $i$ 's, then:

$$\begin{aligned} \text{Average} \left( \frac{a}{b_1}, \frac{a}{b_2}, \dots, \frac{a}{b_n} \right) &= \frac{n}{\frac{b_1}{a} + \frac{b_2}{a} + \dots + \frac{b_n}{a}} \\ &= \frac{na}{\sum_{i=1}^n b_i} \end{aligned}$$

Example: MIPS using the same benchmark

## Mean of Ratios (Cont)

2. *If the numerator and the denominator are expected to follow a multiplicative property such that  $a_i = c b_i$ , where  $c$  is approximately a constant that is being estimated, then  $c$  can be estimated by the geometric mean of  $a_i/b_i$ .*

□ Example: Program Optimizer:  $a_i = c b_i$

□ Where,  $b_i$  and  $a_i$  are the sizes before and after the program optimization and  $c$  is the effect of the optimization which is expected to be independent of the code size.

$$\log c = (\log b_i - \log a_i)$$

or  $\log a_i = \log c + \log b_i$

□  $\log c =$  arithmetic mean of  $(\log b_i - \log a_i)$   
 $\Rightarrow c$  geometric mean of  $b_i/a_i$

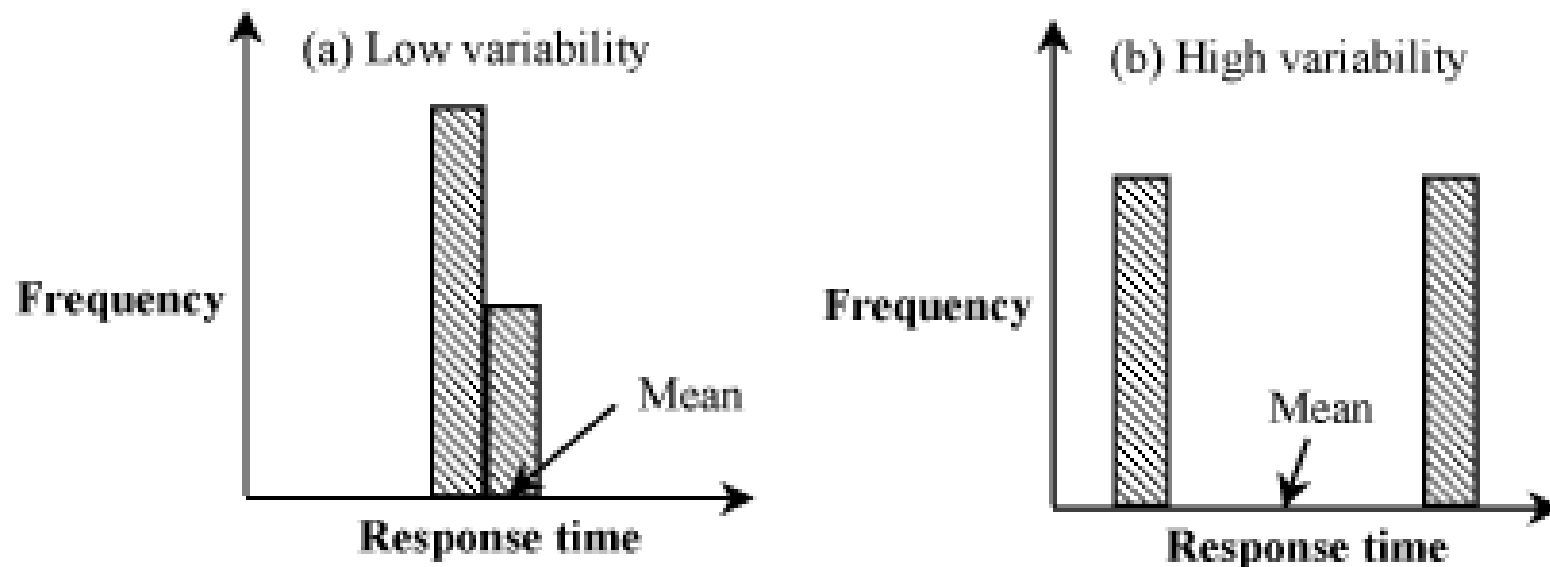
# Program Optimizer Static Size Data

Program	Code Size		Ratio
	Before	After	
BubbleP	119	89	0.75
IntmmP	158	134	0.85
PermP	142	121	0.85
PuzzleP	8612	7579	0.88
QueenP	7133	7062	0.99
QuickP	184	112	0.61
SieveP	2908	2879	0.99
TowersP	433	307	0.71
Geometric Mean			0.79

# Summarizing Variability

- “Then there is the man who drowned crossing a stream with an average depth of six inches.”

