

Probability and Statistics

ishma hafeez
notes
rePSht

1. Walpole
2. Prem S
3. Neil A Weiss

Jishma Hafeez notes

Mid 1	15%
Mid 2	15%
Final	50%
Os/As	20%

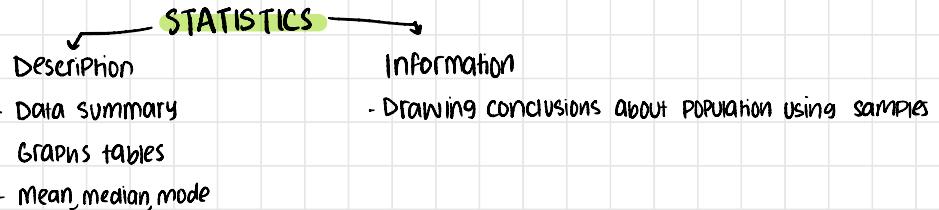
- Introduction, variables, momentum of scalars, types of data
- Mean, median, mode, quartiles, variants, standard deviation, co.efficiency of variant } Mid 1
- Histogram, Bar plot, Box plot, dot plot, frequency, polygon, stan & leaf
- Probability, counting technique, Permutation, combination
- Venn diagram, addition rule, multiplicative rule
- Conditional Probability and Bayes' rule

Statistics

The science of collection, presentation, analysis and interpretation of numerical facts/data

OR

It's the collection of procedures and principals for gathering data and analysing information to help people make decisions when faced with uncertainty



Observational Study

↳ Sample Survey

1. researchers observe characteristics
2. take measurements

Designed Experiment

1. researchers impose treatments and controls

2. then observe characteristics
3. take measures

POPULATION:

The collection of all individuals / objects under consideration in a statistical study

Parameter:

characteristics of population
mean (μ), variant (σ^2) etc

Sample:

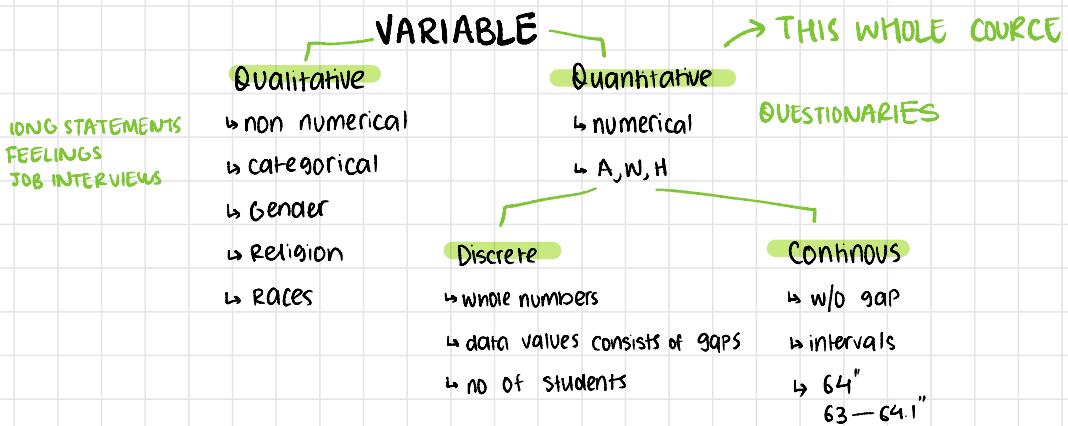
Subset of POPULATION

Statistic:

characteristic of Sample
mean (\bar{x}), variant (s^2) etc

Random sample:

It is a sample selected in such a way that every member of the population has an equal chance of being selected



DATA

PRIMARY

- ↳ raw data
- ↳ first hand information
- ↳ Survey questionnaires

SECONDARY

- ↳ already collected by someone
- ↳ Published reports of research organisations

Cross section data

- data collected on different elements at the same point in time
- ↳ Objects are different
- ↳ Time is constant

University	Annual Cost (dollars)
Harvard University	69,600
Princeton University	66,150
Stanford University	69,100
University of California, Berkeley	65,000
University of Chicago	75,715
University of Pennsylvania	71,715
Yale University	71,290
MIT	67,830

Time series data

- data collected on same element for same element at different intervals in time
- ↳ 1 variable
- ↳ different time intervals

Years	Total Population
2015	320,878,310
2016	323,015,995
2017	325,084,756
2018	327,096,265
2019	329,064,917

Panel

- ↳ different cross sections
- ↳ different time period

Qualitative

↓
String
if you can't take avg

Nominal scales

- ↳ non numerical data
- ↳ categorical data
- ↳ classification
- ↳ Gender / Blood group
- ↳ Religion
- ↳ Profession
- ↳ Zip code

MEASUREMENTS OF SCALES

Ordinal scale

- ↳ ranking
- ↳ ordered data
- ↳ customer rating
- ↳ Likert scale
- ↳ Performance of students

Strongly agree
Agree
Neutral
Disagree
Strongly disagree

Interval scale

- ↳ +, -
- ↳ absolute zero does not exist
- ↳ Temperature
- ↳ IQ scores
- ↳ variables for which '0' isn't meaningful

Quantitative

Ratio scale

- ↳ +, -, *, ÷
- ↳ absolute zero exist
- ↳ Area / Volume / Length
- ↳ Distance / Weight / Height
- ↳ Time
- ↳ Salary
- ↳ Age

Q)	Variable	Scale
	Ranking of golfers	ordinal
	temp	interval
	Weight	ratio
	Salaries	ratio
	ratings	ordinal
	categories	nominal

Frequency Distribution

↳ Histogram

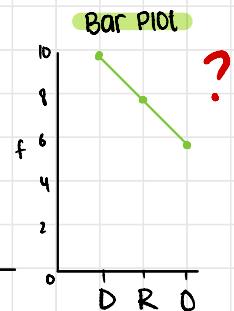
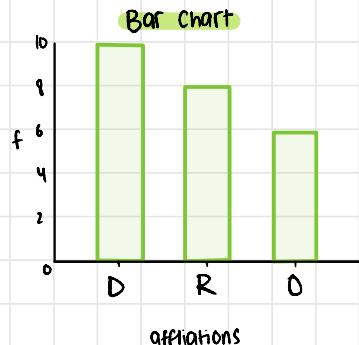
↳ Polygon

↳ Ogive

Ordered Array

↳ Stem and leaf

Affiliations	Tally	f	qualitative
D	## / ##	10	▷ R D R R
R	### / / /	8	▷ ▷ R R R
O	### /	6	○ ○ D ▷



Stem and leaf Plot

↳ leading digits

↳ trailing digits

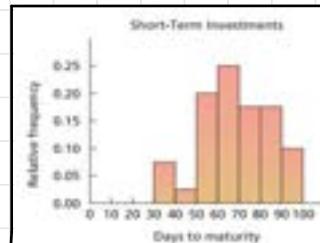
70	64	99	55	64	89	87	65
62	38	67	70	60	69	78	10
75	56	71	51	99	68	95	86
57	53	47	50	55	81	80	98
51	36	63	66	85	79	83	70

SORT IT

3	8, 6, 9
4	7
5	7, 1, 6, 3, 5, 1, 0, 5
6	2, 4, 7, 3, 6, 4, 0, 7, 5
7	0, 5, 1, 0, 9, 8
8	5, 9, 1, 7, 0, 3, 6
9	9, 9, 5, 8

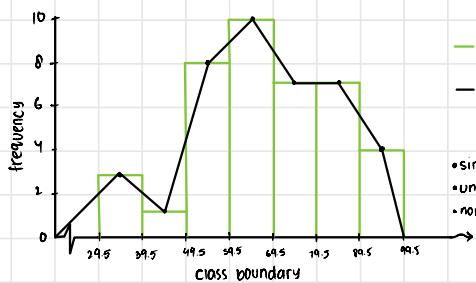
Stem leaf

3	6 8 9
4	7
5	0 1 1 3 5 5 6 7
6	0 2 3 4 4 5 6 7 7
7	0 0 1 5 8 9
8	0 1 3 5 6 7 9
9	5 8 9 9



Class interval C.I	Tally	frequency f	cumulative frequency CF(<)	CF(>)	R.F ($\frac{f}{\text{Total}}$)	-0.5 +0.5	Class boundary
30 - 39	3	3	3	40	$\frac{3}{40} = 0.075$	29.5 - 39.5	
40 - 49	1	4	4	37	$\frac{1}{40}$	39.5 - 49.5	
50 - 59	8	12	12	36	.	49.5 - 59.5	
60 - 69	10	22	22	28	.	59.5 - 69.5	
70 - 79	7	29	29	18	.	69.5 - 79.5	
80 - 89	7	36	36	11	.	79.5 - 89.5	
90 - 99	4	40	40	4	$\frac{4}{40}$	89.5 - 99.5	
Total	40						

Histogram / Polygon

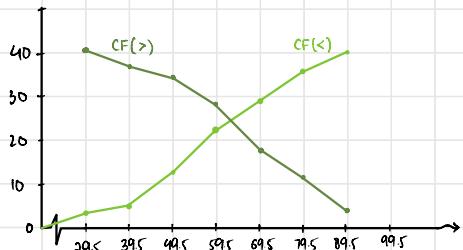


— histogram

— POLYGON
↳ midpoints

- single peak
- uni-modal
- normal

Ogive

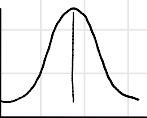


↳ CF

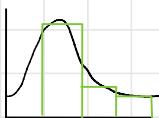
-ve skew



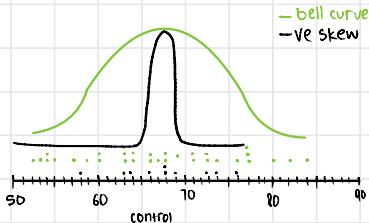
bell curve



+ve skew



Dot Plot



— bell curve
-ve skew

Intervention	Control
68	66
74	58
69	63
68	73
64	76
74	77
60	53
77	76
63	54
60	73
66	66
55	55
71	71
68	68
64	64
82	80

Pareto Charts

When the variable displayed on the horizontal axis is qualitative or categorical, a Pareto chart can often be used to represent the data.

A **Parato chart** is a bar chart representing a frequency distribution for a categorical variable, where frequencies are displayed by the heights of vertical bars, which are arranged in order from highest to lowest.

Histograms (Bar Charts)

The data shows how counts of the number of homeless people for a sample of selected cities (Census) and analysis is Pareto chart for the data.

City	Number
Atlanta	4632
Baltimore	2964
Chicago	6680
St. Louis	1485
Washington	3312

BAR GRAPH
but arranged
by their heights
in descending order

Solution

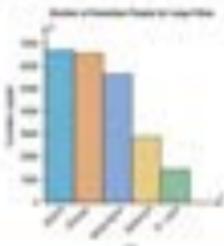
Step 1 Arrange the data from the largest to smallest according to frequency.

City	Number
Atlanta	4632
Chicago	6680
Washington	3312
Baltimore	2964
St. Louis	1485

Step 2 Draw and label the x and y axes.

Step 3 Draw the bars corresponding to the frequencies.

The graph shows that the number of homeless people is about the same for Atlanta and Chicago and a lot less for Baltimore and St. Louis.



Applications for Describing (Parato) Charts

1. Make the bars the same width.
2. Arrange the data from largest to smallest according to frequency.
3. Make the width of the bars equal for the frequency equal or one.

The Time Series Graph

When data are collected over a period of time, they can be represented by a time series graph used to study patterns.

A time series graph represents data that occur over a specific period of time.

Markups (Homeless)

The number of homicides that increased in the workplace for the years 2007 to 2009 is shown. Draw and analyse a time series graph for the data.

Year	2007	2008	2009	2010	2011	2012
Number	612	559	567	540	628	517

Source: Bureau of Justice Statistics.

Solution

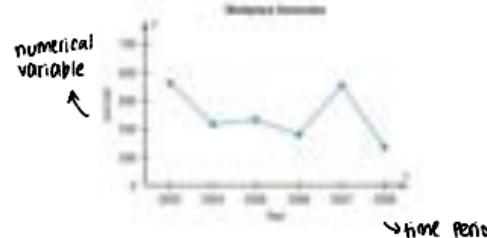
Step 1 Draw and label the x and y axes.

Step 2 Label the x-axis for years and the y-axis for the number.

Step 3 Plot each point according to the ratio.

Step 4 Draw line segments connecting adjacent points. Draw a straight line through the data points.

There was a slight decrease in the years '08, '09, and '10, compared to '07, and again an increase in '11. The largest decrease occurred in '09.



The Pie Graph

The graphs are used commonly in statistics. The purpose of the pie graph is to show the relationship of the parts to the whole by visually comparing the sizes of the sections. Percentages or proportions can be used. The variable is nominal or categorical.

A **pie graph** is a circle that is divided into sections or segments according to the percentage of frequencies in each category of the distribution.

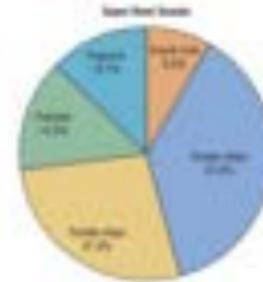
Super Bowl Snack Foods

The frequency distribution shows the number of pounds of each snack food eaten during the Super Bowl. Construct a pie chart for the data.

Snack	Pounds (frequency)	Percentages	Angles
Potato chips	11.2 million	$(11.2/30)^{*}100 = 37.3$	$(11.2/30)^{*}360 = 134^\circ$
Tortilla chips	8.2 million	$(8.2/30)^{*}100 = 27.3$	$(8.2/30)^{*}360 = 98^\circ$
Pretzels	4.3 million	14.3	52
Popcorn	3.8 million	12.7	46
Snack nuts	2.5 million	8.3	30
Total	30.0 million	100.0	360

Step 6

Now, using a protractor and a compass, draw the graph using the appropriate degree measures found in step 5, and label each section with the name and percentage.



C.I	Tally	Frequency f	Mid Point u	f <u>n</u>	c.f
30 - 39	3	$\frac{30+39}{2} = 34.5$	34.5	$f \times n = 103.5$	3
40 - 49	1	44.5			4
50 - 59	8	54.5		12	\rightarrow preceding value
60 - 69	10	64.5		22	\rightarrow closest to 20
70 - 79	7	74.5		29	
80 - 89	7	84.5		36	
90 - 99	4	94.5	37.8	40	
Total	40		2720		

✓

$$\text{mean}(\bar{u}) = \frac{\sum u}{n}$$

$$\text{coeff of variance} = \frac{\sigma}{\bar{u}} \times 100$$

$$\text{median} = \left(\frac{n+1}{2} \right)^{\text{th}}$$

Standard deviation (σ)

$$\tilde{n} = l + \frac{f}{f-f_p}$$

(L.C.B) lower class boundary
(C.F.B) upper class boundary
preceding class
cumulative frequency
of preceding class

$$\text{LCB} = 0.5$$

$$\text{UCB} = 10.5$$

$$\text{mode} = l + \frac{f_m - f_{m-1}}{2f_m - f_{m-1} - f_{m+1}} h$$

(L.C.B) lower class boundary
(M.C.B) median class
preceding class
following class

grouped \leftarrow variance \rightarrow ungrouped

$$\sigma^2 = \frac{\sum [f(u-\bar{u})^2]}{\sum f}$$

$$\sigma^2 = \frac{\sum (u-\bar{u})^2}{n} = \frac{\sum (u-\bar{u})^2}{\text{no. of observations}}$$

GROUP

Q) find mean

$$\bar{u} = \frac{2720}{40} = 68$$

Q) find mode

$$59.5 + \frac{10-8}{2(10)-8-7} \times 10$$

$$= 63.5$$

$$Q) \tilde{n} = 59.5 + \frac{10}{10} \left[\frac{40}{2} - 12 \right]$$

(L.C.B) 59.5
(U.C.B) 69
preceding class
to 32

$$= 67.5$$

Q) find standard deviation

$$\bar{u} = 68 \rightarrow \text{found above}$$

$$\sigma^2 = \frac{3(34.5-68)^2 + 1(44.5-68)^2 + 8(54.5-68)^2 + 10(64.5-68)^2 + 7(74.5-68)^2 + 7(84.5-68)^2 + 4(94.5-68)^2}{40}$$

$$= \sqrt{\frac{10510}{40}} = 16.209$$

$$68 > 67.5 > 63.5$$

\bar{x} mean

\tilde{x} median

\hat{x} mode

$\bar{x} > \tilde{x} > \hat{x}$

+ve skewed

UNGROUP

Q) $u = 2, 4, 6, 8$, find mean

$$\bar{u} = \frac{2+4+6+8}{4} = 5$$

\rightarrow No mode as all values occur exactly once

Q) $1, 2, 3$, find medium

$$= \frac{3+1}{2} = 2^{\text{nd}} \text{ position} \rightarrow \text{has 2}$$

Q) $2, 4, 6, 8$, find medium

$$\frac{4+1}{2} = 2.5^{\text{th}} \rightarrow \frac{2+3^{\text{rd}}}{2} \rightarrow \frac{4+6}{2} = 5$$

Q) $2, 3, 5, 8$, find variance

$$\bar{u} = \frac{18}{4} = 4.5$$

$$\sigma^2 = \frac{(2-4.5)^2 + (3-4.5)^2 + (5-4.5)^2 + (8-4.5)^2}{4}$$

$$= \sqrt{5.25}$$

$$= 2.3$$

APPLICATIONS

3.49 The following table gives the frequency distribution of the number of hours spent last week on cell phones (making phone calls and texting) by all 100 students in the tenth grade at a school.

Hours	Number of Students	f	n	fn	cf
0 to less than 4	14	2	28	14	
4 to less than 8	18	6	108	32	
8 to less than 12	24	10	250	51	
12 to less than 16	18	14	252	75	
16 to less than 20	16	18	288	91	
20 to less than 24	9	22	198	100	
				1124	

Find the mean, variance, and standard deviation.

$$\mu = \frac{1124}{100} = 11.24$$

$$mde = 8.5 \left[\frac{25 - 18}{2(25) - 18 - 18} \right] \times 4$$

$$\sigma^2 =$$

Empirical Rule
68-95-99.7 rule
Normal Distribution

$\bar{x} \pm s \rightarrow$ Approx. 68% observations

$\bar{x} \pm 2s \rightarrow$ Approx. 95% observations

$\bar{x} \pm 3s \rightarrow$ Approx. 99.7% observations

Consider the following sample of exam scores, arranged in increasing order. The sample mean and sample standard deviation of these exam scores were, respectively, 85 & 16.1.

86	87	88	88	89	90
89	90	91	91	92	93
91	92	93	93	94	95
92	93	94	94	95	96
93	94	95	95	96	97
94	95	96	96	97	98
95	96	97	97	98	99

Use the data to obtain the exact percentage of observations that lie within two standard deviations to either side of the mean. Compare your answer.

Empirical rule is fixed if the data is normally distributed or symmetric.

$\bar{x} + s \text{ to } \bar{x} + 3s \rightarrow$ 68% observations will fall in this interval

$\bar{x} + 2s \text{ to } \bar{x} + 2s \rightarrow$ 95% observations will fall in this interval

$\bar{x} + 3s \text{ to } \bar{x} + 3s \rightarrow$ 99.7% observations will fall in this interval.

