

CX2100 Probability & Statistics

Part 1 Instructor: A/P Lau Chiew Tong
Room: N4-2B-58; Email: asctlau@ntu.edu.sg

Part 2 Instructor: A/P Adams Kong
Room: N4-2A-17; Email: adamskong@ntu.edu.sg

~ 24 Sessions of lectures

~ 12 Sessions of tutorials

Assessments (to be confirmed):

- 4 Quizzes – 25% each
- No Final exam

Part 1 (Week 1-7):

Ch1 Introduction to Statistics

Ch2 Presenting Data

Ch3 Summarizing Distributions

Ch4 Bivariate Data

Ch5 Probability Theory

Ch6 Probability Distribution

- Random Variables
- Discrete Distribution
- Continuous Distribution

Part 2 (Week 8-13):

To be taught by Prof Adams Kong

- TEL supporting Materials for Ch1-4 (short video clips and notes) are available in NTULearn. Students are advised to watch the video clips before attending the lectures for Part 1 (Ch1-4).
- Lectures (Week 1 to 7)
 - Discussion & worked examples on topics covered in the TEL materials
 - Additional topics not covered in TEL materials (Ch 5 and 6)

- Recess week (Self-study non-exam)
 - Use of R programming for analysis and presentation of statistical data (Practice materials will be available in NTULearn)
- Tutorial – One session per week,
to begin in Week 3

All Course Materials are available in [NTULearn](#).

Prof Kong will take over the teaching from Wk 8.

Students participation during lectures using Student Response System - Wooclap

A QR code or URL will be presented for
students to join Wooclap.

Scan the QR code with your mobile phone
and access the web page. Account
registration is not required.

Let's try out the polling session by answering this question using your mobile phone:

www.wooclap.com/CX2100

Have you taken a course on Probability and Statistics... in JC or Poly?

1 Yes 0% 0

2 No 0% 0

wooclap 90 % 0 / 0

Ch 1. Introduction to Statistics

- Descriptive & Inferential Statistics
- Types of Variables
- Percentiles
- Types of Measurement Scale
- Distributions
- Linear Transformations

Statistics involve gathering, organizing, analyzing, interpreting and presenting data.

Statistics are being used everywhere:

- Number of students enrolled in this course
- Index measuring stock market
- Singapore household income
- Live data of new Covid-19 cases
- Opinion poll, benchmarks poll, tracking polls, etc.

Example: A toothpaste manufacturer claims that more than 80% of Dentists recommend a particular brand of toothpaste. This was based on surveys of dentists which allow selection of one or more brands.

Is the claim misleading?

Yes.

Because it may be understood that 80% of dentists recommend this brand over the others. It should be noted that other brands were also recommended and may be as much as that particular brand.

Why Statistics are important?

- Predicting the spreading of diseases
- Weather forecasting based on statistics
- Provide informed choice on investment decision
- AI or machine learning based on past statistical data

■ Descriptive & Inferential Statistics

Descriptive statistics – summarize and describe important features of the data collected. Does not generalize beyond the data collected.

Stay
in
group
topic

Inferential statistics - collection of sample to draw inferences about the population, i.e. formal guesses of statistical parameters about the population by looking at the samples.



A teacher wishes to know whether the males in his class have more conservative attitudes than the females. A questionnaire is distributed assessing attitudes and the males and the females are compared.

Is this an example of descriptive or inferential statistics?

1

Descriptiv
e

0%

0

2

Inferential

0%

0

A cognitive psychologist is interested in comparing two ways of presenting stimuli on subsequent memory. Twelve subjects are presented with each method and a memory test is given.

What would be the roles of descriptive and inferential statistics in the analysis of these data?

Descriptive statistics – we describe and analyze the data from the sample.

Inferential statistics – we use the data from the sample to generalize to a larger population of people.

■ Types of Variables

In statistic, we can broadly group variable into two categories: Qualitative and Quantitative.

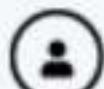
Express in name, etc etc

Express in number

Categorize the following variables as being qualitative or quantitative.

www.wooclap.com/CX2100

Enter A for qualitative and B for quantitative.



1. Rating of the quality of a service on a 10-point scale - 1.



2. Favourite colour - 2.

3 best timing to complete task

B
A
B

wooclap



90

%



0 / 0



■ Percentile

In certain experiments, it is more meaningful to compare the outcomes obtained.

Eg: if you know that your quiz marks is 65 out of 100, you may not know how well you have done compared to others in your class.

A **percentile** is a comparison score between a particular score and the scores of the rest of a group.

■ Percentile - calculation

Calculation of P^{th} Percentile for a set of N data:

Data arranged in the order of magnitude

1. Compute the rank $R = \frac{P}{100} \times (N + 1)$
2. Let I_R = Integer part of R and F_R = Fractional part of R
3. P^{th} Percentile = Data at rank I_R +
(Data at rank (I_R+1) – Data at rank I_R) $\times F_R$

Eg: Given data: [3, 5, 7, 8, 9, 11, 13, 15], compute the
 25^{th} and 75^{th} percentile.

$$25^{\text{th}} = \frac{1}{4}(9) = 2.25$$
$$5 + (7-5)0.25 = 5.5$$

$$75^{\text{th}} = \frac{3}{4}(9) = 6.75$$
$$11 + (13-11)0.75$$
$$11 + 1.5 = 12.5$$

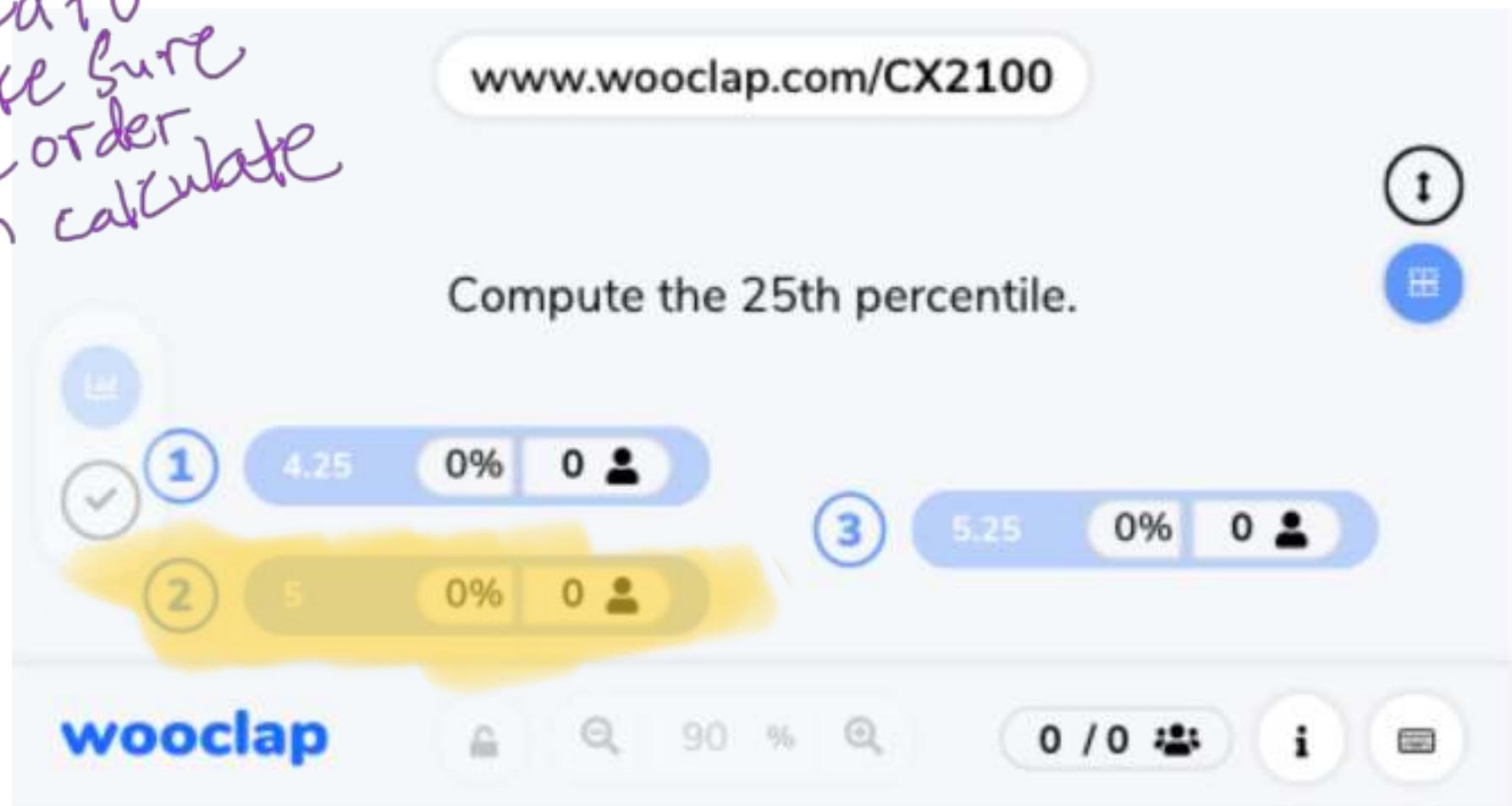
Given a data set of 20 values:

1	2	3	4	5	5	6	6	6	7
7	8	8	8	9	9	9	10	10	10

need to
make sure
the order
then calculate

www.wooclap.com/CX2100

Compute the 25th percentile.



Given a data set of 20 values:

1	2	3	4	5	5	10	10	10	6
6	9	9	9	7	7	7	8	8	8

www.wooclap.com/CX2100

Compute the 85th percentile.

The image shows a Wooclap poll interface. At the top, the URL "www.wooclap.com/CX2100" is displayed. Below it, the instruction "Compute the 85th percentile." is shown. Three options are listed: 1 (9.85) in yellow, 2 (10) in blue, and 3 (7.85) in light blue. Each option has a progress bar at 0% and a count of 0 people. On the left, there are icons for a list, a checkmark, and a grid. On the right, there are icons for help, a person, and a barcode. The Wooclap logo is at the bottom left, and a navigation bar with icons for lock, search, percentage, and refresh is at the bottom right. A status bar at the bottom shows "0 / 0" and other icons.

Practical example:

Graduate Management Admissions Test - a standardized test for application to graduate-level business programs.

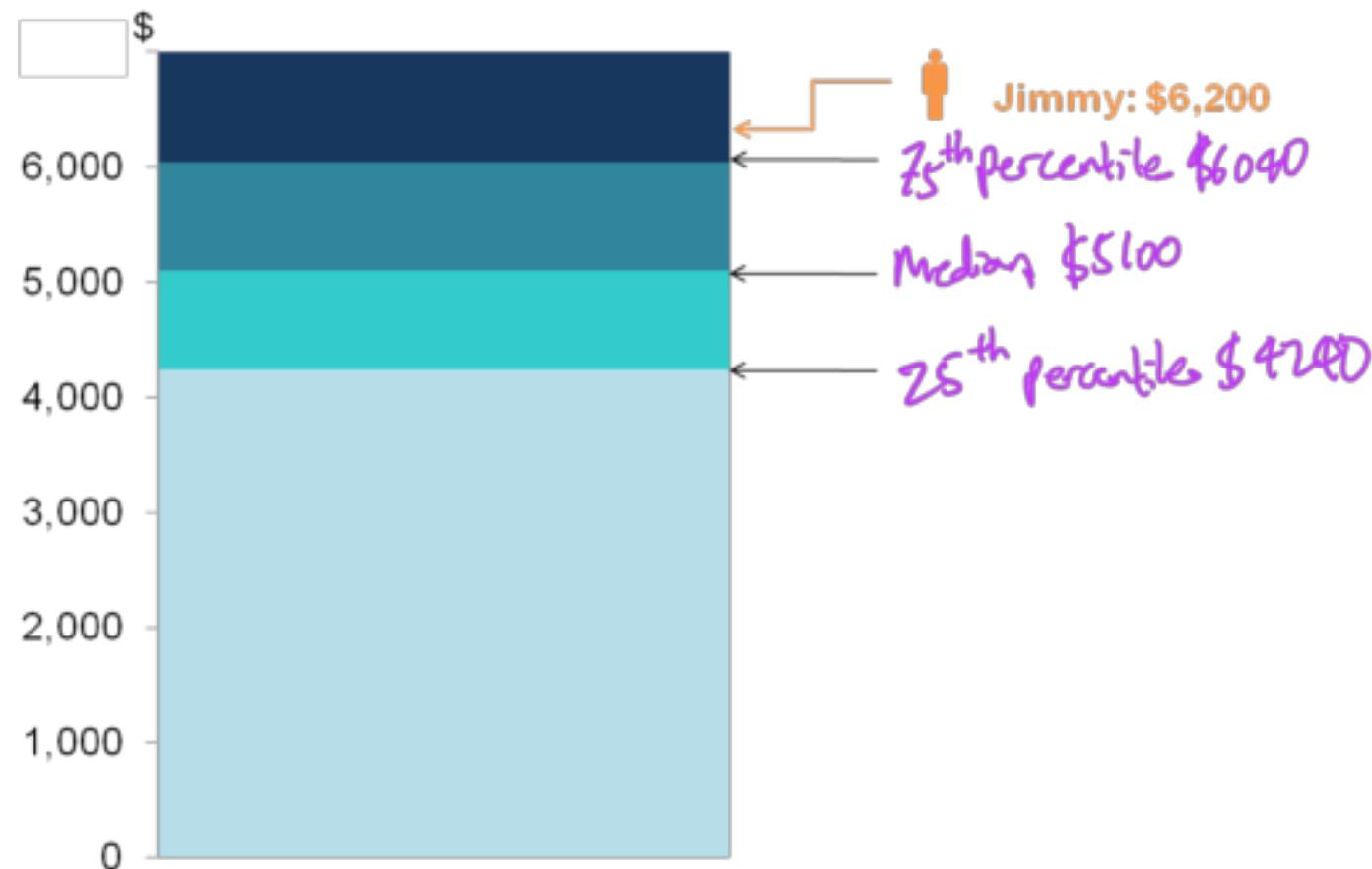
UNOFFICIAL GMAT® SCORE REPORT		
	Scaled Score	Percentile
Integrated Reasoning	6	67
Quantitative	47	70
Verbal	47	99
Total*	750	98

otal score is derived from your Verbal and Quantitative scores.

ficial Score Report, including your Analytical Writing Assessment score
ing score and percentile ranking.

Example: Percentile for benchmarking.

Suppose Jimmy earns a monthly wage of \$6,200. He would like to compare his wages with those working in the same occupation and industry.





■ Types of Measurement Scales

Four basic levels of measurement scales:

Nominal – names or labels with no specific order

Ordinal – variables in a specific order

Interval – numerical scales in which intervals have
the same interpretation throughout, but
no true zero

Numerical

Ratio – include all the characteristics of interval
scale, plus it has zero position indicating the
absence of the quantity being measured

Specify the level of measurement used for the following variables:

www.wooclap.com/CX2100

Specify the level of measurement for the following:
Enter O for Nominal, O for Ordinal, I for Interval and ...
for Ratio.

I

temperature

N

Cities

R

N

H

O

1. SAT scores - ①

2. Favourite colour - ②

3. Time to complete a task - ③

4. Rating of the quality of a service on a 10-point scale - ④

I

Calendar year

R

amount of money in pocket

O

Ranking of army officers

O

Movie ranking

I

N

rainbow color O

N

I?
O?