- W. I. E. Gates



# Indices of Dispersion

1. Range: Minimum and maximum of the values observed
2. Variance or standard deviation
3. 10- and 90- percentiles
4. Semi inter-quantile range
5. Mean absolute deviation

# Range

- ❑ Range = Max-Min
- ❑ Larger range => higher variability
- ❑ In most cases, range is not very useful.
- ❑ The minimum often comes out to be zero and the maximum comes out to be an ``outlier" far from typical values.
- ❑ Unless the variable is bounded, the maximum goes on increasing with the number of observations, the minimum goes on decreasing with the number of observations, and there is no ``stable" point that gives a good indication of the actual range.
- ❑ Range is useful if, and only if, there is a reason to believe that the variable is bounded.

# Variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- ❑ The divisor for  $s^2$  is  $n-1$  and not  $n$ .
- ❑ This is because only  $n-1$  of the  $n$  differences  $(x_i - \bar{x})$  are independent.
- ❑ Given  $n-1$  differences,  $n^{\text{th}}$  difference can be computed since the sum of all  $n$  differences must be zero.
- ❑ The number of independent terms in a sum is also called its **degrees of freedom**.



## Variance (Cont)

- Variance is expressed in units which are square of the units of the observations.  
=> It is preferable to use standard deviation.
- Ratio of standard deviation to the mean, or the **coefficient of variation (COV)**, is even better because it takes the scale of measurement (unit of measurement) out of variability consideration.

# Percentiles

- ❑ Specifying the 5-percentile and the 95-percentile of a variable has the same impact as specifying its minimum and maximum.
- ❑ It can be done for any variable, even for variables without bounds.
- ❑ When expressed as a fraction between 0 and 1 (instead of a percent), the percentiles are also called **quantiles**.  
=> 0.9-quantile is the same as 90-percentile.
- ❑ **Fractile**= quantile.
- ❑ The percentiles at multiples of 10% are called **deciles**. Thus, the first decile is 10-percentile, the second decile is 20-percentile, and so on.

# Quartiles

- ❑ **Quartiles** divide the data into four parts at 25%, 50%, and 75%.  
=> 25% of the observations are less than or equal to the first quartile  $Q_1$ ,  
50% of the observations are less than or equal to the second quartile  $Q_2$ , and  
75% are less than the third quartile  $Q_3$ .
- ❑ Notice that the second quartile  $Q_2$  is also the median.
- ❑ The  $\alpha$ -quantiles can be estimated by sorting the observations and taking the  $[(n-1)\alpha+1]$ th element in the ordered set. Here,  $[.]$  is used to denote rounding to the nearest integer.
- ❑ For quantities exactly half way between two integers use the lower integer.

# Semi Inter-Quartile Range

- Inter-quartile range =  $Q_3 - Q_1$
- Semi inter-quartile range (SIQR)

$$\text{SIQR} = \frac{Q_3 - Q_1}{2} = \frac{x_{0.75} - x_{0.25}}{2}$$

# Mean Absolute Deviation

$$\text{Mean absolute deviation} = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

- No multiplication or square root is required

# Comparison of Variation Measures

- ❑ Range is affected considerably by outliers.
- ❑ Sample variance is also affected by outliers but the affect is less
- ❑ Mean absolute deviation is next in resistance to outliers.
- ❑ Semi inter-quantile range is very resistant to outliers.
- ❑ If the distribution is highly skewed, outliers are highly likely and SIQR is preferred over standard deviation
- ❑ In general, SIQR is used as an index of dispersion whenever median is used as an index of central tendency.
- ❑ For qualitative (categorical) data, the dispersion can be specified by giving the number of most frequent categories that comprise the given percentile, for instance, top 90%.