[use this if you are asked to check the upper empirical rule otherwise not needed]

Rule	Lower to Upper	Frequency
85 ± 16.1	68.9 to 101.1	
$85 \pm 2 \cdot 16.1$	52.8 to 117.2	
$85 \pm 3 \cdot 16.1$	36.7 to 133.3	

Quartiles

5, 15, 16, 20, 21, 25, 26, 27, 30, 30
 31, 32, 32, 34, 35, 38, 38, 41, 43, 66 → outlier outside max

data set represent no. of observations

$$25\%: Q_1 = \left(\frac{n+1}{4} \right)^{th} \text{ lower data set} \quad \left(\frac{20+1}{4} \right)^{th} \Rightarrow 5.25^{th} \rightarrow \frac{5^m + 6^m}{2} = 23$$

$$50\%: Q_2 = \left(\frac{2(n+1)}{4} \right)^{th} \text{ median} \quad 2(5.25) \Rightarrow 10.5^{th} \rightarrow \frac{10^m + 11^m}{2} = 30.5$$

$$75\%: Q_3 = \left(\frac{3(n+1)}{4} \right)^{th} \text{ upper data set} \quad 3(5.25) \Rightarrow 15.75 \rightarrow \frac{15^m + 16^m}{2} = 36.5$$

Box and Whiskers Plot

↳ Data shape (skew)

$$IQR = Q_3 - Q_1$$

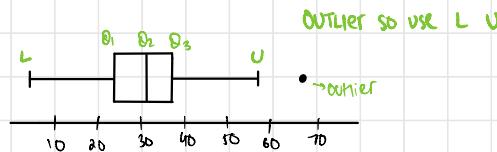
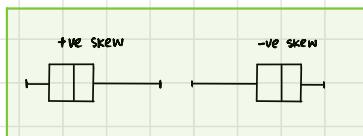
↳ Outliers

Min ↓
L ↓
U ↓ Max

↳ Variation (Box length)(IQR)

$$Q_1 - 1.5(IQR) = 2.75$$

$$56.75 = Q_3 + 1.5(IQR)$$

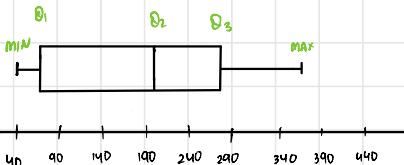


NO outlier so use max min

1 2 3 4 5 6 7 8

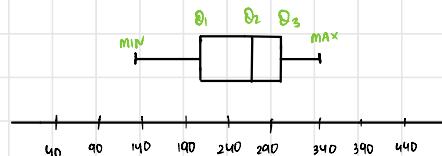
Real cheese $x_i \rightarrow 40, 45, 90, 180, 220, 240, 310, 420$

2.25	4.50	6.75	no outliers		
Q ₁	Q ₂	Q ₃	IQR	L	U
67.5	200	275	207.5	-243.75	586.25
215	265	300	85	87.5	431.25



REAL CHEESE +ve skew

has more sodium content
greater variation



CHEESE SUBSTITUTE -ve skew

Probability

Statistical Experiment

Any activity/process that generates data

↳ Tossing a coin

↳ Rolling a dice

↳ Playing cards

↳ Survey/opinion of voters

Sample Space

List of all possible outcomes of a statistical experiment

↳ Tossing a coin 2 times: { HH, HT, TH, TT }

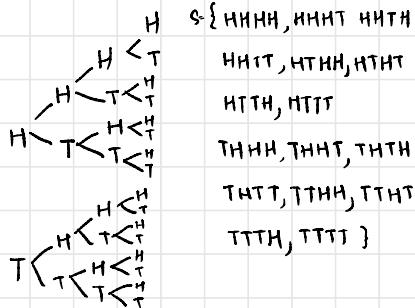
↳ Rolling a device: { 1, 2, 3, 4, 5, 6 }

event: subset of a sample space

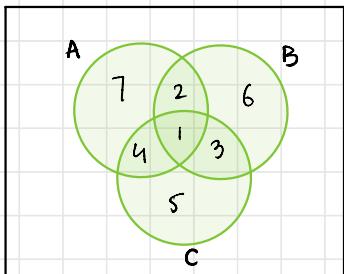
* Toss a coin 4 times / 4 coins tossed together

$$\text{possibilities} = 2^4 = 16$$

list of elements



VENN DIAGRAM



$$\begin{aligned}A \cap B &= 1, 2 \\B \cap C &= 1, 3 \\A \cup C &= 1, 2, 3, 4, 5, 7 \\B' \cap A &= 4, 7 \\A \cap B \cap C &= 1 \\(A \cup B) \cap C' &= 7, 2, 6\end{aligned}$$

Permutations

↳ order matters

↳ organise

↳ ways

$${}^n P_r = \frac{n!}{(n-r)!}$$

Combinations

↳ order doesn't matter

↳ different types

↳ select

$${}^n C_r = \frac{n!}{(n-r)!r!}$$

Empirical
Formulae

Probability

It is a measure of chance that an uncertain event will occur

Subjective

↳ personal experience

↳ judgement

Objective

↳ classical approach

$P(\text{Event} = A) = \frac{\text{favourable case of } A}{\text{all possible cases}}$

$$\leq \text{Prob} = 1, 0 \leq P(A) \leq 1$$

- 1) A coin is tossed twice. Prob atleast 1 head

$$S = \{HH, TH, HT, TT\} \rightarrow 3/4$$

$$Q) 10 \quad E \quad 10 \quad M$$

$$25 \quad I \quad 8 \quad C$$

- 3) A dice drawn such that an even no is twice likely to occur than odd

$$\left\{ \begin{array}{ccccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ w & +2w & +w & +2w & +w & +2w \end{array} \right\}$$

$$P(X < 4) = \frac{1}{9} + \frac{2}{9} + \frac{1}{9} = \frac{4}{9}$$

$$a) P(I) = 25/53$$

$$b) P(C \text{ or } E) = P(C) + P(E) = \frac{8}{53} + \frac{10}{53}$$

- Q) Dice rolled. Sum 7 or 11 → exclusive

$$P(7) = \frac{16}{36}, \frac{25}{36}, \frac{36}{36} \quad P(11) = \frac{6}{36}, \frac{5}{36}$$

- 5) A card is drawn at random from 52 cards

$$a) P(I) = \frac{4C_1}{52C_1} = \frac{4}{52}$$

$$P(7 \text{ or } 11) = \frac{6}{36} + \frac{2}{36}$$

- b) In a poker hand consisting of 3 cards
Prob of 2 Aces and 3 Jacks

$$\frac{\binom{4}{2} \binom{4}{3}}{\binom{52}{5}} =$$

2 Aces 3 Jacks
 Total + 52 C 5
 > need
 5 cards

OR has + → only one event from choice

MUTUALLY EXCLUSIVE EVENTS

↳ nothing in common

$$P(A \cup B) = P(A) + P(B)$$

$$P(A \cap B) = \emptyset$$

AND has * → more than 1 events

NON MUTUALLY EXCLUSIVE EVENTS

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap C) - P(A \cap B) - P(B \cap C) + P(A \cap B \cap C)$$

Independent events

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

if $P(A|B) = P(A)$ then independent

else dependent

Dependent events

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

BAYES RULE

$$P(A|B) = \frac{P(B|A_i) P(A_i)}{\sum_{i=1}^n P(B|A_i) P(A_i)}$$

Conditional Probability

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

A will occur
given that B already

$$\begin{aligned} &\rightarrow P(A|B) = \frac{P(B|A) P(A)}{P(B)} \\ &\hookrightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(B|A) = \frac{P(A \cap B)}{P(A)} \\ &\hookrightarrow P(A|B) P(B) = P(B|A) P(A) \end{aligned}$$

CONDITIONAL PROBABILITY

Q)

	non smokers	moderate smokers	heavy smokers	Total
H	21	36	30	87
NH	48	26	19	93
Total	69	62	49	180

a) $P(H/H_3) = \frac{P(H \cap H_3)}{P(H_3)} = \frac{30}{49}$

b) $P(NH/NH) = \frac{P(NH \cap NH)}{P(NH)} = \frac{48}{93}$

Q) $P(\text{of choosing } M) \quad P(\text{of } D \text{ of } M) \quad P(\text{of it being } M \text{ and } D)$

$P(M_1) = 0.30 \times P(D/M_1) = 0.02 = 0.006$

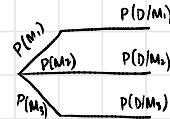
$P(M_2) = 0.45 \times P(D/M_2) = 0.03 = 0.0135$

$P(M_3) = 0.25 \times P(D/M_3) = 0.02 = 0.005$

$P(D) = ? \quad 0.0245$

$$\begin{aligned} P(D) &= P(D \cap M_1) + P(D \cap M_2) + P(D \cap M_3) \\ &= P(M_1)P(D|M_1) + P(M_2)P(D|M_2) + P(M_3)P(D|M_3) \\ &= 0.0245 \end{aligned}$$

$$P(M_1|D) = \frac{P(M_1 \cap D)}{P(D)} = \frac{0.006}{0.0245}$$



1) Station	A	B	C	Total
Electric	2	1	1	
Computers	4	3	2	
Equipment	5	4	2	
Human Error	7	7	5	
Total				

2)	C	D	
Male	10	2	
Female	5	3	
		20	

P(diploma holder is a male)

$$P(M/D) = \frac{P(M \cap D)}{P(D)} = \frac{2}{5}$$

3)	1	2	3	Total
G	1	3	4	8
Y	2	2	3	2
R	4	5	1	10
Total	7	10	8	25

$$a) P(Y) = P(Y \cap B_1) + (Y \cap B_2) + P(Y \cap B_3)$$

$$\frac{1}{3} \times \frac{2}{7} + \frac{1}{3} \times \frac{2}{10} + \frac{1}{3} \times \frac{3}{8} = \frac{241}{840}$$

$$b) P(B_2 \cap Y) = \frac{P(B_2 \cap Y)}{P(Y)}$$

$$\frac{\frac{1}{3} \times \frac{2}{10}}{\frac{241}{840}}$$

- 4) $M_1, M_2, M_3 \rightarrow$ chances M hits
 $\frac{1}{6}, \frac{1}{4}, \frac{1}{3}$

P(target was hit by M_1) =

$$P(M_1/H) = \frac{P(M_1 \cap H)}{P(H)} = \frac{\frac{1}{3} \times \frac{1}{6}}{\frac{1}{3}(\frac{1}{6}) + \frac{1}{3}(\frac{1}{4}) + \frac{1}{3}(\frac{1}{3})} = \frac{2}{9} ?$$

- 5) $H_1, H_2, H_3 \rightarrow$ chances
 $\frac{1}{6}, \frac{1}{4}, \frac{1}{3}$

$$P(M_1) = \frac{1}{6}(1 - \frac{2}{6})(1 - \frac{1}{3})$$

$$P(M_2) = (1 - \frac{1}{6})^2 \cdot \frac{1}{4} (1 - \frac{1}{3})$$

$$P(M_3) = (1 - \frac{1}{6})(1 - \frac{2}{4}) \cdot \frac{1}{3}$$

P(target was hit by M_1) =

$$P() =$$

a) How many diff ways can student check of

1 ans to each question. 5 questions 4 option

$$4^5 =$$

b) Get all wrong

$$3^5 \rightarrow 5 \text{ questions}$$

$$3^5 \rightarrow 3 \text{ wrongs}$$

2.37) 4B, 5G

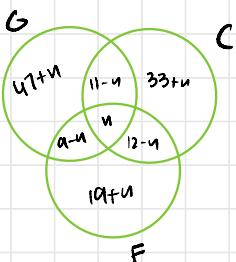
if B, G must alternate

$$\frac{5 \times 4 \times 4 \times 3 \times 3 \times 2 \times 2 \times 1 \times 1}{6 \ B \ G \ B \ G \ B \ G \ B} = 2880$$

8) How many ways no 2 students have same bday. 60 students

$$\text{days} \leftarrow 365 P_{60} \text{ students}$$

9) In a survey, B=136, G=67, C=56, F=40, G\cap C=11, C\cap F=12, G\cap F=9, each boy played atleast
How many played all games



0, 1 2 3 4 5 6 NO REPETITION, in 3 DIGITS

a) $\boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad}$
 $6 \times 6 \times 5 = 180$

for 0

b) How many of these are odd numbers

$$5 \times 5 \times 3 = 75$$

c) How many > 330

$$\begin{array}{r} 3 \times 6 \times 5 + 1 \times 3 \times 5 = 105 \\ 405 \\ \hline 330 \end{array}$$

a) How many no each of 5 digits constructed

$$0 \quad \underline{4 \ 5} \quad 6 \ 7$$

4 and 5 next to each other

$$3 \times 3! \times 2! = 36$$

b) Not to be together 4 and 5

$$5 \text{ digits. } 4 \times 4! = 96$$

$$96 - 36 = 60$$

APPENDIX D

- 2.11 A local gas station collected data from the day's receipts, recording the gallons of gasoline each customer purchased. The following table lists the frequency distribution of the gallons of gas purchased by all customers on this one day at this gas station.

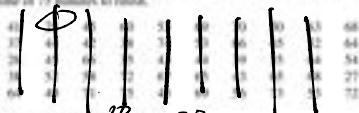
Gallons of Gas	Number of Customers	mid point \bar{x}	R.f.	P.D.
0 to less than 4	31	2	0.084	0.4
4 to less than 8	78	6	0.211	21.1
8 to less than 12	49	10	0.133	13.3
12 to less than 16	81	14	0.226	22.6
16 to less than 20	137	18	0.311	31.1
20 to less than 24	13	22	0.035	3.5
	369			

369

- How many customers were served on this day at this gas station?
- Find the class midpoints. Do all of the classes have the same width? If so, what is this width? If not, what are the different class widths? all have same widths $n=4$
- Prepare the relative frequency and percentage distribution columns.
- What percentage of the customers purchased 12 gallons or more? 57.2
- Explain why you cannot determine exactly how many customers purchased 10 gallons or less not specified
- Prepare the cumulative frequency, cumulative relative frequency, and cumulative percentage distributions using the given table.

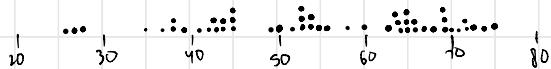
Make a dotplot for these data.

- 2.25 The following data give the times (in minutes) taken by 50 students to complete a statistics examination that was given a maximum time of 75 minutes to finish.



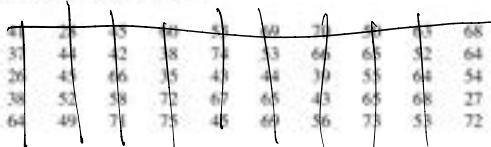
Create a dotplot for these data.

- 2.26 The following data give the cumulative examination times (in



each stem in increasing order.

- 2.27 The following data give the times (in minutes) taken by 50 students to complete a statistics examination that was given a maximum time of 75 minutes to finish.



Use
for
2.31
book
unit

EMPERICAL?
SETS DIGRAM?
PERMUTATIONS for strings?

- Prepare a stem-and-leaf display for these data. Arrange the leaves for each stem in increasing order.
- Prepare a split stem-and-leaf display for the data. Split each stem into two parts. The first part should contain the leaves 0,

2 8 6 7
3 7 8 8 5 9
4 1 5 4 5 9 2 3 5 4
5 3 0 2 8 3 3 6 5 2 3 4
6 0 9 3 8 4 6 1 5 9 6 5 6 4 8 4
7 0 1 2 5 4 3 2

ishma hafeez
notes

repsht
tree

RANDOM VARIABLE (RV)

A RV is a function that associates a real no. with each element in the sample space

Discrete random variable: Possible outcomes countable

Continuous random variable: Values on a continuous scale

Q) Let $X = \text{No. of Heads in toss of 3 coins}$

x	$2^3 = 8$	$P(x)$	$F(x)$
0	$\leftarrow \text{TTT}$	$\frac{1}{8}$	$\frac{1}{8}$
1	$\leftarrow \text{HTT, THT, TTH}$	$\frac{3}{8}$	$\frac{4}{8}$
2	$\leftarrow \text{HHT, THH, HTH}$	$\frac{3}{8}$	$\frac{7}{8}$
3	$\leftarrow \text{HHH}$	$\frac{1}{8}$	$\frac{8}{8}$

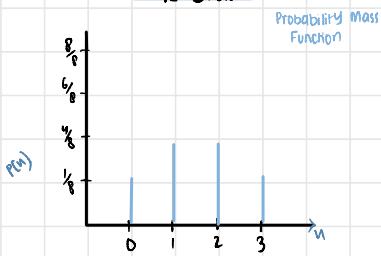
Cumulative Distribution Function (CDF)

$$F(x) = P(X \leq n)$$

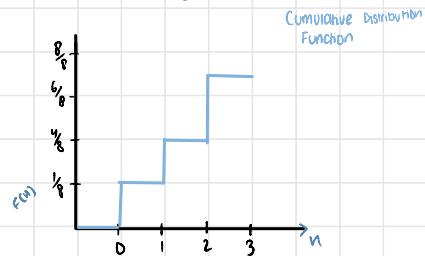
$$F(x) = \begin{cases} 0 & n < 0 \\ \frac{1}{8} & 0 \leq n < \frac{1}{8} \\ \frac{4}{8} & \frac{1}{8} \leq n < \frac{4}{8} \\ \frac{7}{8} & \frac{4}{8} \leq n < \frac{7}{8} \\ 1 & \frac{7}{8} \leq n < 1 \end{cases}$$

$P_{\leq 0} = P(n=0)$
 $P_{(0, \frac{1}{8})} = P(n=1)$
 $P_{(\frac{1}{8}, \frac{4}{8})} = P(n=2)$
 $P_{(\frac{4}{8}, \frac{7}{8})} = P(n=3)$
 $P_{\geq 1} = P(n=0) + P(n=1) + P(n=2) + P(n=3)$

Line Graph \rightarrow PMF \rightarrow uses $P(x)$



Line Graph \rightarrow CDF \rightarrow uses $F(x)$



Q) Let $X = \text{No. of Heads in toss of 4 coins}$

x	$2^4 = 16$	$P(x)$	$F(x)$
0	$\leftarrow \text{TTTT}$	$\frac{1}{16}$	$\frac{1}{16}$
1	$\leftarrow \text{HTTT, THTT, TTHT, TTTT}$	$\frac{4}{16}$	$\frac{5}{16}$
2	$\leftarrow \text{HHTT, THHH, HTHT, THHT}$	$\frac{6}{16}$	$\frac{9}{16}$
3	$\leftarrow \text{THHH, HTHH, HHTH, HHHT}$	$\frac{4}{16}$	$\frac{13}{16}$
4	$\leftarrow \text{HHHH}$	$\frac{1}{16}$	$\frac{14}{16}$

Cumulative Distribution Function (CDF)

$$\hookrightarrow F(n) = P(X \leq n)$$

for continuous RV

$$\rightarrow \int_{-\infty}^n f(t) dt$$

$$\hookrightarrow F(x) = P(X \leq x) = \sum P(x)$$

$$\hookrightarrow P(a < x < b) = F(b) - F(a)$$

$$f(n) = \frac{dF(n)}{dn}$$

$$F(n) = \sum_{n=1}^{\infty} (n+1)^{-n}$$

EXAMPLE 1

* Consider the following pmf: $f(x) = (x/6)$, $x = 1, 2, 3$, zero elsewhere.

- (i) Find distribution function and its graph.
- (ii) Calculate $P(1.5 < x \leq 4.5)$.

$$i) f(n) = \frac{n}{6}$$

n	$f(n)$	$F(n)$
1	$\frac{1}{6}$	$\frac{1}{6}$
2	$\frac{2}{6}$	$\frac{3}{6}$
3	$\frac{3}{6}$	$\frac{6}{6}$

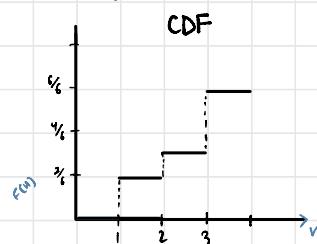
→ Distribution function

$F(n)$	0	$n < 1$
$\frac{1}{6}$	$1 \leq n < 2$	
$\frac{3}{6}$	$2 \leq n < 3$	
$\frac{6}{6}$	$n \geq 3$	

$$ii) P(1.5 < n \leq 4.5)$$

$$\begin{aligned} & \cdot F(4.5) - F(1.5) \\ & \cdot n \geq 3 - (1 \leq n < 2) \\ & \cdot 1 - \frac{1}{6} \\ & \cdot \frac{5}{6} \end{aligned}$$

graph



EXAMPLE 3

(3) Consider the following functions.

- (i) $f(x) = (x+2) / 5$ for $x = 1, 2, 3, 4, 5$.
- (ii) $f(x) = (4Cx) / (2^5)$, for $x = 0, 1, 2, 3$, and 4 ,

and check whether the functions can serve as a pmf?

$$i) f(n) = \frac{(n+2)}{5}$$

n	$f(n)$
1	$\frac{3}{5}$
2	$\frac{4}{5}$
3	1
4	$\frac{6}{5}$
5	$\frac{7}{5}$

Since $0 \leq n \leq 1$ → what

$$\leq f(x) \neq 1$$

So not PMF

$$\leq f(x) > 5$$

$$ii) f(n) = \frac{4Cx}{2^5}$$

n	$f(n)$
1	$\frac{1}{32}$
2	$\frac{1}{8}$
3	$\frac{3}{16}$
4	$\frac{1}{8}$
5	$\frac{1}{32}$

not PMF as
 $\leq f(n) \neq 1$

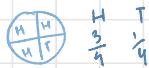
$$\leq f(n) = \frac{1}{2}$$

Example 4

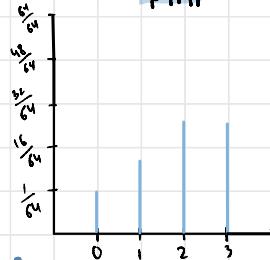
(4) A coin is biased so that a head occurs 3 times of tail. If the coin is tossed 3 times, find the probability distribution for the number of heads and also find $P(1 \leq n \leq 3)$.

$$\begin{aligned} 0 & \left\{ \text{TTT} \quad \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \right\} \\ 1 & \left\{ \begin{array}{l} \text{HTT} \quad \frac{3}{4} \times \frac{1}{4} \times \frac{1}{4} \times 3 \\ \text{THT} \\ \text{THH} \end{array} \right. \\ 2 & \left\{ \begin{array}{l} \text{HHT} \quad \frac{3}{4} \times \frac{3}{4} \times \frac{1}{4} \times 3 \\ \text{HTH} \\ \text{TTH} \end{array} \right. \\ 3 & \left\{ \text{HHH} \quad \frac{3}{4} \times \frac{3}{4} \times \frac{3}{4} \right\} \end{aligned}$$

X	$P(n)$
0	$\frac{1}{64}$
1	$\frac{9}{64}$
2	$\frac{27}{64}$
3	$\frac{27}{64}$



PMF



$$\begin{aligned} P(1 \leq n \leq 3) &= \frac{9}{64} + \frac{27}{64} + \frac{27}{64} \\ &= \frac{63}{64} \end{aligned}$$

Example 7

The distribution function for a discrete random variable x is given as:
 $F(x) = 1 - (1/2)^x (x+1)$, for $x = 0, 1, 2, \dots$

Find:

- $P(X=3)$
- $P(7 \leq x < 10)$
- Probability Mass Function.

$$f(u) = 1 - \left(\frac{1}{2}\right)^{u+1}$$

$$a) P(u=3) = F(3) - F(2)$$

$$\therefore P(X=x) = F(x) - F(x-1) = \left[1 - \left(\frac{1}{2}\right)^{3+1} \right] - \left[1 - \left(\frac{1}{2}\right)^{2+1} \right] = \frac{1}{16}$$

$$b) P(7 \leq u < 10) = F(9) - F(6) \quad \text{w.r.t } F(10) - F(6) \\ = \frac{1}{1024} \quad = \frac{15}{2048}$$

c) PMF

$$\begin{aligned} & F(u) - F(u-1) \\ & = 1 - \frac{1}{2}^{u+1} - 1 + \frac{1}{2}^{u+1} \\ & = \left(-\frac{1}{2}\right)^{u+1} + \left(\frac{1}{2}\right)^{u+1} \\ & = \frac{1}{2}^{u+1} \rightarrow \text{w.r.t } \end{aligned}$$

Example 8

3.5 Determine the value c so that each of the following functions can serve as a probability distribution of the discrete random variable X:

$$(a) f(x) = c(x^2 + 4), \text{ for } x = 0, 1, 2, 3$$

$$a) f(u) = c(u^2 + u)$$

$$\sum_{n=0}^3 f(u) = 1$$

$$c(0) + c(1) + c(2) + c(3) = 1$$

$$30c = 1$$

$$c = \frac{1}{30}$$

Q) find $P(n=2) = f(n=2)$ using CDF

$$P(2) = F(2) - F(1)$$

$$= P(n \leq 2) - P(n \leq 1)$$

$$= \frac{1}{16} - \frac{5}{16}$$

$$= \frac{3}{16}$$

Q1) A fair coin is tossed until H appears for the first time. Find

a) PMF: $(\frac{1}{2})^n \quad n = 1, 2, 3, \dots$

b) CDF: $\sum_{n=1}^{n=\infty} (\frac{1}{2})^n$

c) $F(4) = \sum_{n=1}^{n=4} (\frac{1}{2})^n$

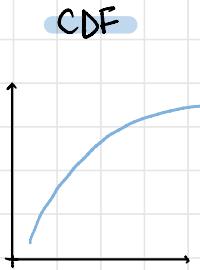
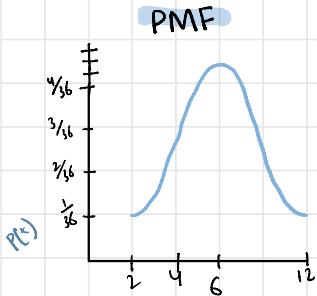
Q) If two dice are rolled once, find the PMF of the sum of points on dice, CDF and their graph

$$\text{let } t = x+y$$

z	2	3	4	5	6	7	8	9	...	12
$f(t) = P(t)$:	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{7}{36}$	$\frac{8}{36}$	

$$F(t) = \frac{1}{36}, \frac{3}{36}, \frac{6}{36}, \dots$$

$$F(t) = \begin{cases} 0 & z < 2 \\ \frac{1}{36} & 2 \leq z < 3 \\ \frac{3}{36} & 3 \leq z < 4 \\ \vdots & \vdots \\ \frac{36}{36} & z \geq 12 \end{cases}$$