I

R

O

H

O

O

wooclap



90

%



0 / 0

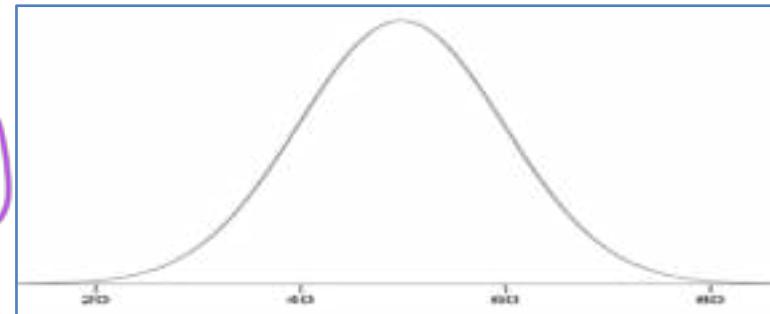


Distributions

Frequency distribution

- Discrete variables
- Continuous variables

(display in
bar chart)



Symmetric

Probability distribution

- Probability mass function
- Probability density function

(discrete random variable
is exactly = some value)
(continuous random Variable)

Positive skew (to the right)

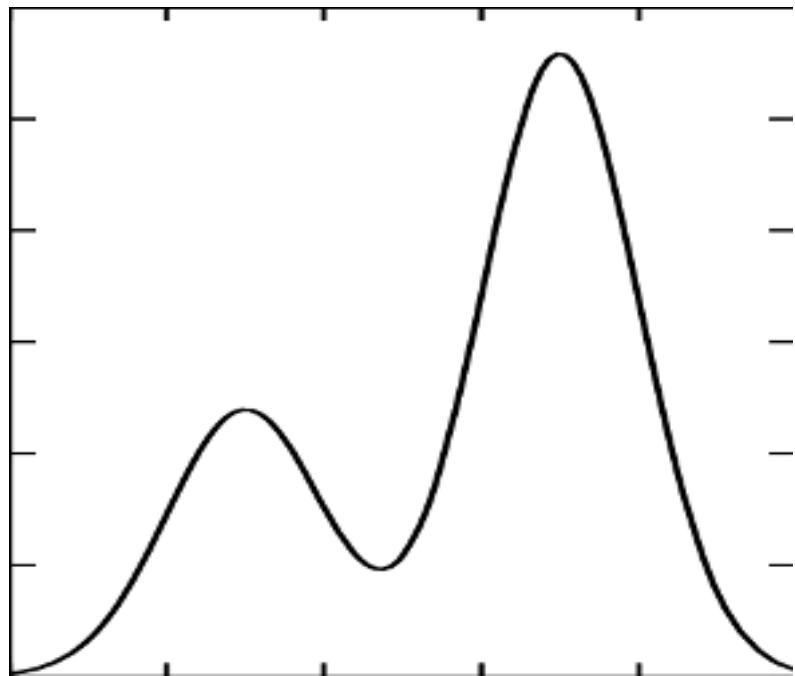
Shapes of distributions

- symmetric
- skewed to the right
- skewed to the left

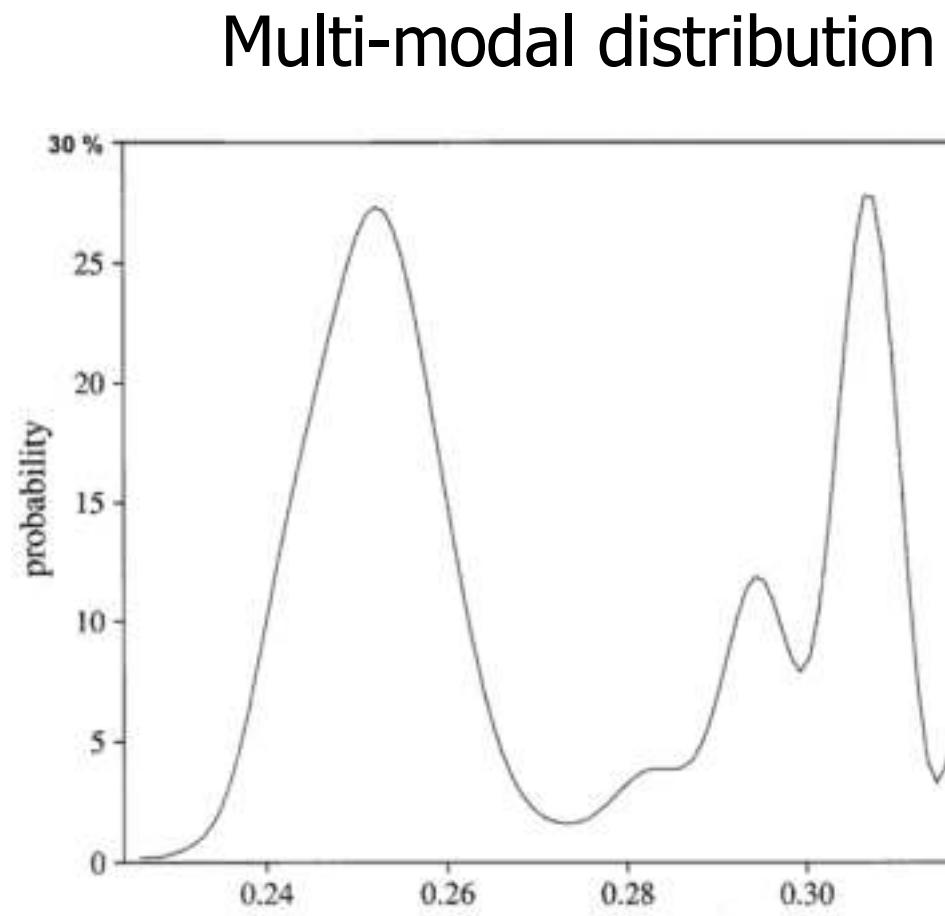


Negative skew (to the left)

- Distributions – other shapes:



Bimodal distribution

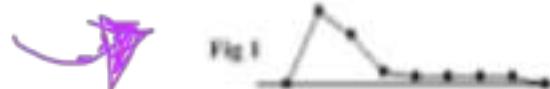


Multi-modal distribution

Which of the following frequency polygons has a large positive skew?

www.wooclap.com/CX2100

Which of these frequency polygons has a large ... *positive Skew*



1

Fig 1

0%

0



2

Fig 2

0%

0



3

Fig 3

0%

0



wooclap



90

%



0 / 0



■ Linear Transformations

Transform data from one measurement scale to another.

Examples:

Convert length measured in X feet to measurement in Y metre, i.e.

$$Y = 0.3048 X$$

Convert temperature in Fahrenheit to Centigrade:

$$C = 0.5556 F - 17.778$$

■ Linear Transformations

Which of the following are linear transformations?

www.wooclap.com/CX2100

Which of the following are linear transformations?

(answer: Y for Yes and N for No)



1. Converting from pounds to kilograms. 1
2. Squaring each side to find the area. 2
3. Dividing all numbers by 2 and then adding 3. 3
4. Computing the square root of each person's weight. 4

Y

N

Y

N

wooclap



90

%



0 / 0



Ch 2. Presenting Data

- Frequency tables and Charts
- Bar Charts
- Stem and Leaf Displays
- Histograms
- Box Plots
- etc

- Presenting qualitative or discrete data:

Previous Ownership	Frequency	Relative Frequency
None	85	0.17
Windows	60	0.12
Macintosh	355	0.71
Total	500	1.00

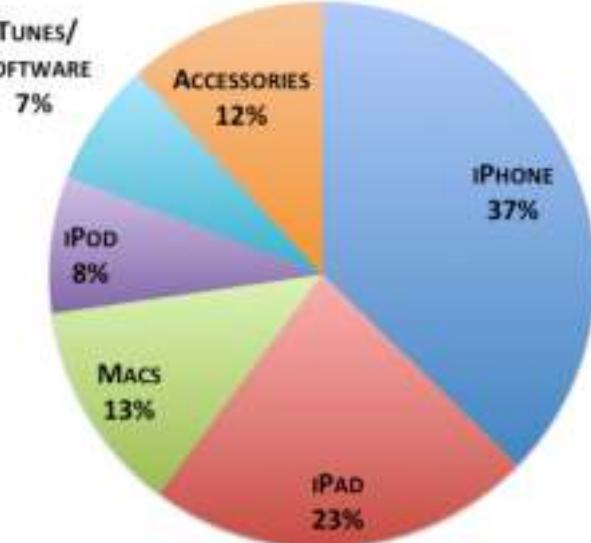
Frequency table

Singapore Average Monthly Wages



Bar chart

How Apple Makes Its Money:
Q1 2013



Pie chart

Examples of misleading presentations:



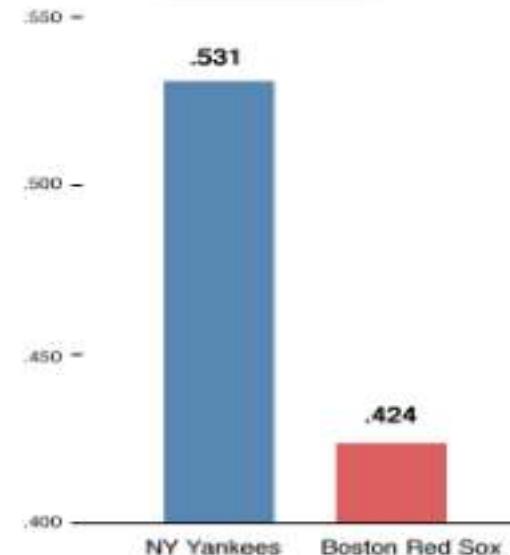
The battery on the right has 70% more capacity than the one on the left.

Misleading display?

Yes

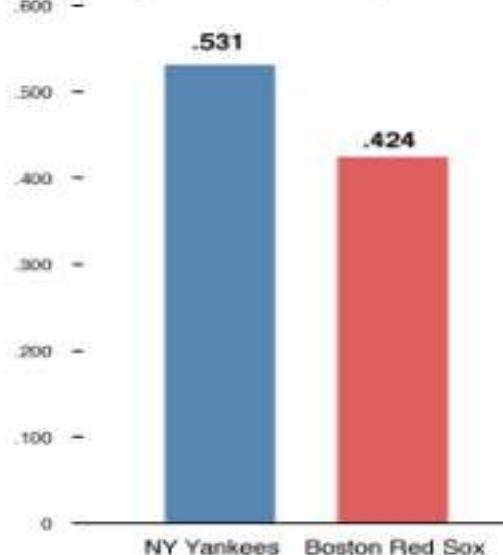
Percentage of victories

WRONG



Percentage of victories

RIGHT



- Stem-and-leaf Presentation
- Useful when data are not too numerous

Eg: No. of touchdown passes by each of the 31 football teams.

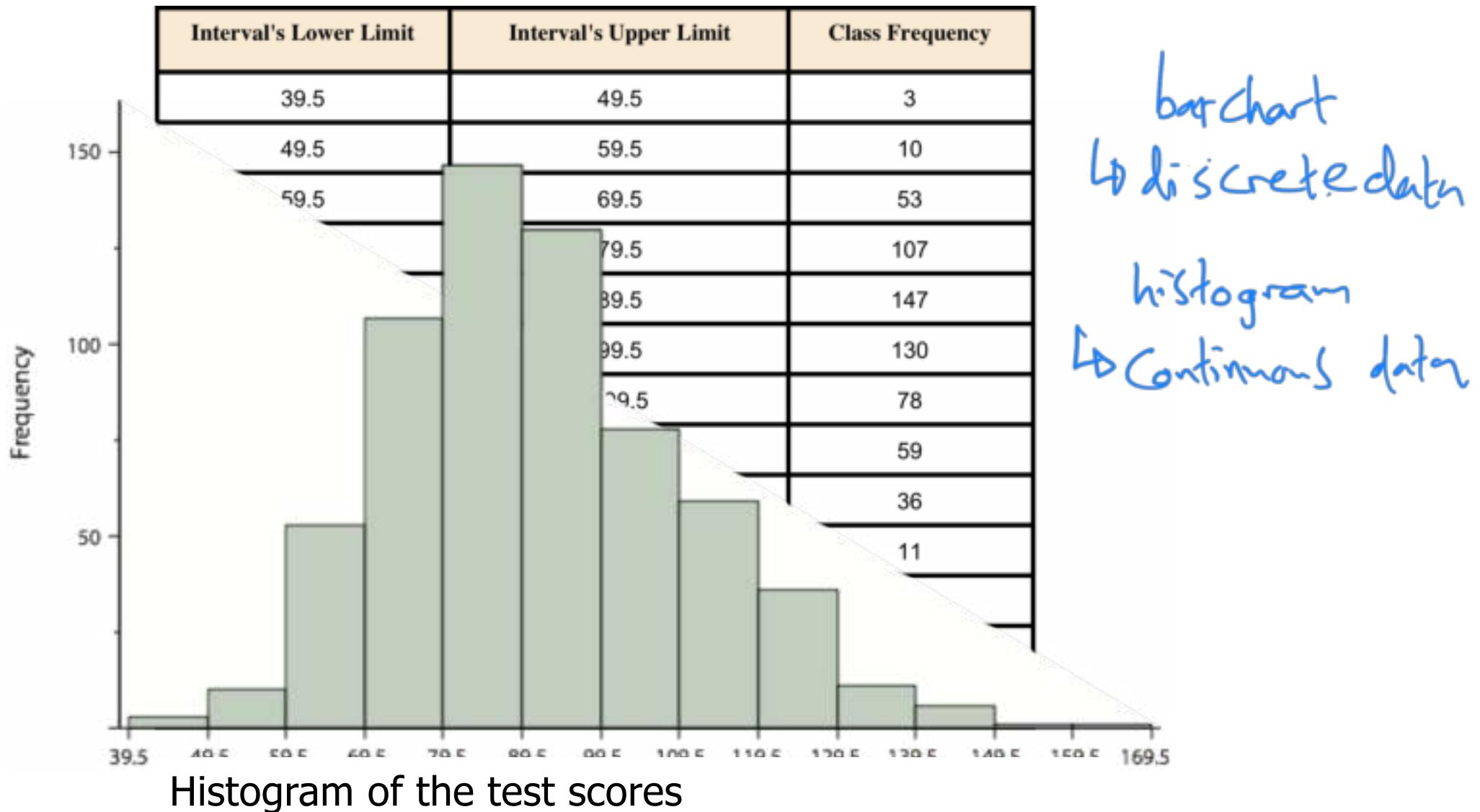
37, 33, 33, 32, 29, 28, 28, 23, 22, 22, 22, 21, 21,
21, 20, 20, 19, 19, 18, 18, 18, 18, 16, 15, 14, 14,
14, 12, 12, 9, 6

(~~xx~~) ↗
↗ ~~(xx)~~

3		2337
2		001112223889
1		2244456888899
0		69

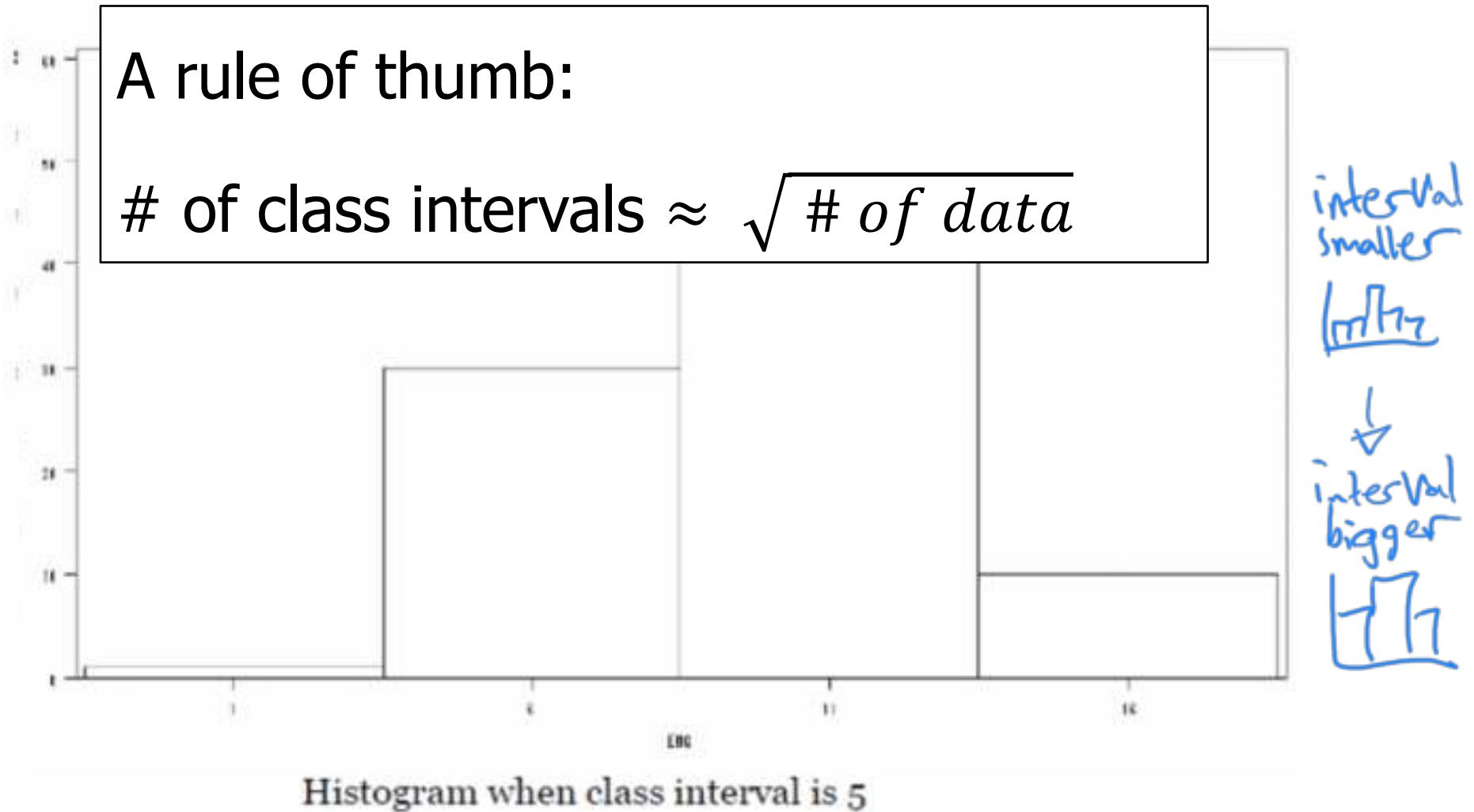
Stem Leaf

- Presenting continuous data:
Data range is divided in class intervals or bins



Presenting continuous data - Histogram:

Examples: Class interval size affects the visual presentation





You have to decide between displaying data with a histogram or a stem-and-leaf display. What factor would affect your choice?