# Measures of Variation: Example

In an experiment, which was repeated 32 times, the measured CPU time was found to be {3.1, 4.2, 2.8, 5.1, 2.8, 4.4, 5.6, 3.9, 3.9, 2.7, 4.1, 3.6, 3.1, 4.5, 3.8, 2.9, 3.4, 3.3, 2.8, 4.5, 4.9, 5.3, 1.9, 3.7, 3.2, 4.1, 5.1, 3.2, 3.9, 4.8, 5.9, 4.2}.

□ The sorted set is {1.9, 2.7, 2.8, 2.8, 2.8, 2.9, 3.1, 3.1, 3.2, 3.2, 3.3, 3.4, 3.6, 3.7, 3.8, 3.9, 3.9, 3.9, 4.1, 4.1, 4.2, 4.2, 4.4, 4.5, 4.5, 4.8, 4.9, 5.1, 5.1, 5.3, 5.6, 5.9}.

□ 10-percentile =  $[1+(31)(0.10)] = 4\text{th element} = 2.8$

□ 90-percentile =  $[1+(31)(0.90)] = 29\text{th element} = 5.1$

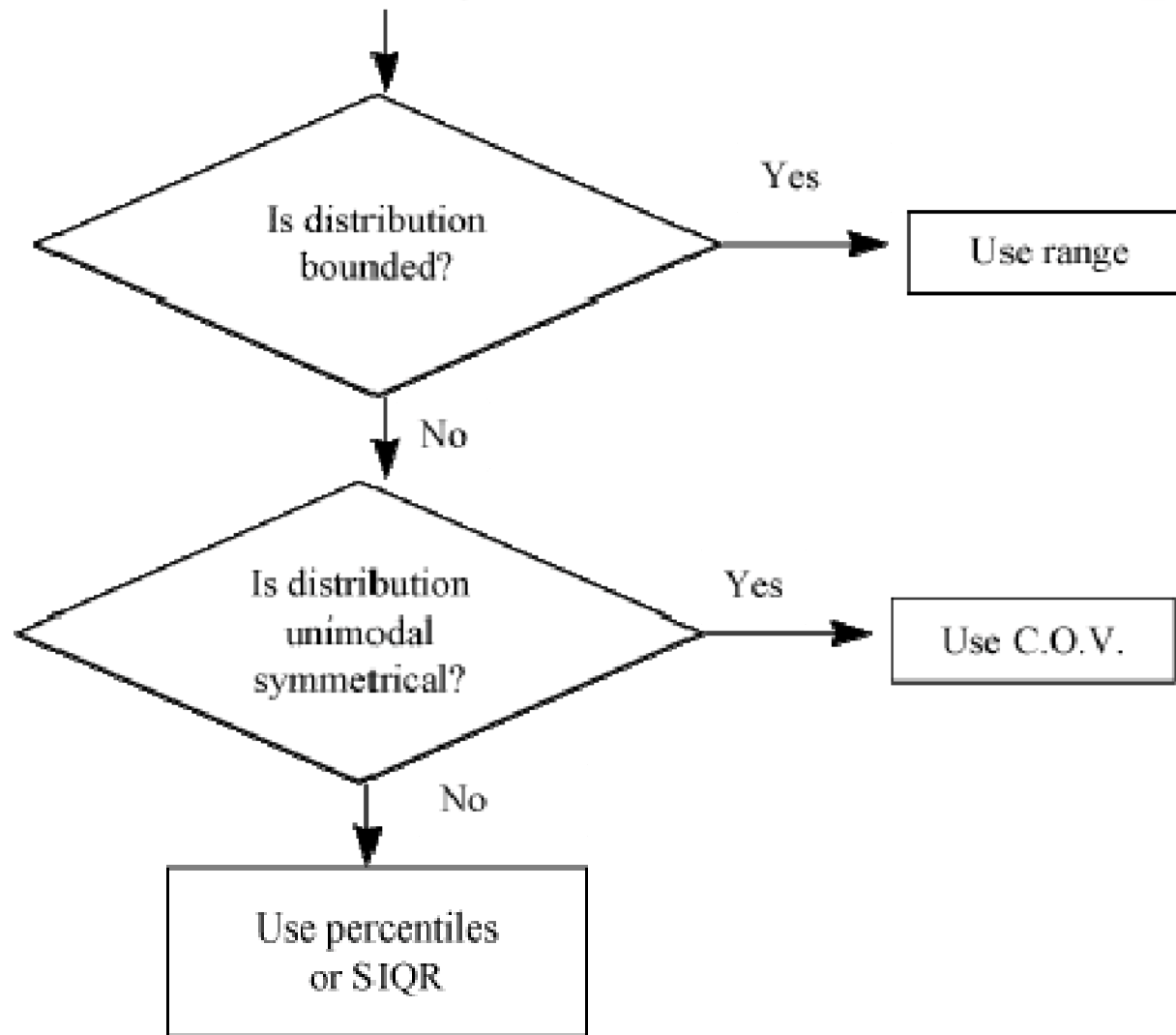
□ First quartile  $Q_1 = [1+(31)(0.25)] = 9\text{th element} = 3.2$

