


Q : sample v/s population meaning?

Variability (spread / dispersion)

- shows how 'spread out' data is
- measures of variability are :
 - range
 - interquartile range
 - variance
 - standard deviation

Range

range = (highest score - lowest score)

Interquartile Range

- $IQR = (75^{th}\% - 25^{th}\%) \approx H\text{-spread}$

Semi-IQR

- $SIR = (IQR) / 2$
- in a symmetric distribution, the median \pm the SIR contains half the scores in the distribution

Variance

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

mean of deviation squared
 X : data
 μ : mean
↖ variance in a population

$$s^2 = \frac{\sum (X - M)^2}{N - 1}$$

↖ variance est. from a sample

Standard Deviation

$$\begin{aligned} \text{sd}(\sigma) &= \sqrt{\sigma^2} && \text{(population standard deviation)} \\ \text{sd}(s) &= \sqrt{s^2} && \text{(sample standard deviation)} \end{aligned}$$

Normal Distributions

can calculate the proportion of the distribution within a given number of SD from the mean

68% w/in 1 SD (+/-)

95% w/in 2 SD (+/-)

one that is not skewed.

Central Tendency

for symmetric distributions: mean, median, trimean, trimmed mean are

↑ equal
(mode is also equal except in bimodal dist)
↳ in bell shaped normal dist

in positive skew, mean is usually higher than the median and geometric mean is lower than all measures except the mode.

with skewed distributions, the geometric mean, trimean and trimmed mean are between the median and the mean.

what report? : mean, median and trimean / mean trimmed 50%

media reports the median for skewed distributions

trimean

weighted average of the 25th%, 50th% and 75th%.

$$\text{trimean} = (P_{25} + 2P_{50} + P_{75}) / 4$$

geometric mean

$$\text{geometric mean} = (\prod x)^{1/N}$$

multiply all numbers and take the n^{th} root

↑ related to logs

trimmed mean

mean calcd after removing some of the higher and lower scores

trimmed $x\%$ $\Rightarrow \frac{x}{2}\%$ scores from the bottom and $\frac{x}{2}\%$ scores from the top are removed

less influenced by extreme scores

median is basically mean trimmed 98+%

measure of CT ①

mean

balance point of the distribution

minimises sum of squared deviation

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

$$\sum (x - \bar{x})^2$$

(basically means that if deviation is calcd and squared wrt mean, will be the smallest sum)

median

minimises the sum of absolute deviations

↓

Same as value that minimises avg absolute deviation.

midpoint of dist (same no of scores above & below)
↑ 50th percentile value

mode

- most frequently occurring value
- for continuous data, the mode is computed as the midpoint of the most frequent interval

measure of
CT ②

- smallest absolute difference (when calcd wrt a
22.7 → mid value)

measure
of CT ③

Variance of the sum of 2 variables

- select a number each from 2 populations, add them. repeat that.
- What is the variance of the sums?

↓

$$\sigma_{\text{sum}}^2 = \sigma_{p_1}^2 + \sigma_{p_2}^2$$

(or diff)

$$\text{mean}(\text{sum}) = \text{mean}(p_1) + \text{mean}(p_2)$$

← regardless of sign

variance sum law: $\sigma_{x+y}^2 = \sigma_x^2 + \sigma_y^2 \pm 2\rho\sigma_x\sigma_y$

ONLY WHEN x, y are independent variables

← for independent this term is 0

(randomly paired \Rightarrow independence)
(for 2 people or something)

effects of linear transformation

- mean \rightarrow changed (same)
- std dev \rightarrow changed (coeff)
- variance \rightarrow changed (coeff²)
- corr \rightarrow unchanged.

if X has mean μ_x
std dev σ_x
var σ_x^2

$$Y = bX + A$$

$$\mu_y = b\mu_x + A$$

$$\sigma_y = b\sigma_x$$

$$\sigma_y^2 = b^2\sigma_x^2$$

measuring distribution skew

pearson's measure of skew : $\frac{3(\text{mean} - \text{median})}{\sigma}$

third moment, measure of skew : $\sum \left\{ \frac{(X - \mu)^3}{N\sigma^3} \right\}$
(more common) about the mean

estimating the skew of samples : $\frac{n}{(n-1)(n-2)} \sum \frac{(X - M)^3}{s^3}$
sample size sample mean sample std dev

measure of kurtosis : $\sum \frac{(X - \mu)^4}{N\sigma^4} - 3$
fourth moment about the mean kurtosis of a normal dist

by hand : $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum \frac{(X - M)^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$

lecture notes

eg 1

1 3 4 4 4 5 5 7 8 9 31

$$\bar{x} = \frac{\sum x}{n} = \frac{81}{11} = 7.36 \dots$$

$$\text{median (50th \%)} = 5$$

$$\text{mode} = 5$$

$$\text{trimean} = \frac{4 + 2 \cdot 5 + 8}{4} = 5.5$$

$$\text{geometric mean} = (1 \times 3 \times 4 \times \dots \times 31)^{1/11} = 5.2$$

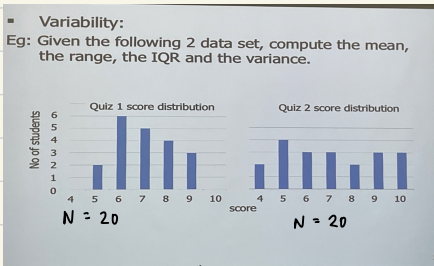
$$\text{mean trimmed 18.2\%} = \frac{49}{9} = 5.4$$