1

Data ...

0%

0

2

Amount o...

0%

0

3

Data-type

0%

0

wooclap



90

%



0 / 0



For a given set of data, the choice is to select between a histogram or a stem-and-leaf display.

With more data, a histogram can be very useful since it shows the overall shape of the distribution.

Stem-and-leaf display is better for smaller sets of data.

Question:

Suppose you are constructing a histogram for describing the distribution of salaries for individuals who are 40 years or older, but are not yet retired. (a) What is on the Y-axis? (b) What is on the X-axis? (c) What would be the probable shape of the salary distribution? Explain why.

- a) The Y-axis would be of the frequency of individuals.
- b) Salary would be on the X-axis because this is the variable whose distribution is of interest.
- c) The distribution is expected to be positively skewed with a few wealthy CEOs earning well above the average salary.

! Box Plot:

Draw the Box Plot for the following data set (31 values).

14	15	16	16	17	17	17	17	17	18	
18	18	18	18	18	19	19	19	20	20	
20	20	20	20	21	21	22	23	24	24	29

Upper Hinge = 75th Percentile (3 rd quartile)	20
Lower Hinge = 25th Percentile (1 st quartile)	17
H-Spread = Upper Hinge - Lower Hinge (IQR)	3
Step = 1.5 x H-Spread	4.5
Upper Inner Fence = Upper Hinge + 1 Step	24.5
Lower Inner Fence = Lower Hinge - 1 Step	12.5

IQR : Upper Hinge - Lower Hinge 3

Step

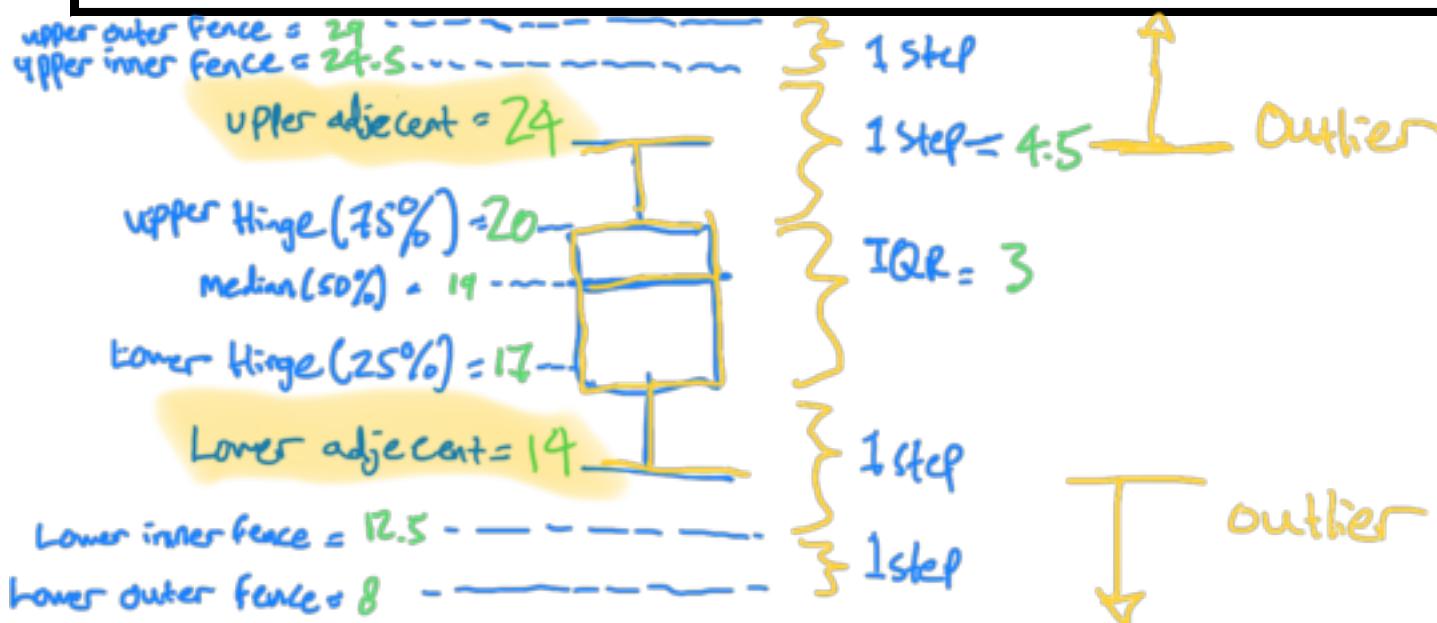
①

②

③

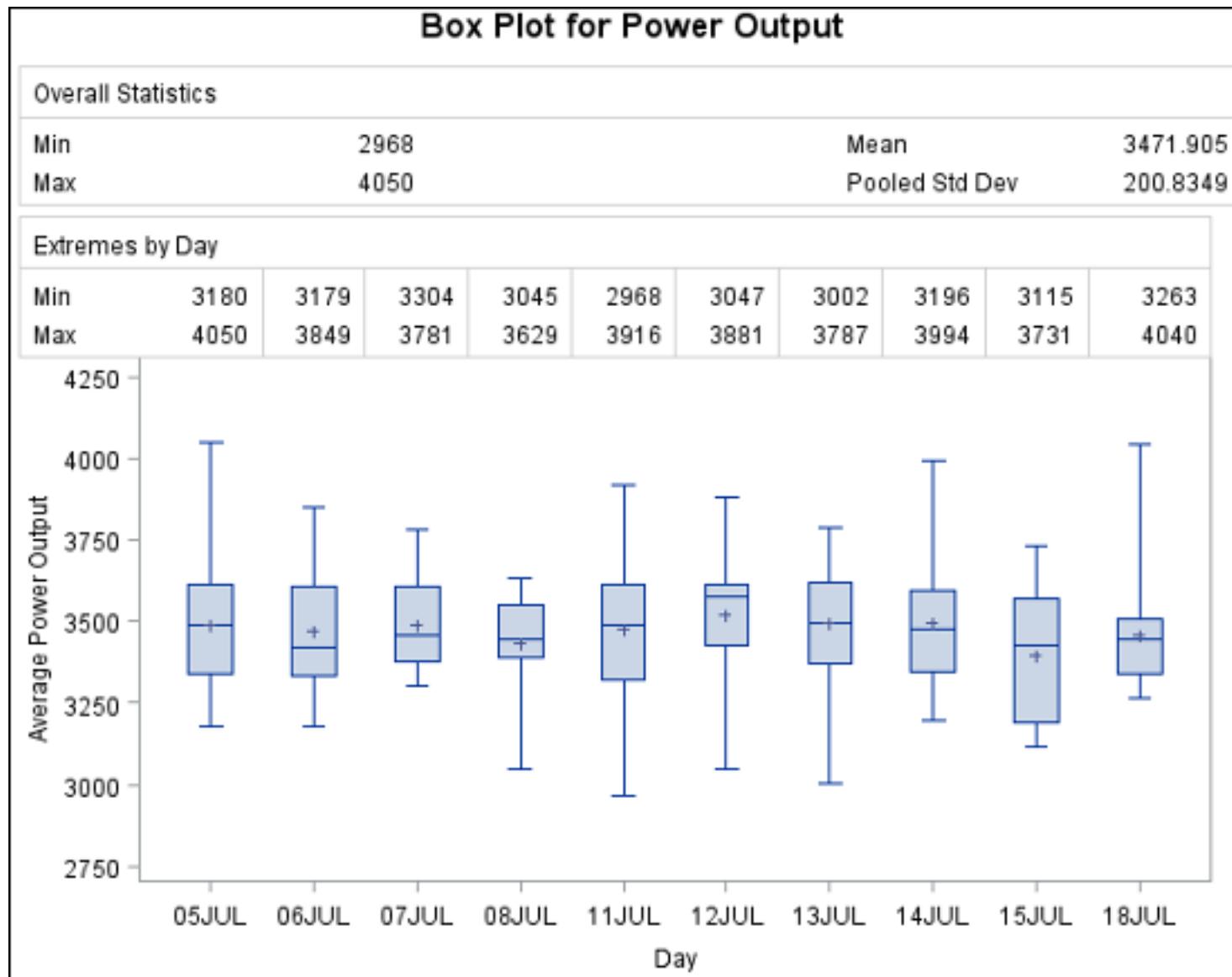
④

Upper Outer Fence = Upper Hinge + 2 Steps	29
Lower Outer Fence = Lower Hinge - 2 Steps	8
Upper Adjacent = Largest \leq Upper Inner Fence	24
Lower Adjacent = Smallest value \geq Lower Inner Fence	14
Outside Value Far Out Value	Outlier = Value $>$ upper adjacent or $<$ lower adjacent



Find 25 - 75 % + 50%
 Find IQR / H-Spread
 Find Step
 Find upper/lower inner
 Find \leq adjacent
 find upper/lower outer

Practical example of Box Plot: Daily Power Output of an Electric Generator



Another practical example of Box Plot: Daily share prices

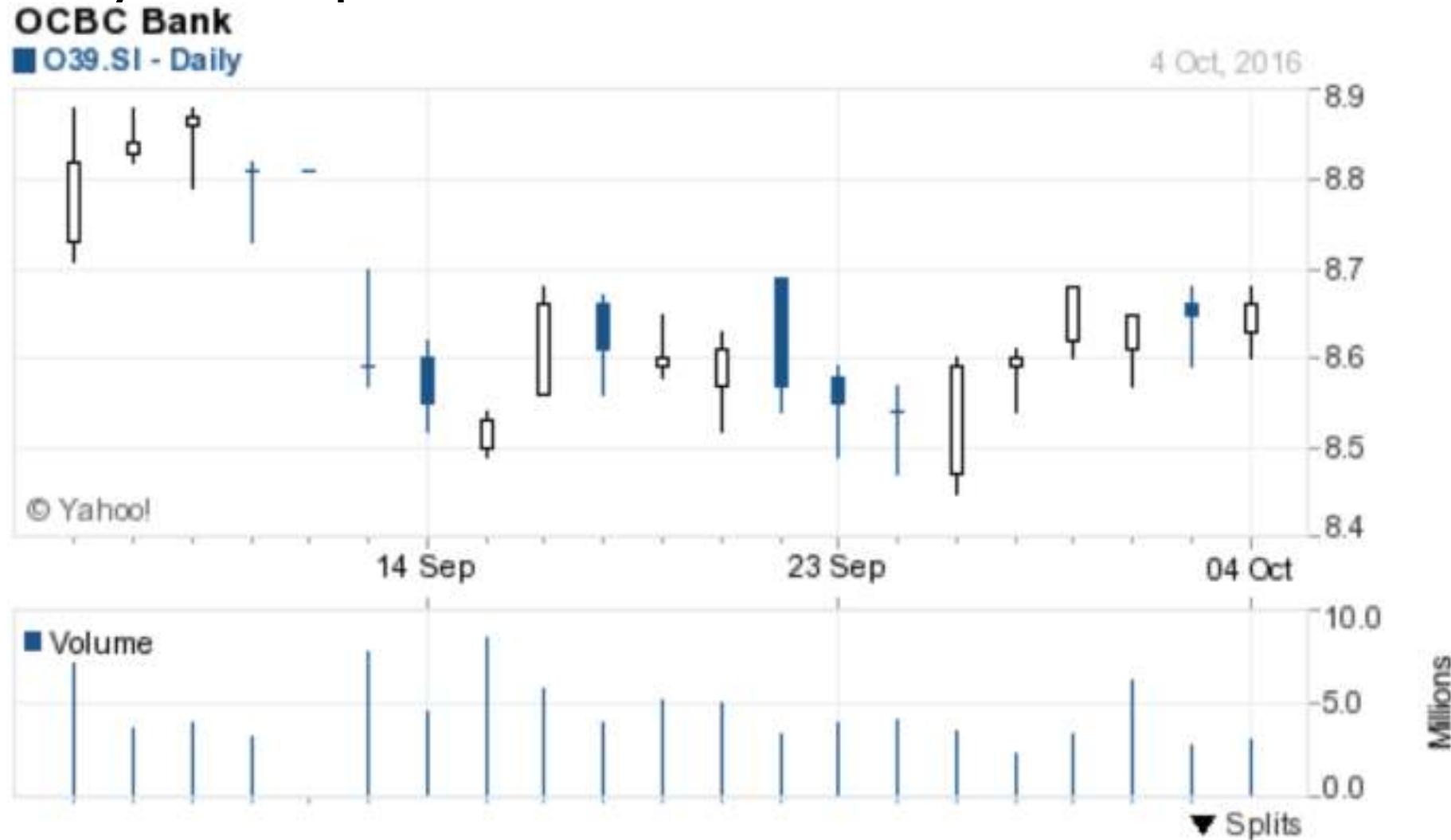


Figure shows the boxplots of the Anger Expression Index by sports participation. Does it look like there are any outliers? Which group reported expressing more anger?

www.wooclap.com/CX2100

Figure shows the boxplots of the Anger Expression index by Sport participation

1) does it look like there are any outliers?
2) which group reported expressing more anger?

The boxplot displays the distribution of the Anger Expression Index for two groups: Sports 1 and Sports 2. The y-axis represents the Anger Expression Index, ranging from 0 to 70. Sports 1 has a median around 30, with whiskers extending from approximately 15 to 50. Sports 2 has a higher median around 40, with whiskers extending from approximately 10 to 65. Both boxplots show some internal variability, with Sports 2 appearing slightly more spread out.

Sport	Median	Q1	Q3	Min	Max
Sports 1	~30	~25	~35	~15	~50
Sports 2	~40	~30	~50	~10	~65

1 (1) Yes. ... 0% 0

3 (1) No. ... 0% 0

2 (1) Yes. ... 0% 0

4 (1) No. ... 0% 0

NO group 2

wooclap

90 %

0 / 0

i

40

Ch 3. Summarizing Distributions

- Central Tendency: mean, median & mode
- Other Measures of Central Tendency
- Comparing Central Tendency
- Measures of Variability: Range, IQR, Variance
- Linear Transformation of variable
- Variance Sum Law I

- Central Tendency

Summarizes a distribution by its central location

Different ways to define central tendency:

- Mean $\mu = \frac{\sum X}{N}$ *Sum of all value*
- Median = 50th Percentile
- Mode = the value with the highest frequency
- Trimean = $\frac{25^{\text{th}} \text{ Percentile} + 2 * \text{Median} + 75^{\text{th}} \text{ Percentile}}{4}$
Multiply all value
- Geometric mean = $(\prod X)^{\frac{1}{N}}$, where \prod means to multiply
Eg: $\prod_{i=1}^5 x_i = x_1 \times x_2 \times x_3 \times x_4 \times x_5$
- Trimmed mean = mean for data with some higher and lower values removed

$$\prod_{i=1}^5 x_i = x_1 \times x_2 \times x_3 \times x_4 \times x_5$$

- Central Tendency:

Eg: Given the following data set, compute the mean, the median, the mode, the trimean, the geometric mean and the mean trimmed 20%.

