

STATISTICS

- ✓ Final
 - 6 - Randomly selected from homework
 - 4 from ~~part 1~~
 - 2 from ~~part 2~~
- Text book - Basic business statistics.

BASIC PROBABILITY

→ probability - the extent to which something is likely to occur
 - the most likely cause of something.

→ sample space - the set of all possible outcomes.

→ event - subset of sample space, the set of all outcomes that produce a specific result.
 Simple (2 simple outcomes) joint (compound) result.

→ A complement (A') - subset of outcomes that are not part of event.

→ mutually exclusive events
 - a set of events that can't occur at the same time eg. heads and tails

→ collectively exhaustive events
 - if one of the events must occur
 eg. in a die roll
 event 1 - getting odd num.
 event 2 - getting even num.
 - they may be mutually exclusive or not.

Types of probability

1. Prior probability / classical
 is based on having a prior knowledge of the outcomes that can occur.
 eg. tossing a coin.

2. Empirical probability

- based on the observed data
 eg. football odds
- mainly identified by survey.

3. Subjective probability

- a probability that differs from person to person.

eg. probability of you marrying.

4. Arithmetic probability

→ simple probability

- probability of occurrence of single event

$$P = \frac{x}{n}$$

← no. of outcomes of event
 ← all possible outcomes.

→ joint probability

- The probability of occurrence of one two or more events.

→ marginal probability

- is an event consisting a set of joint probabilities.

eg.

	A	B
H		
T		

$$P(A) = P(A \text{ and } H) + P(A \text{ and } T)$$

→ General addition rule

$$P(A \text{ or } B) = P(A) + P(B) - P(A \cap B)$$

Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

← joint prob.
 ← marginal prob.

- The probability of occurrence of event A given that event B has already occurred.

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- If conditional probability of event B is known.

⊗ $P(A|B) = P(A)$ or $P(A|B) = P(A) \cdot P(B)$
 if they are independent of one another.

⊗ $P(A|B) = 0$
 if they are mutually exclusive

$$P(A \cap B) = P(A|B) \cdot P(B)$$

→ Marginal probability using multiplication rule.

$$P(A) = (A \cap B_1) + (A \cap B_2) + (A \cap B_n)$$

$$\rightarrow \{P(A) = (P(A|B_1) \cdot P(B_1)) + P(A|B_2) \cdot P(B_2) \dots\}$$

Where $B_1, B_2 \dots B_n$ are mutually exclusive
- collectively exhaustive

→ Bayes' theorem

$$P(B_i|A) = \frac{P(A|B_i) \cdot P(B_i)}{P(A|B_1)P(B_1) + P(A|B_2)P(B_2) \dots P(A|B_k)P(B_k)}$$

where B_i is i th event out of k mutually exclusive and collectively exhaustive events.

Counting rules

rule 1

- no. of possible outcomes if any one of k different mutually exclusive and collectively exhaustive events. in n trials is k^n

eg. variation of licence plates having 3 nos and 3 letters is

$$10^3 \cdot 26^3 = \dots$$

rule 2

- if there are k_1 events in 1st trial k_2 events in 2nd trial and k_n in n th trial

then no. of possible outcomes

$$(k_1)(k_2) \dots (k_n)$$

rule 3

- The no. of ways that all n items can be ordered.

$$n!$$

rule 4 - permutation

- no. of ways of arranging x items selected from n items.

$${}_n P_x = \frac{n!}{(n-x)!}$$

rule 5 - combination

- no. of ways of selecting x items from n items irrespective of order.

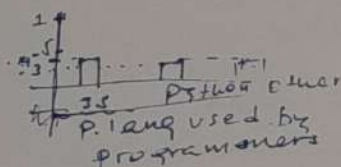
$${}_n C_x = \frac{n!}{x!(n-x)!}$$

CH-5

Discrete probability distribution

- is mutually exclusive list of all possible numerical outcomes along with the probability of each outcome.

eg	No of meals per day	probability
	0	1%
	1	15%
	2	29%
	3	40%
	4	15%



- * Expected value of discrete variable
- mean based on probability

$$\mu = E(X) = \sum_{i=1}^n x_i P(X=x_i)$$

- * Variance

$$\sigma^2 = \sum_{i=1}^n [x_i - E(X)]^2 P(X=x_i)$$

- * Standard deviation

$$\sigma = \sqrt{\sigma^2}$$

Binomial distribution

- n is the probability of event of interest.

$$P(X=x|n, \pi) = \frac{n!}{x!(n-x)!} \pi^x (1-\pi)^{n-x}$$

- n - no of observation
- x - no of event of interest in the sample
- π - probability of an event of interest.

