


Statistics of data


Topics: Mean, Median, Mode, Variance, Std. Deviation, Normal distribution

Measuring the Central Tendency: (1) Mean


- Mean (algebraic measure) (sample vs. population):

Note: n is sample size and N is population size.


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


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

- Weighted arithmetic mean:


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

- Trimmed mean:

- Chopping extreme values (e.g., Olympics gymnastics score computation)

Measuring the Central Tendency: (2) Median

- Median:

- Middle value if odd number of values, or average of the middle two values otherwise

- Estimated by interpolation (for *grouped data*):

<i>age</i>	<i>frequency</i>
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44

Approximate median

Sum before the median interval

Interval width ($L_2 - L_1$)

Low interval limit

$$median = L_1 + \left(\frac{n / 2 - (\sum freq)_l}{freq_{median}} \right) width$$

The diagram illustrates the formula for the approximate median for grouped data. A yellow box labeled 'Approximate median' has a curved arrow pointing to the formula. A green box labeled 'Sum before the median interval' has a downward arrow pointing to the term $(\sum freq)_l$ in the numerator. A green box labeled 'Interval width ($L_2 - L_1$)' has a downward arrow pointing to the variable 'width' in the formula. A green box labeled 'Low interval limit' has an upward arrow pointing to the variable L_1 in the formula.

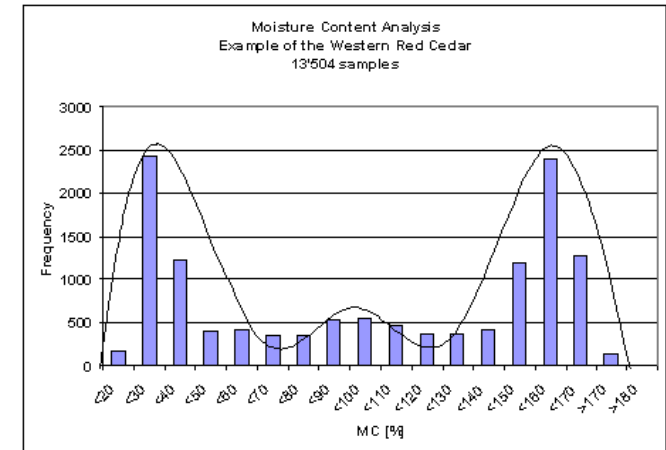
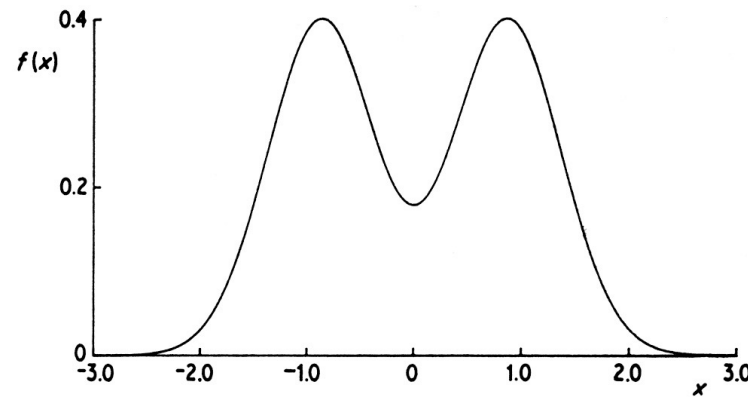
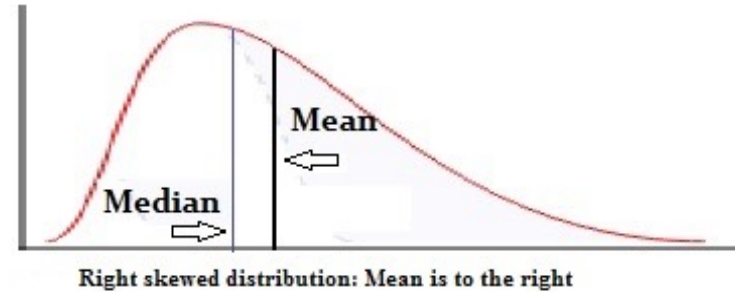
Measuring the Central Tendency: (3) Mode