□ Median  $Q_2 = [1+(31)(0.50)] = 16\text{th element} = 3.9$

□ Third quartile  $Q_3 = [1+(31)(0.75)] = 24\text{th element} = 4.5$

□ 
$$\text{SIQR} = \frac{Q_3 - Q_1}{2} = \frac{4.5 - 3.2}{2} = 0.65$$

# Selecting the Index of Dispersion





## Selecting the Index of Dispersion (Cont)

- ❑ The decision rules given above are not hard and fast.
- ❑ Network designed for average traffic is grossly under-designed.  
The network load is highly skewed  
=> Networks are designed to carry 95 to 99-percentile of the observed load levels  
=>Dispersion of the load should be specified via range or percentiles.
- ❑ Power supplies are similarly designed to sustain peak demand rather than average demand.
- ❑ Finding a percentile requires several passes through the data, and therefore, the observations have to be stored.
- ❑ Heuristic algorithms, e.g.,  $P^2$  allows dynamic calculation of percentiles as the observations are generated.
- ❑ See Box 12.1 in the book for a summary of formulas for various indices of central tendencies and dispersion

# Determining Distribution of Data

- ❑ The simplest way to determine the distribution is to plot a histogram
- ❑ Count observations that fall into each cell or bucket
- ❑ The key problem is determining the cell size.
- ❑ Small cells => large variation in the number of observations per cell
- ❑ Large cells => details of the distribution are completely lost.
- ❑ It is possible to reach very different conclusions about the distribution shape
- ❑ One guideline: if any cell has less than five observations, the cell size should be increased or a variable cell histogram should be used.

# Quantile-Quantile plots

- $y_{(i)}$  is the observed  $q_i$ th quantile  
 $x_i =$  theoretical  $q_i$ th quantile
- $(x_i, y_{(i)})$  plot should be a straight line
- To determine the  $q_i$ th quantile  $x_i$ , need to invert the cumulative distribution function.

$$q_i = F(x_i)$$

- or  $x_i = F^{-1}(q_i)$

- Table 28.1 lists the inverse of CDF for a number of distributions.