(remove 1, 31)

$$\text{rank} = \frac{1}{N} (n^{\text{th}} \% + 1)$$

$$\%_{\text{th}} = \frac{\text{data at rank}_R}{N} +$$

eg 2



μ

7

7

range

$$9 - 5 = 4$$

$$10 - 4 = 6$$

IQR

$$8 - 6 = 2$$

$$9 - 5 = 4$$

σ^2

$$\frac{\sum (x - \mu)^2}{N} = 1.5$$

$$3.9$$

Population

mean $\mu = E[X] = \frac{\sum X}{N}$

variance $\sigma^2 = E[(X - \mu)]^2 = \frac{\sum (X - \mu)^2}{N} = \frac{\sum X^2}{N} - \frac{(\sum X)^2}{N^2}$

$$E(Y^2) = \sigma_y^2 + \mu_y^2 \qquad = E[X^2] - (E[X])^2 \\ = E[X^2] - \mu^2$$

Sample

mean $\bar{x} = \frac{\sum X}{n}$

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

only for linear, so if graph is curved, nope!

↖ Pearson Correlation

population

$$\rho = \frac{E[(X - \mu_x)(Y - \mu_y)]}{\sigma_x \sigma_y} \xleftarrow{\text{cov}(X,Y)} = \frac{E[XY] - \mu_x \mu_y}{\sqrt{E[(X - \mu_x)^2] E[(Y - \mu_y)^2]}}$$

$$= \frac{\frac{\sum XY}{N} - \frac{\sum x \sum y}{N^2}}{\sqrt{E[X^2] - \mu_x^2} \sqrt{E[Y^2] - \mu_y^2}} = \frac{\sum XY - \frac{\sum x \sum y}{N}}{\sqrt{(\sum x^2 - \frac{(\sum x)^2}{N})(\sum y^2 - \frac{(\sum y)^2}{N})}}$$

sample!!

sample

$$r = \frac{E[(X - \bar{x})(Y - \bar{y})]}{s_x s_y} \xleftarrow{\text{cov}(X,Y)} = \frac{1}{n-1} \sum (X - \bar{x})(Y - \bar{y})$$

$$= \frac{\sum XY - \frac{\sum x \sum y}{n}}{\sqrt{\sum x^2 - \frac{(\sum x)^2}{n}} \sqrt{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

Chpt 4

Bivariate Data

↳ dataset w/ pair of variables which may be correlated to each other

eg. ice cream sales + temperature

Pearson Correlation ρ

$$\rho = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} \quad \text{co-variance of X and Y denoted as } \text{cov}(XY)$$

$$= \frac{E[XY] - \mu_X \mu_Y}{\sqrt{(E[X^2] - (\mu_X)^2)(E[Y^2] - (\mu_Y)^2)}}$$

$$\rho = \frac{\sum XY - (\sum X \sum Y) / N}{\sqrt{\sum X^2 - \frac{(\sum X)^2}{N}} \sqrt{\sum Y^2 - \frac{(\sum Y)^2}{N}}}$$

$$\mu_X = \mu_Y = 0$$

$$\rho = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}}$$

Pearson Correlation
Given the data set, calculate \bar{X} , \bar{Y} , $\text{cov}(X, Y)$ and the Pearson Correlation

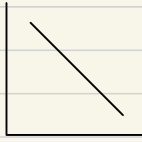
X	2	5	6	8	9
Y	8	5	2	4	1

$\mu_X = \frac{\sum X}{N} = 6$ $\sigma_X^2 = \frac{\sum (X - \mu_X)^2}{N} = 6$
 $\mu_Y = \frac{\sum Y}{N} = 4$ $\sigma_Y^2 = 6$
 $\text{cov}(X, Y) = \frac{\sum (X - \mu_X)(Y - \mu_Y)}{N} = -5.2$
 $\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = -0.867$

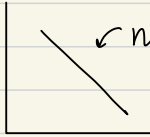
Eg. 4 boys and 2 girls sit in a row. Find the following:
 (i) No. of ways of putting these 6 people in a row.

$$\text{cov}(y, x) = \frac{1}{n-1} \sum (x - \bar{x})(y - \bar{y})$$

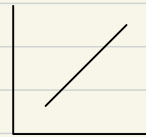
↑
sample



points along a line : linear relationship



negative association



positive association

Pearson Product-Moment correlation

- Strength of linear relationship b/w 2 variables
- valid only for linear relationships
- $\rho \rightarrow$ population
 $r \rightarrow$ sample
- linear transformation of a variable does not change its correlation with other variable
- $$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$
- variance sum law for dependent variables :
$$\sigma_{x \pm y}^2 = \sigma_x^2 + \sigma_y^2 \pm 2\rho\sigma_x\sigma_y$$

estimate mean and variance of a population of size N

population mean $\mu = E[X] = \frac{\sum X}{N}$

population variance $\sigma^2 = E[(X - \mu)^2] = \frac{\sum (X - \mu)^2}{N}$
or $\frac{\sum X^2 - \frac{(\sum X)^2}{N}}{N}$

$$\approx \sigma^2 = E[X^2] - \mu^2$$
$$\therefore E[X^2] = \sigma^2 + \mu^2$$