1	3	4	4	4	5	5	7	8	9	31
---	---	---	---	---	---	---	---	---	---	----

mean = $\frac{81}{11} = 7.36$

median = 5

mode

trimean = $\frac{4 + (2 \times 5) + 8}{4} = 5.5$

$$25^{\text{th}} = \frac{2+1}{100} \times 12 = 3^{\text{rd}}$$

$$4 + 4 = 8$$

$$75^{\text{th}} = \frac{7+3}{100} \times 12 = 9^{\text{th}}$$

$$+ 10 = 18$$

geometric mean
 $(\prod x_i)^{1/11} = (74995200)^{1/11} = 5.2$

mean trimmed by 20%

$$\frac{3 + 4 + \dots + 9}{8}$$

$$8$$

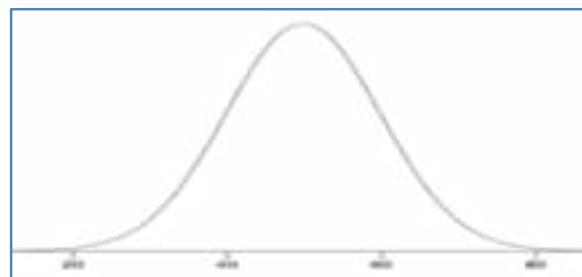
$$10\% \text{ of } 11 = 1.1 \approx 1$$

remove 1 data each side

- Central Tendency – comparing various measures

For symmetric distributions:

Mean = Median= Trimean = Trimmed mean
= Mode (except bimodal distr)



Symmetric

For skewed distributions:

Differences among the measures.

Example – the mean is typically higher than
the median for a positive skewed distribution

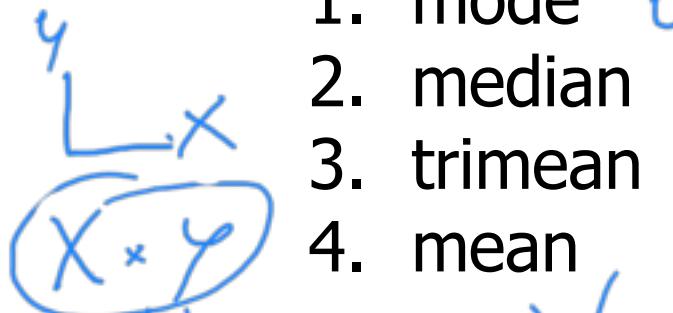
- Central Tendency – comparing various measures

Example:

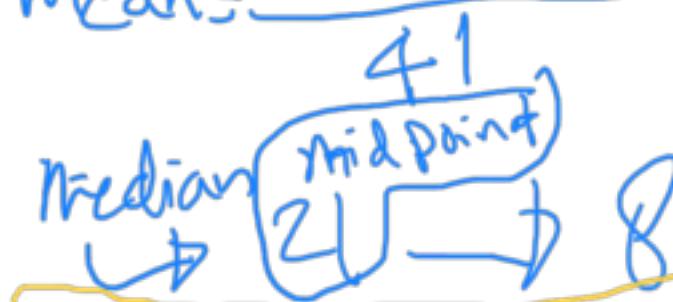
Negative skewed

Calculate:

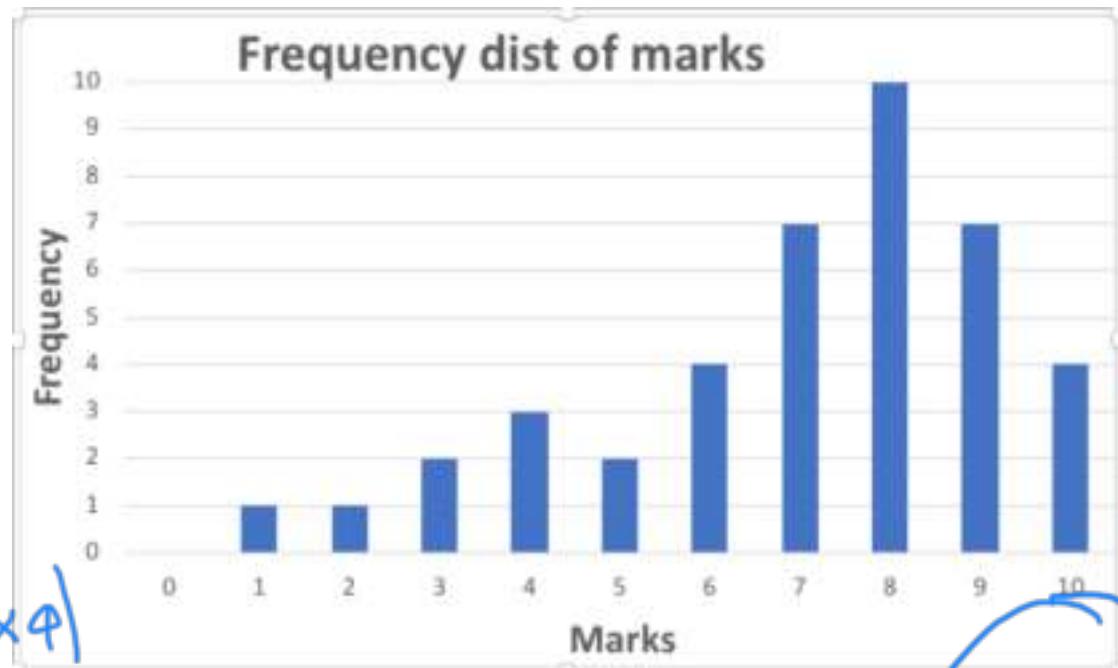
1. mode 8
2. median
3. trimean
4. mean



$$\text{mean} = \frac{(1 \times 1) + (2 \times 1) + (3 \times 2) + (4 \times 3) + (5 \times 2) + (6 \times 4) + (7 \times 7) + (8 \times 9) + (9 \times 7) + (10 \times 4)}{41}$$



$$\begin{aligned} 25^{\text{th}} \text{ percentile} &= 10.5 \rightarrow 6 \\ 75^{\text{th}} \text{ percentile} &= 31.5 \rightarrow 9 \end{aligned}$$



$$\text{trimean} = \frac{6 + 1.8 + 9}{4} = \frac{31}{4} = 7.75$$

- Measures of Variability

An indication of how spread out is the distribution

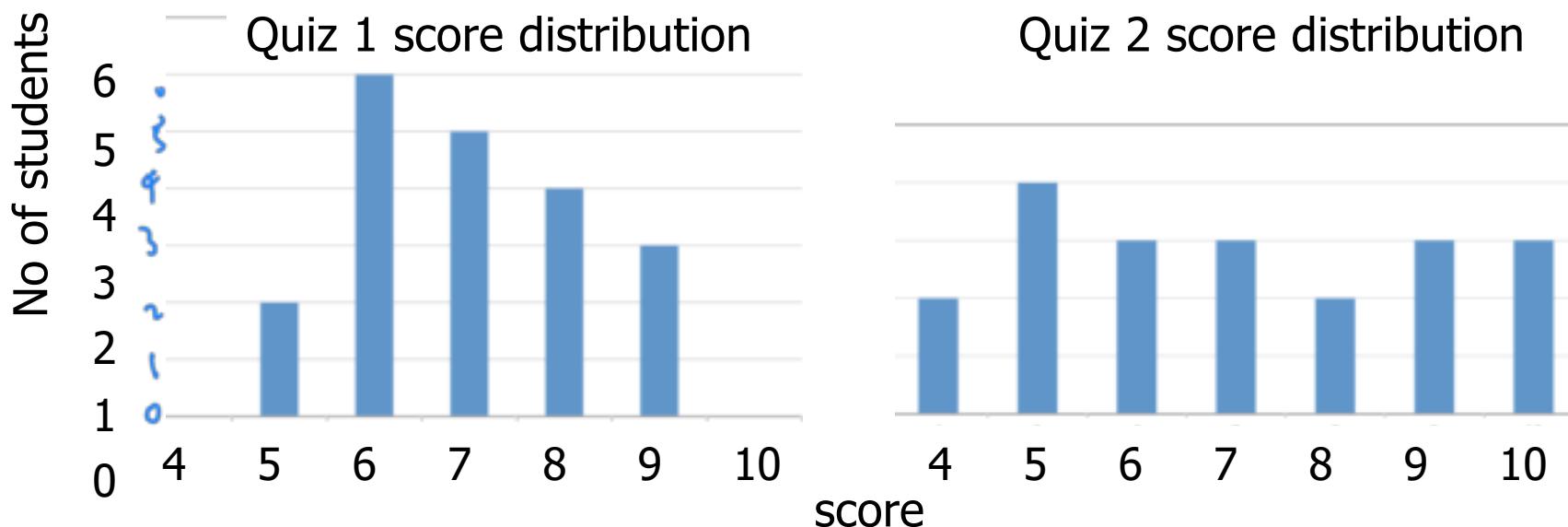
Frequently used measures of variability:

- Range = Highest value – Lowest value
- Interquartile Range IQR = 75th – 25th Percentile
- Variance $\sigma^2 = \frac{\sum(X - \mu)^2}{N}$
- Standard deviation = $\sqrt{\text{Variance}}$

Variance must be +

■ Variability:

Eg: Given the following 2 data set, compute the mean, the range, the IQR and the variance.



Mean μ	7	7
Range	$9 - 5 = 4$	$10 - 4 = 6$
IQR	$8 - 6 = 2$	4
Var σ^2	$\sigma^2 = \frac{\sum (x - \mu)^2}{n} = 1.5$ $2(5-7)^2 + 6(6-7)^2 + 5(7-7)^2 + 4(8-7)^2 + 3(9-7)^2 / 20$ $8+6+0+4+12 = 20 = 1.5$	3.9

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

$$\begin{aligned} n\mu &= E[\sum x] = \sum E[x] \\ n\sigma^2 &= \text{Var}[\sum x] = \sum \text{Var}[x] \end{aligned}$$

Compute mean and variance from the **population** of size N :

Population Mean

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

Population Variance

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

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

$$\text{or } \frac{\sum X^2 - \frac{(\sum X)^2}{N}}{N}$$

Estimate mean and variance from a **sample** of size n :

Sample Mean $\bar{x} = \frac{\sum X}{n}$

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

$$\text{or } \frac{\sum X^2 - \frac{(\sum X)^2}{n}}{n-1}$$

Why $n-1$?

Suppose the denominator of s^2 is n , instead of $(n - 1)$:

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

For unbiased estimate, we expect the mean of s^2 to be equal to σ^2 :

$$\begin{aligned} E[s^2] &= E \left[\frac{1}{n} \left(\sum X^2 - \frac{(\sum X)^2}{n} \right) \right] \\ &= \frac{1}{n} \left(\sum E[X^2] - \frac{E[(\sum X)^2]}{n} \right) \\ &\quad \text{$\sigma^2 + \mu^2$ from the previous slide} \\ &= \frac{1}{n} \left(n \sigma^2 + n \mu^2 - \frac{E[(\sum X)^2]}{n} \right) \end{aligned}$$

$$\begin{aligned}
E[s^2] &= \frac{1}{n} \left(n \sigma^2 + n \mu^2 - \frac{E[(\sum X)^2]}{n} \right) \quad \text{Let } Y = \sum X \\
&= \frac{1}{n} \left(n \sigma^2 + n \mu^2 - \frac{Var[\sum X] + (E[\sum X])^2}{n} \right) \\
&= \frac{1}{n} \left(n \sigma^2 + n \mu^2 - \frac{\sum Var[X] + (\sum E[X])^2}{n} \right) \\
&= \frac{1}{n} \left(n \sigma^2 + n \mu^2 - \frac{n \sigma^2 + (n \mu)^2}{n} \right) \\
&= \frac{1}{n} (n \sigma^2 + n \mu^2 - \sigma^2 - n \mu^2) = \frac{1}{n} (n - 1) \sigma^2
\end{aligned}$$

Let $Y = \sum X$
 $E[Y^2] = \sigma_Y^2 + \mu_Y^2 \rightarrow E[\sum X]^2$
 \downarrow
 $Var[\sum X]$

If the denominator is $(n - 1)$,
then mean of s^2 is equal to σ^2 , i.e. $E[s^2] = \sigma^2$

True or False questions:

www.wooclap.com/CX2100

Answer by entering T for true and F for false.



1. A bimodal distribution has two modes and two medians. ✓ ✗ 1
2. The best way to describe a skewed distribution is to report the mean. 2
3. When plotted on the same graph, a distribution with a mean of 50 and a standard deviation of 10 will look more spread out than with a distribution with a mean of 60 and a standard deviation of 5.

F (only 1 median)
F (only Mean Cannot determine Skew)

T (Standard deviation)
higher = more spread out

- Linear transformation of variable
 - Transform data from one measurement scale to another
- Eg. Consider a taxi trip from point A to B. The taxi service initial charge is \$3 and additional \$0.50 per km for the trip. Let y be the cost of the taxi ride and x be the distance travelled. We have:

$$y = 0.5x + 3$$

The mean of $y = 0.5$ (mean of x) + 3

The variance of $y = 0.5^2$ (variance of x)

! ■ Variance Sum Law I

- Linear combination of 2 independent variables

Eg. Consider a couple A and B. A works x hrs/day and earns \$5/hr while B works y hrs/day and earns \$10/hr. Total daily income t is: Daily income difference d :

$$t = 5x + 10y$$

$$d = 5x - 10y$$

Mean = 5 (mean of x) $-$ 10 (mean of y)

$$\text{Variance} = 5^2 \sigma_x^2 + 10^2 \sigma_y^2$$

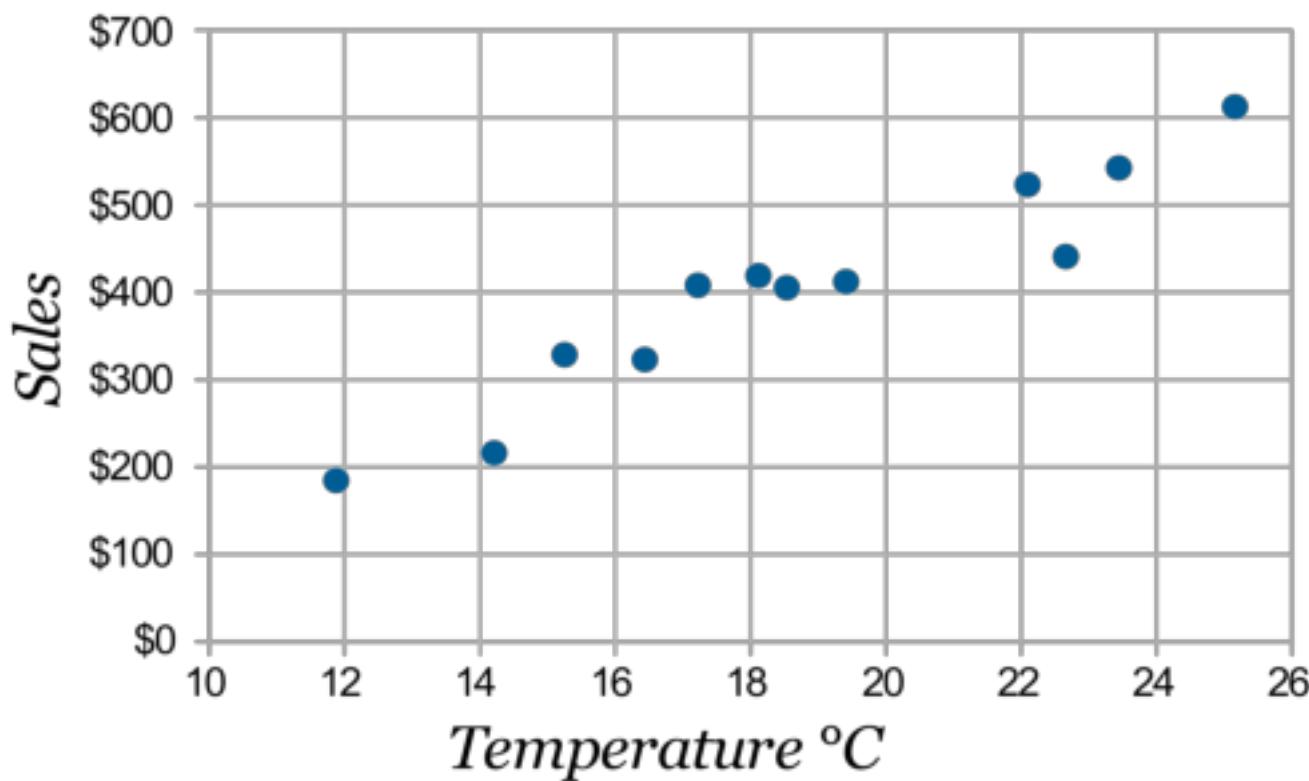
Ch 4. Bivariate Data

- Introduction to Bivariate Data
- Pearson Correlation and Covariance
- Properties of ~~Pearson~~ Correlation
- Variance Sum Law II

- Bivariate Data

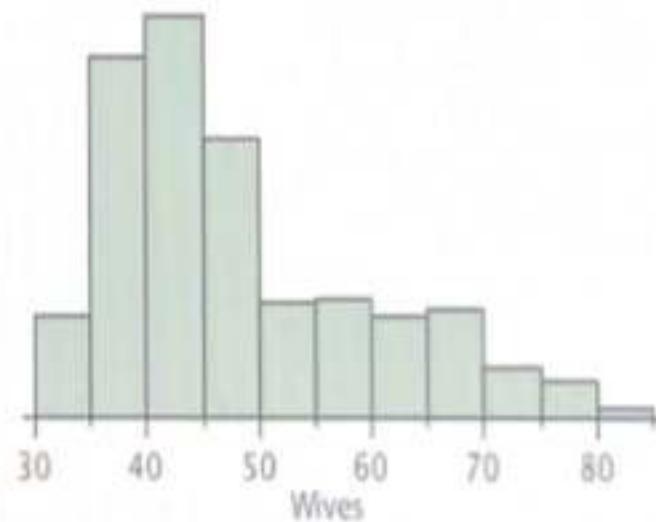
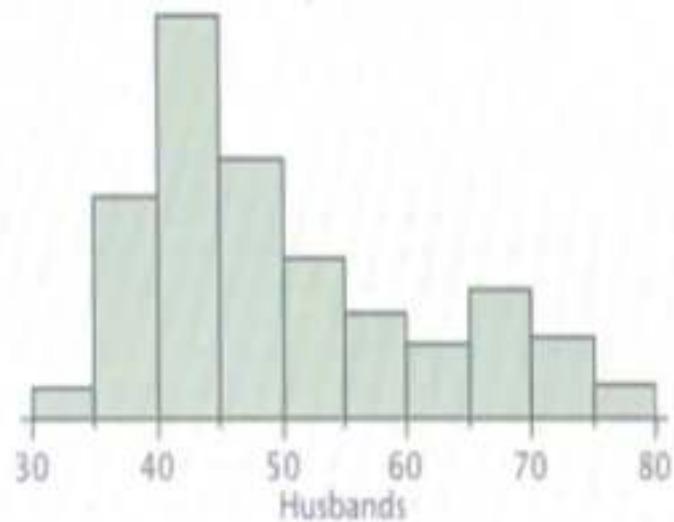
- A dataset with a pair of variables which may be correlated to one another.