- * mean of binomial distribution.

$$\mu = E(X) = n\pi$$

- * standard deviation of b dist

$$\sigma = \sqrt{\sigma^2} = \sqrt{\text{var}(X)} = \sqrt{n\pi(1-\pi)}$$

Poisson distribution

- probability of an event in area or interval of time given λ (expected no of event per unit)

$$P(X=x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$e = 2.7$$

λ = no of events.

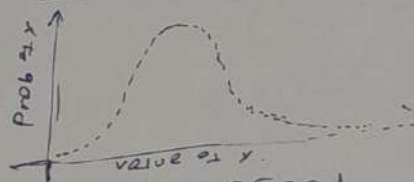
CH-6

The normal distribution and other continuous distribution

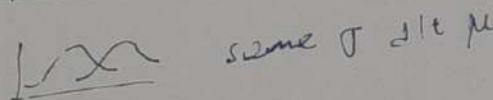
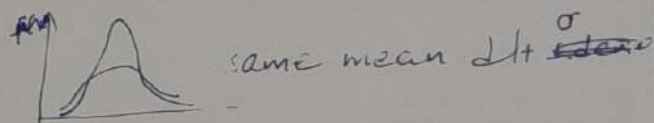
- it doesn't have finite no of variables or events

- Area under the curve is probability.
- The whole A. under the curve = 1
- single value can't have a prob because do can't have area

Normal distribution



- Symmetrical
- Its mean and median are equal
- Its range is infinite $(-\infty, \infty)$



- * Normal probability density function.

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

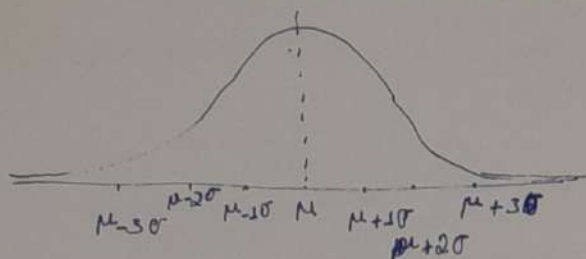
x - any value of the variable

$$Z = \frac{x - \mu}{\sigma}$$

Converts normally distributed variable into standardized normal variables

μ of $Z = 0$
 σ of $Z = 1$
 but normally distributed variables have their own μ and σ

if the ave load time (μ) for a website is 7 sec and standard deviation of 2



x scale 1 3 5 7 9 11 13
 z scale -3 -2 -1 0 +1 +2 +3

so load time of 1 sec is -3 standardized unit (3 standard deviations) below the mean value.

to find the probability in a normal distribution you can find the z index and find prob from the cumulative distribution table. and if x is asked given probability you can do the same.

⇒ Uniform distribution

values are evenly distributed in the range b/w the smallest value A and the largest value B.

uniform prob density fun. $f(x) = \frac{1}{(b-a)}$ if $a \leq x \leq b$

mean $\frac{a+b}{2}$

$$\sigma^2 = \frac{(b-a)^2}{12}$$

$$\sigma = \sqrt{\frac{(b-a)^2}{12}}$$

CH-1

Sampling the distribution of the mean

sampling distribution is the distribution of the result if you actually selected all possible samples.

sampling distribution of mean

- is distribution of results if you actually selected mean of all possible samples
- the mean of all possible sample means is equal to the pop. mean.
- it's always a normal distribution (becomes less normal if n of sample is less)
- standard error of the mean (standard deviation of sample)