* Expected Value

Let suppose a coin tossed 2 times

X = Head calculate

$E(X)$ = Mass of a row X

X	$P(X)$	$X \cdot P(X)$	$X^2 \cdot P(X)$
0	$\frac{1}{4}$	0	0
1	$\frac{3}{4}$	$\frac{3}{4}$	$\frac{3}{4}$
2	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{4}{4}$
Sum	1	$E(X) = 1$	$E(X^2) = \frac{3}{2}$

$$E(x) = \sum x \cdot P(x) \rightarrow \text{Expected value}$$

$$\sigma^2 = V(X) = E(X^2) - [E(X)]^2 \rightarrow \text{Variance}$$

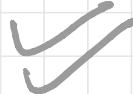
Joint Probability Distribution

$$1. f(u, y) \geq 0$$

$$2. \sum_{u} \sum_y f(u, y) = 1$$

$$3. P(X=u, Y=y) = f(u, y)$$

$$\hookrightarrow P[(X, Y) \in A] = \sum_u \sum_y f(u, y)$$



EXAMPLE 9

Two ballpoint pens are selected at random from a box that contains 8 blue pens, 2 red pens, and 3 green pens. If X is the number of blue pens selected and Y is the number of red pens selected, find:
 (a) the joint probability function $f(x, y)$
 (b) $P[X \leq 1, Y \leq 1]$, where it is the region $(x, y) | x \leq 1, y \leq 1$

→ WHAT?

$$B=3, R=2, G=3 \Rightarrow 8 \rightarrow \text{select 2 w/o replacement}$$

JOINT PROB MASS FUNCTION

		pred			marginal dist $g(u)$		
		0	1	2	$g(u)$	$xg(x)$	$x^2 g(x)$
Blue	0	$\frac{C_0 C_0}{C_0 C_0} = \frac{3}{28}$	$\frac{C_0 C_1}{C_0 C_1} = \frac{6}{28}$	$\frac{C_0 C_2}{C_0 C_1} = \frac{3}{28}$	$\frac{10}{28}$	0	0
	1	$\frac{C_1 C_0}{C_0 C_1} = \frac{6}{28}$	$\frac{C_1 C_1}{C_0 C_1} = \frac{6}{28}$	—	$\frac{15}{28}$	$\frac{15}{28}$	$\frac{15}{28}$
	2	$\frac{C_2 C_0}{C_0 C_1} = \frac{3}{28}$	—	—	$\frac{3}{28}$	$\frac{6}{28}$	$\frac{12}{28}$
$m(y)$		$\frac{15}{28}$	$\frac{12}{28}$	$\frac{1}{28}$	1	$\frac{21}{28}$	$\frac{42}{28}$
$y g(y)$		0	$\frac{12}{28}$	$\frac{1}{28}$	$E(Y)$	$E(Y^2)$	
marginal dist of y		0	$\frac{12}{28}$	$\frac{1}{28}$	$E(Y)$	$E(Y^2)$	
$y g(y)$		0	$\frac{12}{28}$	$\frac{1}{28}$	$E(Y)$	$E(Y^2)$	

$$E(x) = 6^2 x - V(x) = E(x^2) - [E(x)]^2 \rightarrow \text{Variance}$$

$$E(\sqrt{x}) = \sqrt{\frac{21}{28}} = \sqrt{\frac{21}{28}}$$

$$E[3x] = 3E(x)$$

$$a) P[X \leq 2, Y=1] = f(0,1) + f(1,1) + f(2,1) \\ = \frac{6}{28} + \frac{6}{28} + \frac{3}{28} \\ = \frac{15}{28}$$

$$b) P[X \geq 2, Y \leq 1] = f(2,1) + f(2,0)$$

$$c) P[X > Y] = f(1,0) + f(2,0)$$

MATHEMATICAL EXPECTATION

$$\text{variance } V(x) = E(x^2) - [E(x)]^2$$

$$V(y) = E(y^2) - [E(y)]^2$$

$$E(xy) = \sum xy f(x, y)$$

$$\text{Covariance}(x, y) = E(xy) - E(x) E(y)$$

$$\text{Correlation}(x, y) = \frac{\text{Covariance}(x, y)}{\sqrt{V(x)} \sqrt{V(y)}}$$

mean

$$\mu = E(x) = \sum x f(x)$$

$\pm 1 \rightarrow \text{strong}$

$\pm 0.5 \rightarrow \text{moderate}$

0 → negligible

$$d) P[X+Y=4] = 0$$

conditional probability

$$e) P[X=0 | Y=1] = \frac{f(u=0, y=1)}{h(y=1)}$$

$$= \frac{6/28}{12/28} = \frac{6}{12}$$

$$f) P[Y=1 | X=0] = \frac{f(y=1, x=0)}{g(x=0)} = \frac{6/28}{10/28} = \frac{6/28}{10/28}$$

$$g) E[XY] = \sum \sum nyf(u, y) = \frac{6}{28}$$

0 0 0	\times	$\frac{3}{28}$
0 1 0	\times	$\frac{6}{28}$
0 2 0	\times	$\frac{3}{28}$
1 0 0	\times	$\frac{6}{28}$
1 1 1	\times	$\frac{6}{28}$
1 2 1	\times	$\frac{3}{28}$
2 0 0	\times	$\frac{3}{28}$
2 1 0	\times	$\frac{6}{28}$
2 2 1	\times	$\frac{3}{28}$

$$\frac{6}{28}$$

$$h) \text{covariance}(x, y) = E(XY) - E(X) E(Y)$$

$$= \frac{6}{28} - \left(\frac{21}{28}\right)\left(\frac{15}{28}\right) \\ = -\frac{9}{56}$$

$$-1 \leq \text{correlation}(x, y) = \frac{\text{covariance}(x, y)}{\sqrt{V(x)} \sqrt{V(y)}} \leq 1$$

$$= \frac{-9/56}{\sqrt{143}/\sqrt{112} \sqrt{143}/\sqrt{28}} \\ = -0.441$$

moderate → negative relationship

JPMF

$$f(u, v) = \frac{u+v}{30} \quad \text{for } u=0, 1, 2, 3 \\ v=0, 1, 2$$

Calculated

$$a) P[X \leq 2, Y=1] = F(0,1) + F(1,1) + F(2,1) \\ = \frac{1}{30} + \frac{2}{30} + \frac{3}{30} = \frac{6}{30}$$

$$b) P[X > 2, Y \leq 1] = F(3,1) + F(3,0) \\ = \frac{4}{30} + \frac{3}{30} = \frac{7}{30}$$

$$c) P[X > Y] = F(3,2) + F(3,1) + F(3,0) + F(2,2) + F(2,1) + F(1,0) \\ = \frac{5}{30} + \frac{4}{30} + \frac{3}{30} + \frac{3}{30} + \frac{2}{30} + \frac{1}{30} = \frac{18}{30}$$

$$d) P[X+Y=4] = F(3,1) + F(2,2) \\ = \frac{4}{30} + \frac{4}{30} = \frac{8}{30}$$

e) Calculate

$$E(x) = \frac{60}{30} = 2 \\ E(x^2) = \frac{150}{30} = 5$$

$$(X-\mu)^2 f(x) \\ (X-2)^2 f(x)$$

$$V(x) = E(x^2) - [E(x)]^2 = 5 - (2)^2 = 1$$

$$V(y) = E(y^2) - [E(y)]^2 = \frac{66}{30} - \left(\frac{38}{30}\right)^2 = 0.59$$

$$E(xy) = \sum \sum xy f(x,y) = 2.4 \longrightarrow E(xy) = \sum \sum xy f(x,y) = 2.4$$

$$\text{Covariance}(x,y) = E(xy) - E(x)E(y) = 2.4 - 2\left(\frac{38}{30}\right) \rightarrow -\frac{2}{15} \rightarrow -0.133$$

$$\text{Correlation}(x,y) = \frac{\text{Covariance}(x,y)}{\sqrt{V(x)} \sqrt{V(y)}} = -0.17277 \\ \hookrightarrow \text{weak relation} \\ \hookrightarrow \text{inverse relation as -ve sign}$$

x \ y	0	1	2	$E(x)$	$E(x^2)$	$E(y)$
0	0	$\frac{1}{30}$	$\frac{2}{30}$	$\frac{3}{30}$	$\frac{4}{30}$	0
1	$\frac{1}{30}$	$\frac{2}{30}$	$\frac{3}{30}$	$\frac{4}{30}$	$\frac{5}{30}$	$\frac{1}{30}$
2	$\frac{2}{30}$	$\frac{3}{30}$	$\frac{4}{30}$	$\frac{5}{30}$	$\frac{6}{30}$	$\frac{2}{30}$
3	$\frac{3}{30}$	$\frac{4}{30}$	$\frac{5}{30}$	$\frac{6}{30}$	$\frac{7}{30}$	$\frac{3}{30}$
$E(y)$	$\frac{6}{10}$	$\frac{10}{30}$	$\frac{14}{30}$	1	2	5
$E(y^2)$	0	$\frac{10}{30}$	$\frac{28}{30}$	$\frac{38}{30}$		
$E(xy)$	0	$\frac{10}{30}$	$\frac{56}{30}$	$\frac{156}{30}$		

00	0	$\frac{1}{30} \times 0$
01	0	$\frac{1}{30} \times 0$
02	0	$\frac{1}{30} \times 0$
10	0	$\frac{1}{30} \times 0$
11	1	$\frac{1}{30} \times 1$
12	2	$\frac{1}{30} \times 2$
20	0	$\frac{1}{30} \times 0$
21	2	$\frac{1}{30} \times 2$
22	4	$\frac{1}{30} \times 4$
30	0	$\frac{1}{30} \times 0$
31	3	$\frac{1}{30} \times 3$
32	6	$\frac{1}{30} \times 6$

Conditional Distribution of X

$$f(u|y) = \frac{f(u,y)}{h(y)} \quad \therefore h(y) > 0$$

Slide 20

3.51 Three cards are drawn without replacement from the 12 face cards (jacks, queens, and kings) of an ordinary deck of 52 playing cards. Let X be the number of kings selected and Y the number of jacks. Find:

- (a) the joint probability distribution of X and Y;
- (b) $P\{(X,Y) \in A\}$, where A is the region given by $\{(x,y) \mid x+y \geq 2\}$.

X = no of Kings Y = no of Jacks total facecards = 12

total cards = 52

marginal dist
of x

y \ x	0	1	2	3	Total g(n)
0	$\frac{14C_0 \cdot 12C_0}{52C_3} \cdot \frac{1}{55}$	$\frac{14C_1 \cdot 12C_1}{52C_3} \cdot \frac{6}{55}$	$\frac{14C_2 \cdot 12C_2}{52C_3} \cdot \frac{15}{55}$	$\frac{14C_3 \cdot 12C_3}{52C_3}$	$\frac{14}{55}$
1	$\frac{14C_1 \cdot 12C_1}{52C_3} \cdot \frac{6}{55}$	$\frac{14C_2 \cdot 12C_2}{52C_3} \cdot \frac{15}{55}$	$\frac{14C_3 \cdot 12C_3}{52C_3} \cdot \frac{6}{55}$	-	$\frac{28}{55}$
2	$\frac{14C_2 \cdot 12C_2}{52C_3} \cdot \frac{6}{55}$	$\frac{14C_3 \cdot 12C_3}{52C_3} \cdot \frac{6}{55}$	-	-	$\frac{12}{55}$
3	-	-	-	-	$\frac{6}{55}$
Total n(y)	$\frac{14}{55}$	$\frac{28}{55}$	$\frac{12}{55}$	$\frac{6}{55}$	1

marginal
dist of y

Conditional Distribution of Y

$$f(y|u) = \frac{f(u,y)}{g(u)} \quad \therefore g(u) > 0$$

$$b) P(X, Y \in A) = P(X+Y \geq 2)$$

$$P = 1 - P(X+Y \leq 1)$$

$$= 1 - [P(0,0) + P(0,1) + P(1,0)]$$

$$= 1 - \left[\frac{1}{55} + \frac{6}{55} + \frac{6}{55} \right]$$

$$P(X+Y \geq 2) = \frac{42}{55}$$

Statistical independence

$$f(u,y) = g(u) h(y)$$

$$f(u,y) = f(u|y) h(y)$$

$$g(u) = f(u|y) = \int_{-\infty}^{\infty} f(u,y) dy = 1$$

J: Joint Probability Distribution for Example 3.14

f(x,y)	x			Row Totals	
	0	1	2		
y	0	$\frac{3}{28}$	$\frac{9}{28}$	$\frac{3}{28}$	$\frac{15}{28}$
	1	$\frac{3}{14}$	$\frac{3}{14}$	0	$\frac{6}{28}$
	2	$\frac{1}{28}$	0	0	$\frac{1}{28}$
Column Totals		$\frac{5}{28}$	$\frac{15}{28}$	$\frac{3}{28}$	1

Show nor statistically independent

$$\text{taking } f(0,1) = \frac{3}{14}$$

$$g(0) = \frac{5}{14}$$

$$h(1) = \frac{15}{28}$$

$$f(0,1) = g(0) h(1)$$

$$\frac{3}{14} = \frac{3}{14} \times \frac{15}{28}$$

$$\frac{3}{14} \neq \frac{15}{392}$$

Slide 20

- 2.32 A coin is tossed twice. Let Z denote the number of heads on the first toss and H the total number of heads on the 2 tosses. If the coin is unbiased and a head has a 40% chance of occurring, find
 (a) the joint probability distribution of W and Z .
 (b) the marginal distribution of H .
 (c) the marginal distribution of Z .
 (d) the probability that at least 1 head occurs.

a)

$y \times$	no. of heads on first toss		marginal dist of y Total $g(y)$
z	$z=0$	$z=1$	
$w=0$	$(0.6)(0.6) = 0.36$	$(0.4)(0.6) = 0.24$	0.36
$w=1$	$(0.4)(0.6) = 0.24$	$(0.4)(0.4) = 0.16$	0.48
$w=2$	$(0.6)(0.4) = 0$	$(0.4)(0.4) = 0.16$	0.16
Total $h(w)$	0.6	0.4	1
Marginal dist of w			

$$\begin{aligned} H &= 0.6 \\ T &= 0.4 \end{aligned}$$

marginal Distribution \rightarrow Discrete

$$g(u) = \sum_y f(u,y) dy$$



$$h(y) = \sum_u f(u,y) du$$

b) Marginal distribution of w

$$\begin{aligned} w=0 &= 0.36 \\ w=1 &= 0.48 \\ w=2 &= 0.16 \end{aligned}$$

c) Marginal distribution of z

$$\begin{aligned} z=0 &= 0.6 \\ z=1 &= 0.4 \end{aligned}$$

$$\begin{aligned} d) P(\text{at least 1 } H) &= P(w=1 \text{ or } w=2) \\ &= P(w=1) + P(w=2) \\ &= 0.48 + 0.16 \\ &= 0.64 \end{aligned}$$

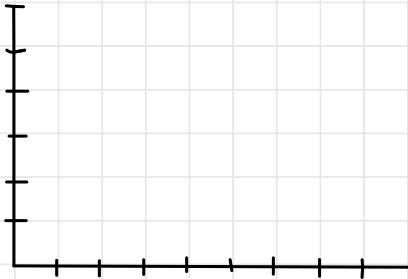
Slide 22

Let the number of phone calls received by a switchboard during a 5-minute interval be a random variable X with probability function

$$f(x) = \frac{e^{-3} 3^x}{x!}, \quad \text{for } x = 0, 1, 2, \dots$$

- (a) Determine the probability that X equals 0, 1, 2, 3, 4, 5, and 6.
- (b) Graph the probability mass function for these values of x .
- (c) Determine the cumulative distribution function for these values of X .

x	$f(x)$	$F(x) \rightarrow \text{CDF}$
0	0.1353	0.1353
1	0.2707	0.4060
2	0.2707	0.6767
3	0.1804	0.8571
4	0.0902	0.9473
5	0.0361	0.9834
6	0.01202	0.9954



Probability Density Function (PDF)

1. $f(n) \geq 0$, for all $n \in \mathbb{R}$
2. $\int_{-\infty}^{\infty} f(n) dn = 1 \rightarrow \text{verify PDF}$
3. $P(a < X < b) = \int_a^b f(n) dn$

EXAMPLE 1

Suppose that the error in the reaction temperature is "n", for a controlled laboratory experiment is a continuous random variable N having the probability density function

$$f(n) = \begin{cases} \frac{8}{9}, & -1 < n < 2, \\ 0, & \text{otherwise.} \end{cases}$$

- (a) Verify that $f(n)$ is a density function.
- (b) Find $P(0 < X \leq 1)$.

$$\begin{aligned} \text{a)} \int_{-1}^2 \frac{8}{9} n^2 dn &= \left| \frac{8}{9} \frac{n^3}{3} \right|_{-1}^2 = \frac{8}{9} + \frac{1}{9} = 1 \xrightarrow{\substack{\text{to equal} \\ \text{Hence verified}}} \end{aligned}$$

$$\begin{aligned} \text{b)} \int_0^1 \frac{8}{9} n^2 dn &= \frac{1}{9} \end{aligned}$$

EXAMPLE 4

The shelf life, in days, for bottles of a certain prescribed medicine is a random variable having the density function:

$$f(x) = \begin{cases} \frac{20000}{(x+100)^3}, & x > 0, \\ 0, & \text{elsewhere.} \end{cases}$$

Find the probability that a bottle of this medicine will have a shelf life of

- (a) at least 200 days;
- (b) anywhere from 80 to 120 days.

Example #07

(a) at least 200 days

$$\int_{200}^{\infty} \frac{20000}{(u+100)^3} du = \int_{200}^{\infty} \frac{20000}{u^3} du$$

$$= \left[-\frac{20000}{2(u+100)^2} \right]_{200}^{\infty}$$

$$= \left[\frac{10000}{(u+100)^2} \right]_{200}^{\infty}$$

$$= 1$$

(b) anywhere from 80 to 120 days

$$\int_{80}^{120} \frac{20000}{(u+100)^3} du = \int_{80}^{120} \frac{20000}{u^3} du$$

$$= \left[-\frac{20000}{2(u+100)^2} \right]_{80}^{120}$$

$$= \left[\frac{10000}{(u+100)^2} \right]_{80}^{120}$$

$$= 25 - 25$$

$$= 0$$

$$= 0.1540$$

continuous Probability Distribution

- ↳ $f(n) \rightarrow \text{probability density function (PDF)}$
- ↳ Areas used to represent probabilities
- ↳ can not be given in tabular form



→ A probability density function is constructed so that the area under its curve bounded by the x axis is equal to 1 when computed over the range of X .