- Mode: Value that occurs most frequently in the data

- Unimodal

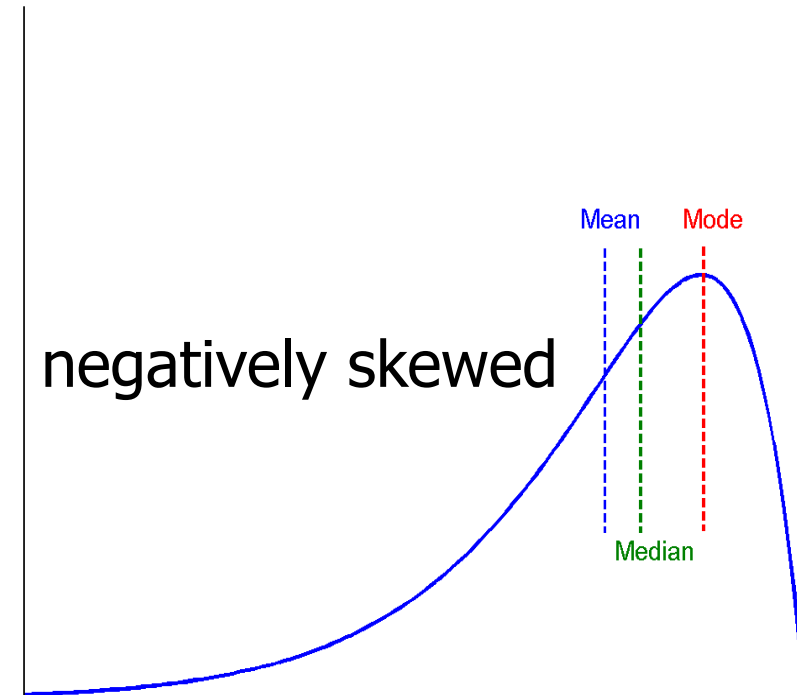
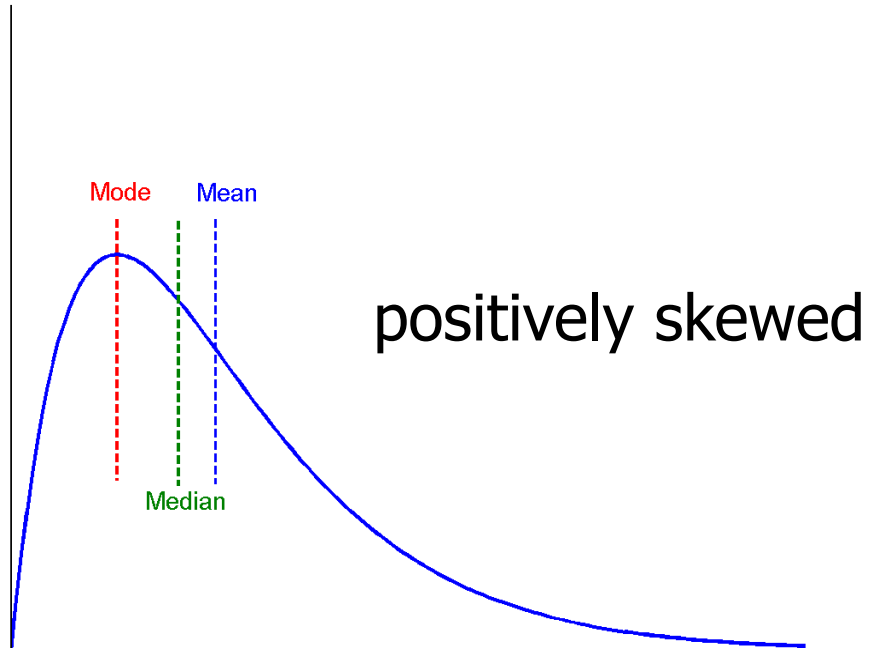
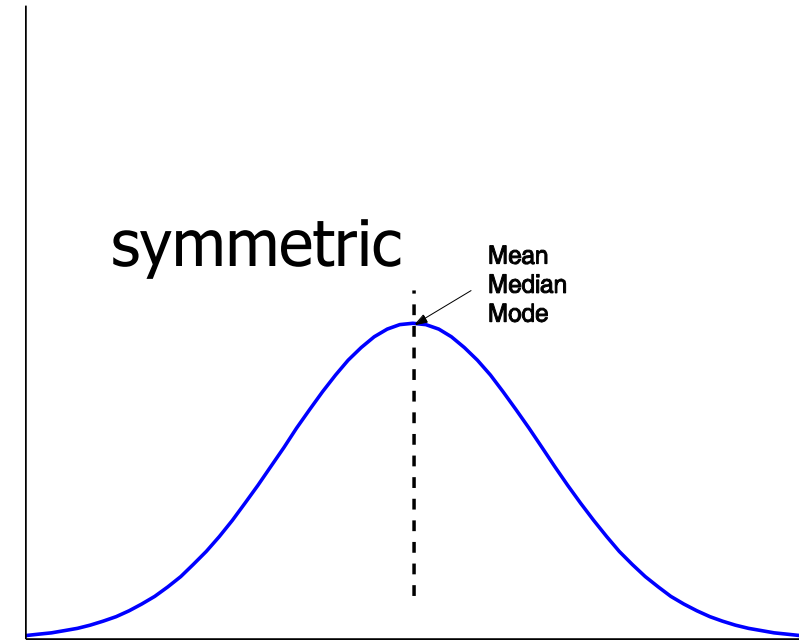
- Multi-modal

- Bimodal
- Trimodal



Symmetric vs. Skewed Data

- Median, mean and mode of symmetric, positively and negatively skewed data



Measures Data Distribution: Variance and Standard Deviation

- Variance and standard deviation (*sample: s, population: σ*)

- **Variance:**

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

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

Note: The subtle difference of formulae for sample vs. population

- n : the size of the sample
- N : the size of the population

- **Standard deviation** s (*or σ*) is the square root of variance s^2 (*or σ^2*)

Properties of Normal Distribution Curve

Normal distribution is

- a distribution of data that has roughly the same amount of data on either side of the middle and
- has its most common values around the middle of the data