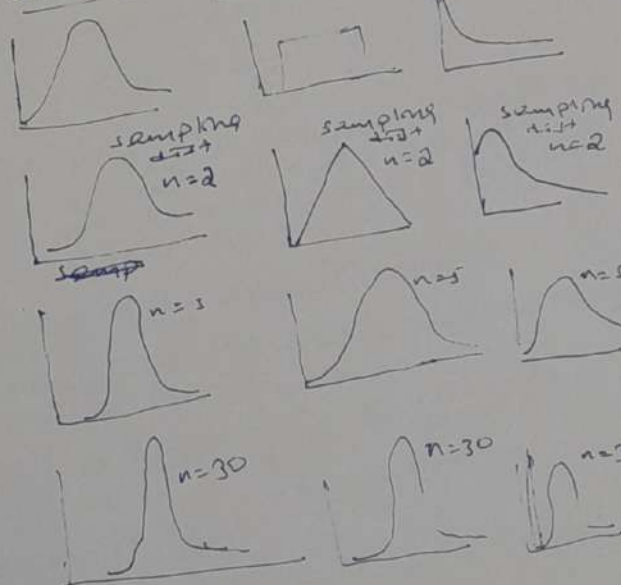
$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \quad n = \text{no of elem in the sample.}$$

z-score of the sampling dist of the mean

$$z = \frac{(\bar{x} - \mu)}{\sigma_{\bar{x}}} = \frac{(\bar{x} - \mu)}{\frac{\sigma}{\sqrt{n}}}$$

$$\bar{x} = \mu + z \frac{\sigma}{\sqrt{n}}$$

Normal dist Uniform dist Exponential dist



variance ... are called
parameters of pop'n.

Sampling proportion of the proportion

sampling proportion (p) = $\frac{\text{No of elem of variable of interest}}{\text{sample size}}$

$$p = \frac{x}{n}$$

- used to estimate the pop's proportion π
- use when dealing with a categorical variable.

standard error of the proportion

$$\sigma_p = \sqrt{\frac{\pi(1-\pi)}{n}}$$

$$Z = \frac{p - \pi}{\sigma_p}$$

CH-8

Confidence Interval Estimation.

- is a range of vals that a probab. given parameter of a pop'n is true at a given probability. (most of the times 95%)
- The variation of sampling statistics from sample to sample is called sampling error.

Conf interval for mean (σ known)

$$\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$Z_{\alpha/2}$ - critical value
 $\alpha = 1 - \text{given prob}$
 $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ - sampling error.

Confidence interval of mean (σ unknown)

$$\bar{x} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right)$$

$v = n - 1$
then find $t_{\alpha/2}$ using v & $\alpha/2$ from t -score table.

Confidence interval for proportion.

$$p \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

Determining sample size.

$$\text{sampling error } e = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$\text{so } n = \frac{Z_{\alpha/2}^2 \sigma^2}{e^2}$$

for proportion replace σ with $\sqrt{\pi(1-\pi)}$

$$n = \frac{Z_{\alpha/2}^2 (\pi(1-\pi))}{e^2}$$

CH-9.

Fundamentals of hypothesis testing - one sample tests.

Null hypothesis (H_0) - states the status quo claim.

collectively exhaustive & mutually exclusive
refers to pop'n parameters such as μ , σ
always includes equal sign.

Alternative hypothesis (H_1)
states a claim that is contrary to the null hypothesis.

Risks

TYPE 1 Error - if you reject the null... when H_0 is true and should not be rejected.

- is a false alarm.
- its probability is α

TYPE 2 Error - if you don't reject the H_0 when H_0 is false and should be rejected.

- is missed opportunity
- its prob is β risk

- α - level of significance
- Confidence coefficient ($1 - \alpha$)
- β - risk.
- power of statistical test ($1 - \beta$)

Z-test to the mean (σ known)

step 1 - find critical values $\pm Z_{\alpha/2}$

step 2 - evaluate Z_{STAT}

$$Z_{STAT} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

step - if Z_{STAT} is within the range of critical value then accept the null hypothesis. and if it is not reject the H_0 .

t-test of the mean (σ known)

same with the above procedure but replace the formula by

$$t_{STAT} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

And used t -dist instead of Z -dist
 s - sometimes called s .
use $v = n - 1$.