EXAMPLE 3

The Department of Energy (DOE) puts projects out on bid and generally estimates what a reasonable bid should be. Call the estimate b . The DOE has determined that the density function of the winning (low) bid is:

$$f(y) = \begin{cases} \frac{54}{8b}, & 0 \leq y \leq 2b, \\ 0, & \text{elsewhere.} \end{cases}$$

Find $F(y)$ and use it to determine the probability that the winning bid is less than the DOE's preliminary estimate b .

$$F(y) = \int_{\frac{y}{2b}}^y \frac{54}{8b} dy = \left[\frac{54}{8b} y \right]_{\frac{y}{2b}}^y = \frac{54}{8b} y - \frac{1}{4} y^2$$

$$P(Y \leq b) = F(b) = \frac{54}{8b} b - \frac{1}{4} b^2 = \frac{3}{8}$$

$$\begin{aligned} &\int_{200}^{\infty} \frac{1}{(u+100)^3} du \\ &= \int_{200}^{\infty} u^3 du \\ &= \left[\frac{u^2}{2} \right]_{200}^{\infty} \\ &= \left[\frac{(u+100)^2}{2} \right]_{200}^{\infty} \\ &= \left[-\frac{10000}{(u+100)^2} \right]_{200}^{\infty} \end{aligned}$$

$u = u+100$
 $\frac{du}{du} = 1 \Rightarrow du = \frac{du}{1}$

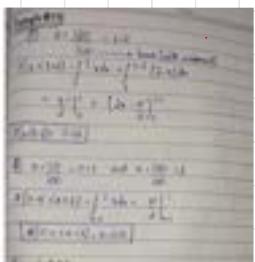
Example 5

The total number of hours, measured in units of 100 hours, that a family uses a vacuum cleaner over a period of one year is a continuous random variable X that has the density function

$$f(x) = \begin{cases} x, & 0 < x < 3, \\ 2-x, & 1 \leq x < 2, \\ 0, & \text{elsewhere.} \end{cases}$$

Find the probability that over a period of one year, a family uses their vacuum cleaner:

- (a) less than 120 hours;
- (b) between 50 and 100 hours.



$$02) f(u) = \begin{cases} 3u^4 & ; u > 1 \\ 0 & ; \text{elsewhere} \end{cases}$$

a) Verify that $f(u)$ is a PDF

$$\begin{aligned} & \int 3u^4 \\ & \cdot \left[\frac{3}{5}u^5 \right]_1^\infty \\ & = \left[-\frac{1}{5}u^5 \right]_1^\infty \\ & \cdot \left(-\frac{1}{5} \right) - \left(\frac{1}{5} \right) \\ & -0 - (-1) = 1 \end{aligned}$$

$$b) F(u) = \text{CDF}$$

$$\begin{aligned} & \text{use } u \text{ as upper limit} \quad \leftarrow u \\ & \int \frac{3}{u^4} du \\ & \cdot \left[-\frac{3}{3}u^{-3} \right]^u \\ & = -u^3 + 1 \\ & = 1 - u^3 \end{aligned}$$

$$c) P(u > 4)$$

$$\begin{aligned} & \rightarrow \int_u^\infty 3u^4 du \\ & \left| -u^5 \right|_u^\infty \\ & -\frac{1}{64} - \frac{1}{\infty^5} \\ & -\frac{1}{64} \end{aligned}$$

Example 7

Measurements of scientific systems are always subject to variation, more than others. There are many structures for measurement error, and statisticians spend a great deal of time modeling these errors. Suppose the measurement error X of a certain physical quantity is decided by the density function

$$f(x) = \begin{cases} k(3-x^2), & -1 \leq x \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

- (a) Determine k that renders $f(x)$ a valid density function.
- (b) Find the probability that a random error in measurement is less than $1/2$.
- (c) For this particular measurement, it is undesirable if the magnitude of the error (i.e., $|x|$) exceeds 0.8. What is the probability that this occurs?

$$a) \int_{-1}^1 k(3-u^2) du = 1 \rightarrow \text{equate to 1 as valid PDF evaluates to 1}$$

$$K(3u - u^3/3) = 1$$

$$K \left[3 - \frac{1}{3} - \left(3(-1) - \frac{(-1)^3}{3} \right) \right] = 1$$

$$K = \frac{3}{16}$$

$$b) P(u < 1/2) = \int_{-1}^{1/2} \frac{3}{16} (3-u^2) du \Rightarrow \frac{99}{128}$$

$$c) P(|u| < 0.8) = P(u < -0.8) + P(u > 0.8)$$

↓
WHAT

$$= F(-0.8) + [1 - F(0.8)]$$

$$= 0.164$$

$$-0.8 < u < 0.8$$

Joint Density Function (JDF)

$$1. f(u, y) \geq 0$$

↳ for continuous variable

$$2. \int_{-\infty}^{\infty} \int f(u, y) dy du = 1$$

$$3. P[(X, Y) \in A] = \int_A \int f(u, y) du dy$$

Marginal Distribution

→ for continuous variable

$$g(u) = \int_{-\infty}^{\infty} f(u, y) dy \rightarrow \text{w.r.t. } y$$

$$h(y) = \int_{-\infty}^{\infty} f(u, y) du \rightarrow \text{w.r.t. } u$$

$$\text{Q3)} f(u, y) = \begin{cases} \frac{2}{5} (2u + 3y) & 0 \leq u \leq 1 \\ 0 & 0 \leq y \leq 1 \\ \text{elsewhere} & \end{cases}$$

a) Verify that $f(u, y)$ is a valid JPDF

$$\int_0^1 \int_0^1 \frac{2}{5} (2u + 3y) dy du$$

$$\frac{2}{5} \int_0^1 \int_0^1 (u + \frac{3y^2}{2}) dy du$$

$$\frac{2}{5} \left[u + \frac{3y^3}{2} \right]_0^1$$

$$\frac{2}{5} \left[u + \frac{3}{2} \right]$$

$$\frac{2}{5} [\frac{5}{2}]$$

1

c) find $g(u) \& h(y)$

$$\text{Marginal } g(u) = \frac{2}{5} \int_0^1 (2u + 3y) dy$$

$$= \frac{2}{5} \left(2uy + \frac{3y^2}{2} \right)_0^1$$

$$= \frac{2}{5} \left(2u + \frac{3}{2} \right)$$

$$\text{Marginal } h(y) = \frac{2}{5} \int_0^1 (2u + 3y) du$$

$$\text{b) } P[0 \leq u \leq \frac{1}{2}, \frac{1}{4} \leq y \leq \frac{1}{2}]$$

$$= \frac{2}{5} \int_{\frac{1}{4}}^{\frac{1}{2}} \int_0^1 (2u + 3y) du dy$$

$$= \frac{2}{5} \int_{\frac{1}{4}}^{\frac{1}{2}} (u^2 + 3uy) \Big|_0^1 dy$$

$$= \frac{2}{5} \int_{\frac{1}{4}}^{\frac{1}{2}} \left(\frac{1}{4} + \frac{3y^2}{2} \right) dy$$

$$= \frac{2}{5} \left(\frac{y}{4} + \frac{3y^3}{6} \right) \Big|_{\frac{1}{4}}^{\frac{1}{2}}$$

$$= \frac{2}{5} \left[\frac{1}{8} + \frac{3}{16} - \left(\frac{1}{16} + \frac{3}{64} \right) \right]$$

$$= 0.08$$

Example 9

* A privately owned business operates both a drive-in facility and a walk-in facility. On a randomly selected day, let X and Y , respectively, be the proportions of the time that the drive-in and the walk-in facilities are in use, and suppose that the joint density function of these random variables is

$$f(x, y) = \begin{cases} \frac{2}{5}(2x + 3y), & 0 \leq x \leq 1, 0 \leq y \leq 1, \\ 0, & \text{otherwise} \end{cases}$$

(a) Verify for PDF?

(b) Find $P[X, Y \in A]$, where $A = \{(x, y) / 0 \leq x \leq 1/2, 1/4 \leq y \leq 1/2\}$.

$$\begin{aligned} \text{Example 9C} \\ \text{(a)} \int_0^{1/2} \int_{1/4}^{1/2} \frac{2}{5} (2x + 3y) dy dx \\ = \int_0^{1/2} \frac{2}{5} \left(2x + 3y^2 \right) \Big|_{1/4}^{1/2} dx \\ = \frac{2}{5} \left(x + \frac{3}{2} \right) \Big|_0^{1/2} \\ = \frac{2}{5} \left(\frac{1}{2} + \frac{3}{4} \right) \\ = 0.25 \end{aligned}$$

$$\begin{aligned} \text{(b)} \text{Find } P[(X, Y) \in A], \text{ where } A = \{(x, y) / \\ = \int_0^{1/2} \int_{1/4}^{1/2} \frac{2}{5} (2x + 3y) dy dx \\ = \frac{2}{5} \int_0^{1/2} \left(2x + 3y^2 \right) \Big|_{1/4}^{1/2} dx \\ = \frac{2}{5} \left(\frac{1}{2} + \frac{3}{4} \right) \\ = \frac{2}{5} \left(\frac{5}{4} \right) \\ = 0.5 \end{aligned}$$

Mean of a Random Variable



i) Let X be a random variable with probability distribution $f(x)$. The expected value of the random variable $g(X)$ is

$$\mu_{g(X)} = E[g(X)] = \sum_x g(x)f(x)$$

if X is discrete, and

$$\mu_{g(X)} = E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x) dx$$

if X is continuous.

Example 4.1: Suppose that the number of cars X that pass through a car wash between 4:00 pm and 5:00 pm on any given Friday has the following probability distribution:

x	4	5	6	7	8	9	10
$P(X = x)$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{12}$

Let $g(X) = (X - 1)$ represent the amount of money, in dollars, paid to the attendant by the customer. Find the attendant's expected earnings for this particular time period.

Solution: By Theorem 4.1, the attendant can expect to receive

$$\begin{aligned} E[g(X)] &= E[(X - 1)] = \sum_{x=1}^{10} (x - 1)f(x) \\ &= (1)(\frac{1}{12}) + (2)(\frac{1}{12}) + (3)(\frac{1}{6}) + (4)(\frac{1}{6}) \\ &\quad + (5)(\frac{1}{6}) + (6)(\frac{1}{6}) = \$12.50. \end{aligned}$$

Let X and Y be random variables with joint probability distribution $f(x,y)$. The mean, or expected value, of the random variable $g(X,Y)$ is

$$\mu_{g(X,Y)} = E[g(X,Y)] = \sum_x \sum_y g(x,y)f(x,y)$$

if X and Y are discrete, and

$$\mu_{g(X,Y)} = E[g(X,Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y)f(x,y) dx dy$$

if X and Y are continuous.

Example 4.2: Let X and Y be the random variables with joint probability distribution indicated in Table 11 on page 96. Find the expected value of $g(X,Y) = XY$. The table is reproduced here for convenience.

		Row Total			
		0	1	2	3
0	0	0.0	0.0	0.0	0.0
	1	0.0	0.0	0.0	0.0
	2	0.0	0.0	0.0	0.0
	3	0.0	0.0	0.0	0.0
Column Total		0.0	0.0	0.0	0.0

Solution: By Definition 4.2, we write

$$\begin{aligned} E(XY) &= \sum_{x=0}^3 \sum_{y=0}^3 xyf(x,y) \\ &= (0)(0)f(0,0) + (0)(1)f(0,1) \\ &\quad + (0)(2)f(0,2) + (0)(3)f(0,3) \\ &\quad + (1)(0)f(1,0) + (1)(1)f(1,1) \\ &\quad + (1)(2)f(1,2) + (1)(3)f(1,3) \\ &\quad + (2)(0)f(2,0) + (2)(1)f(2,1) \\ &\quad + (2)(2)f(2,2) + (2)(3)f(2,3) \\ &\quad + (3)(0)f(3,0) + (3)(1)f(3,1) \\ &\quad + (3)(2)f(3,2) + (3)(3)f(3,3) \\ &= \mu_{XY} = \frac{3}{12}. \end{aligned}$$

when it is clear to which random variable we refer.

Let X be a random variable with probability distribution $f(x)$ and mean μ . The variance of X is

$$\sigma^2 = E[(X - \mu)^2] = \sum_x (x - \mu)^2 f(x), \quad \text{if } X \text{ is discrete, and}$$

$$\sigma^2 = E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx, \quad \text{if } X \text{ is continuous.}$$

The positive square root of the variance, σ , is called the standard deviation of X .

Let X be a random variable with probability distribution $f(x)$. The mean, or expected value, of X is

$$\mu = E(X) = \sum_x xf(x)$$

if X is discrete, and

$$\mu = E(X) = \int_{-\infty}^{\infty} xf(x) dx$$

if X is continuous.

The reader should note that the way to calculate the numerical value of μ depends

Variance & Covariance of Random Variables

$$\mu = E(x) = \sum x f(x) \rightarrow \text{mean}$$

$$\sigma^2 = (x - \mu)^2 \rightarrow \text{variance}$$

$$\sigma^2 = E(x^2) - [E(x)]^2$$

$$= \sum x^2 f(x) - \mu^2$$

Let X be a random variable with probability distribution $f(x)$ and mean μ . The variance of X is

$$\sigma^2 = E[(X - \mu)^2] = \sum (x - \mu)^2 f(x), \quad \text{if } X \text{ is discrete, and}$$

$$\sigma^2 = E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx, \quad \text{if } X \text{ is continuous.}$$

The positive square root of the variance, σ , is called the standard deviation of X .

MATHEMATICAL EXPECTATION

\rightarrow $\sigma^2 \rightarrow$ variance

$$V(x) = E(x^2) - [E(x)]^2$$

mean \rightarrow expected
 $\mu = E(x) = \sum x f(x)$

$$E(xy) = \sum \sum xy f(x,y)$$

$$\text{Covariance}(x,y) = E(xy) - E(x)E(y)$$

$$\text{Correlation}(x,y) = \frac{\text{Covariance}(x,y)}{\sqrt{V(x)} \sqrt{V(y)}}$$

$\pm 1 \rightarrow$ strong
 $\pm 0.5 \rightarrow$ moderate
 $0 \rightarrow$ negligible

Example 4.16

The fraction X of male runners and the fraction Y of female runners who compete in marathon runs are described by the joint density function

$$f(x,y) = \begin{cases} 8xy, & 0 \leq y \leq x \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

Find the covariance of X and Y .

Solution: We first compute the marginal density functions. They are

$$g(x) = \begin{cases} 4x^2, & 0 \leq x \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

and

$$h(y) = \begin{cases} 8y(1-y^2), & 0 \leq y \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

From these marginal density functions, we compute

$$\mu_x = E(X) = \int_0^1 4x^3 dx = \frac{4}{5} \quad \text{and} \quad \mu_y = \int_0^1 8y^2(1-y^2) dy = \frac{8}{15}.$$

From the joint density function given above, we have

$$E(XY) = \int_0^1 \int_y^1 8x^2y^2 dx dy = \frac{4}{9}.$$

Then

$$\sigma_{xy} = E(XY) - \mu_x \mu_y = \frac{4}{9} - \left(\frac{4}{5}\right)\left(\frac{8}{15}\right) = \frac{1}{225}.$$

DISCRETE PROBABILITY DISTRIBUTION

1. Binomial Distribution

- The number X of successes in n Bernoulli trials is called a binomial random variable.
- The probability distribution of this discrete random variable is called the binomial distribution.
- The probability of a success in a binomial experiment can be computed with this formula:

$$f(x; n, p) = \binom{n}{x} p^x q^{n-x}, \quad x = 0, 1, 2, \dots, n.$$

2. The Bernoulli Process

↳ repeated trials } TWO POSSIBLE OUTCOMES

↳ the $P(\text{success})$ remains constant

↳ repeated trials are independent

$$P + Q = 1$$

Mean = NP Probability Success
→ no of trials Probability fail
 $q=1-p$

Variance = $NPQ \rightarrow \sigma^2$

Standard Deviation = $\sqrt{\text{Variance}}$

PMF = $f(x) = {}^n C_x P^x q^{n-x}$

Q) A coin is tossed 4 times. Find mean, variance, SD of heads that will be obtained

$$n=4, P=\frac{1}{2}, Q=\frac{1}{2}$$

1. Mean = NP

$$4\left(\frac{1}{2}\right) = 2$$

2. Variance = NPQ

$$4\left(\frac{1}{2}\right)\left(1-\frac{1}{2}\right) = 1$$

3. $SD = \sqrt{\text{Variance}}$

$$= \sqrt{1} = 1$$

$$2^4 = 16$$

x	$P(x)$	$xP(x)$	$x^2 P(x)$
0	$\frac{1}{16}$		
1			
2			
3			
4			

$$V(X) = E(X^2) - [E(X)]^2$$
$$= 5 - 4$$
$$= 1$$

Q) A dice is rolled 480 times. Find the mean, variance of the no of 3s that will be rolled

$$n=480 \quad P=\frac{1}{6} \quad Q=\frac{5}{6}$$

Mean = NP

$$480\left(\frac{1}{6}\right) = 80$$

$SD = \sqrt{NPQ}$

$$\sqrt{480\left(\frac{1}{6}\right)\left(\frac{5}{6}\right)} = 8.16$$

3. Hypergeometric Experiment

↳ Trial Outcomes $\begin{matrix} \text{FAIL} \\ \text{SUCCESS} \end{matrix}$

↳ $P(\text{success})$ changes each trial

↳ successive trials are dependent

↳ experiment repeated a fix no. of times

$$P(x) = \frac{^a C_x \cdot ^b C_{n-x}}{^{a+b} C_n}$$

↑ same size

] selecting with Out replacement

$a+b = \text{total population}$

one kind of items
of
another kind of items

Q)

- Assistant Manager Applicants: Ten people apply for a job as assistant manager of a restaurant. Five have completed college and five have not. If the manager selects 3 applicants at random, find the probability that all 3 are college graduates.

$$\begin{matrix} c & n-c \\ a & b \\ {}^a C_x & {}^b C_{n-x} \end{matrix}$$

$$\begin{matrix} a+b \\ = 10 \\ C_3 \end{matrix}$$

$$P(\text{all 3 college graduates}) = \frac{{}^5 C_3 \cdot {}^5 C_0}{{}^{10} C_3}$$

Q)

- House Insurance: A recent study found that 2 out of every 10 houses in a neighborhood have no insurance. If 5 houses are selected from 10 houses, find the probability that exactly 1 will be uninsured.

Select 5

$P(1 \text{ will be insured})$

2 → not insured

$$\frac{{}^8 C_u \cdot {}^2 C_1}{{}^{10} C_5} = \frac{5}{9}$$

8 → insured

$$\begin{matrix} I & NI \\ a & b \\ {}^a C_x \cdot {}^b C_{n-x} & C_5 \end{matrix} = 10$$

4. Geometric Experiment

- ↳ Trial Outcomes 
- ↳ $P(\text{success})$ remain constant for each trial
- ↳ each trial is independent
- ↳ experiment repeated a variable no. of times, until first success is obtained



$$g(x) = pq^{x-1}$$

$$\text{mean} \rightarrow \mu = \frac{1}{p}$$

$$\text{variance} \rightarrow \sigma^2 = \frac{1-p}{p^2}$$

POISSON Distribution



↳ avg no of success

↳ time interval

↳ mean

Q)

• **Radioactive Particles:** During a laboratory experiment, the average number of radioactive particles passing through a counter in 1 millisecond is 4. What is the probability that 6 particles enter the counter in a given millisecond?

$$\lambda = \overset{\text{ms}}{\underset{1}{\uparrow}} \times 4 = 4$$

$$P(X=6)$$

$$P(X=6, \lambda=4)$$

if 6 Particles in 2ms

$$\lambda = \overset{\text{ms}}{\underset{2}{\uparrow}} (4) = 8$$

$$e^{-4} \left(\frac{4^6}{6!} \right) = 0.1042$$

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

no of times event occurred
mean of x

Discrete Random Variable

- ↳ PMF
- ↳ CDF
- ↳ Mean
- ↳ Variance

Continuous Random Variable

- ↳ PDF
- ↳ CDF
- ↳ Mean
- ↳ Variance

ishma hafeez
notes

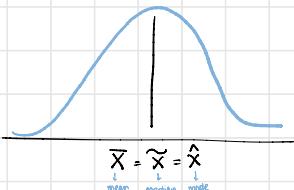
reprt
tree

Normal Distribution

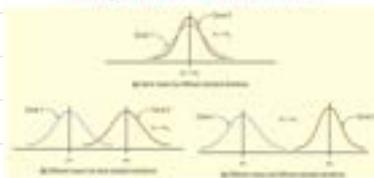
↳ continuous type Distribution

$$\text{PDF: } f(x) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}}, \quad -\infty < x < \infty$$

↑ Population mean
↓ Population std



Shapes of Normal Distribution



Standard Normal Distribution

$\mapsto \mu = 0$ and $\sigma = 1$

$$Z\text{-score} = \frac{X - \mu}{\sigma}, \quad y = \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}}$$

$$Z \sim SN(\mu=0, \sigma=1)$$

Areas under the standard Normal Curve

- The figure consists of four separate normal distribution curves, each with a mean labeled μ . The area under the curve represents the total probability, which is 1 for a normal distribution.

 - Left of a value:** The first curve shows a yellow shaded area to the left of a vertical line at $x = a$. This represents the probability $P(X < a)$.
 - Right of a value:** The second curve shows a yellow shaded area to the right of a vertical line at $x = a$. This represents the probability $P(X > a)$.
 - Between two values:** The third curve shows a yellow shaded area between two vertical lines at $x = a$ and $x = b$, where $a < b$. This represents the probability $P(a < X < b)$.
 - Outside two values:** The fourth curve shows a yellow shaded area outside two vertical lines at $x = a$ and $x = b$, where $a < b$. This represents the probability $P(X < a \text{ or } X > b) = P(X < a) + P(X > b)$.

Example 6t

- Find the area to the left of $z = -2.08$.
 - Find the area to the right of $z = -1.39$.
 - Find the area between $z = 1.68$ and $z = -1.37$.

Example # 12

*Gauges are used to reject all components for which a certain dimension is not within the specification (1.50 ± 0.1). It is known that this measurement is normally distributed with mean 1.50 and standard deviation 0.2. Determine the value of such that the specifications "cover" 95% of the measurements.

$$\frac{d+d-\mu}{\mu+100} = \frac{2}{5} = 2$$

$$\mu = 1.5 \quad \sigma = 0.2 \quad \alpha = 0.02$$

$$AL = 0.95 + 0.025 = 0.975$$

CHP 6



9.17

6.5 Given a standard normal distribution, find the area under the curve that lies

- to the left of $z = -1.39$; $1 - 0.9177 = 0.0823$
- to the right of $z = 1.96$; $1 - 0.950 = 0.050$
- between $z = -2.16$ and $z = -0.65$; $1 - 0.9846 = 1 - 0.7422$
 $0.0154 - 0.2518 = 0.2424$
- to the left of $z = 1.43$; 0.9239
- to the right of $z = -0.89$; 0.8078
- between $z = -0.48$ and $z = 1.74$. 0.5326

6.6 Find the value of z if the area under a standard normal curve

- to the right of z is 0.3622; $1 - 0.3622 = 0.6378 \rightarrow 0.35$
- to the left of z is 0.1131; $1 - 0.1131 = 0.8869 \rightarrow 1.21$
- between 0 and z , with $z > 0$, is 0.4838; $0.5 + 0.4838 = 0.9838 \rightarrow 2.14$
- between $-z$ and z , with $z > 0$, is 0.9500. $0.25 + 0.95 = 0.975 \rightarrow 1.96$

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

6.8 Given a normal distribution with $\mu = 30$ and $\sigma = 6$, find

- the normal curve area to the right of $x = 17$; $\frac{17-30}{6} = -0.5$
- the normal curve area to the left of $x = 22$; $\frac{22-30}{6} = -0.67$
- the normal curve area between $x = 32$ and $x = 41$; $\frac{32-30}{6} = \frac{2}{3}$ and $\frac{41-30}{6} = \frac{11}{6}$
- the value of x that has 80% of the normal curve area to the left; $x = 24.8$
- the two values of x that contain the middle 75% of the normal curve area.

6.9 Given the normally distributed variable X with mean 18 and standard deviation 2.5, find

- $P(X < 15)$; $\frac{15-18}{2.5} = -1.2$
- the value of k such that $P(X < k) = 0.2236$; $k = (2.5)(-0.76) + 18 = 16.1$
- the value of k such that $P(X > k) = 0.1814$; $k = (2.5)(0.91) + 18 = 20.275$
- $P(17 < X < 21)$.

$$\frac{x - 30}{6}$$

$$\frac{32-30}{6} - \frac{41-30}{6}$$

$$x = Z \times \sigma + \mu = 35.04$$

$$0.8 \times 6 + 30$$

$$\mu = 18$$

$$8.25 \quad \frac{x - 18}{2.5}$$

- 6.9 (a) $z = (15 - 18)/2.5 = -1.2$; $P(X < 15) = P(Z < -1.2) = 0.1151$.
- (b) $z = -0.76$, $k = (2.5)(-0.76) + 18 = 16.1$.
- (c) $z = 0.91$, $k = (2.5)(0.91) + 18 = 20.275$.
- (d) $z_1 = (17 - 18)/2.5 = -0.4$, $z_2 = (21 - 18)/2.5 = 1.2$;
 $P(17 < X < 21) = P(-0.4 < Z < 1.2) = 0.8849 - 0.3446 = 0.5403$.