Eg: two variables – ice cream sales and temperature



- Bivariate Data

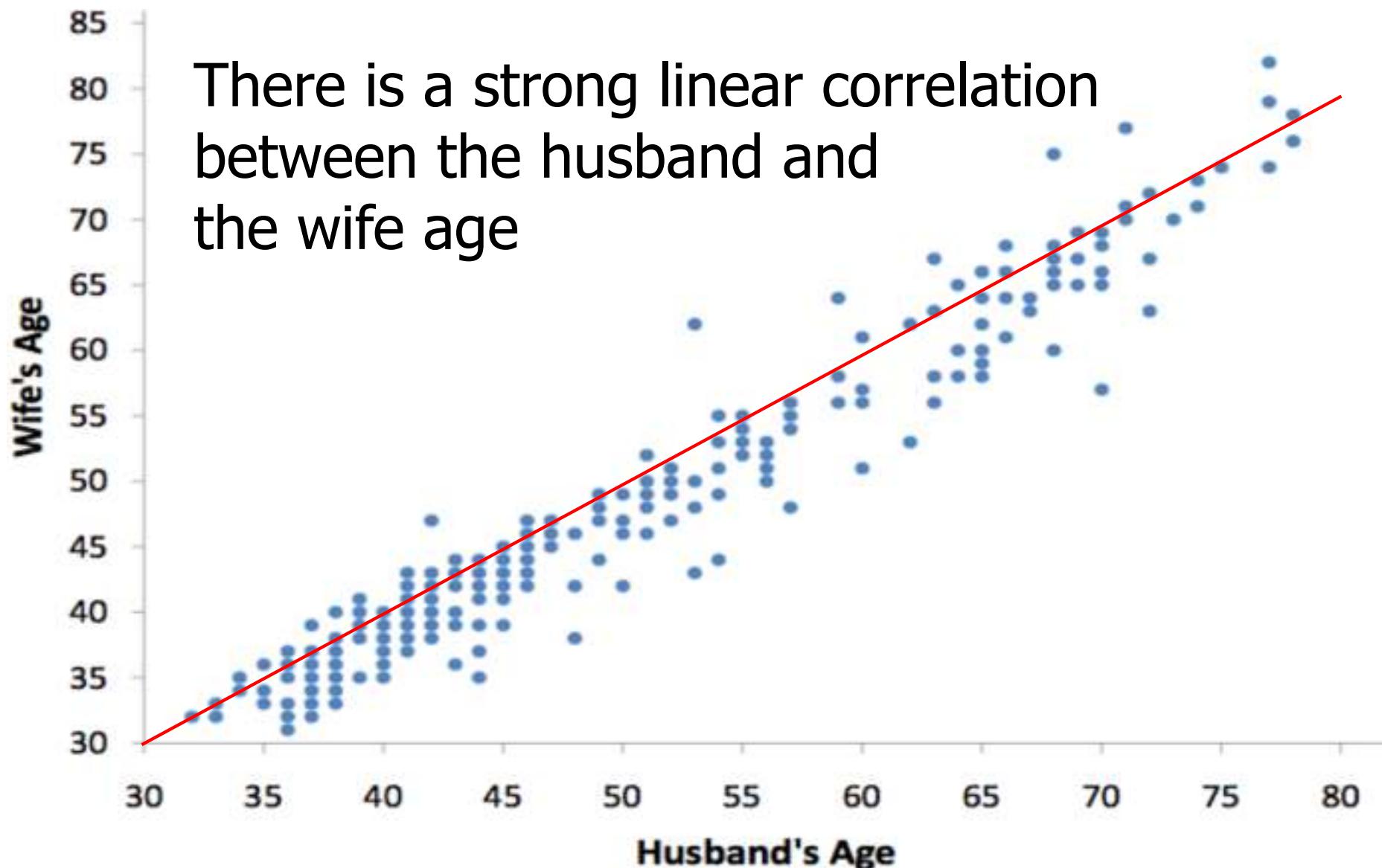
Eg: Figures below show the frequency distribution of the age of husband and wife for 280 couples.



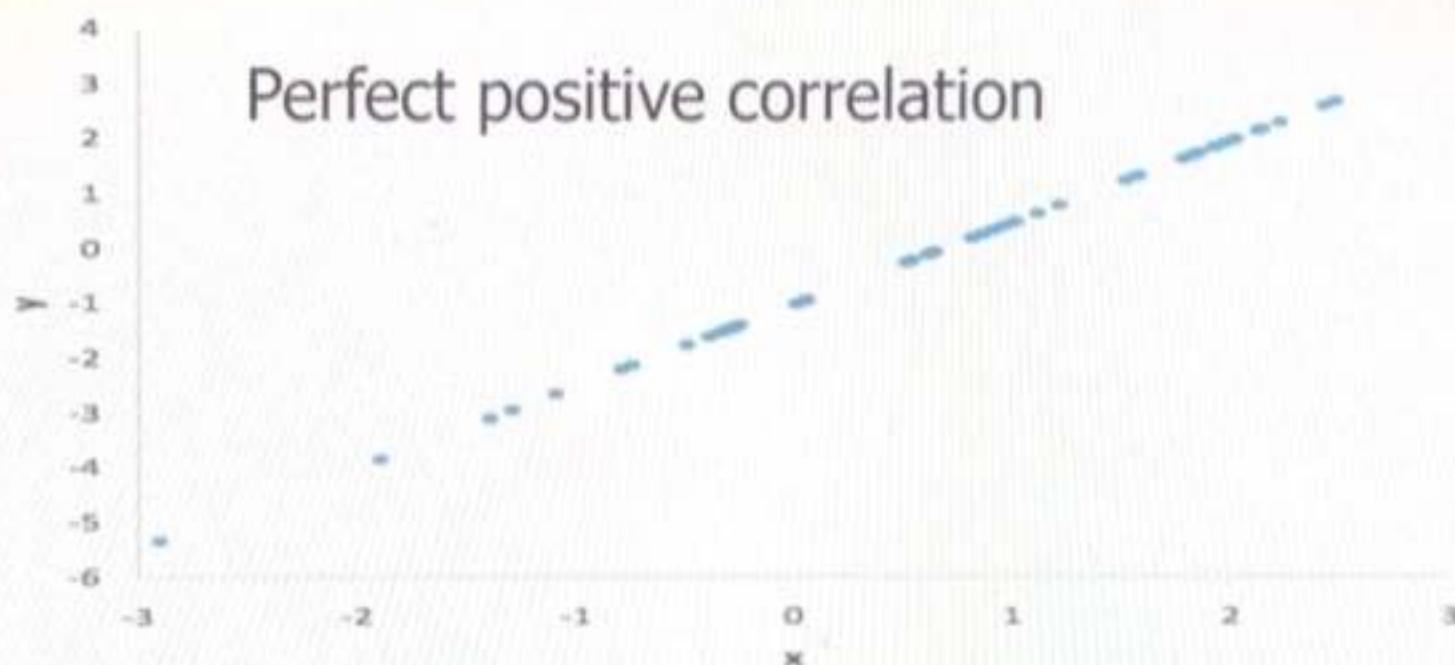
Indication on the central tendency, variability and the shape of distribution of husband and wife age.
No indication on the correlation of the 2 variables.

We
this one

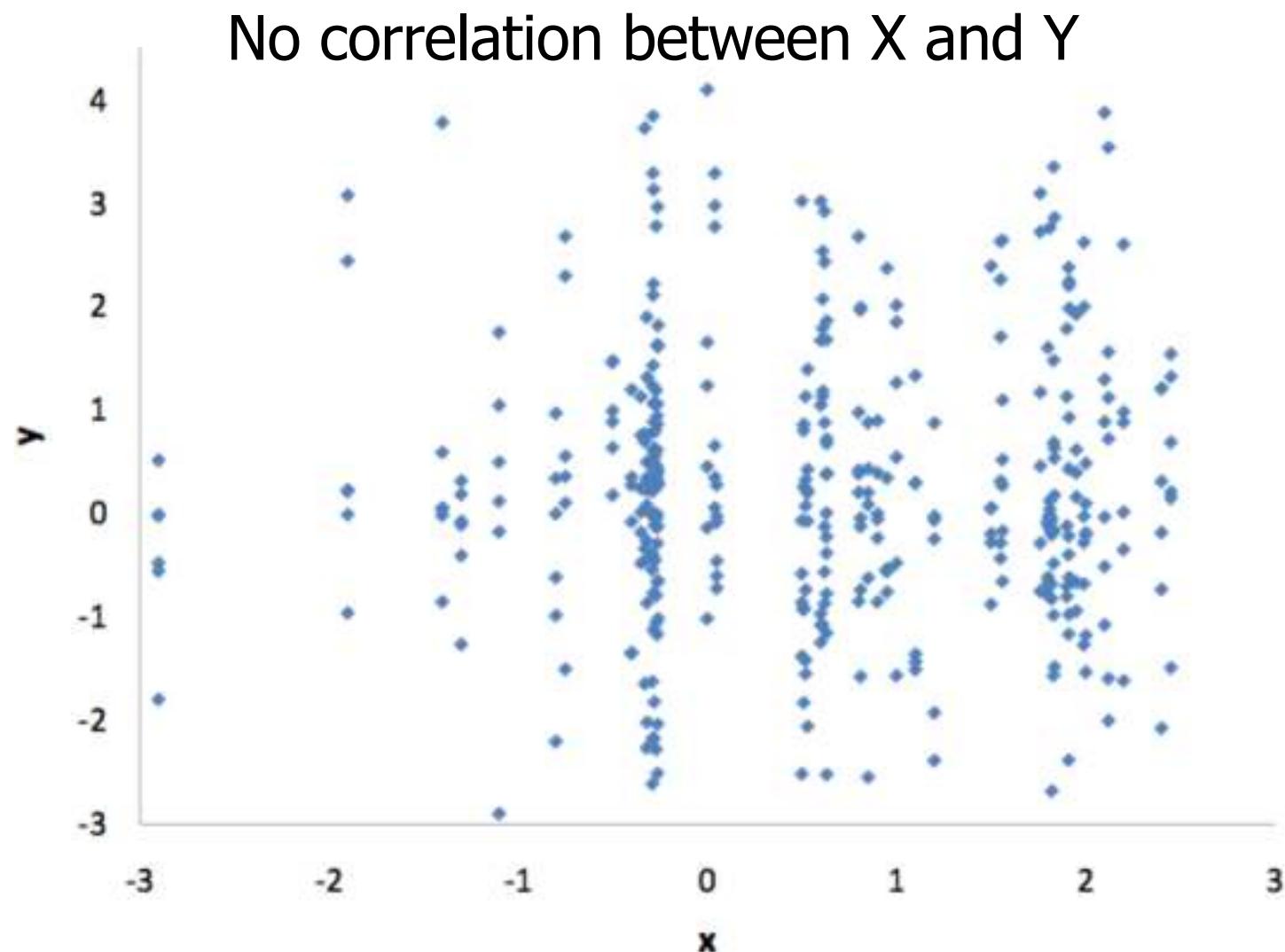
- Bivariate Data



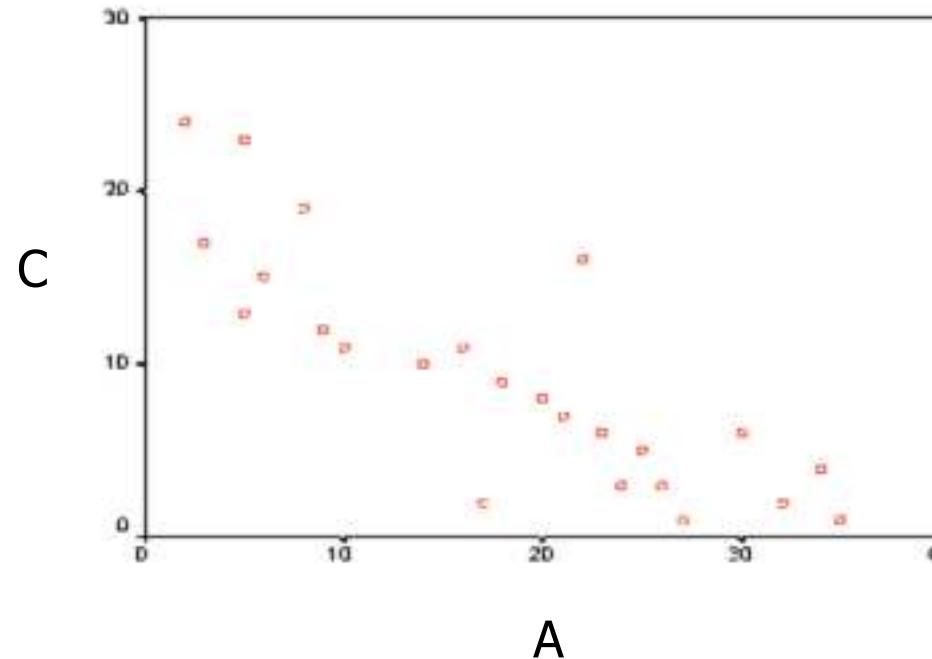
Eg: Correlation of variables X and Y



Eg: Correlation of variables X and Y



Question: Describe the relationship between variables A and C. Think of things these variables could represent in real life.



Negative relationship between A and C.
There is a negative relationship between price and quantity of the products that we buy.