One tail test

when you fall $< > \geq \leq$ probs

- if the rejection area is in the upper tail check $t_{STAT} > t_{\alpha/2}$ if false don't reject (use Z if σ is known)
- if the rejection area is in the lower tail check $t_{STAT} < -t_{\alpha/2}$ if false don't reject.

Z-Test for hypothesis for Proportion

The use of the above steps then use..

$$Z_{STAT} = \frac{p - \pi}{\sqrt{\frac{\pi(1-\pi)}{n}}} \quad \text{②}$$

CH-10

TWO SAMPLE TESTS

① pooled variance^t test for diff b/w means.

+ find critical values using $\alpha/2$ and $v = n_1 + n_2 - 2$.

+ identify H_0, H_a and rejection region.

+ evaluate

$$t_{STAT} = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$S_p^2 = \frac{(n_1 - 2)S_1^2 + (n_2 - 2)S_2^2}{(n_1 - 2) + (n_2 - 2)}$$

(pooled variance)

→ Reject H_0 if test stat falls on the rejection area.

Population consists of all the ~~things~~ individuals which you want to reach conclusion.

Variable \rightarrow a characteristic of an item or individual

Data are d/t values associated with elements

Defining and collecting data

Classifying variables by type

Numerical variable \rightarrow whose data represent a counted or measured quantity.

Categorical variable \rightarrow whose data represent categories (Eg: Gender)

sample - is a portion of a population selected for analysis.

Numerical variable

Discrete

Have data that arise from a counting process (represent a number of sth.)

Continuous

Have data that arise from a measuring process (Eg: The time spent waiting on a check-out line)

If spent timing measurements

Measurement scales

Numerical variables

Interval \rightarrow expresses a difference b/n measurements that do not include a true zero point.

Ratio \rightarrow an ordered scale that includes a scale true zero point.

\Rightarrow For both interval and ratio scales, what the d/c of 1 unit represents twice the height of of values, represents remains the same among pairs

weakest form of measurement (because you cannot specify any ranking across the various categories)

Categorical variable

nominal scale \rightarrow category values express no order or ranking

ordinal scale \rightarrow an ordering or ranking of category values is implied. (good, better, best)

A parameter is a measure that describes a characteristic of a pop.

A statistic is a measure that describes a characteristic of a sample.

Online

Parameter summarizes the value of a popn for a specific variable.
Collecting data

μ, σ, \dots

Statistic - summarizes the value of a specific variable for sample data.
 $A_1 = A_2 + 1$

A sample contains ^{only a} portion of interest

Types of sampling methods

Frame - is a complete or partial listing of items

⇒ Non-probability sample: you select items (individuals) without knowing their probabilities of selection.

- Convenience sample: you select items that are easy, inexpensive or convenient to sample.

- Judgment sample: you collect the opinions of preselected experts in the subject matter.

⇒ Probability sample: you select items based on known probabilities.

⇒ Simple random sample → every item from a frame has the same chance of selection as every other item and every sample of a fixed ~~2~~ size has the same chance of selection as every other sample of that size.

[You use "n" to represent the sample size.
"N" to represent the Frame size.]

⇒ Sampling with replacement → After you select an item you return it to the frame; where it has the prob of being selected again.

⇒ Sampling without replacement → Once you selected an item you cannot select it again.

[On the first selection $\frac{1}{N}$
On the 2nd " $\frac{1}{N-1}$]

$N=40$ $N=800$ $= 400$ partition the frame of 800 into 40
 each of which contains 20 employees.
 $k = \frac{800}{40} = 20$
 002, 0143, 060, 080, 108, 122

⇒ Systematic sample: you partition the N -items in
 the frame into n -groups of k -items, where $k = \frac{N}{n}$
 • Assign number to every pop. sample
 • Select a random number
 • Select samples at regular interval
 you round k into the next integer

⇒ Stratified sample

↳ You first subdivide the N -items in the
 frame into separate sub-populations, or strata

↳ More efficient than either single random
 & systematic sampling

⇒ Cluster sample

↳ You divide N -items in the frame into clusters
 that contain several items.