Central limit theorem

$$z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

standard error of mean

Formula	Use
1. $z = \frac{X - \mu}{\sigma}$	Used to gain information about an individual data value when the variable is normally distributed.
2. $z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$	Used to gain information when applying the central limit theorem about a sample mean when the variable is normally distributed or when the sample size is 30 or more.

Example 3

An upgrade of a certain software package requires the installation of 60 new files. Each file is installed independently. On average, it takes 15 seconds to install one file, with a variance of 11 seconds.

- (a) How likely is it that the whole package will be updated in less than 35 minutes?
 (b) There is a later version of the package. It requires only N new files to be installed, and it is presumed that 99% of the time upgrading takes less than 10 minutes. Based on this information, calculate N.

Solution:

$$(a) n = 60, \mu = 15 \text{ sec}, \text{Var}(X) = 11 \text{ sec}^2, \sigma = \sqrt{11} \approx 3.316$$

$$P(X_n \leq 600) \approx P\left(\frac{X_n - 60}{\sqrt{11}} \leq \frac{600 - 60}{\sqrt{11}}\right) = P(Z \leq -0.20) = 0.4844$$

$$(b) X \sim N(\mu, \sigma^2), \mu = 15 \text{ sec}, \sigma = \sqrt{11} \text{ sec}$$

$$P(X_n \leq 600) = P\left(\frac{X_n - 60}{\sqrt{11}} \leq \frac{600 - 60}{\sqrt{11}}\right) = P(Z \leq 4.54) = 0.99$$

So for $N = 17$ we get the correct probability more than 99%.

Finite Population Correction Factor

$$\sqrt{\frac{N-n}{N-1}}$$

population size
sample size

ishma hafeez
notes

reversht
heez

- ↪ Estimation ↪ point estimate: \bar{x}, s^2, P (representative point)
- ↪ interval estimate: $16 < \mu < 25$ (margin of error)
- ↪ Hypothesis testing
- ↪ ANOVA
- ↪ Regression and Correlation ↪ ^{Gallup} pulse

Properties of good estimators

1. Unbiased estimator

↪ \bar{x} = parameter being estimated

2. Consistent estimator

3. Relatively efficient estimator

↪ has the smallest variance

$$\bar{X} \stackrel{\text{unbiased estimator of population mean } (\mu)}{\approx} \mu$$

↓ ↓
statistic parameter

sampling error: $\bar{X} - \mu$

Q) 90% $\rightarrow Z \frac{(1-0.9)}{2} = 1.65$

95% $\rightarrow Z \frac{(1-0.95)}{2} = 1.96$

99% $\rightarrow Z \frac{(1-0.99)}{2} = 2.58$

eg. $n=50, \bar{x}=54, \sigma=6$

95% of CI mean?

$54 \pm (1.96) \left(\frac{6}{\sqrt{50}} \right)$

$52.3 < \mu < 55.7$

95% chance that the above interval contains true population mean

Some Important Notations

For Sample Data

- n = Sample Size
- \bar{X} = Sample Mean
- s^2 = Biased Sample Variance
- s^2 = Unbiased Sample Variance
- S = Sample Standard Deviation
- P = Sample Proportion
- $S_x = \frac{s}{\sqrt{n}}$ (Sample Standard Error)

For Population Data

- N = Population Size
- μ = Population Mean
- σ^2 = Population Variance
- σ = Population Standard Deviation
- π = Population Proportion
- α = Level Of Significance
- $1-\alpha$ = Level Of Confidence
(or)
Confidence Interval

PROB ASSIGNMENT

K214688

Q1)

$$a) f(u, y, z) = \begin{cases} Ky^2z, & 0 < y \\ 0, & y \leq 1 \\ 0 < z < 2 \end{cases}$$

$$\int_0^2 \int_0^1 \int_0^y Ky^2z \, du \, dy \, dz = 1$$

$$\int_0^2 \int_0^1 \left| \frac{Ky^3z^2}{2} \right|_0^y = 1$$

$$\int_0^2 \int_0^1 \frac{Ky^3z^2}{2} = 1$$

$$\int_0^2 \left| \frac{Ky^3z^2}{6} \right|_0^y = 1$$

$$\int_0^2 \frac{Kz^2}{6} = 1$$

$$\left| \frac{Kz^3}{18} \right|_0^2 = 1$$

$$\frac{K}{3} = 1$$

$$\boxed{K=3}$$

$$b) P(u < \frac{1}{4}, y > \frac{1}{2}, 1 < z < 2)$$

$$\int_1^2 \int_{\frac{1}{2}}^1 \int_0^y Ky^2z \, du \, dy \, dz$$

$$\int_1^2 \int_{\frac{1}{2}}^1 \left| \frac{Ky^3z^2}{2} \right|_0^y$$

$$\int_1^2 \int_{\frac{1}{2}}^1 \frac{Ky^3z^2}{32}$$

$$\int_1^2 \left| \frac{Ky^3z^2}{96} \right|_{\frac{1}{2}}$$

$$\int_1^2 \frac{7z}{768}$$

$$\int_1^2 \left| \frac{7z^2}{1536} \right|_1^2$$

$$\frac{21}{512}$$

c) Marginal

$$f(y) = 3 \int_0^1 \int_0^y Ky^2z \, du \, dz$$

$$E(y) = \int_0^1 3y^3$$

$$= \frac{3}{4}$$

$$E(x^2) = \int_0^1 3y^6$$

$$\left| \frac{3y^7}{7} \right|_0^1$$

$$= \frac{3}{7}$$

$$\text{var}(y) = \frac{3}{5} - \left(\frac{3}{4} \right)^2$$

$$= \frac{3}{80}$$

$$\int_0^1 \int_0^y \frac{Ky^2z^2}{2} \, du \, dz$$

$$E(z) = \int_0^1 \frac{z^2}{2}$$

$$= \frac{4}{3}$$

$$\text{var}(z) = 2 - \left(\frac{4}{3} \right)^2$$

$$= \frac{2}{9}$$

$$\int_0^1 \int_0^y \frac{Ky^2z^2}{4} \, du \, dz$$

$$E(z^2) = \int_0^1 \frac{z^2}{2}$$

$$\left| \frac{z^3}{6} \right|_0^1$$

$$= \frac{1}{6}$$

$$\text{covariance} = \frac{1}{2} - \frac{3}{80} \left(\frac{2}{3} \right)$$

$$= \frac{59}{120}$$

$$f(z) = 3 \int_0^1 \int_0^y Ky^2z \, dy \, dz$$

$$\int_0^1 \int_0^y \frac{Ky^2z}{2} \, dy \, dz$$

$$\int_0^1 \left| \frac{Ky^3z}{6} \right|_0^y$$

$$= \frac{Kz}{2}$$

$$E(yz) = \frac{3}{2} \int_0^1 \int_0^y yz \, dy \, dz$$

$$= \frac{3}{2} \int_0^1 \int_0^y y^2z^2 \, dy \, dz$$

$$= \frac{3}{2} \int_0^1 \left| \frac{y^3z^2}{9} \right|_0^y$$

$$= \frac{3}{2} \int_0^1 \frac{z^3}{9} \, dz$$

$$= \frac{3}{2} \left| \frac{z^4}{36} \right|_0^1$$

$$= \frac{1}{2}$$

$$\text{correlation} = \frac{59/120}{\sqrt{3/80} \sqrt{2/9}}$$

$$= 11.758$$

ESTIMATION

Point estimate

Interval estimate

1. $\bar{X} = \frac{\sum x}{n}$ sample mean
2. $S = \sqrt{\frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2}$ $\rightarrow SD$
3. $S_{\bar{X}} = \frac{s}{\sqrt{n}}$ standard error of mean
4. $E = Z \frac{\sigma}{\sqrt{n}}$ margin of error \rightarrow standard error of mean
5. $n = \left(\frac{Z \frac{\sigma}{E}}{2} \right)^2$ \rightarrow critical value of z-distribution

Example # 04

* Depth of a River: A scientist wishes to estimate the average depth of a river. He wants to be 99% confident that the estimate is accurate within 2 feet. From a previous study, the standard deviation of the depths measured was 4.35 feet.

$$n = \left(\frac{Z_{0.99} \cdot \sigma}{E} \right)^2 = \left[\frac{2.576 \cdot 4.35}{2} \right]^2 \approx 51.7$$

confidence interval for μ

$$1. \bar{X} \pm Z \frac{\sigma}{\sqrt{n}}$$

population SD

$\hookrightarrow \sigma$ known

$$2. \bar{X} \pm Z \frac{s}{\sqrt{n}}$$

sample SD

$\hookrightarrow \sigma$ unknown

$\hookrightarrow n \geq 30$

$$3. \bar{X} \pm t \frac{s}{\sqrt{n}}$$

sample SD

critical value of t-distribution

$\hookrightarrow \sigma$ unknown

$\hookrightarrow n < 30$

confidence interval for $\mu_1 - \mu_2$

$$1. (\bar{X}_1 - \bar{X}_2) \pm Z \frac{\sqrt{\sigma_1^2 + \sigma_2^2}}{\sqrt{n_1 + n_2}}$$

$\hookrightarrow \sigma_1, \sigma_2$ known

$$2. (\bar{X}_1 - \bar{X}_2) \pm Z \frac{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}{\sqrt{n_1 + n_2}}$$

$\hookrightarrow \sigma_1, \sigma_2$ unknown

$\hookrightarrow n \geq 30$

$$3. (\bar{X}_1 - \bar{X}_2) \pm t \frac{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}{\sqrt{n_1 + n_2}}$$

$\hookrightarrow \sigma_1, \sigma_2$ unknown

$\hookrightarrow n < 30$

?

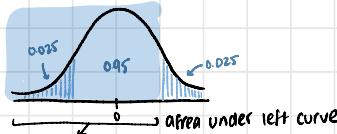
Z-distribution

Question # 1

An electrical firm manufactures light bulbs that have a length of life with mean μ and a standard deviation of 40 hours. If a sample of 100 bulbs has an average life of 780 hours, find a 95% confidence interval for the population mean of all bulbs produced by this firm.

$$\sigma = 40 \text{ hours} \quad n = 100 \quad \bar{x} = 780 \text{ hours} \quad CI = 95\% \quad \mu = ?$$

① as $n > 30$ we use z test



A_L - confidence level + α

$$A_L = 0.95 + 0.025 = 0.975 \rightarrow \text{Find on Table}$$

$$Z \frac{\alpha}{2} = 1.96$$

$$\textcircled{2} \quad \bar{x} \pm Z \frac{\sigma}{\sqrt{n}}$$

$$= 780 \pm 1.96 \left(\frac{40}{\sqrt{80}} \right)$$

$$= 780 \pm 7.84$$

$$= 772.16, 787.84$$

$$772.16 < m \leq 787.84$$

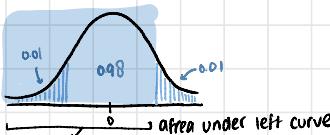
--

Question #2

The height of a random sample of 50 college students showed a mean of 174.5 cm and a standard deviation of 6.9cm. Construct a 98% Confidence interval for the mean height of all college students.

$$n=50, \bar{X}=174.5\text{cm}, s=6.9\text{cm}, CI=98\%, \mu=?$$

① $n > 50$ so z test



A_L = confidence level + α

$$A_L = 0.98 + 0.01 = 0.99 \rightarrow$$

$$Z_{\frac{\alpha}{2}} = 2.3263$$

$$\textcircled{2} \quad \bar{x} \pm Z \frac{\alpha}{2} \frac{s}{\sqrt{n}}$$

$$= 174.5 \pm 2.3263 \left(\frac{6.9}{\sqrt{50}} \right)$$

$$= 174.5 \pm 2.2$$

• 172.23 - 176.77

$$172.33 < \mu < 176.71$$

$$\begin{aligned} &\rightarrow \text{CLOSEST TO 0.99} \\ &\rightarrow \text{TAKE AVG} \\ &\rightarrow \frac{0.02 + 0.03}{2} = 0.025 \\ &\rightarrow 2.325 \end{aligned}$$

Question #3

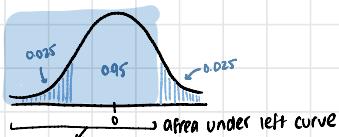
A random sample of 100 farms in a certain year gives an average yield of barley of 2100 lbs. per acre. A random sample of 100 farms in the following years given an average yield of 2000 lbs. per acre. The S.D for two populations are 224 and 192 respectively. Compute a 95% CI for the difference between two population means. Assuming data follows normal distribution.

$$n_1 = 100 \quad \bar{x}_1 = 2100 \quad CI = 95\%$$

$$n_2 = 100 \quad \bar{x}_2 = 2000 \quad \mu_1 - \mu_2 = ?$$

$$\sigma_1 = 224 \quad \sigma_2 = 192$$

① As $n \geq 30$ so z test



A_L = confidence level + α

$$A_L = 0.95 + 0.025 = 0.975 \rightarrow \text{Find on Table}$$

$$Z_{\frac{\alpha}{2}} = 1.96$$

$$② (\bar{x}_1 - \bar{x}_2) \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$= (2100 - 2000) \pm 1.96 \sqrt{\frac{224^2}{100} + \frac{192^2}{100}}$$

$$= 100 \pm 51.82$$

$$= 42.18, 157.82$$

$$= 42.18 < \mu < 157.82$$

Confidence Interval for Proportions

$$1. \hat{P} \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}\hat{q}}{n}}$$

Sample proportion \hat{P} → no. of sample units
 $\frac{x}{n}$ → x → 1 - P

Example #09

Ensuring College Costs: A survey conducted by Sallie Mae and Gallup of 1404 respondents found that 323 students paid for their education by student loans. Find the 90% confidence of the true proportion of students who paid for their education by student loans.

$$n = 1404, x = 323, CI = 90\%$$

$$P = \frac{323}{1404} = 0.23 \quad q = 1 - 0.23 = 0.77$$

$$0.23 - 1.65 \sqrt{\frac{(0.23)(0.77)}{1404}}$$

$$= 21.1\% < p < 24.9\%$$

Example #10

Religious Books: A survey of 1721 people found that 25.9% of individuals purchase religious books at a Christian bookstore. Find the 95% confidence interval of the true proportion of people who purchase their religious books at a Christian bookstore.

$$n = 1721, \hat{P} = 0.259, CI = 95\%$$

$$\hat{q} = 1 - 0.259 = 0.741$$

$$= 0.259 \pm 1.96 \sqrt{\frac{(0.259)(0.741)}{1721}}$$

$$= 0.142 < p < 0.316$$

Minimum 'n' needed Interval Estimate of Population Proportion

$$n = \hat{p}\hat{q} \left(\frac{Z_{\alpha/2}}{E} \right)^2$$

→ rearrange confidence interval
for proportions

Example # 11

*Home Computers: A researcher wishes to estimate, with 95% confidence, the proportion of people who own a home computer. A previous study shows that 40% of those interviewed had a computer at home. The researcher wishes to be accurate within 2% of the true proportion. Find the minimum sample size necessary.

$$CI=95\%, \hat{p}=40\%, E=2\% \\ \hat{q}=1-40=60\%$$

$$n = (0.4)(0.6) \left(\frac{1.96}{0.02} \right)^2 \\ = 2304.96$$

Example # 12

*M&M Colors: A researcher wishes to estimate the percentage of M&M's that are brown. He wants to be 95% confident and be accurate within 3% of the true proportion. How large a sample size would be necessary?

$$CI=95\%, E=3\% \\ \hat{p}=5\% \quad \hat{q}=5\%$$

$$n = (0.5)(0.5) \left(\frac{1.96}{0.03} \right)^2 \\ = 1067.1$$

Chi Square Distribution

- ↳ random sample
- ↳ normal distribution

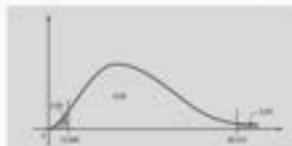
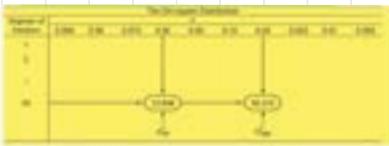
$$1. \frac{(n-1)s^2}{\chi^2_{right}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{left}}$$

$$2. df = n - 1$$

Example # 14

Find the values for χ^2_{right} and χ^2_{left} for a 90% confidence interval when $n = 25$.

$$\chi^2_{right}: 1 - 0.90 = 0.05, \chi^2_{left}: 1 - 0.05 = 0.95, df: 25 - 1 = 24$$



$$\chi^2_{left} = 13.848 \quad \chi^2_{right} = 36.415$$

Example # 16

Cost of Ski Lift Tickets: Find the 90% confidence interval for the variance and standard deviation for the price in dollars of an adult single-day ski lift ticket. The data represent a selected sample of nationwide ski resorts. Assume the variable is normally distributed.

59	54	53	52	51
39	49	46	49	48

Degrees of Freedom	Area to the Right of Critical Value									
	0.990	0.980	0.970	0.960	0.950	0.940	0.930	0.920	0.910	0.900
1	0.455	0.833	0.983	0.990	0.994	0.997	0.998	0.999	0.9995	0.9999
2	0.872	1.323	1.475	1.532	1.567	1.591	1.610	1.626	1.640	1.653
3	1.323	1.735	1.922	1.992	2.053	2.103	2.143	2.173	2.203	2.233
4	1.773	2.179	2.361	2.441	2.507	2.563	2.613	2.653	2.693	2.733
5	2.223	2.616	2.792	2.862	2.923	2.973	3.013	3.053	3.093	3.133
6	2.673	3.063	3.238	3.308	3.368	3.418	3.463	3.503	3.543	3.583
7	3.123	3.508	3.683	3.753	3.813	3.863	3.908	3.953	3.993	4.033
8	3.573	3.953	4.123	4.193	4.253	4.303	4.353	4.403	4.453	4.503
9	4.023	4.393	4.558	4.623	4.683	4.733	4.783	4.823	4.873	4.923
10	4.473	4.833	5.003	5.063	5.123	5.173	5.223	5.273	5.323	5.373

$$n=10, CI: 90$$

$$\chi^2_{right}: 1 - 0.90 = 0.05, \chi^2_{left}: 1 - 0.05 = 0.95, df = 9$$

$$\chi^2_{left} = 16.919 \quad \chi^2_{right} = 3.825$$

$$\frac{(10-1)(28.2)^2}{16.919} < \sigma^2 < \frac{(10-1)(28.2)^2}{3.825}$$

$$15.0 < \sigma^2 < 76.3$$

$$S = \sqrt{\frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2}$$

$$= \sqrt{\frac{25254}{10} - \left(\frac{500}{10}\right)^2}$$

$$= 5.039$$

Example # 15

→ Nicotine Content: Find the 95% confidence interval for the variance and standard deviation of the nicotine content of cigarettes manufactured if a sample of 20 cigarettes has a standard deviation of 1.6 milligrams.

