# 1   Data

> **Definition**
>
> **Categorical Variables** represent data that can be divided into different groups.

Common examples of categorical data include ethnicity, income level, education, age group, and gender. Categories are described by words or letters. Categorical data is much harder to analyze mathematically than numerical data.

> **Definition**
>
> **Quantitative Variables** represents numerical data from population measurement. Quantitative data can be either **discrete** or **continuous**.

Common examples of quantitative data include height, weight, income, age, and cost. If data set can only take on specific values (e.g. integer values), then the data set is discrete. If the data is not restricted to specific values, then the data set is continuous.

# 2   Sampling

Gathering the data on an entire population may be virtually impossible due to limitations in time and cost. Instead, a **sample** of the population will be studied. The sample must be representative of the population being sampled in order for any statistical inference to be made. Random sampling is used to to create samples that minimize any bias from the sampling process.

> **Definition**
>
> A **simple random sample (SRS)** is a sampling method in which each member of a population has an equal chance of being selected.

In practice, it may be difficult to achieve a random sample. Care must be taken to eliminate biases that arise in a sampling process. Many data selection methods are not truly random. Random number generators can be powerful tools to add randomization to the sampling process. Two other common types of sampling methods are stratified sampling and cluster sampling. Introductory statistics assumes a simple random sampling process.

# 3   Histograms

> **Definition**
>
> A **Histogram** is a bar chart representing how many data points are present in each specified interval.

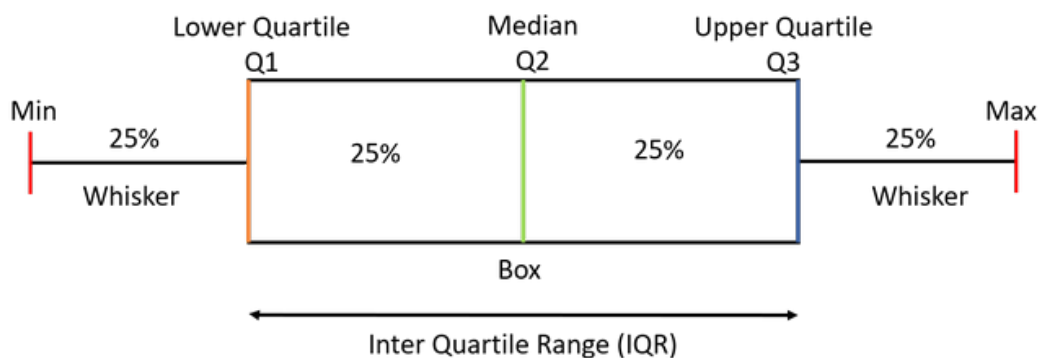To construct a histogram from a data set:

1. Determine number of bins.

2. Find the bin width and determine the bin intervals.

3. Count the number of data points per bin.

4. Plot data point counts vs bins on a bar graph.

# 4 Percentiles

> **Definition**
>
> **Percentiles** measure the location of data relative to the entire data set. The $k$th percentile is a score below which $k$ percent of the data falls below.
>
> ---
>
> The 25th percentiles is called the **first quartile** $Q_1$. The 50th percentile is called the **median** or **second quartile** $Q_2$. The 75th percentile is called the **third quartile** $Q_3$.



> **Definition**
>
> A **box plot** is a diagram representing five values: the minimum, maximum, and the three quartiles.
>
> ---
>
> To calculate a box plot on a TI-84 calculator:
>
> 1. STAT EDIT. Fill a list with values.
>
> 2. Press STAT and arrow to CALC
>
> 3. 1-VarStats. Choose list with values entered. ENTER

# 5 Statistical Mean

> **Definition**
>
> The **mean (arithmetic mean)** is one of the most common measures of center of a data set in statistics. The mean of a data set representing an entire population is called a **population mean** and is denoted $\mu$.

The mean can be calculated by the formula

$$\mu = \frac{1}{N} \sum x$$

where the sum is taken over the entire data set and $N$ represents the population size. The mean of a data set representing a sample from a population is called a **sample mean** and is denoted $\bar{x}$. The formuala for sample mean

is identical to the formula for population mean

$$\bar{x} = \frac{1}{n} \sum x$$

except that the sum is taken over the sample data set and $n$ represents the sample size. For datasets described by frequency tables or histograms,

$$\mu = \frac{\sum fm}{\sum f}$$

where $f$ is the frequency of the interval and $m$ is the midpoint of the interval.

# 6  Standard Deviation

> **Definition**
>
> The **standard deviation** provides a measure of the overall variation in a data set. The standard deviation can be used to determine whether a data value is close to or far away from the mean.

The **population standard deviation** $\sigma$ is computed by the formula

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

where $\mu$ is the population mean and $N$ is the population size. The **sample standard deviation** $s$ has a similar formula

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$

where $\bar{x}$ is the sample mean and $n$ is the sample size. Observe that the denominator is $n - 1$ not $n$.