Pearson Correlation ρ

- An indicator on the strength of the linear relationship between two variables.

$$\text{Definition: } \rho = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

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

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

If $\mu_X = \mu_Y = 0$, then $\rho = \frac{\sum XY}{\sqrt{\sum X^2} \sqrt{\sum Y^2}}$

- Computation of Correlation based on a sample of size n

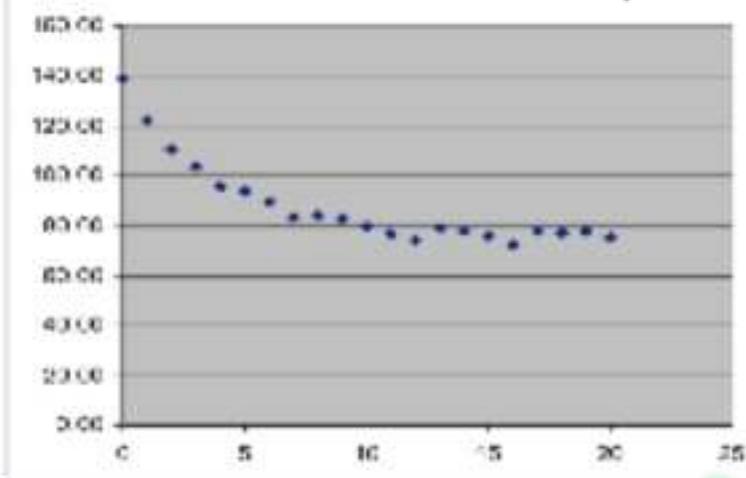
$$\text{cov}(XY) = \frac{1}{n-1} \sum (X - \bar{X})(Y - \bar{Y})$$

$$\text{Correlation } r = \frac{E[(X-\bar{X})(Y-\bar{Y})]}{S_X S_Y}$$

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

$$\text{If } \bar{X} = \bar{Y} = 0, \text{ then } r = \frac{\sum XY}{\sqrt{\sum X^2} \sqrt{\sum Y^2}}$$

Given the data shown in the figure, is it appropriate to...
Pearson correlation to describe the relationship between x and y?



1

Yes

0%

0

2

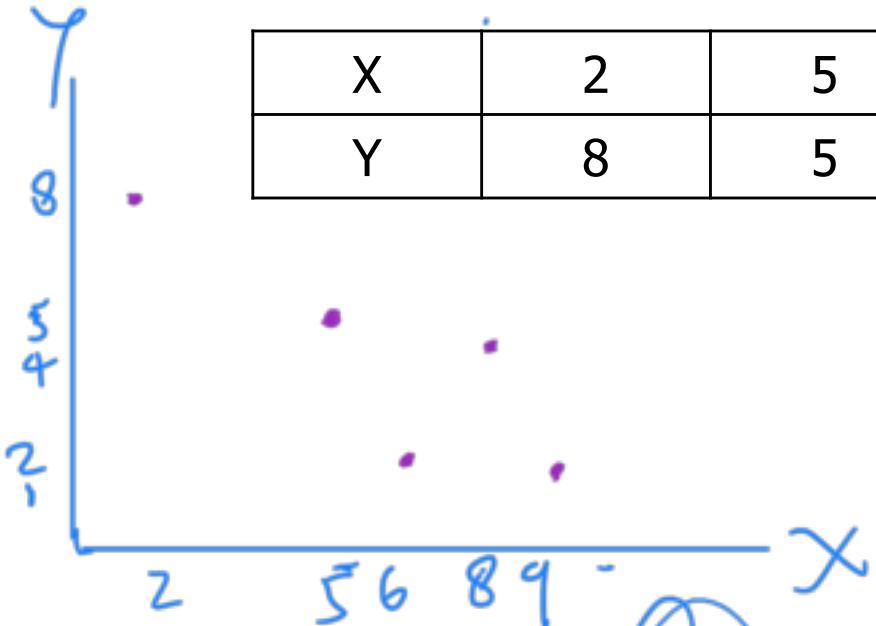
No

0%

0

because it
is not a
linear curve

Eg: Given the data set, calculate σ_x , σ_y , $\text{cov}(X,Y)$ and the Pearson Correlation.



$$\mu_x = \frac{\sum x}{N} = \frac{30}{5} = 6$$

$$\mu_y = \frac{\sum y}{N} = \frac{20}{5} = 4$$

$$\sigma_x^2 = \frac{\sum (x - \mu_x)^2}{N} = \frac{16 + 1 + 0 + 4 + 9}{5} = 6$$

$$\sigma_y^2 = \frac{\sum (y - \mu_y)^2}{N} = \frac{16 + 1 + 4 + 0 + 9}{5} = 6$$

$$\text{cov}(x,y) = E[(x - \mu_x)(y - \mu_y)]$$

$$\frac{\sum (x - \mu_x)(y - \mu_y)}{N} = \frac{-16 - 1 + 0 + 9}{5} = -2.6$$

$$P = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} = \frac{-5.2}{\sqrt{6} \cdot \sqrt{6}} = -0.867$$

Properties of Correlation

- Value in the range of $[-1, +1]$
- Symmetric: correlation of X with Y
= correlation of Y with X
- Unaffected by linear transformations:
Correlation of Y with X
= correlation of Y with $A X + B$
where A and B are constants

Eg: If the correlation between weight (in pounds) and height (in feet) is 0.58, find:

- (a) the correlation between weight (in pounds) and height (in yards)
- (b) the correlation between weight (in kilograms) and height (in meters).

The correlation for both (a) and (b) is still 0.58 because linear transformations do not affect the value of Pearson's correlation, and both of the above instances are linear transformations.

- Variance Sum Law II

- Linear combination of 2 independent variables X and Y

$$\text{Variance of } X \pm Y: \sigma_{x \pm y}^2 = \sigma_x^2 + \sigma_y^2$$

- If the variables X and Y are correlated

$$\text{Variance of } X \pm Y: \sigma_{x \pm y}^2 = \sigma_x^2 + \sigma_y^2 \pm 2\rho\sigma_x\sigma_y$$

- For computation based on a sample

$$s_{x \pm y}^2 = s_x^2 + s_y^2 \pm 2rs_xs_y$$

Eg: Students took 2 parts of a test, each worth 50 points. Part A has a variance of 25, and Part B has a variance of 49. The correlation between the test scores is 0.6.

- (a) If the teacher adds the grades of the two parts together to form a final test grade, what would the variance of the final test grades be?
- (b) What would the variance of Part A - Part B be?

$$\sigma_a^2 = 25$$

$$\sigma_b^2 = 49$$

$$\rho = 0.6$$

a) $\sigma_{x+y}^2 = \sigma_x^2 + \sigma_y^2 + 2\rho\sigma_x\sigma_y$
 [adds] $= 25 + 49 + 2 \cdot 0.6 \cdot 5 \cdot 7$
 $= 74 + 42 = 116$

b) $\sigma_{x-y}^2 = 74 - 42 = 32$
 (Part A - Part B)

Ch 5. Probability

- Basic Concepts – Review of Set Theory and Venn Diagram
- Counting – Ordered, Unordered, Sampling With and Without Replacement.
- Probability Theory
- Base Rate: Bayes' Theorem

- Set Theory - Review

A set is a collection of elements.

Eg: $A = \{\text{Head, Tail}\}$, $\text{Head} \in A$, $\text{Tail} \in A$

$B = \{1, 2, 3, 4, 5, 6\}$

$C = \{1, 3, 5\}$

Set C is a subset of Set B since all elements in C are also in B. Notation: $C \subset B$

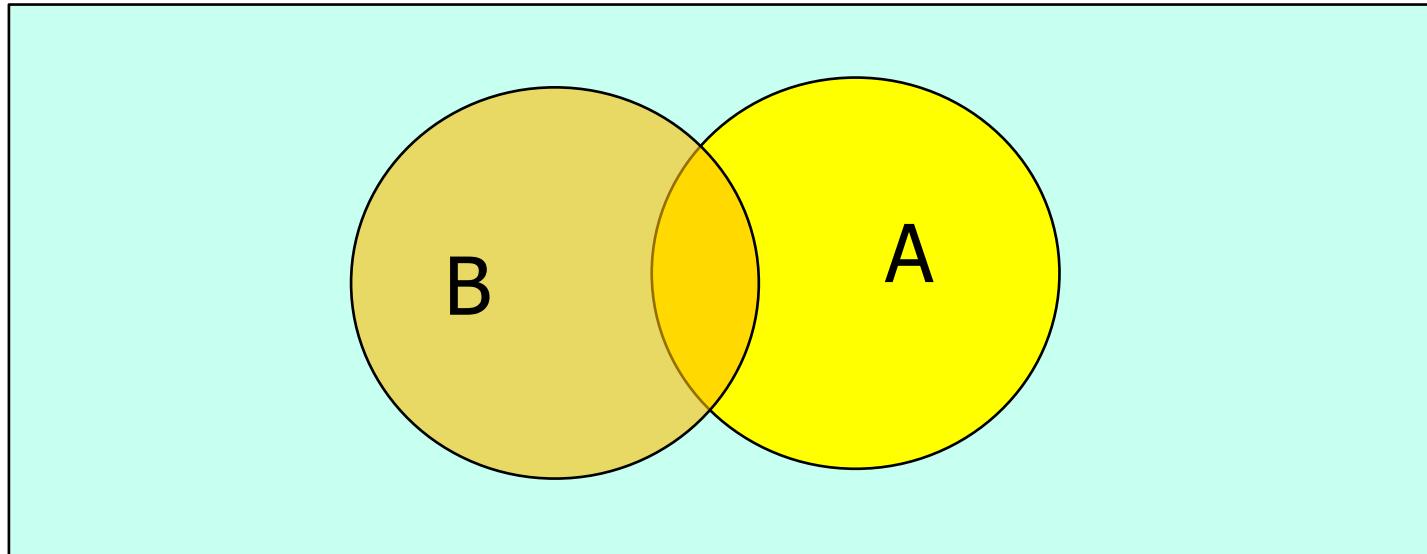
The set of all elements in a given context is known as an universal set Ω . $=$ *Set of all Elements*

Null set \emptyset is a set with no element.

A' : Complement of A

- Venn Diagram

universal set Ω



Set $A \subset \Omega$. Complement of A , denoted by A' , consists of all elements in Ω that are not in A .

Eg: $\Omega = \{1,2,3,4,5,6\}$, $A=\{1,2\}$, $A'=\{3,4,5,6\}$

Union of A and B consists of all elements in A or B . Notation: $A \cup B$.

Eg: $\{1,2,3\} \cup \{2,3,4\} = \{1,2,3,4\}$

Intersection of A and B consists of all elements in both A and B . Notation: $A \cap B$.

Eg: $\{1,2,3\} \cap \{2,3,4\} = \{2,3\}$

Consider rolling a die. Let the outcome set $A=\{1,2,3\}$ and set $B=\{2,4,6\}$. Match the following:

www.wooclap.com/CX2100

Consider rolling a die. Let the outcome set $A=\{1,2,\underline{3}\}$ and set $B=\{2,4,6\}$. Match the following:

 Let's vote!

0%
of participants have already answered

$A \cup B$ — $\{1,2,3,4,6\}$
 $A \cap B$ — $\{2\}$
 $A' \cap B$ — $\{4,6\}$
 $A' \cap B'$ — $\{5\}$




wooclap   90 %  0 / 0   

- Counting

When solving probability problems, we may involve counting.

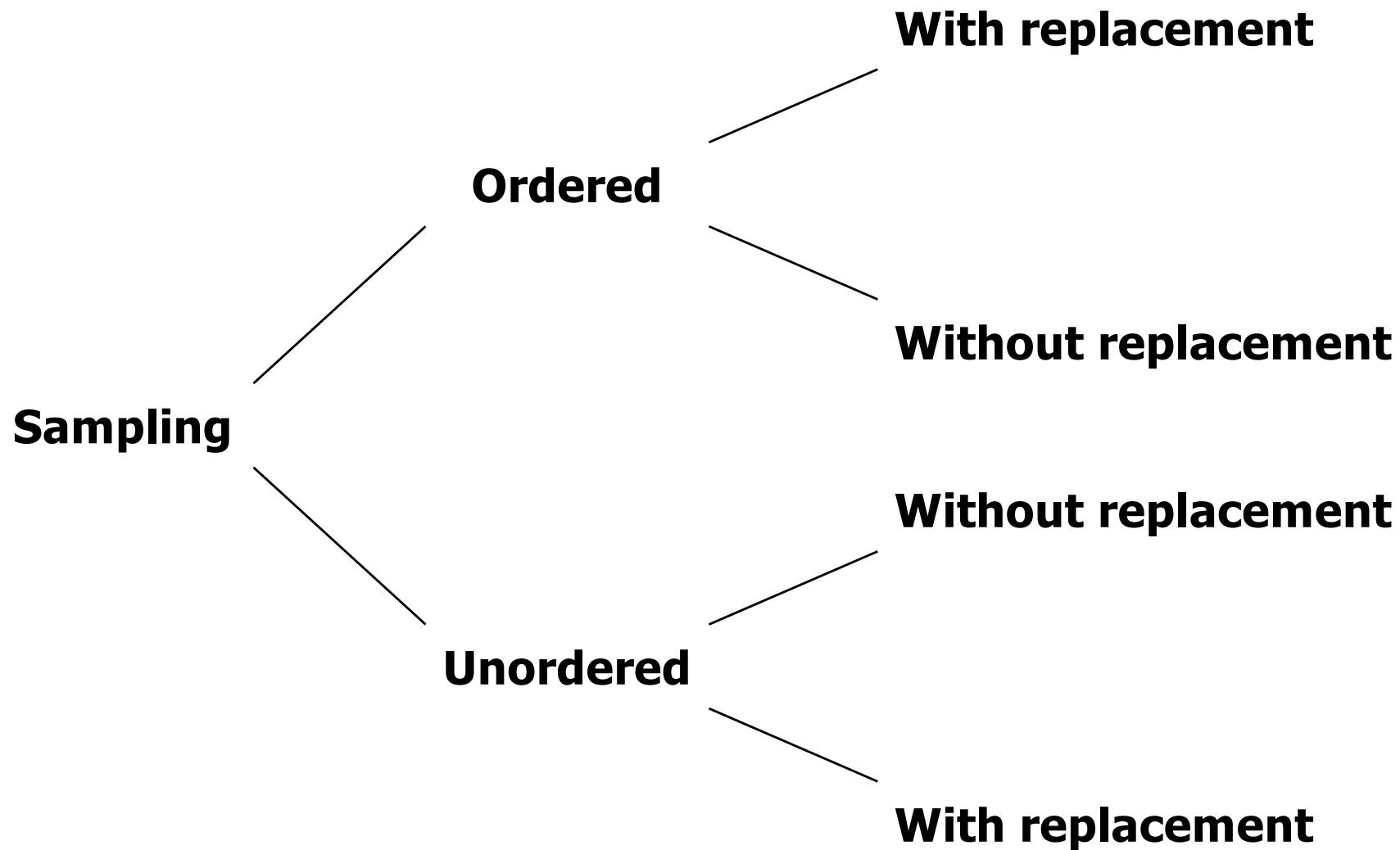
Eg: Calculating the probability of guessing the correct 4 digits numeric code.

Each digit is sampled from $\{0,1,2,\dots,9\}$.

Different scenarios for consideration:

- sampling with replacement
- sampling without replacement
- digits drawn are ordered
- digits drawn are unordered

- Counting – general scenario



- Counting – ordered sampling with replacement

Ordered sampling k elements from a set of n elements with replacement.

! No. of possible arrangements = n^k

Eg: Calculate the total number of different 4-digit numeric codes. Each digit is independently selected from $\{0,1,2,\dots,9\}$.

1st digit: 10 possible numbers

2nd digit: also 10 possible numbers.

3rd and 4th digit: each also 10 possible numbers.

Total arrangements = $10 \times 10 \times 10 \times 10 = 10^4$

- Counting – ordered sampling without replacement

Ordered sampling k elements from a set of n elements without replacement.

No. of possible arrangements =

$$n \times (n - 1) \times (n - 2) \times (n - 3) \times \dots \times (n - k + 1)$$

Number of k permutations of n elements:

$${}^n P_k = \frac{n!}{(n-k)!}$$

$n = 10$

$k = 4$

Eg: Calculate the total number of different 4-digit numeric codes. Each digit is sampled from $\{0,1,2,\dots,9\}$ without replacement.

1st digit: 10 possible numbers

2nd digit: 9 possible numbers

3rd digit: 8 possible numbers and so on.

Total arrangements = $10 \times 9 \times 8 \times 7$

i.e. there are ${}^{10}P_4 = \frac{n!}{(n-k)!} = \frac{10!}{(10-4)!}$ arrangements

Eg: 4 boys and 2 girls sit in a row. Find the following:

www.wooclap.com/CX2100

4 boys and 2 girls sit in a row. Find the following:



(1) No. of ways of putting these 6 people in a row.

1

$$6! = 720$$



(2) No. of ways such that each girl has a boy to her left and to her right.

b b g b g b
b g b g b b

b g b b q b
possible pattern

3 × ways of placing boys → ways of placing girl

$$3 \times 4! \times 2! \\ = 3 \times 4! = 144$$

wooclap



90 %



0 / 0



- Counting – unordered sampling without replacement

We choose k elements from a set of n elements and the ordering does not matter.