$$CI: 95\%, n=20, s.d.: 1.6$$

$$X_{\text{right}} = \frac{1-0.95}{2} = 0.025, X_{\text{left}} = 1 - 0.025 \cdot 0.975, df = 19$$

$$X_{\text{right}}^2 = 32.852, X_{\text{left}}^2 = 8.907$$

$$\frac{(20-1)(1.6)^2}{32.852} < \sigma^2 < \frac{(20-1)(1.6)^2}{8.907}$$

$$1.5 < \sigma^2 < 5.5$$

Degrees of Freedom	Chi-Square (χ^2) Distribution								
	0.995	0.99	0.875	0.90	0.95	0.10	0.05	0.025	0.01
1	—	0.001	0.000	0.000	2.706	3.890	6.631	9.821	13.821
2	0.000	0.000	0.000	0.000	5.991	9.238	12.877	16.917	20.977
3	0.002	0.001	0.000	0.000	7.814	11.347	15.787	20.277	24.788
4	0.007	0.002	0.001	0.000	10.821	15.779	19.488	23.497	27.371
5	0.021	0.014	0.007	0.000	14.811	20.536	25.885	31.271	34.960
6	0.057	0.034	0.017	0.000	18.856	25.645	32.592	39.869	44.852
7	0.098	0.060	0.030	0.000	22.879	30.807	38.852	47.870	52.376
8	0.160	0.110	0.050	0.000	27.878	37.753	47.387	57.208	63.980
9	0.230	0.160	0.080	0.000	32.877	43.773	54.807	66.953	73.993
10	0.304	0.214	0.115	0.000	38.876	51.753	64.889	78.053	85.986
11	0.380	0.367	0.214	0.000	45.875	60.733	75.907	90.187	100.396
12	0.457	0.521	0.321	0.000	53.874	70.732	87.906	105.186	123.395
13	0.530	0.675	0.427	0.000	62.873	81.731	101.905	122.185	142.394
14	0.600	0.729	0.531	0.000	72.872	93.730	116.904	140.184	162.393
15	0.667	0.882	0.635	0.000	82.871	105.729	132.903	160.183	182.392
16	0.730	0.935	0.738	0.000	92.870	117.728	146.902	176.182	198.391
17	0.787	0.988	0.841	0.000	102.869	127.727	157.901	187.181	209.390
18	0.840	1.041	0.944	0.000	112.868	142.726	172.895	203.180	225.389
19	0.887	1.094	1.047	0.000	122.867	152.725	182.894	213.179	235.388
20	0.930	1.147	1.149	0.000	132.866	162.724	192.893	223.178	245.387
21	0.967	1.199	1.252	0.000	142.865	172.723	202.892	233.177	255.386
22	1.000	1.251	1.354	0.000	152.864	182.722	212.891	243.176	265.385
23	1.030	1.294	1.456	0.000	162.863	192.721	222.890	253.175	275.384
24	1.057	1.342	1.558	0.000	172.862	202.720	232.889	263.174	285.383
25	1.083	1.385	1.660	0.000	182.861	212.719	242.888	273.173	295.382
26	1.100	1.427	1.762	0.000	192.860	222.718	252.887	283.172	305.381
27	1.120	1.469	1.864	0.000	202.859	232.717	262.886	293.171	315.380
28	1.136	1.511	1.966	0.000	212.858	242.716	272.885	303.170	325.379
29	1.153	1.553	2.068	0.000	222.857	252.715	282.884	313.169	335.378
30	1.168	1.595	2.169	0.000	232.856	262.714	292.883	323.168	345.377
31	1.183	1.637	2.271	0.000	242.855	272.713	302.882	333.167	355.376
32	1.197	1.679	2.373	0.000	252.854	282.712	312.881	343.166	365.375
33	1.210	1.721	2.475	0.000	262.853	292.711	322.880	353.165	375.374
34	1.223	1.763	2.576	0.000	272.852	302.710	332.879	363.164	385.373
35	1.236	1.805	2.678	0.000	282.851	312.709	342.878	373.163	395.372
36	1.248	1.847	2.779	0.000	292.850	322.708	352.877	383.162	405.371
37	1.260	1.889	2.881	0.000	302.849	332.707	362.876	393.161	415.370
38	1.272	1.931	2.983	0.000	312.848	342.706	372.875	403.160	425.369
39	1.284	1.973	3.084	0.000	322.847	352.705	382.874	413.159	435.368
40	1.295	2.015	3.186	0.000	332.846	362.704	392.873	423.158	445.367
41	1.306	2.057	3.287	0.000	342.845	372.703	402.872	433.157	455.366
42	1.317	2.099	3.388	0.000	352.844	382.702	412.871	443.156	465.365
43	1.328	2.141	3.489	0.000	362.843	392.701	422.870	453.155	475.364
44	1.338	2.183	3.590	0.000	372.842	402.700	432.869	463.154	485.363
45	1.348	2.225	3.691	0.000	382.841	412.700	442.868	473.153	495.362
46	1.358	2.267	3.792	0.000	392.840	422.700	452.867	483.152	505.361
47	1.368	2.309	3.893	0.000	402.839	432.700	462.866	493.151	515.360
48	1.378	2.351	3.994	0.000	412.838	442.700	472.865	503.150	525.359
49	1.388	2.393	4.095	0.000	422.837	452.700	482.864	513.149	535.358
50	1.397	2.435	4.196	0.000	432.836	462.700	492.863	523.148	545.357
51	1.406	2.477	4.297	0.000	442.835	472.700	502.862	533.147	555.356
52	1.415	2.519	4.398	0.000	452.834	482.700	512.861	543.146	565.355
53	1.424	2.561	4.499	0.000	462.833	492.700	522.860	553.145	575.354
54	1.433	2.603	4.599	0.000	472.832	502.700	532.859	563.144	585.353
55	1.442	2.645	4.699	0.000	482.831	512.700	542.858	573.143	595.352
56	1.451	2.687	4.799	0.000	492.830	522.700	552.857	583.142	595.351
57	1.460	2.729	4.899	0.000	502.829	532.700	562.856	593.141	605.350
58	1.469	2.771	4.999	0.000	512.828	542.700	572.855	603.140	615.349
59	1.478	2.813	5.099	0.000	522.827	552.700	582.854	613.139	625.348
60	1.486	2.855	5.199	0.000	532.826	562.700	592.853	623.138	635.347
61	1.495	2.897	5.299	0.000	542.825	572.700	602.852	633.137	645.346
62	1.504	2.939	5.399	0.000	552.824	582.700	612.851	643.136	655.345
63	1.513	2.981	5.499	0.000	562.823	592.700	622.850	653.135	665.344
64	1.522	3.023	5.599	0.000	572.822	602.700	632.849	663.134	675.343
65	1.530	3.065	5.699	0.000	582.821	612.700	642.848	673.133	685.342
66	1.539	3.107	5.799	0.000	592.820	622.700	652.847	683.132	695.341
67	1.548	3.149	5.899	0.000	602.819	632.700	662.846	693.131	705.340
68	1.556	3.191	5.999	0.000	612.818	642.700	672.845	703.130	715.339
69	1.565	3.233	6.099	0.000	622.817	652.700	682.844	713.129	725.338
70	1.573	3.275	6.199	0.000	632.816	662.700	692.843	723.128	735.337
71	1.582	3.317	6.299	0.000	642.815	672.700	702.842	733.127	745.336
72	1.590	3.359	6.399	0.000	652.814	682.700	712.841	743.126	755.335
73	1.598	3.401	6.499	0.000	662.813	692.700	722.840	753.125	765.334
74	1.606	3.443	6.599	0.000	672.812	702.700	732.839	763.124	775.333
75	1.614	3.485	6.699	0.000	682.811	712.700	742.838	773.123	785.332
76	1.622	3.527	6.799	0.000	692.810	722.700	752.837	783.122	795.331
77	1.630	3.569	6.899	0.000	702.809	732.700	762.836	793.121	805.330
78	1.638	3.611	6.999	0.000	712.808	742.700	772.835	803.120	815.329
79	1.646	3.653	7.099	0.000	722.807	752.700	782.834	813.119	825.328
80	1.654	3.695	7.199	0.000	732.806	762.700	792.833	823.118	835.327
81	1.662	3.737	7.299	0.000	742.805	772.700	802.832	833.117	845.326
82	1.670	3.779	7.399	0.000	752.804	782.700	812.831	843.116	855.325
83	1.678	3.821	7.499	0.000	762.803	792.700	822.830	853.115	865.324
84	1.686	3.863	7.599	0.000	772.802	802.700	832.829	863.114	875.323
85	1.694	3.905	7.699	0.000	782.801	812.700	842.828	873.113	885.322
86	1.702	3.947	7.799	0.000	792.800	822.700	852.827	883.112	895.321
87	1.710	3.989	7.899	0.000	802.800	832.700	862.826	893.111	905.320
88	1.718	4.031	7.999	0.000	812.800	842.700	872.825	903.110	915.319
89	1.726	4.073	8.099	0.000	822.800	852.700	882.824	913.109	925.318
90	1.734	4.115	8.199	0.000	832.800	862.700	892.823	923.108	935.317
91	1.742	4.157	8.299	0.000	842.800	872.700	902.822	933.107	945.316
92	1.750	4.199	8.399	0.000	852.800	882.700	912.821	943.106	955.315
93	1.758	4.241	8.499	0.000	862.800	892.700	922.820	953.105	965.314
94	1.766	4.283	8.599	0.000	872.800	902.700	932.819	963.104	975.313
95	1.774	4.325	8.699	0.000	882.800	912.700	942.818	973.103	985.312
96	1.782	4.367	8.799	0.000	892.800	922.700	952.817	983.102	995.311
97	1.790	4.409	8.899	0.000	902.800	932.700	962.816	993.101	1005.310
98	1.798	4.451	8.999	0.000	912.800	942.700	972.815	1003.100	1015.309
99	1.806	4.493	9.099	0.000	922.800	952.700	982.814	1013.100	1025.308
100	1.814	4.535	9.199	0.000	932.800	962.700	992.813	1023.100	1035.307

CHP 9

9.3 Many cardiac patients wear an implanted pacemaker to control their heartbeat. A plastic connector module mounts on the top of the pacemaker. Assuming a standard deviation of 0.0015 inch and an approximately normal distribution, find a 95% confidence interval for the mean of the depths of all connector modules made by a certain manufacturing company. A random sample of 75 modules has an average depth of 0.310 inch.

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

$$\delta = 0.0015 \quad CI = 95\%$$

$$n = 75 \quad \bar{x} = 0.310$$

$$\frac{1 - 0.95}{2} = 0.025$$

$$A_L = 0.025 + 0.95$$

$$= 0.975$$

$$P = 1.96$$

$$0.309 < \mu < 0.3103$$

9.6 How large a sample is needed in Exercise 9.2 if we wish to be 96% confident that our sample mean will be within 10 hours of the true mean?

$$n = ? \quad CI = 96\% \quad \bar{x} = 10 \quad \delta = 40 \quad P = 10$$

$$\frac{1 - 0.96}{2} = 0.02$$

$$E = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$A_L = 0.98$$

$$P = 2.055$$

$$n = \left(\frac{Z_{\alpha/2} \sigma}{E} \right)^2$$

$$n = \left(\frac{2.055 \times 40}{10} \right)^2 = 68 \text{ round up}$$

9.35 A random sample of size $n_1 = 25$, taken from a normal population with a standard deviation $\sigma_1 = 5$, has a mean $\bar{x}_1 = 80$. A second random sample of size $n_2 = 36$, taken from a different normal population with a standard deviation $\sigma_2 = 3$, has a mean $\bar{x}_2 = 75$. Find a 94% confidence interval for $\mu_1 - \mu_2$.

$$n_1 = 25 \quad \bar{x}_1 = 80 \quad \sigma_1 = 5 \quad CI = 94\%$$

$$n_2 = 36 \quad \bar{x}_2 = 75 \quad \sigma_2 = 3$$

$$\frac{1 - 0.94}{2} = 0.03$$

$$A_L = 0.91$$

$$P = 1.885$$

$$80 - 75 \pm 1.885 \sqrt{\frac{5^2}{25} + \frac{3^2}{36}}$$

$$2.9 < \mu_1 - \mu_2 < 7.1$$

HYPOTHESIS TESTING

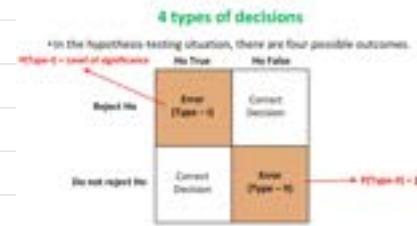
Null hypothesis: $H_0: \mu = 20$

Alternate Hypothesis: $H_1: \mu \neq 20$ or $\mu > 20$ or $\mu < 20$

level of significance = $\alpha: 0.05$ or 0.01 or 0.10

$$\Leftrightarrow P[\text{rejecting } H_0 \text{ when } H_0 \text{ is true}] = P[\text{Type I Error}]$$

β : Accepting false H_0 : Type II error



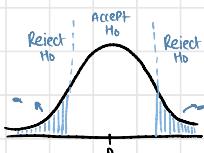
Statistical Tests

- ↳ z-test
- ↳ t-test
- ↳ f-test

TWO TAILED TEST

$$H_0: \mu = K$$

$$H_1: \mu \neq K$$



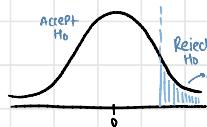
$$\alpha: \frac{1 - \text{confidence level}}{2}$$

one Tailed Test

Right tailed test

$$H_0: \mu = K$$

$$H_1: \mu > K$$

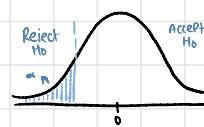


$$\alpha: 1 - \text{confidence level}$$

Left-tailed test

$$H_0: \mu = K$$

$$H_1: \mu < K$$



Z-Test for mean

↪ $n \geq 30$

↪ σ is known

↪ independent sample

↪ population is normally distributed

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

↑ sample mean ↑ H_0 mean
 ↓ standard deviation ↓ sample size

1 State H_0, H_1

$n \geq 30$

2 calculate z formula

3 find rejection region

4 if calculated z inside rejection region

↪ reject H_0

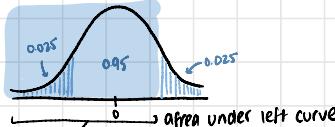
1. A factory has a machine that dispenses 80 ml. of fluid in a bottle. An employee believes the average amount of fluid is not 80 ml. Using 40 samples, he measures the average amount dispensed by the machine to be 78 ml. with a standard deviation of 2.5 ml. State the null and alternative hypothesis. ($\alpha = 0.05$, 95% confidence level) → there enough evidence to support the claim that the machine is not working properly?

$\bar{X} = 78, S = 2.5, n = 40$
 ↪ as $n > 30$ so z-test

① null hypotheses alternate hypotheses
 $H_0: \mu = 80$ $H_a: \mu \neq 80$

② $Z = \frac{\bar{X} - \mu_0}{S / \sqrt{n}} = \frac{78 - 80}{2.5 / \sqrt{40}} = -5.06$

③ $\alpha = \frac{1 - 0.95}{2} = 0.025$



$A_L = \text{confidence level} + \alpha$

$A_L = 0.95 + 0.025 = 0.975 \rightarrow$ find on Table

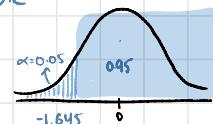
$Z_{\alpha/2} = 1.96$

0.005 + 0.99

③ area under right curve

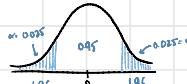
$A_R = 0.95 \rightarrow$ check table

$Z_{\alpha/2} = -1.645$



④ Z falls in rejection region

Reject H_0



④ Z falls in rejection region
 Reject H_0

so with 95% confidence we believe that
 machine is not working properly

Close to 10.045
 late outcome
 $\rightarrow 0.95 - 0.045$
 $\rightarrow -0.045$
 $\rightarrow -1.645$

Z-test for two means

L n > 30

↳ σ known

↳ independent sample

An admission test was administered to incoming freshmen in the College of Medical Laboratory Sciences and College of Radiologic Technology with 100 students each college randomly selected. The mean scores of the given samples were $\bar{x}_1 = 90$ and $\bar{x}_2 = 85$ and the variances of the test scores were 40 and 36 respectively. Is there a significant difference between the two groups? Use .01 level of significance.

$$\bar{X}_1 = 90 \quad \bar{X}_2 = 85 \quad n = 10$$

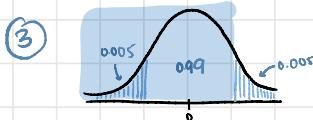
$$S^2 = 40 \quad S^2 = 35 \quad \alpha = 0.01$$

$$\textcircled{1} \quad H_0: \mu = 0$$

$$H_1: \mu \neq 0$$

$$\textcircled{2} \quad Z = \frac{(90 - 85) - (0 - 0)}{=} = 5.714$$

$$\frac{40^2}{100} + \frac{35^2}{100}$$



$$AL = 0.005 + 0.99 = 0.995$$

$$Z_{\frac{\alpha}{2}}^{\alpha} = 2.575$$

(4) Z falls in rejected region
 H_0 rejected

$$= \frac{0.07+0.08}{2} = 0.075$$

T-Test for mean

- ↳ variance > 1
- ↳ $n < 30$
- ↳ σ is known
- ↳ independent sample
- ↳ population is normally distributed

$$t = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

df: n-1
positive of freedom

sample mean
 H_0 mean
standard deviation
sample size

- 1 State H_0, H_1 , $n < 30$
- 2 calculated t formula
- 3 find rejection region
- 4 if calculated t inside rejection region
↳ reject H_0

2. A company manufactures car batteries with an average life span of 2 or more years. An engineer believes this value to be less. Using 10 samples, he measures the average life span to be 1.8 years with a standard deviation of 0.15. (a) State the null and alternative hypotheses.
(b) At a 99% confidence level, is there enough evidence to discard the null hypothesis?

① $H_0: \mu \geq 2$

$H_a: \mu < 2$

②

$\mu_0 = 2, \bar{x} = 1.8, s = 0.15, n = 10 \rightarrow$ as $n < 30$
so we use t test

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}} = \frac{1.8 - 2}{0.15/\sqrt{10}} = -4.22$$

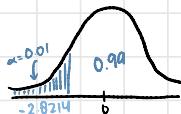
③ df: n-1

$= 10-1 = 9$

$\alpha = 1 - 0.99 = 0.01$

find on table

$t_{\frac{\alpha}{2}} = -2.8214$



④ t falls in rejection region

H_0 rejected

so with 99% confidence we believe that H_a

	0.25	0.20	0.15	0.10	0.08	0.05	0.03	0.02	0.01	0.005	0.001
1	1.378	1.380	1.378	1.374	1.371	1.368	1.362	1.359	1.353	1.349	1.313
2	0.878	1.081	1.386	1.386	2.600	4.200	4.848	8.865	9.923	22.33	
3	0.768	0.878	1.250	1.608	2.350	3.162	3.482	4.541	5.841	10.21	
4	0.741	0.941	1.180	1.533	2.132	2.778	2.998	3.747	4.804	7.173	
5	0.737	0.939	1.156	1.476	2.015	2.571	2.757	3.585	4.632	5.880	
6	0.718	0.905	1.134	1.440	1.945	2.487	2.612	3.403	3.707	5.208	
7	0.711	0.896	1.119	1.415	1.869	2.389	2.517	3.269	3.498	4.786	
8	0.706	0.886	1.104	1.387	1.800	2.306	2.449	3.186	3.338	4.301	
9	0.700	0.883	1.090	1.363	1.831	2.329	2.388	3.129	3.250	4.297	
10	0.705	0.879	1.080	1.352	1.812	2.328	2.386	3.129	3.249	4.148	
11	0.697	0.871	1.088	1.363	1.796	2.201	2.328	2.716	3.106	4.026	
12	0.690	0.873	1.083	1.368	1.782	2.179	2.300	2.681	3.098	3.900	
13	0.684	0.878	1.079	1.382	1.771	2.160	2.282	2.668	3.012	3.803	
14	0.680	0.868	1.078	1.345	1.761	2.145	2.264	2.646	2.977	3.787	

T-test for two means

↳ $n < 30$

↳ σ known

↳ independent samples

Karla grows tomatoes in two separate fields. When the tomatoes are ready to be picked, she is curious as to whether the sizes of her tomato plants differ between the two fields. She takes a random sample of plants from each field and measures the heights of the plants. Here is a summary of the results:

	Field A	Field B
Mean	1.3 m	1.6 m
Standard deviation	0.3 m	0.2 m
Number of plants	22	24

↳ $\sigma_1 \neq \sigma_2$ ^{variance}

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

sample mean H_0 mean
 \bar{x}_1 \bar{x}_2
 σ_1^2 σ_2^2
 standard deviation standard deviation
 sample size sample size

$$df_2 = \frac{(\delta_1 + \delta_2)^2}{\frac{\delta_1^2}{n-1} + \frac{\delta_2^2}{n-1}}$$

↳ $\sigma_1 = \sigma_2$ ^{variance}

$$\delta = \frac{(n-1)(\delta_1)^2 + (n-1)(\delta_2)^2}{(n+n)-2}$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\delta \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\textcircled{1} H_0: \mu = 0 \rightarrow \mu_1 = \mu_2$$

$$H_1: \mu \neq 0 \rightarrow \mu_1 \neq \mu_2$$

$$\textcircled{2} n_1 = 22 \quad \sigma_1 = 0.5 \quad \bar{x}_1 = 1.3$$

assume (as not mentioned)
 $\alpha = 0.05$

$$n_2 = 24 \quad \sigma_2 = 0.3 \quad \bar{x}_2 = 1.6$$

$$t = \frac{(1.3 - 1.6)(0-0)}{\sqrt{\frac{0.5^2}{22} + \frac{0.3^2}{24}}} = -2.4402$$

$$\textcircled{3} df_1 = 21 \quad df_2 = 23 \quad \alpha = 0.05/2 = 0.025$$

take $t_{\alpha/2}$ for the smallest df

$$t_{\alpha/2} = 2.080$$

t lies in rejection region

H_0 rejected

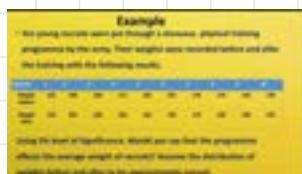


T-test for two means

$\hookrightarrow n < 30$

$\hookrightarrow \sigma$ unknown

\hookrightarrow dependent samples



$$\textcircled{1} \quad H_0: \mu = 0 \quad H_1: \mu \neq 0$$

$$\textcircled{2} \quad n=10, \alpha=0.05$$

$$\sum D = (125-136) + (145-201) + (160-158) + (171-184) + (140-145) + (201-195) + (170-175) + (176-190) + (195-190) + (139-145) = -47$$

$$\bar{x} = \frac{-47}{10} = 4.7$$

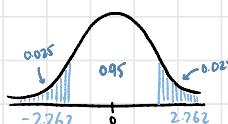
$$\sum D^2 = (125-136)^2 + (145-201)^2 + (160-158)^2 + (171-184)^2 + (140-145)^2 + (201-195)^2 + (170-175)^2 + (176-190)^2 + (195-190)^2 + (139-145)^2 = 673$$

$$S.D. = \sqrt{\frac{10(673) - (47)^2}{10(10-1)}} = 7.09$$

$$\textcircled{3} \quad t = \frac{4.7 - 0}{7.09 / \sqrt{10}} = 2.09$$

$$\textcircled{4} \quad df = 9, \alpha = \frac{0.05}{2} = 0.025 \\ t_{\frac{\alpha}{2}} = 2.262$$

$\textcircled{5}$ t lies in acceptance region
H₀ accept



$$1. D = X_1 - X_2 \quad 2. D^2 = (X_1 - X_2)^2$$

$$3. \bar{X} = \frac{D}{n} \quad 4. S_D = \sqrt{\frac{n \sum D^2 - (\sum D)^2}{n(n-1)}}$$

$$5. t = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \quad 6. df = n-1$$

$$\hookrightarrow \sigma_1 = \sigma_2 \xrightarrow{\text{variance}}$$

$$8. \frac{(n-1)(s_1)^2 + (n-1)(s_2)^2}{(n+n)-2}$$

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\delta \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\delta = \sqrt{\frac{n \sum x^2 - (\sum x)^2}{n(n-1)}}$$

	0.25	0.20	0.15	0.10	0.05	0.025	0.01	0.005	0.001	
1	1.331	1.333	1.337	1.344	1.371	1.388	1.414	1.436	1.466	1.513
2	0.818	1.261	1.386	1.486	2.020	4.003	4.889	8.865	9.823	29.33
3	0.760	0.871	1.250	1.609	2.253	3.982	3.482	4.541	5.841	10.21
4	0.741	0.844	1.180	1.533	2.133	2.778	2.988	3.747	4.854	7.173
5	0.727	0.823	1.156	1.476	2.015	2.571	2.757	3.365	4.032	5.883
6	0.718	0.805	1.134	1.446	1.945	2.447	2.612	3.143	3.707	5.208
7	0.711	0.800	1.118	1.415	1.889	2.317	2.497	3.049	3.699	4.799
8	0.706	0.789	1.104	1.387	1.860	2.206	2.449	2.886	3.338	4.301
9	0.700	0.783	1.090	1.363	1.833	2.169	2.382	2.825	3.250	4.287
10	0.700	0.779	1.080	1.352	1.812	2.159	2.358	2.764	3.169	4.144
11	0.697	0.776	1.088	1.363	1.798	2.201	2.331	2.716	3.108	4.025
12	0.690	0.773	1.083	1.362	1.782	2.179	2.303	2.681	3.098	3.900
13	0.684	0.770	1.079	1.350	1.771	2.160	2.282	2.660	2.912	3.802
14	0.692	0.768	1.078	1.345	1.761	2.146	2.264	2.636	2.877	3.787