## Quantile-Quantile plots (Cont)

- Approximation for normal distribution  $N(0,1)$

$$x_i = 4.91[q_i^{0.14} - (1 - q_i)^{0.14}]$$

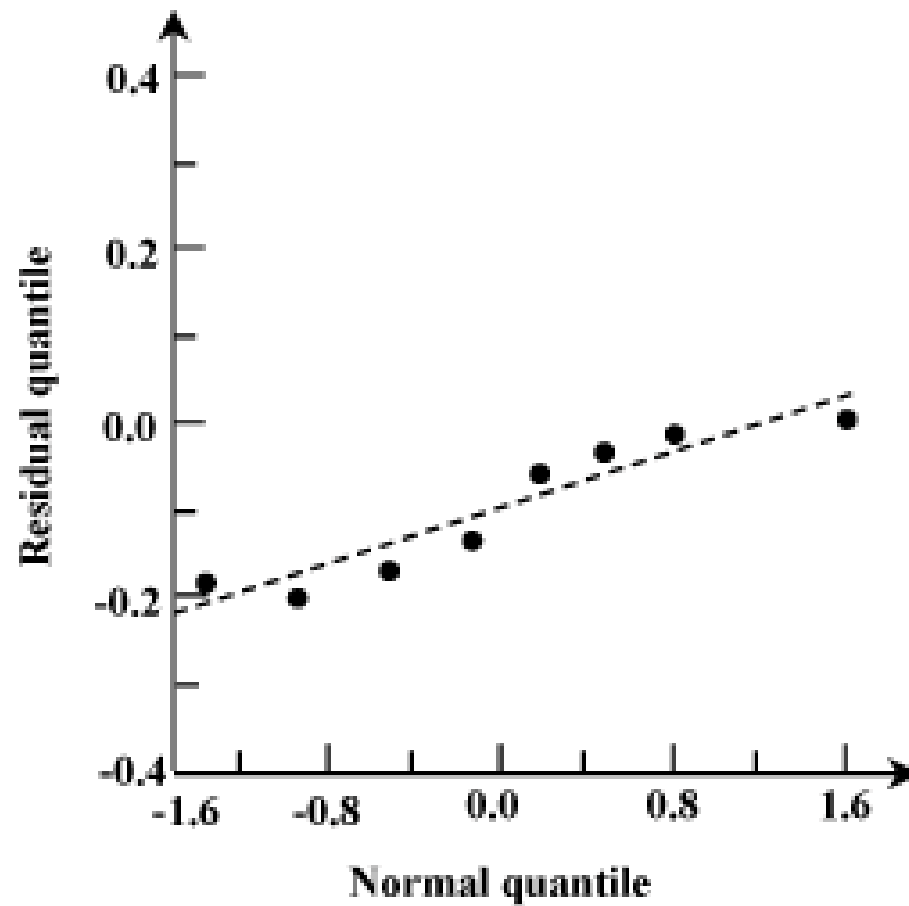
- For  $N(\mu, \sigma)$ , the  $x_i$  values computed above are scaled to  $\mu + \sigma x_i$  before plotting.

# Quantile-Quantile Plots: Example

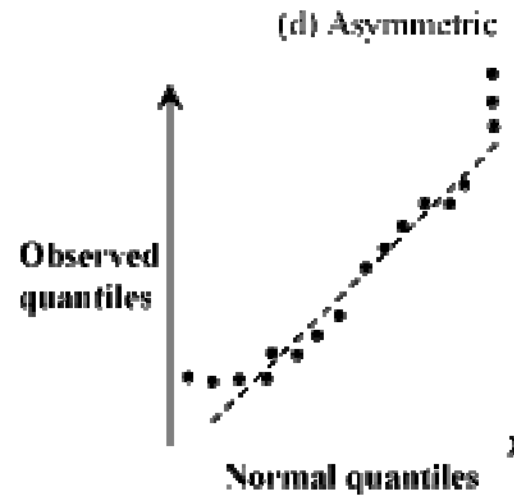
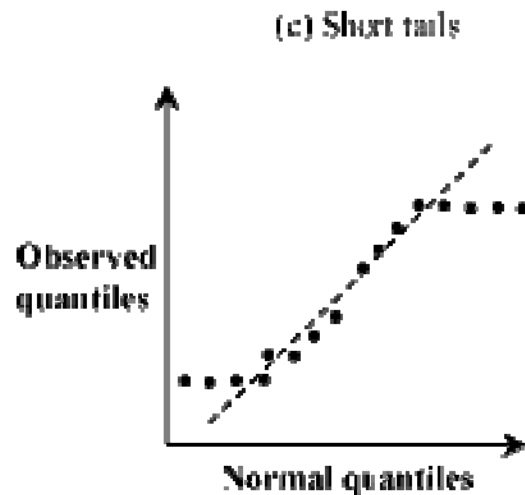
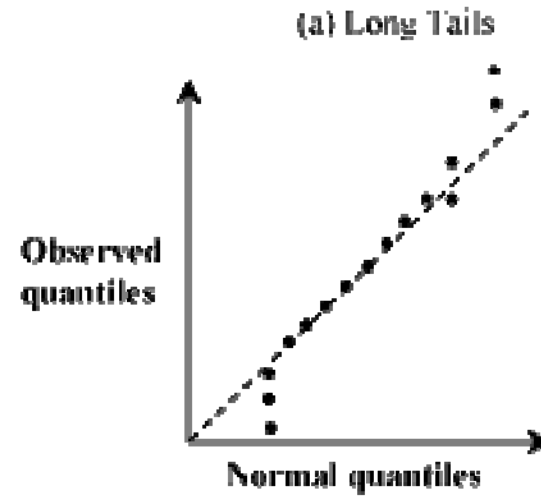
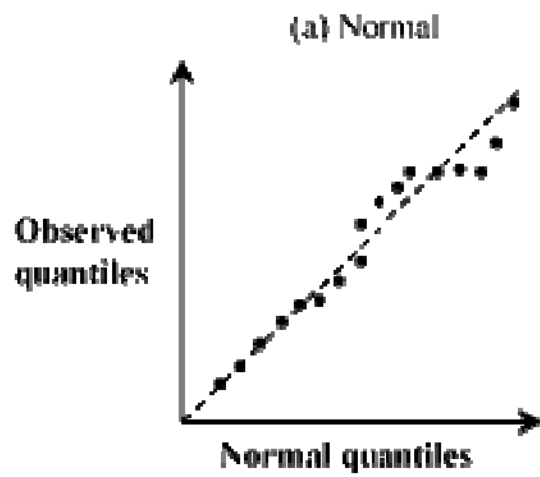
- The difference between the values measured on a system and those predicted by a model is called modeling error. The modeling error for eight predictions of a model were found to be -0.04, -0.19, 0.14, -0.09, -0.14, 0.19, 0.04, and 0.09.