Number of arrangements equals to k combinations of n elements

$$\binom{n}{n-k} = \binom{n}{k} = \frac{n!}{k! (n-k)!}$$

use for
→ binomial
too

Note: $\binom{n}{k}$

$$= \frac{n!}{k! (n-k)!}$$
$$= \binom{n}{n-k}$$

Eg: How many combinations of two numbers between 1 and 6 are there:

$$\text{Ans: } \frac{6!}{2! (6-2)!} = \frac{6 \times 5}{2 \times 1} = 15$$

- Counting – unordered sampling with replacement

Pick k elements from a set of n elements, one at a time with replacement and the ordering does not matter.

Eg: Pick two numbers from $\{1, 2, 3\}$, we have

6 arrangements with $n=3$ and $k=2$.

i.e. $(1,1)$, $(1,2)$, $(1,3)$, $(2,2)$, $(2,3)$ and $(3,3)$

In general:

$$\text{No. of arrangements} = \binom{n+k-1}{k}$$

Eg: A 6-bit code is made up of 4 1's and 2 0's.
Calculate:

www.wooclap.com/CX2100

A 6-bit code is made up of four '1's and two '0's.

Find the following:

$$\binom{6}{2} {}^6C_4 = {}^6C_2 = \frac{6!}{2!4!} = 15$$

(1) No. of distinct codewords in the code.

(2) No. of codewords such that a '0' has a '1' to its left and to its right.

101011
110101
101101

wooclap



90 %



0 / 0



- Probability Theory

In a random experiment, the outcome of the experiment occurs with a certain probability.

Eg: We toss a coin. If the coin is unbiased, then the chances of getting a Head is 50%.

Eg: We roll a fair die two times. The outcome $\in \{11, 12, 13, 14, 15, 16, 21, 22, \dots, 65, 66\}$.

Total of 36 possible outcomes.

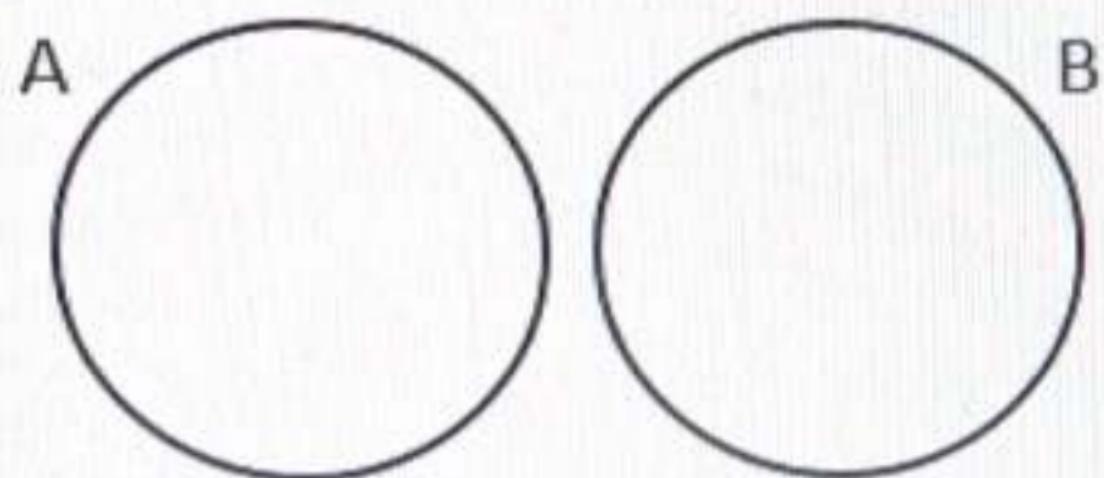
Probability of getting a “1” followed by another “1” is $1/36$.

Probability of getting 1st no. = 2nd no. is $6/36$.

Definitions

- A sample space is the set S of all possible outcomes of an experiment.
- An event is a set of one or more (favorable) outcomes in the sample space.
- Two events are mutually exclusive if they have no outcomes in common.
- They are exhaustive if they cover all possible outcomes.
- Two events are independent if the probability that one occurs is not affected by whether or not the other has occurred.

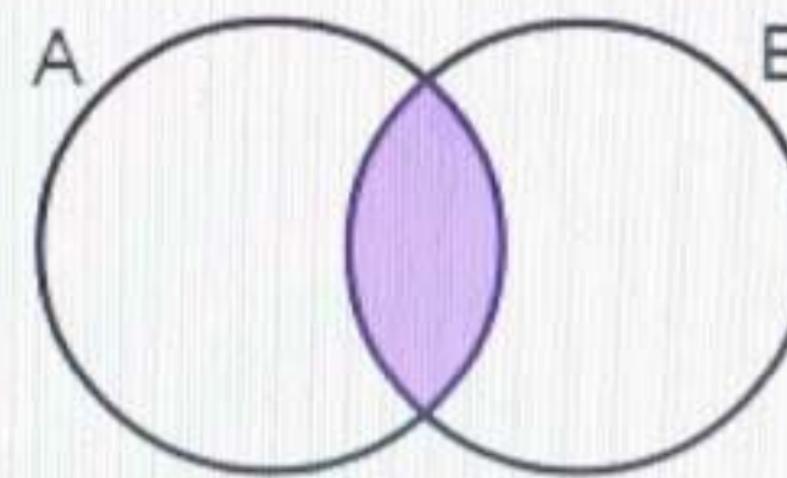
Mutually Exclusive Events



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

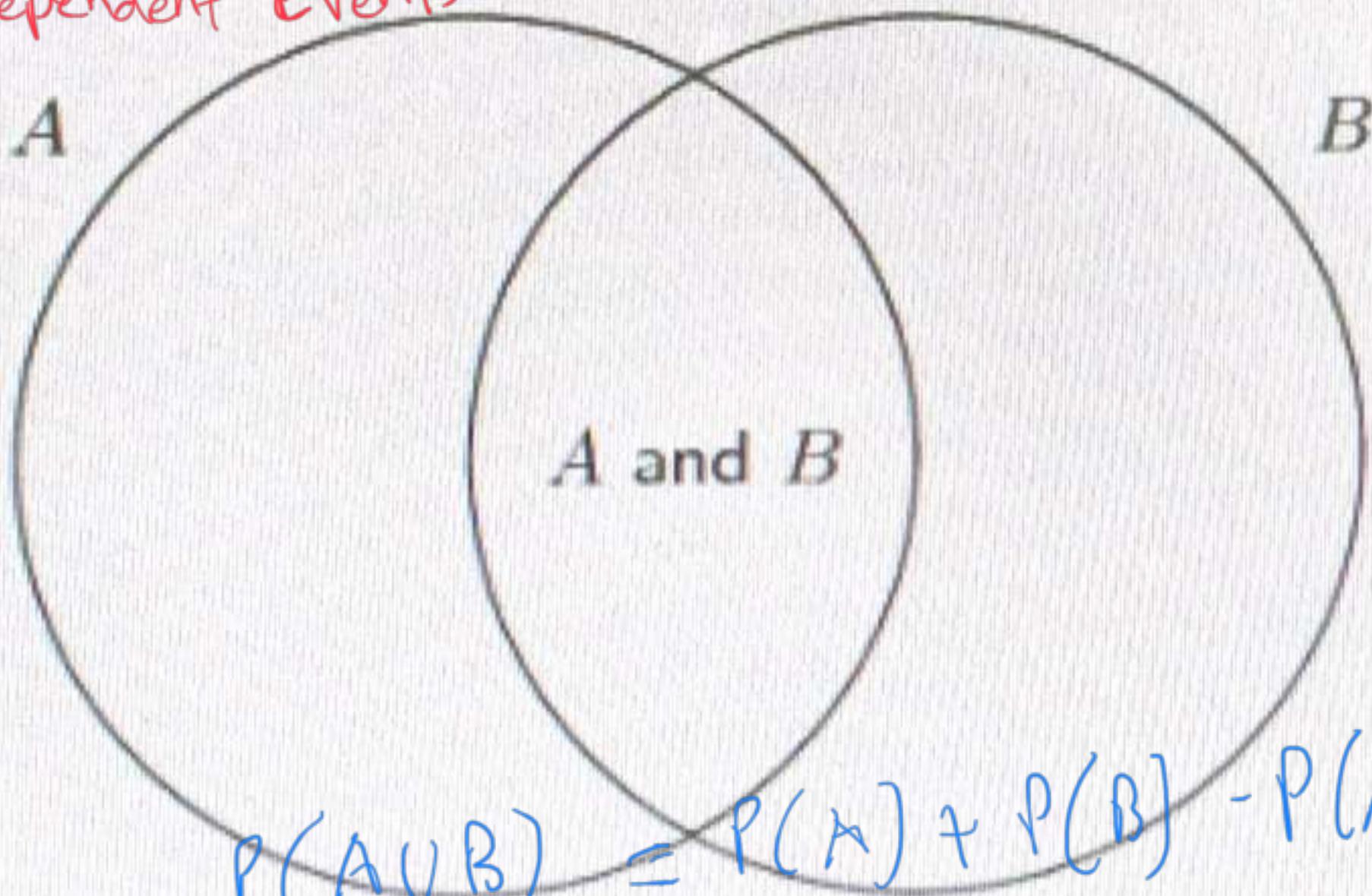
$$P(A \cap B) = 0$$

Non-Mutually Exclusive Events



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

Independent Events



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

Axioms of Probability:

- (1) Probability of any outcome or event X is a non-negative:

$$P(X) \geq 0$$

- (2) Probability of the sample space S is 1:

$$P(S) = 1$$

- (3) If X_1, X_2, X_3, \dots are mutually exclusive events, then:

$$P(X_1 \text{ or } X_2 \text{ or } X_3 \dots) = P(X_1) + P(X_2) + P(X_3) + \dots$$

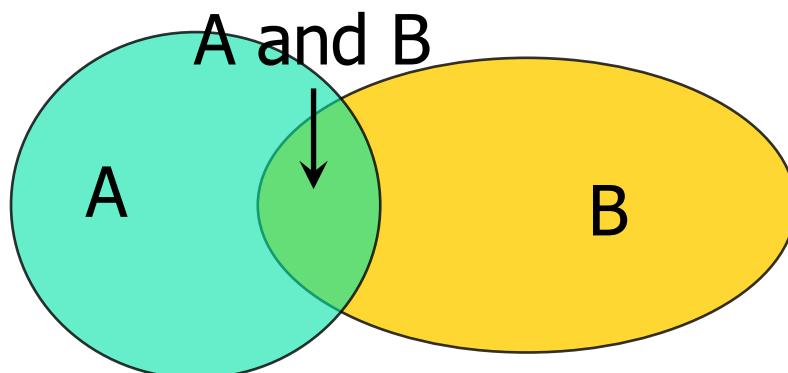
Probabilities and Events (mutually exclusive)

Eg: A company has decided that in the next 5 years, 40% of their new employees will be men, 30% will be Singaporean, and 35% will be foreigner women. What percentage of new employees will be Singaporean men?

	Men	Women	Total
Foreigner	0.35	0.35	0.7
Singaporean	0.05	0.25	0.3
Total	0.4	0.6	1

Probabilities and Events (non-mutually exclusive)

Eg: A certain kind of fruit is grown in 2 districts, A and B. Both areas sometimes get fruitflies. Suppose the probabilities are $P(A)=0.1$, $P(B)=0.05$ and $P(A \text{ and } B)=0.02$, what is the probability that one or other (or both) districts are infected at a given time?



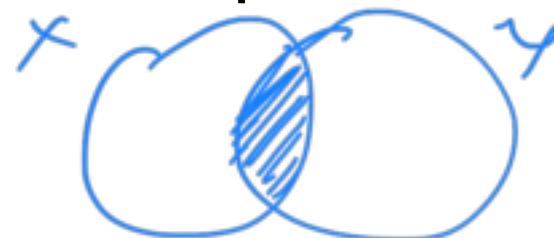
$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ &= 0.1 + 0.05 - 0.02 \\ &= 0.13 \end{aligned}$$

Probabilities and Events (independent)

Eg: A company has 2 guards. Each carries a pager activated by sensors. Guard 1 and guard 2 respond to pager alert 80% and 50% of the time respectively. They independently report any alert. What is the probability that at least one will report an alert?

$$\begin{cases} x = G_1 = 0.8 = P(x) \\ y = G_2 = 0.5 = P(y) \end{cases}$$

↓ events



probability of 2 independent event = $P(x) \cdot P(y)$
= 0.4

Probability of at least one reports = $P(x) + P(y) - P(x \cap y)$
= 0.8 + 0.5 - 0.4
= 0.9 ↗

Conditional Probabilities

$$P(x | y)$$

↳ given that

Assuming that 10 percent of the days are rainy in a certain city.

$$P(\text{rain}) = 0.1$$

The probability that it rains given that it is cloudy might be, say,

$$P(\text{rain}|\text{cloudy}) = 0.8$$

This is known as conditional probability: the probability of event A given that event B has occurred.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \text{ where } P(B) > 0$$

Conditional Probabilities

www.wooclap.com/CX2100

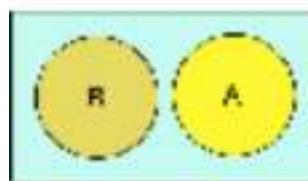


Fig 1.



Fig 2.

Given $P(A)=0.5$ and $P(B)=0.8$, determine the following.

- (1) $P(A|B)$ for Fig 1.
- (2) $P(A|B)$ for Fig 2.
- (3) $P(B|A)$ for Fig 2.

The most frequent answers are

(1)



Conditional Probabilities

Eg: A jar contains 10 blue, 5 red, 4 green and 1 yellow marbles. Two marbles are randomly picked. What is the probability that one will be blue and the other yellow?

$$P(1^{\text{st}} B \cap 2^{\text{nd}} Y) = P(1^{\text{st}} B) \cdot P(2^{\text{nd}} Y \mid 1^{\text{st}} B)$$
$$\approx \frac{10}{20} \cdot \frac{1}{19} = 0.0263$$

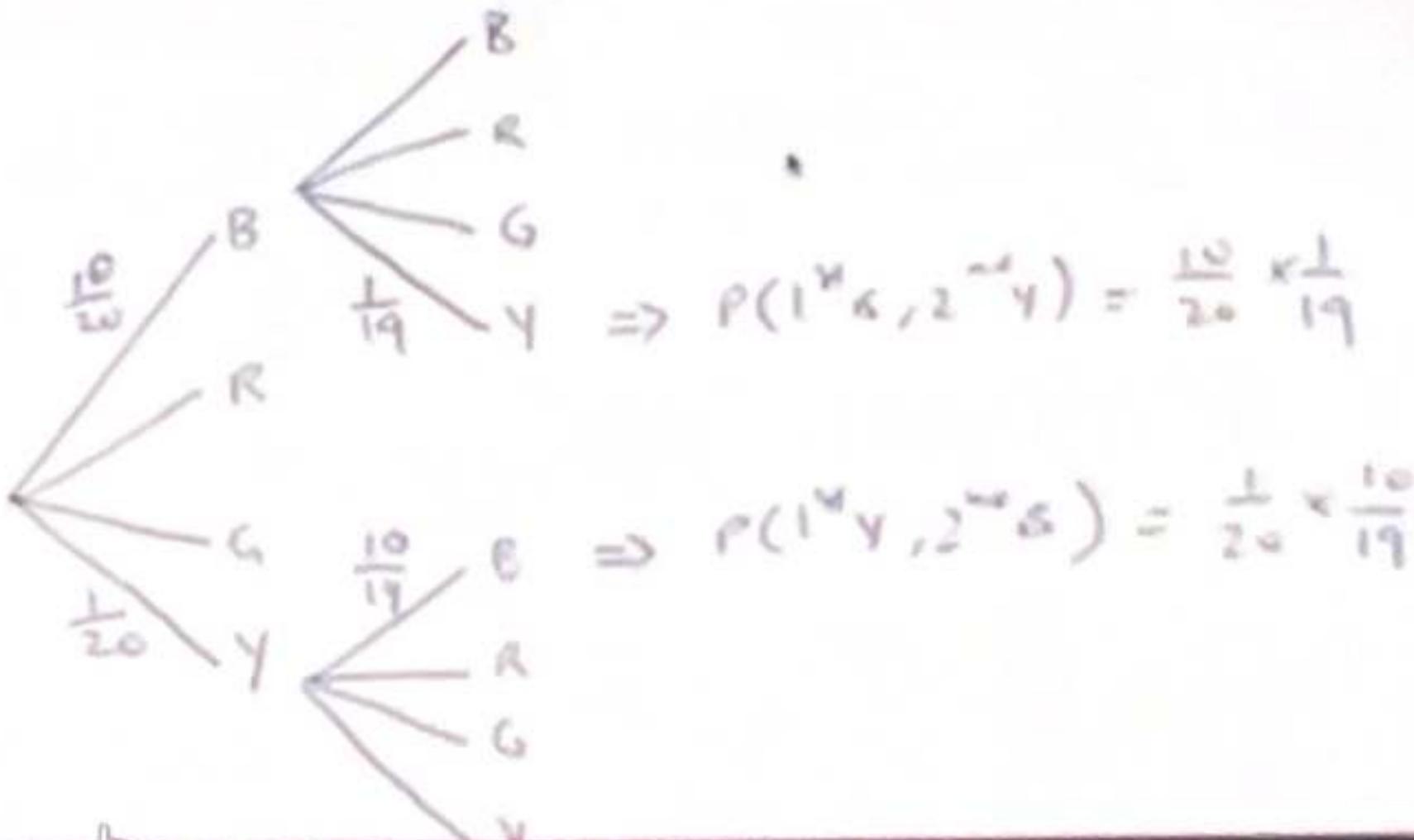
$$P(1^{\text{st}} Y \cap 2^{\text{nd}} B) = P(1^{\text{st}} Y) \cdot P(2^{\text{nd}} B \mid 1^{\text{st}} Y)$$
$$\approx \frac{1}{20} \cdot \frac{19}{19} = 0.0263$$

$$P(\text{1 blue} \cap \text{1 yellow}) = P(1^{\text{st}} B \cap 2^{\text{nd}} Y) + P(1^{\text{st}} Y \cap 2^{\text{nd}} B)$$
$$\approx 0.0526$$