CHP 10

10.20 A random sample of 64 bags of white cheddar popcorn weighed, on average, 5.23 ounces with a standard deviation of 0.24 ounce. Test the hypothesis that $\mu = 5.5$ ounces against the alternative hypothesis, $\mu < 5.5$ ounces, at the 0.05 level of significance.

$$n=64 \quad \bar{x}=5.23 \quad \sigma=0.24 \quad \alpha=0.05$$

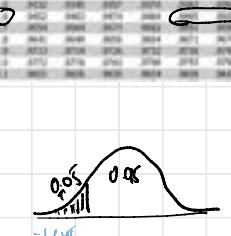
$$H_0: \mu = 5.5$$

$$H_1: \mu < 5.5$$

$$Z = \frac{5.23 - 5.5}{0.24/\sqrt{64}} = -9$$

H_0 reject $\mu < 5.5$

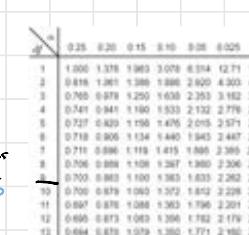
0.0	3000	3049	3098	3137	3176	3215	3254	3293	3332	3371	3410	3449
0.1	3000	3049	3098	3137	3176	3215	3254	3293	3332	3371	3410	3449
0.2	3180	3140	3101	3061	3021	2981	2941	2901	2861	2821	2781	2741
0.3	3270	3237	3205	3163	3123	3081	3040	2998	2957	2916	2874	2833
0.4	3360	3326	3293	3251	3209	3178	3135	3093	3052	3010	2968	2926
0.5	3450	3416	3382	3349	3307	3276	3233	3191	3149	3107	3065	3023
0.6	3540	3505	3471	3437	3395	3362	3320	3278	3236	3194	3152	3110
0.7	3630	3593	3559	3524	3481	3447	3404	3361	3318	3275	3232	3189
0.8	3720	3690	3655	3620	3576	3532	3488	3444	3399	3355	3311	3267
0.9	3810	3786	3750	3714	3669	3624	3579	3534	3488	3443	3398	3353
1.0	3900	3880	3843	3806	3759	3714	3668	3622	3575	3529	3483	3437
1.1	3990	3980	3939	3892	3842	3792	3742	3691	3641	3590	3539	3489
1.2	4080	4070	4028	3978	3928	3878	3827	3776	3725	3674	3623	3572
1.3	4170	4160	4118	4067	4017	3967	3916	3865	3814	3763	3712	3661
1.4	4260	4250	4208	4157	4107	4056	4005	3954	3903	3852	3801	3750
1.5	4350	4340	4298	4247	4196	4145	4094	4043	3992	3941	3890	3839
1.6	4440	4430	4388	4337	4286	4235	4184	4133	4082	4031	3980	3929
1.7	4530	4520	4478	4427	4376	4325	4274	4223	4172	4121	4070	4019
1.8	4620	4610	4568	4517	4466	4415	4364	4313	4262	4211	4160	4109
1.9	4710	4700	4658	4607	4556	4505	4454	4403	4352	4291	4240	4189
2.0	4800	4790	4748	4697	4646	4595	4544	4493	4442	4391	4340	4289
2.1	4890	4880	4838	4787	4736	4685	4634	4583	4532	4481	4430	4379
2.2	4980	4970	4928	4877	4826	4775	4724	4673	4622	4571	4520	4469
2.3	5070	5060	5018	4967	4916	4865	4814	4763	4712	4661	4610	4559
2.4	5160	5150	5108	5057	5006	4955	4904	4853	4802	4751	4699	4648
2.5	5250	5240	5198	5147	5096	5045	5094	5043	5092	5041	5090	5039
2.6	5340	5330	5288	5237	5186	5135	5084	5033	5082	5031	5080	5029
2.7	5430	5420	5378	5327	5276	5225	5174	5123	5072	5021	5070	5019
2.8	5520	5510	5468	5417	5366	5315	5264	5213	5162	5111	5060	5009
2.9	5610	5600	5558	5507	5456	5405	5354	5303	5252	5201	5150	5099
3.0	5700	5690	5648	5597	5546	5495	5444	5393	5342	5291	5240	5189
3.1	5790	5780	5738	5687	5636	5585	5534	5483	5432	5381	5330	5279
3.2	5880	5870	5828	5777	5726	5675	5624	5573	5522	5471	5420	5369
3.3	5970	5960	5918	5867	5816	5765	5714	5663	5612	5561	5510	5459
3.4	6060	6050	6008	5957	5906	5855	5804	5753	5702	5651	5600	5549
3.5	6150	6140	6098	6047	6096	6045	6094	6043	6092	6041	6090	6039
3.6	6240	6230	6188	6137	6086	6035	6084	6033	6082	6031	6080	6029
3.7	6330	6320	6278	6227	6176	6125	6074	6023	6072	6021	6070	6059
3.8	6420	6410	6368	6317	6266	6215	6164	6113	6162	6111	6160	6149
3.9	6510	6500	6458	6407	6356	6305	6254	6203	6252	6201	6250	6239
4.0	6600	6590	6548	6497	6446	6395	6344	6293	6342	6291	6340	6329
4.1	6690	6680	6638	6587	6536	6485	6434	6383	6432	6381	6430	6419
4.2	6780	6770	6728	6677	6626	6575	6524	6473	6522	6471	6520	6509
4.3	6870	6860	6818	6767	6716	6665	6614	6563	6612	6561	6610	6599
4.4	6960	6950	6908	6857	6806	6755	6704	6653	6702	6651	6700	6689
4.5	7050	7040	7098	7047	7096	7045	7094	7043	7092	7041	7090	7079
4.6	7140	7130	7188	7137	7185	7133	7182	7131	7180	7139	7178	7167
4.7	7230	7220	7278	7227	7275	7225	7273	7223	7272	7221	7270	7259
4.8	7320	7310	7368	7317	7365	7315	7363	7313	7362	7311	7360	7349
4.9	7410	7400	7458	7447	7445	7443	7441	7440	7439	7438	7437	7436
5.0	7500	7490	7548	7537	7535	7533	7531	7530	7529	7528	7527	7526
5.1	7590	7580	7638	7627	7625	7623	7621	7620	7619	7618	7617	7616
5.2	7680	7670	7728	7717	7715	7713	7711	7710	7709	7708	7707	7706
5.3	7770	7760	7818	7807	7805	7803	7801	7800	7809	7808	7807	7806
5.4	7860	7850	7908	7897	7895	7893	7891	7890	7889	7888	7887	7886
5.5	7950	7940	7998	7987	7985	7983	7981	7980	7979	7978	7977	7976
5.6	8040	8030	8088	8077	8075	8073	8071	8070	8069	8068	8067	8066
5.7	8130	8120	8178	8167	8165	8163	8161	8160	8159	8158	8157	8156
5.8	8220	8210	8268	8257	8255	8253	8251	8250	8249	8248	8247	8246
5.9	8310	8300	8358	8347	8345	8343	8341	8340	8339	8338	8337	8336
6.0	8390	8380	8438	8427	8425	8423	8421	8420	8419	8418	8417	8416
6.1	8480	8470	8528	8517	8515	8513	8511	8510	8509	8508	8507	8506
6.2	8570	8560	8618	8607	8605	8603	8601	8600	8609	8608	8607	8606
6.3	8660	8650	8708	8697	8695	8693	8691	8690	8689	8688	8687	8686
6.4	8750	8740	8798	8787	8785	8783	8781	8780	8779	8778	8777	8776
6.5	8840	8830	8888	8877	8875	8873	8871	8870	8869	8868	8867	8866
6.6	8930	8920	8978	8967	8965	8963	8961	8960	8959	8958	8957	8956
6.7	9020	9010	9068	9057	9055	9053	9051	9050	9049	9048	9047	9046
6.8	9110	9100	9158	9147	9145	9143	9141	9140	9139	9138	9137	9136
6.9	9200	9190	9248	9237	9235	9233	9231	9230	9229	9228	9227	9226
7.0	9290	9280	9338	9327	9325	9323	9321	9320	9319	9318	9317	9316
7.1	9380	9370	9428	9417	9415	9413	9411	9410	9409	9408	9407	9406
7.2	9470	9460	9518	9507	9505	9503	9501	9500	9499	9498	9497	9496
7.3	9560	9550	9608	9597	9595	9593	9591	9590	9589	9588	9587	9586
7.4	9650	9640	9698	9687	9685	9683	9681	9680	9679	9678	9677	9676
7.5	9740	9730	9788	9777	9775	9773	9771	9770	9769	9768	9767	9766
7.6	9830	9820	9878	9867	9865	9863	9861	9860	9859	9858	9857	9856
7.7	9920	9910	9968	9957	9955	9953	9951	9950	9949	9948	9947	9946
7.8	10010	1000	10058	10047	10045	10043	10041	10040	10039	10038	10037	10036



$$\alpha = 0.05$$

$$\text{from table } = 1.645$$

$$P = 1.645$$



$$P = 1.96$$

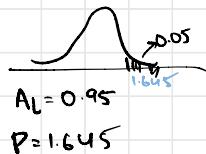
H_0 accept

- H₁
10.25 It is claimed that automobiles are driven on average more than 20,000 kilometers per year. To test this claim, 100 randomly selected automobile owners are asked to keep a record of the kilometers they travel. Would you agree with this claim if the random sample showed an average of 23,500 kilometers and a standard deviation of 3900 kilometers? Use a *P*-value in your conclusion.

$$H_0: \mu = 20,000 \quad n=100 \quad \bar{x} = 23500 \quad \delta = 3900$$

$$H_1: \mu > 20,000 \quad \alpha = 0.05 \text{ as not mentioned}$$

$$z = \frac{23500 - 20000}{3900 / \sqrt{100}} = 8.974$$



H_0 rejected
 $\mu > 20,000$

- 10.27** A study at the University of Colorado at Boulder shows that running increases the percent resting metabolic rate (RMR) in older women. The average RMR of 30 elderly women runners was 34.0% higher than the average RMR of 30 sedentary elderly women, and the standard deviations were reported to be 10.5 and 10.2%, respectively. Was there a significant increase in RMR of the women runners over the sedentary women? Assume the populations to be approximately normally distributed with equal variances. Use a *P*-value in your conclusions.

Decision: Reject H_0 and claim $\mu > 220$ milligrams.

30.27 The hypothesis are:

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 > \mu_2$$

Since $s_p = \sqrt{\frac{(20.0)(10.5)^2 + (20.0)(10.2)^2}{59}} = 10.25$, then

$$P\left[Z > \frac{34.0}{10.25 \sqrt{1/30 + 1/30}}\right] = P(Z > 12.72) \approx 0.$$

Hence, the conclusion is that running increases the mean RMR in older women.

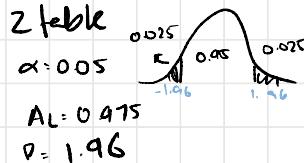
$$H_0: \mu_1 = \mu_2 \quad n=30 \quad \bar{x}_1 - \bar{x}_2 = 34 \quad \delta_1 = 10.5$$

$$H_1: \mu_1 \neq \mu_2 \quad \delta_2 = 10.2$$

$$\delta = \frac{(29)(10.5)^2 - (29)(10.2)^2}{30+30-2} = 10.35$$

$$z = \frac{34}{10.35 / \sqrt{\frac{1}{30} + \frac{1}{30}}} = 0.8481$$

H_0 accept



10.35 To find out whether a new serum will arrest leukemia, 9 mice, all with an advanced stage of the disease, are selected. Five mice receive the treatment and 4 do not. Survival times, in years, from the time the experiment commenced are as follows:

Treatment	2.1	5.3	1.4	4.6	0.9
No Treatment	1.9	0.5	2.8	3.1	

At the 0.05 level of significance, can the serum be said to be effective? Assume the two populations to be normally distributed with equal variances.

T test & unknown

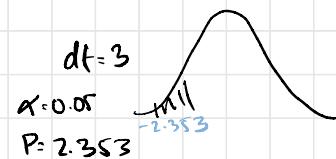
$$n_1 = 5 \quad \bar{X}_1 = \frac{14.3}{5} = 2.86 \quad S_1 = \frac{\sqrt{(56.43) - (14.3)^2}}{\sqrt{4}} = 3.883$$

$$n_2 = 4 \quad \bar{X}_2 = \frac{8.3}{4} = 2.075 \quad S_2 = \frac{\sqrt{(21.91) - (8.3)^2}}{\sqrt{3}} = 1.3625$$

$$X_1 - X_2 = 0.785$$

$$S = \frac{3(1.3625)^2 + 4(3.883)^2}{5+4-2} = 1.674$$

$$t = \frac{0.785}{1.674 / \sqrt{\frac{1}{5} + \frac{1}{4}}} = -0.70$$



$$H_0: \mu_1 - \mu_2 = 0$$

$$H_1: \mu_1 - \mu_2 < 0$$

H_0 accept

normal with equal variances.

10.39 The following data represent the running times of films produced by two motion-picture companies:

Company	Time (minutes)					
1	102	86	98	109	92	
2	81	165	97	134	92	87

Test the hypothesis that the average running time of films produced by company 2 exceeds the average running time of films produced by company 1 by 10 minutes against the one-sided alternative that the difference is less than 10 minutes. Use a 0.1 level of significance and assume the distributions of times to be approximately normal with unequal variances.

10.39 The hypotheses are

$$H_0: \mu_2 = \mu_1 + 10, \quad H_1: \mu_2 > \mu_1 + 10$$

$$\alpha = 0.1$$

Degrees of freedom is calculated as

$$v = \frac{(78.6/3 + 94.3/7)^2}{(78.6/3)^2/4 + (94.3/7)^2/6} = 7.38$$

hence we use 2 degrees of freedom with the critical region $t > 1.31$.

Computation: $t = \frac{(81 - 102) - 10}{\sqrt{78.6/3 + 94.3/7}} = 0.22$

Decision: Fail to reject H_0 .

ANALYSIS OF VARIANCE (ANOVA)

F test

- ↳ independent samples
- ↳ normal distribution
- ↳ σ are same

not t-test as the more means, the more t-tests needed

1. State H_0, H_a, α
2. Degree's of freedom
↳ df between groups = $K - 1$
↳ df within groups = $N - K$ <small>→ no of groups total no of observations</small>
↳ diff total = $N - 1$
↳ F_{critical} <small>→ using table</small>
3. Sum of square deviations from mean
↳ each groups mean $\rightarrow \bar{x}_1, \bar{x}_2, \bar{x}_3$
↳ Grandmean $\rightarrow \bar{\bar{x}}$
↳ $SS_{\text{total}} = \sum (x - \bar{\bar{x}})^2$
↳ $SS_{\text{within}} = \sum (x - \bar{x})^2$
↳ $SS_{\text{between}} = SS_{\text{total}} - SS_{\text{within}}$
4. Mean Square
↳ $MS_{\text{between}} = \frac{SS_{\text{between}}}{df_{\text{between}}}$
↳ $MS_{\text{within}} = \frac{SS_{\text{within}}}{df_{\text{within}}}$
5. Calculate F statistic
↳ $F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$
↳ conclusion
reject H_0
F ratio > F critical
accept H_0
F ratio < F critical
F ratio > 1
F ratio ≤ 1

ANALYSIS OF VARIANCE (ANOVA)

E.G

	A	B	C
1	2	2	2
2	4	3	
3	5	2	4
mean	2.67	2.67	3

$\alpha = 0.5 \rightarrow$ generally

(1)

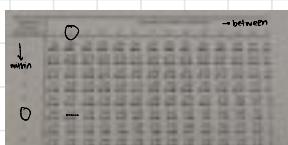
$$H_0 \rightarrow \mu_1 = \mu_2 = \mu_3 \rightarrow \text{no diff in means}$$

$$H_a \rightarrow \text{at least 1 diff in the means}$$

$$(2) df_{\text{between}} = K-1 = 3-1 = 2$$

$$df_{\text{within}} = N - K = 9 - 3 = 6$$

$$\text{diff total} = 6 + 2 = 8$$



$$F_{\text{critical}} = 5.14$$

$$(3) \bar{x}_1 = 2.67, \bar{x}_2 = 2.67, \bar{x}_3 = 3, \bar{\bar{x}} = \frac{25}{9} = 2.78$$

$$\begin{aligned} SS_{\text{total}} &= \sum (x - \bar{\bar{x}})^2 \\ &= (1-2.78)^2 + (2-2.78)^2 + (5-2.78)^2 + \\ &\quad (2-2.78)^2 + (4-2.78)^2 + (2-2.78)^2 + \\ &\quad (2-2.78)^2 + (3-2.78)^2 + (4-2.78)^2 = 13.6 \end{aligned}$$

$$\begin{aligned} SS_{\text{within}} &= \sum (x_n - \bar{x}_n)^2 \\ &= (1-2.67)^2 + (2-2.67)^2 + (5-2.67)^2 + \\ &\quad (2-2.67)^2 + (4-2.67)^2 + (2-2.67)^2 + \\ &\quad (2-3)^2 + (3-3)^2 + (4-3)^2 = 13.34 \end{aligned}$$

$$SS_{\text{between}} = SS_{\text{total}} - SS_{\text{within}}$$

$$= 13.6 - 13.34 = 0.23$$

$$(4) MS_{\text{between}} = \frac{SS_{\text{between}}}{df_{\text{between}}} = \frac{0.23}{2} = 0.12$$

$$MS_{\text{within}} = \frac{SS_{\text{within}}}{df_{\text{within}}} = \frac{13.34}{6} = 2.22$$

1. State H_0, H_a, α

2. Degree's of freedom

$$\hookrightarrow df_{\text{between groups}} = K-1$$

$$\hookrightarrow df_{\text{within groups}} = N - K \xrightarrow{\text{no. of groups}}$$

$$\hookrightarrow \text{diff total} = N - 1$$

$$\hookrightarrow F_{\text{critical}} \rightarrow \text{using table}$$

3. Sum of square Deviations from mean

$$\hookrightarrow \text{each groups mean} \rightarrow \bar{x}_1, \bar{x}_2, \bar{x}_3$$

$$\hookrightarrow \text{Grand mean} \rightarrow \bar{\bar{x}}$$

$$\hookrightarrow SS_{\text{total}} = \sum (x - \bar{\bar{x}})^2$$

$$\hookrightarrow SS_{\text{within}} = \sum (x - \bar{x})^2$$

$$\hookrightarrow SS_{\text{between}} = SS_{\text{total}} - SS_{\text{within}}$$

4. Mean Square

$$\hookrightarrow MS_{\text{between}} = \frac{SS_{\text{between}}}{df_{\text{between}}}$$

$$\hookrightarrow MS_{\text{within}} = \frac{SS_{\text{within}}}{df_{\text{within}}}$$

5. Calculate F statistic

$$\hookrightarrow F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$

conclusion

reject H_0

$$\text{F ratio} > \text{F critical}$$

$$\text{F ratio} > 1$$

accept H_0

$$\text{F ratio} < \text{F critical}$$

$$\text{F ratio} \leq 1$$

(5)

$$F = \frac{MS_{\text{between}}}{MS_{\text{within}}} = \frac{0.12}{2.22} = 0.05 \rightarrow <$$

$$F_{\text{critical}} = 5.14$$

$$0.05 < 5.14 \text{ and } < 1$$

hence accept H_0

CHP 13

13.1 Six different machines are being considered for use in manufacturing rubber seals. The machines are being compared with respect to tensile strength of the product. A random sample of four seals from each machine is used to determine whether the mean tensile strength varies from machine to machine. The following are the tensile-strength measurements in kilograms per square centimeter $\times 10^{-2}$:

Machine					
1	2	3	4	5	6
17.5	16.4	20.3	14.6	17.5	18.3
16.9	19.2	15.7	16.7	19.2	16.2
15.8	17.7	17.8	20.8	16.5	17.5
18.6	15.4	18.9	18.9	20.5	20.1

Perform the analysis of variance at the 0.05 level of significance and indicate whether or not the mean tensile strengths differ significantly for the six machines.

$$\textcircled{1} \quad H_0 \rightarrow \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5 = \mu_6$$

$$H_1 \rightarrow \text{at least 1 } \mu \text{ is different}$$

$$\alpha = 0.05 \quad df_b = 5 \quad P = 2.77$$

$$df_w = 18$$

$$\textcircled{2} \quad df_b = k - 1 = 5$$

$$df_w = N - k = 24 - 6 = 18$$

$$df_t = 23$$

$$\textcircled{3} \quad \begin{array}{ccccccc} \mu_1 & \mu_2 & \mu_3 & \mu_4 & \mu_5 & \mu_6 & \mu_T \\ 17.2 & 17.175 & 18.175 & 17.75 & 18.425 & 18.025 & 17.79 \end{array}$$

$$\textcircled{4} \quad SS_T = 5.5003 + 9.6404 + 11.9004 + 21.6564 + 11.0804 + 8.2084 \\ = 67.9863$$

$$SS_W = 4.1 + 8.1275 + 11.3075 + 21.65 + 9.4675 + 7.9815 \\ = 62.64$$

$$S_b = 67.9863 - 62.64 \\ = 5.3463$$

$$f = \frac{1.06926}{3.48} = 0.3072 < 1 \text{ and } < 2.77 \\ \text{accept } H_0$$

$$\textcircled{5} \quad MS_b = \frac{5.3463}{5} = 1.06926$$

$$MS_w = \frac{62.64}{18} = 3.48$$

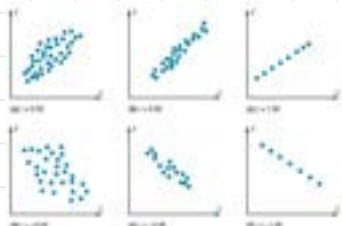
CORRELATION & REGRESSION

does a rs exist between variables

Scatter Plot

x-axis → independent variable can be controlled

y-axis → dependent variable can't be controlled



correlation coefficient

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

no of data pairs

Hypothesis testing for correlation

$$t = r \sqrt{\frac{n-2}{1-r^2}}$$

$$H_0: \rho = 0$$

$$df = n-2$$

Correlation Coefficient (PPMC)

↳ if $r=0$ then no correlation

↳ if $r \neq 0$ then correlation ($r > 0, n \uparrow \rightarrow y \uparrow$)

Type of rs

1. direct cause and effect rs (x causes y)

2. reverse cause and effect rs (y causes x)

3. caused by a third variable rs

4. curvilinear rs



1. find r

2. state H_0, H_1

3. find t

4. find rejection region

5. If t lies in rejection region

↳ H_0 rejected

6. Type of rs

Example #05

- Test the significance of the correlation coefficient found in Example # 6.

