# 1.017/1.010 Class 13
# Populations and Samples

---

## Sampling and Repeated Trials:

Statistics relies on sets of observations of random variables, commonly called **samples**.

A sample can be viewed as a collection of outcomes from **repeated trials**. Random trial $i$ yields an outcome (a particular number) corresponding to an underlying random variable $x_i$. The collection of random variables associated with a sample of $N$ observations is $[x_1, x_2, x_3, \ldots, x_N]$.

Samples are used to estimate distributional properties of the underlying random variables or to make inferences about relationships among different variables.

Samples are sometimes supposed to be drawn from a larger **population** of (which may actually exist or may be hypothesized).

Example:

Define the random variable $x$ to be the height of an MIT student drawn at random from the **population** of all MIT students. Suppose we record the heights of 10 MIT students selected at random. This can be viewed as a composite experiment composed of 10 repeated trials with 10 corresponding random heights $[x_1, x_2, x_3, \ldots, x_{10}]$. The outcomes obtained from this set of trials constitute a **sample**.

The 10 student sample can be used to estimate the unknown mean and variance (or other properties) of the height probability distribution or to make inferences about this distribution (e.g. to compare it to a similar distribution from Stanford).

The following is a sample of 10 student heights (in inches).

[56.25, 59.90, 54.64, 65.73, 56.45, 57.45, 61.26, 57.23, 56.61, 53.67]

What are reasonable estimates of the mean and variance of the student height probability distribution?

## Random samples:

$[x_1, x_2, x_3, \ldots, x_N]$ is a **random sample** if:

1. The $x_i$'s are **independent** $F_{x1, x2, \ldots xN} = F_{x1}(x_1)F_{x2}(x_1)\ldots F_{xN}(x_N)$
2. The $x_i$'s all have the **same distribution** $F_{xi}(x_i) = F_x(x)$

## Functions of many random variables

A function $y = g(x_1, x_2, \ldots x_N)$ of a sample of $N$ random variables is a random variable with CDF $F_y(y)$.

If $[x_1, x_2, \ldots x_N]$ is a random sample $F_y(y)$ can be derived from $g(x_1, x_2, \ldots x_N)$ and $F_x(x)$. In general, this derived distribution problem can be solved with stochastic simulation.

## Mean and variance of a linear function (prove)

Suppose that $x_1$ and $x_2$ are random variables (not necessarily independent) and $a_1, a_2, b$ are constants (not random):

$$E[a_1x_1 + a_1x_2] = a_1E[x_1] + a_2E[x_2] + b$$

$$Var[a_1x_1 + a_1x_2 + b] = a_1^2Var[x_1] + a_2^2Var[x_2] + 2a_1a_2Cov[x_1, x_2]$$

If $x_1, x_2$ are uncorrelated

$$Var[a_1x_1 + a_1x_2 + b] = a_1^2Var[x_1] + a_2^2Var[x_2]$$

Example:

A reasonable **estimate of the mean** of an unknown probability distribution is the arithmetic average (sample mean) $m_x$ of $N$ observations drawn at random from this distribution.

If the observations are outcomes associated with a random sample $x_1, x_2, \ldots x_N$ the sample mean is:

$$m_x = \frac{1}{N}\sum_{i=1}^{N} x_i$$

This is a linear function of the random sample. The mean and variance of $m_x$ are (prove):

$$E[m_x] = E[x] = \bar{x}$$

$$Var[m_x] = \frac{1}{N}Var[x]$$

A reasonable **estimate of the standard deviation** of the unknown distribution is the mean sum of squares (sample variance) $s_x^2$ computed from the random sample:

$$s_x^2 = \frac{1}{N-1}\sum_{i=1}^{N}[x_i - m_x]^2$$

This is a linear function of the set of variables $[x_i - m_x]^2$, $i = 1...N$. The mean and variance of $s_x^2$ are:

$$E[s_x^2] = \sigma_x^2$$

$$Var[s_x^2] \propto \frac{1}{N}$$