# Descriptive Statistics

| Ratio | Interval | Ordinal | Nominal |
|---|---|---|---|
| • absolute zero | arbitrary zero | $f(\text{old}) = \text{new}$ | |
| • geometric/harmonic mean | arithmetic mean | order preserving monotonic | |
| $ax$ | $ax + b$ | | |
| mean, SD, Pearson's corr., t-test, F-test | median, percentile, rank, run test, sign test | mode, entropy, $\chi^2$ test | |

geometric mean $\left[ \sqrt[n]{x_1, x_2 \cdots x_n} \right.$

$\dfrac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \cdots + \frac{1}{x_n}}$ harmonic mean

**4 moments** →
- mean
- variance
- skewness
- kurtosis

cross sectional: several variables at one instance

time series: one var across several instances

longitudinal/panel: combine above

## measures of shape

### skewness
measure of symmetry (or its lack)

**Pearson's moment coefficient of skewness ($g_1$)**

$$g_1 = \frac{\sum\limits_{i=1}^{n} \frac{(x_i - \bar{x})^3}{n}}{\sigma^3}$$

$g_1 \to 0$   when data symmetrical

$g_1 \to +ve$   positively skewed

$g_1 \to -ve$   negatively skewed

symmetric : skewness $\in (-0.5, 0.5)$

moderate skew : $\in (-1, -0.5) \cup (0.5, 1)$

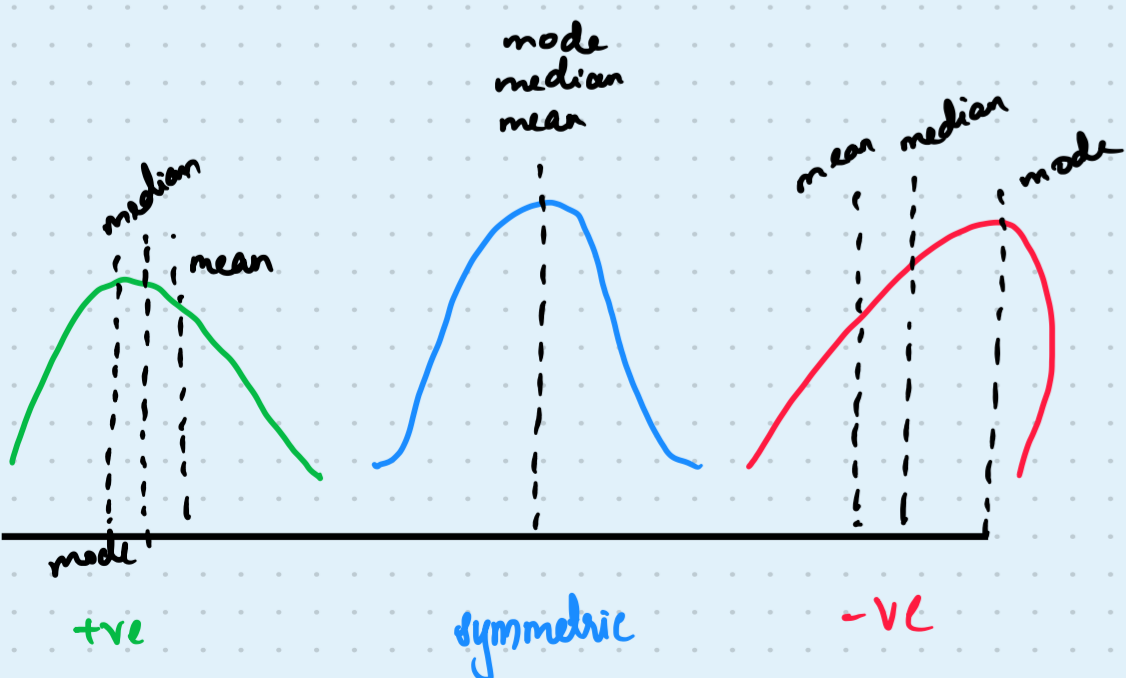highly skew : $\in (-\infty, -1) \cup (1, \infty)$