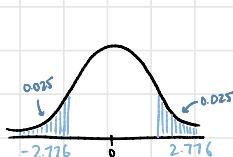
Subject	Age, y	Protein, g	n	χ^2	P
A	1.0-1.5	1.56	3,000	1,040	0.99
B	1.0-1.5	1.50	3,000	2,560	13.98
C	1.0-1.5	1.50	3,000	3,116	0.02
D	1.0-1.5	1.40	4,725	7,771	0.44
E	1.0-1.5	1.40	6,617	6,889	0.08
F	1.0-1.5	1.50	5,000	3,980	0.00
G	1.0-1.5	1.50	5,000	4,000	0.00
H	1.0-1.5	1.50	5,000	4,000	0.00
I	1.0-1.5	1.50	5,000	4,000	0.00
J	1.0-1.5	1.50	5,000	4,000	0.00
K	1.0-1.5	1.50	5,000	4,000	0.00
L	1.0-1.5	1.50	5,000	4,000	0.00
M	1.0-1.5	1.50	5,000	4,000	0.00
N	1.0-1.5	1.50	5,000	4,000	0.00
O	1.0-1.5	1.50	5,000	4,000	0.00
P	1.0-1.5	1.50	5,000	4,000	0.00
Q	1.0-1.5	1.50	5,000	4,000	0.00
R	1.0-1.5	1.50	5,000	4,000	0.00
S	1.0-1.5	1.50	5,000	4,000	0.00
T	1.0-1.5	1.50	5,000	4,000	0.00
U	1.0-1.5	1.50	5,000	4,000	0.00
V	1.0-1.5	1.50	5,000	4,000	0.00
W	1.0-1.5	1.50	5,000	4,000	0.00
X	1.0-1.5	1.50	5,000	4,000	0.00
Y	1.0-1.5	1.50	5,000	4,000	0.00
Z	1.0-1.5	1.50	5,000	4,000	0.00
A'	1.6-2.0	1.50	3,000	1,040	0.99
B'	1.6-2.0	1.50	3,000	2,560	13.98
C'	1.6-2.0	1.50	3,000	3,116	0.02
D'	1.6-2.0	1.40	4,725	7,771	0.44
E'	1.6-2.0	1.40	6,617	6,889	0.08
F'	1.6-2.0	1.50	5,000	3,980	0.00
G'	1.6-2.0	1.50	5,000	4,000	0.00
H'	1.6-2.0	1.50	5,000	4,000	0.00
I'	1.6-2.0	1.50	5,000	4,000	0.00
J'	1.6-2.0	1.50	5,000	4,000	0.00
K'	1.6-2.0	1.50	5,000	4,000	0.00
L'	1.6-2.0	1.50	5,000	4,000	0.00
M'	1.6-2.0	1.50	5,000	4,000	0.00
N'	1.6-2.0	1.50	5,000	4,000	0.00
O'	1.6-2.0	1.50	5,000	4,000	0.00
P'	1.6-2.0	1.50	5,000	4,000	0.00
Q'	1.6-2.0	1.50	5,000	4,000	0.00
R'	1.6-2.0	1.50	5,000	4,000	0.00
S'	1.6-2.0	1.50	5,000	4,000	0.00
T'	1.6-2.0	1.50	5,000	4,000	0.00
U'	1.6-2.0	1.50	5,000	4,000	0.00
V'	1.6-2.0	1.50	5,000	4,000	0.00
W'	1.6-2.0	1.50	5,000	4,000	0.00
X'	1.6-2.0	1.50	5,000	4,000	0.00
Y'	1.6-2.0	1.50	5,000	4,000	0.00
Z'	1.6-2.0	1.50	5,000	4,000	0.00

$$H_0: p = 0 \quad H_1: p \neq 0$$

$$r = \frac{6(47,634) - (345)(819)}{\sqrt{[6(20399) - (345)^2][6(112,443) - (819)^2]}} = 0.897$$

$$t = \frac{6 - 2}{1 - 0.897} = 4059$$

t lies inside rejection region
H₀ rejected



<i>d</i>	0.25	0.30	0.15	0.10	0.05	0.025	0.01	0.005	0.001	
1	1.000	1.378	1.983	3.078	6.354	12.71	19.88	31.62	63.66	318.3
2	0.818	1.261	1.386	1.886	2.620	4.230	4.848	6.865	9.925	22.33
3	0.765	0.979	1.250	1.638	2.253	3.182	3.482	4.541	5.841	10.21
4	0.746	0.849	1.160	1.933	2.120	2.484	2.998	3.747	4.654	7.173
5	0.737	0.829	1.156	1.676	2.015	2.571	2.737	3.262	4.032	5.893
6	0.718	0.805	1.134	1.440	1.845	2.447	2.612	3.143	3.707	5.208
7	0.711	0.806	1.118	1.415	1.848	2.389	2.517	2.988	3.499	4.795
8	0.706	0.809	1.108	1.387	1.860	2.306	2.446	2.886	3.328	4.501
9	0.700	0.803	1.100	1.363	1.823	2.262	2.398	2.821	3.250	4.287
10	0.700	0.807	1.090	1.372	1.852	2.229	2.358	2.704	3.139	4.144
11	0.697	0.810	1.088	1.363	1.819	2.201	2.328	2.718	3.106	4.025
12	0.690	0.813	1.083	1.356	1.812	2.179	2.300	2.681	3.059	3.900
13	0.684	0.816	1.079	1.350	1.771	2.186	2.282	2.688	3.012	3.852
14	0.692	0.808	1.078	1.345	1.781	2.145	2.264	2.639	2.977	3.767

ANSWERING YOUR QUESTIONS

In a study on spread control, it was found that the mean distance for acquisitions must be limited to 10 miles. These values sufficient and to minimize the risk of change. An area that was measured can be the study was that distance required to completely stop a vehicle at warehouse spread. Use the following table to answer the questions.

MPH	Breaking distance (feet)
20	209
30	419
40	619
50	819
60	1019
70	1219

Assume NAPF is going to be used to predict mapping distance

1. Which of the two variables is the independent variable?
 2. Which is the dependent variable?
 3. What type of variable is the independent variable?
 4. What type of variable is the dependent variable?
 5. Construct a scatter plot for the data.
 6. Is there a linear relationship between the two variables?
 7. Is the relationship positive or negative?
 8. Can breaking distance be accurately predicted from MPH?
 9. List some other variables that affect breaking distance.
 10. Compute the value of r.
 11. Is r significant at $\alpha = 0.05$?

Simple Linear Regression

$$y = mn + c$$

↓
regression

Example # 06

- Find the equation of the regression line for the data in Example # 01 (Slide 15), and graph the line on the scatter plot of the data.

Subject	Age, x	Pressure, y	x^2	x^3
A	47	128	5,764	3,449
B	48	129	5,760	3,436
C	56	138	7,796	5,136
D	49	135	5,761	3,721
E	47	131	5,761	4,09
F	76	172	59,840	4,900

$\Sigma x = 367 \quad \Sigma y = 819 \quad \Sigma xy = 47,634 \quad \Sigma x^2 = 26,596 \quad \Sigma x^3 = 122,687$

$$\begin{aligned} \bar{x} &= \frac{\sum x}{n} = \frac{\sum x}{6} = \frac{367}{6} = 61.167 \\ \bar{y} &= \frac{\sum y}{n} = \frac{\sum y}{6} = \frac{819}{6} = 136.5 \\ b_0 &= \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2} = \frac{6(47,634) - (367)(819)}{6(26,596) - (367)^2} = 81.048 \\ b_1 &= \frac{n(\sum x^2) - (\sum x)(\sum x)}{n(\sum x^3) - (\sum x)^2} = \frac{6(26,596) - (367)^2}{6(122,687) - (367)^2} = 0.964 \\ y' &= 81.048 + 0.964x \end{aligned}$$

sign of the correlation coefficient and the sign of the slope of the regression line

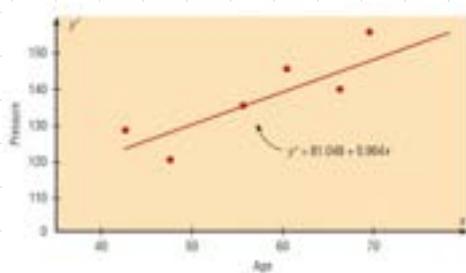
Prediction Using Regression Eq.

- Using the equation of the regression line found in Example # 06, predict the blood pressure for a person who is 50 years old.

Substituting 50 for x in the regression line $y' = 81.048 + 0.964x$ gives:

$$y' = 81.048 + (0.964)(50) = 129.248 \text{ (rounded to 129)}$$

In other words, the predicted systolic blood pressure for a 50-year-old person is 129.



$$y = mn + c$$

Method of Least Sq

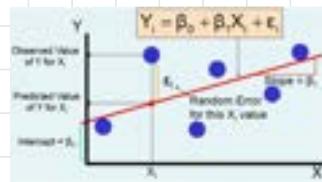
$$b_0 = \frac{(\sum y)(\sum x^2) - (\sum x)(\sum xy)}{n(\sum x^2) - (\sum x)^2}$$

$$b_1 = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

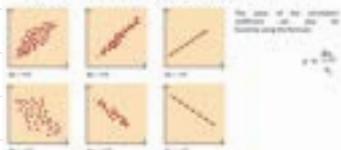
$$Y = \underline{b_0} + \underline{b_1}x + \epsilon$$

Random Error

$$r^2 = \frac{\text{Explained Variation}}{\text{Total Variation}}$$



Relationship b/w Regression & Correlation



CHP 11

- 11.3 The amounts of a chemical compound y that dissolved in 100 grams of water at various temperatures x were recorded as follows:

x ($^{\circ}\text{C}$)	y (grams)
0	8
15	12
30	25
45	31
60	44
75	48

- (a) Find the equation of the regression line.
 (b) Graph the line on a scatter diagram.
 (c) Estimate the amount of chemical that will dissolve in 100 grams of water at 50°C .

$$\sum x = 225 \quad \sum x^2 = 12915 \quad \sum xy =$$

$$\sum y = 488 \quad \sum y^2$$

- 11.9 A study was made by a retail merchant to determine the relation between weekly advertising expenditures and sales.

Advertising Costs (\$)	Sales (\$)
40	385
20	400
25	395
20	365
30	475
50	440
40	490
20	420
50	560
40	525
25	480
50	510

- (a) Plot a scatter diagram.
 (b) Find the equation of the regression line to predict weekly sales from advertising expenditures.
 (c) Estimate the weekly sales when advertising costs are \$35.
 (d) Plot the residuals versus advertising costs. Comment.

b) $\sum x = 410 \quad \sum x^2 = 15,650 \quad \sum xy = 191325$

$$2y = 5445 \quad n = 12$$

$$b_0: \frac{(5445)(15650) - (410)(191325)}{12(15650) - 410^2} = 343.7056$$

$$b_1: \frac{12(191325) - (410)(5445)}{12(15650) - 410^2} = 3.22081$$

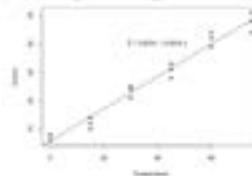
$$y = 343.7056 + 3.22081n$$

- 11.8 (a) $\sum x_i = 675, \sum y_i = 488, \sum x_i^2 = 27,325, \sum x_i y_i = 23,985, n = 18$. Therefore,

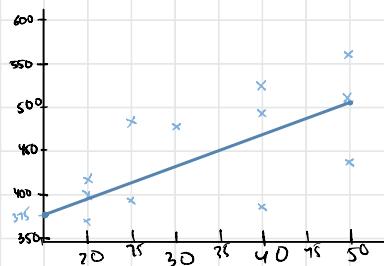
$$\begin{aligned} b &= \frac{(110)(25,985) - (675)(488)}{(110)(27,325) - (675)^2} = 0.3678, \\ a &= \frac{488 - (0.3678)(675)}{18} = 5.8256. \end{aligned}$$

$$\text{Hence } y = 5.8256 + 0.3678x$$

- (b) The scatter plot and the regression line are shown below.



$$(i) \text{ For } x = 35, y = 5.8256 + 0.3678(35) = 34.20 \text{ pesos}$$



c) $y = 343.7056 + 3.22081(35)$
 $y = 456.434$

CHP 11

- 11.43 Compute and interpret the correlation coefficient for the following grades of 6 students selected at random:

	Mathematics grade	70 92 80 74 65 83
	English grade	74 84 63 87 78 90

$$\begin{array}{ccccc} x & x^2 & y & y^2 & xy \\ \hline 464 & 36354 & 476 & 38254 & 36926 \end{array}$$

$$r = \frac{6(36926) - (464)(476)}{\sqrt{[6(36354) - (464)^2][6(38254) - (476)^2]}}$$

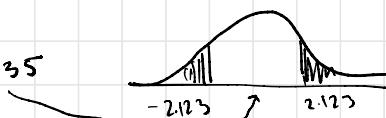
$$r = 0.2396$$

- 11.46 Test the hypothesis that $\rho = 0$ in Exercise 11.43 against the alternative that $\rho \neq 0$. Use a 0.05 level of significance.

$$H_0: \rho = 0 \quad \alpha = 0.05 \quad n = 6$$

$$H_1: \rho \neq 0 \quad r = 0.2396$$

$$t_c = \frac{0.2396 \sqrt{6-2}}{\sqrt{1-0.2396^2}} = 0.4935$$



df = 6-2 = 4, t from table = 2.123 \rightarrow rejection region

DON'T REJECT H_0

- 11.45 With reference to Exercise 11.33 on page 400, assume a bivariate normal distribution for x and y .

- Calculate r .
- Test the null hypothesis that $\rho = -0.5$ against the alternative that $\rho < -0.5$ at the 0.025 level of significance.
- Determine the percentage of the variation in the amount of particulate removed that is due to changes in the daily amount of rainfall.

Daily Rainfall, x (0.01 cm)	Particulate Removed, y ($\mu\text{g}/\text{m}^2$)
4.3	126
4.5	121
5.9	116
5.6	118
6.1	114
5.2	118
3.8	132
2.1	141
7.5	108

a) ... $r = -0.919$

$$t_c = -0.919 \sqrt{\frac{9-2}{1-(0.919)^2}} = -12.71$$

WHY
Z TEST

11.46 (a) From the data of Exercise 11.3 we find $S_{xx} = 344.36 - 47.79 = 396.56$, $S_{yy} = 111.76 - 106.79 = 44.97$ and $S_{xy} = 3000 - (111.76)(106.79) = -111.36$, $S_{yx} = \frac{-111.36}{\sqrt{396.56 \cdot 44.97}} = -0.079$.

(b) The hypotheses are

$$H_0: \rho = -0.5,$$

$$H_1: \rho < -0.5.$$

$$\begin{aligned} &\text{at } \alpha = 0.05, \\ &\text{Critical value } z_c = -1.64, \\ &\text{Computation: } \left| \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} \right| = -4.32, \\ &\text{Decision: Reject } H_0, \rho < -0.5, \\ &\text{(i) } (-0.499)^2/0.4997 > 0.025. \end{aligned}$$

ishma hafeez
notes
reprst
free

BAYES RULE

$$P(A|B) = \frac{P(B|A) P(A)}{\sum_{i=1}^n P(B|A_i) P(A_i)}$$

$$IQR = Q_3 - Q_1$$

$$L = Q_1 - 1.5(IQR)$$

$$U = Q_3 + 1.5(IQR)$$

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

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

$$\hookrightarrow P(A|B) P(B) = P(B|A) P(A)$$

Probability Mass Function (PMF)

1. $f(n) \geq 0$, for all $n \in \mathbb{N}$

$\hookrightarrow \sum f(n) dn = 1 \rightarrow$ verify PMF

$$\hookrightarrow P(X=x) = F(x) - F(x-1)$$

$$\hookrightarrow P(a < X < b) = F(b) - F(a)$$

$$P(n) = (n+1)^n$$

Joint Probability Distribution

$$1. f(n, y) \geq 0$$

$$2. \sum_n \sum_y f(n, y) = 1$$

$$3. P(X=n, Y=y) = f(n, y)$$

$$\hookrightarrow P[(X, Y) \in A] = \sum_n \sum_y f(n, y)$$

Empirical Rule

$$68\% \quad \bar{x} \pm s$$

$$95\% \quad \bar{x} \pm 2s$$

$$99.7\% \quad \bar{x} \pm 3s$$

$V(X) = E(X^2) - [E(X)]^2$ $V(Y) = E(Y^2) - [E(Y)]^2$ $E(XY) = \sum x \sum y f(x, y)$ $Covariance(X, Y) = E(XY) - E(X) E(Y)$ $Correlation(X, Y) = \frac{Covariance(X, Y)}{\sqrt{V(X)} \sqrt{V(Y)}}$	$\xrightarrow{\text{variance}}$ $\mu = E(X) = \sum x f(x)$ $\xrightarrow{\text{mean}}$
---	--

$\pm 1 \longrightarrow \text{strong}$
 $\pm 0.5 \longrightarrow \text{moderate}$
 $0 \longrightarrow \text{negligible}$

Standard Normal Distribution

$$Z\text{-score} = \frac{X - \mu}{\sigma}$$

Central limit theorem

$$z \text{-test}$$

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

standard error of mean

Finite Population Correction Factor

$$\sqrt{\frac{N-n}{N-1}}$$

population size
sample size

Confidence Interval for Proportions

$$1. \hat{P} \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}\hat{q}}{n}}$$

$\frac{x}{n} \rightarrow$ no. of sample units
 $1 - \hat{P}$

$$1. \bar{X} = \frac{\sum x}{n}$$

sample mean
sample size

$$2. S = \sqrt{\frac{\sum x^2}{n} - (\frac{\sum x}{n})^2}$$

SD
standard error of mean

$$3. S_{\bar{X}} = \frac{S}{\sqrt{n}}$$

SD

$$4. E_z = Z_{\frac{\alpha}{2}} \left(\frac{\sigma}{\sqrt{n}} \right)$$

margin of error
standard error of mean

$$5. n = \left(\frac{Z_{\frac{\alpha}{2}} \sigma}{E} \right)^2$$

critical region of z-distribution

Chi Square Distribution

$$1. \frac{(n-1)s^2}{\chi^2_{\text{right}}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{\text{left}}}$$

$$2. df = n-1$$

Confidence Interval for μ

$$2. \bar{X} \pm Z_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$

sample SD

$\hookrightarrow \sigma$ unknown

$\hookrightarrow n \geq 30$

$$3. \bar{X} \pm t_{\frac{\alpha}{2}, \frac{2n}{2n-1}} \frac{s}{\sqrt{n}}$$

sample SD

$\hookrightarrow \sigma$ unknown

$\hookrightarrow n < 30$

Confidence Interval for $\mu_1 - \mu_2$

$$2. (\bar{X}_1 - \bar{X}_2) \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$\hookrightarrow \sigma_1, \sigma_2$ unknown

$\hookrightarrow n \geq 30$

$$3. (\bar{X}_1 - \bar{X}_2) \pm t_{\frac{\alpha}{2}, \frac{2n}{2n-1}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$\hookrightarrow \sigma_1, \sigma_2$ unknown

$\hookrightarrow n < 30$

Z-Test

$\hookrightarrow n \geq 30$

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

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

sample mean H_0 mean
 $\frac{\sigma_1^2}{n_1}$ standard deviation $\frac{\sigma_2^2}{n_2}$
 sample size

$$\delta = \sqrt{\frac{n \sum x^2 - (\sum x)^2}{n(n-1)}}$$

T-Test for mean

$$t = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

$df = n-1$

$\hookrightarrow \sigma_1 \neq \sigma_2$ variance

sample mean H_0 mean

 $t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$

standard deviation standard deviation
 sample size

$$df^2 = \frac{(\delta_1 + \delta_2)^2}{\frac{\delta_1^2}{n-1} + \frac{\delta_2^2}{n-1}}$$

$\hookrightarrow \sigma_1 = \sigma_2$ variance

 $\delta = \frac{(n-1)(\delta_1)^2 + (n-1)(\delta_2)^2}{(n+n)-2}$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\delta \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\delta = \sqrt{\frac{n \sum x^2 - (\sum x)^2}{n(n-1)}}$$

correlation coefficient (PPMC)

correlation coefficient

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

no. of data points

Hypothesis testing for correlation

$$t = r \sqrt{\frac{n-2}{1-r^2}}$$

$H_0: P=0$

$df = n-2$

1. State H_0, H_A, α

2. Degree's of freedom

$\hookrightarrow df$ between groups = $K-1$

$\hookrightarrow df$ within groups = $N-K$ no. of groups

$\hookrightarrow \text{diff total} = N-1$ total no. of observations

$\hookrightarrow F_{\text{critical}}$ using table

3. Sum of square deviations from mean

\hookrightarrow each groups mean $\bar{x}_1, \bar{x}_2, \bar{x}_3$

\hookrightarrow Grandmean \bar{x}

$\hookrightarrow SS_{\text{total}} = \sum (x - \bar{x})^2$

$\hookrightarrow SS_{\text{within}} = \sum (x - \bar{x})^2$

$\hookrightarrow SS_{\text{between}} = SS_{\text{total}} - SS_{\text{within}}$

4. Mean Square

$\hookrightarrow MS_{\text{between}} = \frac{SS_{\text{between}}}{df_{\text{between}}}$

$\hookrightarrow MS_{\text{within}} = \frac{SS_{\text{within}}}{df_{\text{within}}}$

5. Calculate F statistic

$F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$

\hookrightarrow conclusion

reject H_0

F ratio > F critical

F ratio > 1

accept H_0

F ratio < F critical

F ratio ≤ 1

Simple Linear Regression

method of least sq

$$B_0 = \frac{(\sum y)(\sum x^2) - (\sum x)(\sum xy)}{n(\sum x^2) - (\sum x)^2}$$

$$B_1 = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

$$Y = B_0 + B_1 x + \epsilon$$

random error