$i$	$q_i = \frac{i-0.5}{n}$	$y_i$	$x_i$
1	0.0625	-0.19	-1.535
2	0.1875	-0.14	-0.885
3	0.3125	-0.09	-0.487
4	0.4375	-0.04	-0.157
5	0.5625	0.04	0.157
6	0.6875	0.09	0.487
7	0.8125	0.14	0.885
8	0.9375	0.19	1.535

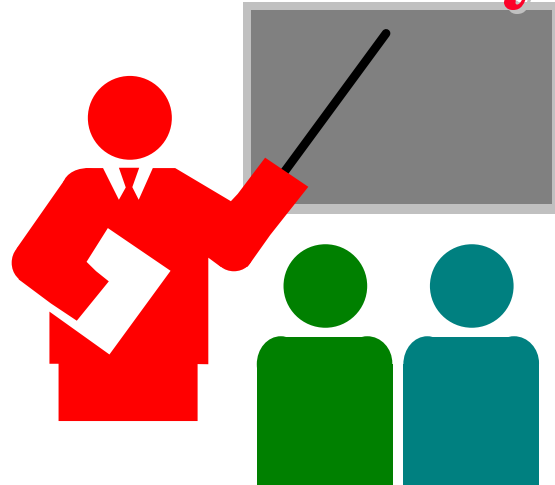
# Quantile-Quantile Plot: Example (Cont)



# Interpretation of Quantile-Quantile Data



# Summary



- ❑ Sum of a large number of random variates is normally distributed.
- ❑ Indices of Central Tendencies: Mean, Median, Mode, Arithmetic, Geometric, Harmonic means
- ❑ Indices of Dispersion: Range, Variance, percentiles, Quartiles, SIQR
- ❑ Determining Distribution of Data: Quantile-Quantile plots



# Homework 12B: Exercise 12.15

- Plot a normal quantile-quantile plot for the following sample of errors:

-0.04444	-0.04439	-0.04165	-0.03268	-0.03235	-0.03182	-0.02771	-0.02650
-0.02569	-0.02358	-0.02330	-0.02305	-0.02213	-0.02128	-0.01793	-0.01668
-0.01565	-0.01509	-0.01432	-0.00978	-0.00889	-0.00687	-0.00543	-0.00084
-0.00083	-0.00048	-0.00024	0.00079	0.00082	0.00106	0.00110	0.00132
0.00162	0.00181	0.00280	0.00379	0.00411	0.00424	0.00553	0.00865
0.01026	0.01085	0.01440	0.01562	0.01975	0.01996	0.02016	0.02078
0.02134	0.02252	0.02414	0.02568	0.02682	0.02855	0.02889	0.03072
0.03259	0.03754	0.04263	0.04276				

- Are the errors normally distributed?