**Skewness with sample of n observations ($G_1$)**

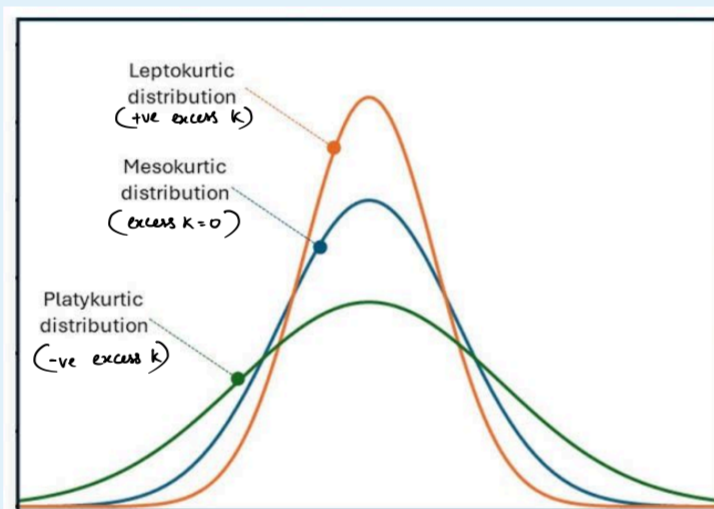$$G_1 = \frac{\sqrt{n(n-1)}}{(n-2)} g_1$$

(Joanes & Gill, 1998)

$G_1 \to 1$ as n increases



+ve   symmetric   -ve

### kurtosis
aimed at the tail; checks whether tail is heavy or light

$$\text{Kurtosis} = \frac{\sum\limits_{i=1}^{n} \frac{(x_i - \bar{x})^4}{n}}{\sigma^{-4}}$$



platykurtic: low frequency of outliers

leptokurtic: high frequency of outliers

**excess kurtosis**
= kurtosis - 3

kurtosis only captures information from outliers thanks to higher power

## exploratory data analytics (EDA)
preliminary exploration

**summary statistics**
part of EDA, summarise data

**quantile**

let $x_1 \cdots x_n$

quantiles : data points that divide the dataset into equal sized parts

$k^{th}$ q-quantile $\Rightarrow$ $\dfrac{k}{q}$ of set before, $\dfrac{(q-k)}{q}$ of set after

$q-1$ quantile points exist

$x$-percentile $\Rightarrow$ $x\%$ of set before $(100 - x)\%$ of set after

$$Q_k = \frac{n+1}{q} \times k$$

Say value at position 5 is 21 and value at position 6 is 22

Value at position 5.1 $\approx$

$21 + 0.1 (22 - 21) = 21 + (0.1 \times 1)$
$= 21.1$

## central tendency

$$\sum (x_i - \bar{x}) = 0$$

mean

median : less influenced by outliers, no need for entire dataset

mode

## measures of variation
identify outliers, how close records are to the mean

feature w/ low variability → unlikely to have statistical significance w/ target variable

**Coefficient of variation**

$$CV = \frac{\sigma}{\bar{x}}$$

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

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

### IQR

$$IQR = Q_3 - Q_1$$

Outliers $\in (-\infty, Q_1 - 1.5(IQR))$
$\cup$
$(Q_3 + 1.5(IQR), \infty)$



Tuesday, February 20, 2024    12:01 AM

**MARKOV'S INEQUALITY**
$X \to$ non-negative random variable with finite mean $\mu$

$$P(X \geq c) \leq \frac{\mu}{c}$$

**CHEBYSHEV'S INEQUALITY**
$X \to$ random variable w/ finite $\mu$ and $\sigma^2$ with any distribution

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

Markov's inequality applied to $(X - \mu)^2$

**Other forms**

① $P(|X - \mu| \geq c) \leq \dfrac{\sigma^2}{c^2}$

$k = \dfrac{c}{\sigma}$

② $P((X - \mu)^2 \geq k^2 \sigma^2) \leq \dfrac{1}{k^2}$

③ $P(\mu - k\sigma < X < \mu + k\sigma) \geq 1 - \dfrac{1}{k^2}$  } complement of original statement

④ $P(X \geq \mu + k\sigma) + P(X \leq \mu - k\sigma) = \dfrac{1}{k^2}$ } $|X - \mu| \geq k\sigma$ → $X \geq \mu + k\sigma$ (2 cases) → $X \leq \mu - k\sigma$

- $P(X \geq \mu + k\sigma) \leq \dfrac{1}{k^2}$
- $P(X \leq \mu - k\sigma) \leq \dfrac{1}{k^2}$

**standard units**: No. of standard deviations that a particular value of a random variable is away from the mean
- Determined by $\dfrac{X - \mu}{\sigma}$
- $-k\sigma \leq X - \mu \leq k\sigma$ for some small $k$
- $\mu - k\sigma \leq X \leq \mu + k\sigma$

describes the percentage of values within a certain $k\sigma$ of the mean; gives bounds for $P(X)$ where $X$ is outside $k\sigma$ limits from the mean. $k$ is strictly +ve real, $k \in \mathbb{R}_{++}$