Data cleaning ~~data~~ into

⇒ Even if you follow proper procedure to collect
 data, it may contain incorrect (inconsistent data) that
 could affect statistical results. Data cleaning corrects
 such defects. & ensure your data contains suitable
 quality for your needs.

⇒ Seeks to correct the five irregularities

• Invalid Variable values, including (being incorrect by simple scanning techniques)

⇒ Non-numerical data for numerical variables

⇒ Invalid categorical values of a categorical variable

⇒ Numeric values outside a defined range

• Coding errors (Poor recording of data)

⇒ Inconsistent categorical values

⇒ Inconsistent code for categorical values

⇒ Extraneous characters

• Data integration error

⇒ Redundant columns

⇒ Duplicated rows

⇒ Differing column lengths

⇒ Different units of

measure or scale for numerical variables

sim/online

stratified members
of strata have
similar characteristics

cluster members of clusters
have d/t characteristics

Missing values → values that were not collected for a variable.

U-2

Organizing and visualizing variables

Organizing categorical variables

Summary table

↳ Helps you see the d/c among the categories by displaying the frequency, amount or percentage of items in a set of categories by displaying the frequency, amount or percentage of

Summary table

↳ helps you see

Contingency table: cross-tabulates, or tallies jointly, the data of two or more categorical variables allowing you to study patterns that may exist b/n the variables.

Frequency distribution → tallies the values of a numerical variable into a set of numerically ordered classes.

Class interval width

$$\text{Interval width} = \frac{\text{highest value} - \text{lowest value}}{\text{number of classes}}$$

Computing the proportion or Relative frequency

$$\text{Proportion} = \text{relative frequency} = \frac{\text{number of values in each class}}{\text{total number of values}}$$

The lower the CV, the less variability of random variable relative to its mean. and vice versa

Sample Variance $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

Sample standard deviation $s = \sqrt{s^2}$

The coefficient of variation (CV)

↳ measures the scatter in the data relative to the mean. (expressed always in percentage)

$CV = \left(\frac{s}{\bar{x}} \right) \times 100\%$

Z-scores

$z = \frac{x - \bar{x}}{s}$

⇒ A z-score of zero indicates that the value is the same as the mean.

⇒ If it's positive or negative, indicates whether the value is above or below the mean and by how many standard deviations.

⇒ Helps identify outliers (the values that seem excessively diff)

⇒ Z-score greater than +3.0 or less than -3.0 indicates outlier values

Shape: skewness

Skewness - measures the extent to which the data values are not symmetrical around the mean

- Mean < median, negative (left skewed dist.)
- mean = median, symmetrical dist. (zero skewness)
- Mean > median, positive (right skewed dist.)

Shape: kurtosis

Kurtosis : measures the peakedness of the curve of the dist., how sharply the curve rises approaching the center of dist.

Q9

Bell-shaped normal distribution \Rightarrow zero kurtosis. ~~at~~ ^{peak}

~~Leptokurtic~~ ^{Leptokurtic}: A dist. that has sharper rising center peak
 ~~flatter~~ ^{flatter} \Rightarrow Positive kurtosis.
 ~~tails tend to be many more values in tails~~ ^(Higher concentration of the values near the mean of the distribution)

~~Platykurtic~~ ^{Platykurtic}: A slower rising (flatter) center peak
 \Rightarrow negative kurtosis.
 ~~(Lower concentration compared to a normal distribution)~~

\Rightarrow Quartiles ($Q_1 = 25\%$ $Q_3 = 75\%$
 $Q_2 = 50\% \Rightarrow$ median)

$$Q_1 = \frac{n+1}{4} \text{ ranked value}$$

(If the ranked value
 $\cdot 75 \rightarrow$ rounded to the next whole number
 $\cdot 25 \rightarrow$ rounded to the previous whole number)

$$Q_3 = \frac{3(n+1)}{4} \text{ ranked value}$$

\Rightarrow Percentiles (split into 100 equal parts)

\Rightarrow Interquartile = $Q_3 - Q_1$
 also called \downarrow range (misread)

\Rightarrow Five number summary

X_{small} Q_1 median Q_3 X_{largest}

Left skewed
 X_{largest} to median \rightarrow median to X_{smallest}
 X_{smallest} to median \rightarrow median to X_{largest}

X_{smallest} Q_1, Q_3
 X_{smallest} to Q_1 Q_3 to X_{largest}

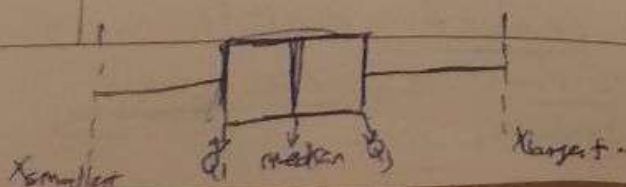
Q_1 median Q_3
 Q_1 to median Q_3 to median

Right skewed
 X_{smallest} to median \leftarrow median to X_{largest}
 X_{smallest} to median \leftarrow median to X_{largest}

X_{smallest} to Q_1 Q_3 to X_{largest}

Q_1 to median Q_3 to median

Box plot



Numerical Descriptive measures of popⁿ

Popⁿ mean $\Rightarrow \mu = \frac{\sum_{i=1}^N x_i}{N}$ *

popⁿ variance & S.D. (\rightarrow measure variation) in the popⁿ

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

The Empirical rule

• Approximately 68% of the values are within ± 1 S.D. of the mean

95% $\rightarrow \pm 2$ S.D. ^{from} the mean.

• 99.7% $\rightarrow \pm 3$ S.D. from the mean.

~~out of 20~~

\Rightarrow Implies about 1 out of 20 values will be beyond 2 standard deviations from the mean in either dirⁿ.

$\mu \pm \sigma \rightarrow 68\%$

\rightarrow values not found in the interval $\mu \pm 2\sigma$ potential outliers. $\Rightarrow 95\%$

\Rightarrow 3 in 1000 will be beyond 3 S.D. from the mean.

$\mu \pm 3\sigma \rightarrow$ values not found in the interval are almost always outliers. $\Rightarrow 99.7\%$

• Chebyshev's theorem \rightarrow For heavily skewed sets ~~that~~ ^{of} data sets that do not appear normally distributed.

$\left(1 - \frac{1}{k^2}\right) \times 100\% \subset$

\rightarrow states that for any data set, regardless of shape, the percentage of values that are found within distances of k S.D. from the ~~mean~~ ^{assumed} must be at least

chebyshev's theorem

\rightarrow General theorem

Empirical

\rightarrow If data approximately

\rightarrow It'll more accurate the greater concentration close to

$\mu \pm \sigma$	chubby at least 0%	Emp. approx. 68%
$\mu \pm 2\sigma$	at least 75%	Approx. 95%
$\mu \pm 3\sigma$	at least 88.89%	Approx. 99.7%

→ You understand how data are distributed around the mean when you have sample data

The covariance & the coefficient of correlation:

Sample

Covariance (measures the strength of the linear relationship between two numerical variables) → X and Y.

$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$$

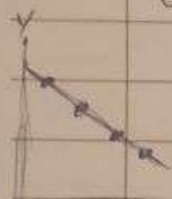
A strong correlation can be produced by chance; by the effect of a lurking variable

Coefficient of correlation

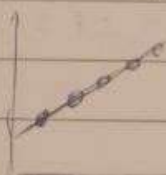
↳ measures the relative strength of linear relationship (strong or weak)

↳ Correlation alone cannot prove that there is a causation effect

(The change in one variable caused the change of the other variable)



Perfect negative (-1)
When X increases, Y decreases



Perfect positive (+1)
(X increases the Y increases)



No correlation
(X changes, Y no change)

⇒ You can say causation implies correlation, but correlation does not imply causation.

$$r = \frac{\text{Cov}(X, Y)}{S_X S_Y}$$

$$S_X = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}}$$

$$S_Y = \sqrt{\frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}}$$

Unit - 4

Basic probability

⇒ mutually exclusive: - can not occur at the same time.

⇒ collectively exhaustive → one of the events must occur.

↓
The two events cover all the possible outcomes

0.40