Illustration using a tree diagram:

Eg: A jar contains 10 blue, 5 red, 4 green and 1 yellow marbles. Probability of getting one blue and one yellow:

Illustration using a tree diagram:



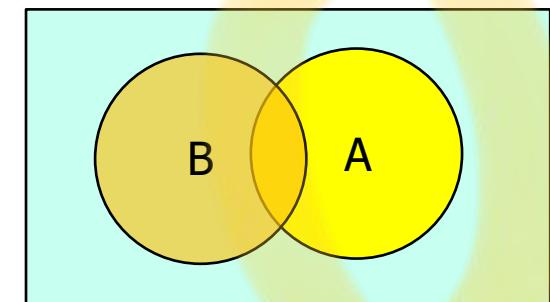
Bayes' Theorem

Given two events A and B, where $P(A) > 0$, we have:

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

Note: $P(A) = P(A \cap B) + P(A \cap B')$

Similarly:



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

Bayes' Theorem

Eg: A test correctly identifies a disease in 95% of people who have it. It correctly identifies no disease in 94% of people who do not have it. In the population, 3% of the people have the disease. What is the probability that one has the disease if tested positive?

$$P(D) = 0.03 \rightarrow P(D') = 0.97 \quad D = \text{one w/ disease}$$
$$P(T|D) = 0.94 \rightarrow P(T|D') = 0.06 \quad T = \text{test is +}$$
$$P(\bar{T}|D) = 0.95 \rightarrow P(\bar{T}|D') = 0.05$$

$$P(D|T) = \frac{P(D \cap T)}{P(T)} = \frac{P(T|D) P(D)}{P(T \cap D) + P(T \cap D')} = \frac{0.95 \cdot 0.03}{P(+|D) P(D) + P(+|D') P(D')} = \frac{0.95 \cdot 0.03}{\frac{0.95 \cdot 0.03 + 0.06 \cdot 0.97}{0.95 \cdot 0.03 + 0.06 \cdot 0.97}} = \frac{0.0285}{0.0867} = 0.329$$

Summary For Probability Calculations



Are X and Y mutually exclusive?



Yes

No

1. $P(X \text{ and } Y) = 0$ because event $(X \text{ and } Y)$ is impossible
2. $P(X \text{ or } Y) = P(X) + P(Y)$

$$1. P(X \text{ and } Y) > 0$$

$$2. P(X \text{ or } Y) = P(X) + P(Y) - P(X \text{ and } Y)$$

Are X and Y independent?

Yes

No

1. $P(X|Y) = P(X)$
2. $P(X \text{ and } Y) = P(X) P(Y)$

$$1. P(X|Y) = \frac{P(X \text{ and } Y)}{P(Y)}$$

$$\begin{aligned}2. P(X \text{ and } Y) &= P(X|Y) P(Y) \\&= P(Y|X) P(X)\end{aligned}$$

Ch 6. Probability Distributions

- Random variables – Discrete
 - Continuous
- Discrete Probability distributions – Binomial, Poisson, Geometric
- Continuous Probability distributions – Uniform, Exponential, Normal

Random Variables

The outcomes of random experiments commonly takes on numerical values.

If the outcomes are not numerical, they can be represented in terms of numbers (this facilitates mathematical analysis).

A random variable (r.v.) is a variable that has a numerical value which is defined on or determined by the outcomes or events of an experiment. The numerical value cannot be predicted exactly but can be described by its probability.

Random Variables

Random variables can be discrete or continuous.

Eg: The no. of seeds germinating from a flower pot. Possible values for X are 0, 1, 2, ... (discrete).

The maximum daily temperature in Singapore.
Possible values are 0 – 50°C, e.g. 26.1276°C
(continuous).

The response to questions with “Yes”, “No”, or
“Don’t know” answers. This are not a r.v. (the
values are not numerical).

Let Y be the number of “Yes”s.
 Y is a discrete r.v.

Random Variables

Expected Value of a r.v.

For a given set of observed data $\{x_1, x_2, \dots, x_n\}$, one of the central tendency measure is the arithmetic average:

$$\bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

\bar{X} is calculated based on a sample of n observed values.

If we know the probability of occurring associated with each value of a r.v., then we can obtain the weighted average which is known as the expectation of the r.v.

Random Variables

(μ_x)

Expected Value of a r.v.

The probability distribution or probability mass function (pmf) of a r.v. X is given as:

Values of X	x_1	x_2	...	x_k
Probabilities	p_1	p_2	...	p_k

Its expected value is defined as:

pmf : $E(X) = x_1 p_1 + x_2 p_2 + \dots + x_k p_k = \sum_{i=1}^k x_i p_i$

Also denoted by μ_x

Random Variables

Eg: (Decision Analysis)

An oil exploration company has a lease for which it must decide on one of the following options:

1. Sell now and can get \$125K.
2. Hold for a year and if oil prices rise (prob = 0.6) it can sell for \$300K or if oil prices fall (prob = 0.4) it can get \$100K. $0.6 \times 300K + 0.4 \times 100K = 180K + 40K = \$220K$
3. Drill now. The cost of drilling is \$200K and drilling will lead to one of the following outcomes: $-100K + 80K + 130K = \$110K$

Well type	Dry	Wet	Gusher
Probability	0.5	0.4	0.1
Profit	\$0 -200	\$400K ²⁰⁰	\$1500K ¹³⁰⁰

What should the company do?

(Z)

Random Variables

Solution: Let X be the financial gain.

Option 1

Option 2

(hold for 1 yr):

X	300	100
Probability	0.6	0.4

$$\therefore E(X) =$$

Option 3, $X = \text{profit} - \text{drilling costs } (\$200K)$

X	0 – 200K	400K – 200K	1,500K – 200K
Probability	0.5	0.4	0.1

$$\therefore E(X) =$$

Decision:

Woodlap Ex:

Given the pmf of X below, if $Y=(X-1)^2$, find $E[Y]$.

x	0	1	2
$P_X(x)$	0.2	0.5	0.3

Random Variables

Variance of a r.v.

Recall that variance is a measure of variability. For a sample of n independent observations the sample variance is obtained using:

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

The variance of a r.v. X is defined as

$$\begin{aligned} \text{var}(X) &= p_1(x_1 - \mu)^2 + p_2(x_2 - \mu)^2 + \cdots + p_n(x_n - \mu)^2 \\ &= \sum_{i=1}^n p_i(x_i - \mu)^2 = E[(x - \mu)^2] \end{aligned}$$

where p_i = probability that x_i occurs.

Random Variables

It represents the theoretical limit of the sample variance s^2 as the sample size n becomes very large. $\text{Var}(X)$ is often denoted by σ^2 .

An alternate formula for $\text{var}(X)$ is as follows:

$$\text{var}(X) = \sum p_i (x_i^2 - 2x_i\mu + \mu^2)$$

$$= (p_1 x_1^2 + \dots + p_n x_n^2) - \mu^2$$

$$= E[X^2] - \mu^2 \quad \text{or} \quad E[X^2] - (E[X])^2$$

Random Variables

Eg

Let the random variable X be the number of male students in a group from a class of size 5. The probability distribution of X is

X	0	1	2	3	4	5
$P(X = x)$	1/32	5/32	10/32	10/32	5/32	1/32

What is the expected number of males $E(X)$ in the group, and what is the standard deviation, σ , of X ?

Random Variables

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

$$\mu = 0 + \frac{5}{32} + \frac{20}{32} + \frac{30}{32} + \frac{20}{32} + \frac{5}{32}$$

$$\mu = \frac{80}{32} = 2.5$$

$$\text{var}(X) = \sum x^2 P(x) - \mu^2$$

$$\left(0^2 \cdot \frac{1}{32} + 1^2 \cdot \frac{5}{32} + 2^2 \cdot \frac{10}{32} + 3^2 \cdot \frac{10}{32} + 4^2 \cdot \frac{5}{32} + 5^2 \cdot \frac{1}{32} \right) - \mu^2$$
$$\left(0 + \frac{5}{32} + \frac{40}{32} + \frac{90}{32} + \frac{80}{32} + \frac{25}{32} \right) - (2.5)^2$$

$$\sigma^2 = 7.5 - 6.25 = 1.25$$

Wooclap Ex:

Given the pmf of X below, if $Y = (X-1)^2$, find $\text{Var}[Y]$.

x	0	1	2
$P_X(x)$	0.2	0.5	0.3

$$\text{Var}[Y] = E[Y^2] - E[Y]^2$$

$$E[Y^2] = 0 \times 0.2 + 1 \times 0.5 + 4 \times 0.3 = 0.5 + 1.2 = 1.7$$

$$! E[Y] = 0.5$$

$$\text{Var}[Y] = 1.7 - 0.5^2 = 1.7 - 0.25 = 1.45$$

Random Variables

Empirical quantity (m observed data)	Theoretical quantity (mathematical d.r.v.)	Remarks
Relative freq of x_i is $\frac{f_i}{m}$	$P[X = x_i] = p_i$	$\frac{f_i}{n} \rightarrow p_i$ as $n \rightarrow \infty$
$\sum_{i=1}^m \frac{f_i}{m} = 1$	$\sum_{i=1}^n p_i = 1$	
Mean: $\bar{X} = \sum_{i=1}^m \frac{f_i}{m} x_i$	Expectation: $E[X] = \mu = \sum_{i=1}^n p_i x_i$	$\bar{X} \rightarrow \mu$ as $n \rightarrow \infty$
Variance: $s^2 = \sum_{i=1}^m \frac{(x_i - \bar{X})^2 f_i}{m-1}$	Var[X]: $\sigma^2 = \sum_{i=1}^n (x_i - \bar{X})^2 p_i$	$s^2 \rightarrow \sigma^2$ as $n \rightarrow \infty$

Random Variables

Expected Value and Variance for a Function of r.v.

(Linear transformation of r.v.)

If $Y = a + bX$ where X is a r.v. and a and b are known constant values, then

$$E(Y) = a + b E(X)$$

and

$$\text{var}(Y) = b^2 \text{var}(X)$$

so

$$\sigma_y = \sqrt{b^2 \text{var}(X)} = \sqrt{b^2 \sigma_x^2} = b\sigma_x$$

Random Variables

Similarly if $T = a + bX + cY$ where X and Y are uncorrelated r.v. and a , b and c are known constants, then

$$E(T) = a + bE(X) + cE(Y)$$

and

$$\text{var}(T) = b^2 \text{ var}(X) + c^2 \text{ var}(Y)$$

Random Variables

Eg: A company makes products for local and export markets. The number of sales for next year are estimated as follows:

local, X units	1,000	3,000	5,000	10,000
probability	0.1	0.3	0.4	0.2

export, Y units	300	500	700
probability	0.4	0.5	0.1

Hence $E(X) = \frac{100 + 900 + 2000 + 2000}{5000}$

Random Variables

$$E(Y) = 120 + 250 + 70 = 440$$

Suppose the company makes a profit of \$2000 on each unit sold on the local market and \$3500 on each exported unit.

Hence the total profit is $2000X + 3500Y$

And the expected profit is

$$E(T) =$$

$$\frac{2000 \cdot 5000 + 3500 \cdot 440}{11540000}$$

Random Variables

Consider a component is made by cutting a piece of metal to length X and trimming it by amount Y . Both processes are somewhat imprecise.

Nett length $T = bX + cY$ with $b = 1$ and $c = -1$.

$$E(T) = bE(X) + cE(Y) = E(X) - E(Y)$$

$$\text{Var}(T) = b^2 \text{var}(X) + c^2 \text{var}(Y)$$

↑
Note

i.e. $\text{var}(T)$ is greater than either $\text{var}(X)$ or $\text{var}(Y)$, even though $T = X - Y$, because both X and Y contribute to the variability in T .

Prob. Distribution of Discrete r.v.

Special Probability Distributions

- Binomial Distribution
- Poisson Distribution
- Geometric Distribution

Prob. Distribution of Discrete r.v.

The list of possible values that a d.r.v. X can take and their probabilities is called the discrete probability distribution (also known as probability mass function, pmf) for X .

Eg:

Let the random variable X be the number of girls in a family of 3 children.

Possible values:

$X = 3$	GGG		
$X = 2$	GGB	GBG	BGG
$X = 1$	BBG	BGB	GBB
$X = 0$	BBB		

Prob. Distribution of Discrete r.v.

Assume that the 8 outcomes are equally likely,

x_i	0	1	2	3
Probability Distribution $P(X=x_i)$	1/8	3/8	3/8	1/8
Cumulative Distribution $P(X \leq x_i)$	1/8	4/8	7/8	1

In general, we write:

$$P(X = x_i) = p_i \text{ for } i = 1, \dots, k, \text{ and } 0 \leq p_i \leq 1$$

Notation convention: we often use capital letters for random variables and small letters for specific values.

Prob. Distribution of Discrete r.v.

Eg: Let X denote the sum of the results of 2 dice thrown.

Outcomes:

1,1	2,1	3,1	...	6,1
1,2	2,2	3,2	...	6,2
:	:	:		:
1,6	2,6	3,6	...	6,6

36 outcomes and the values of X are 2, 3, ..., 12.

Assuming that the outcomes are equally likely, i.e. the probability of each outcome is $1/36$, the probability distribution of X is

x	2	3	4	...	10	11	12
$P(X = x)$	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$...	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$

Exercise: find the cumulative distribution.

$$\frac{1}{36}, \frac{3}{36}, \frac{6}{36}, \dots, \frac{33}{36}, \frac{35}{36}, 1$$

Prob. Distribution of Discrete r.v.

Binomial Distribution

- Consider n Bernoulli trials, where each trial is an “experiment” with exactly 2 possible outcomes, "success" and "failure".
- Assume that the probability of success (S) is the same for all trials, $P(S) = p$, and the probability of failure $P(F) = 1 - p$.
- Assume also that trials are independent, i.e., the probability for any given combination of successes and failures can be obtained by multiplying the probabilities for each trial outcome.

Prob. Distribution of Discrete r.v.

Eg: Consider 5 trials, the probability,

$$\begin{aligned}P(SSFSF) &= p \cdot p \cdot (1-p) \cdot p \cdot (1-p) \\&= p^3(1-p)^2\end{aligned}$$

The probability of obtaining any 3 successes and 2 failures in 5 trials, i.e. *SSSFF*, *SSFSF*, ... etc., is $p^3(1-p)^2$ for each of the different ways this could occur.

Prob. Distribution of Discrete r.v.

The number of distinct "arrangements" of 3 successes and 2 failures can be calculated using the binomial coefficient $\binom{n}{x}$ or nC_x

The binomial **coefficient** is defined as: ${}^nC_x = \frac{n!}{x!(n-x)!}$
where $n! = n \times (n-1) \times \dots \times 2 \times 1$

In this example, $\binom{5}{3} = \frac{5!}{3!2!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{(3 \times 2 \times 1)(2 \times 1)} = 10$

so there are 10 distinct ways to obtain 3 successes in 5 trials, with each arrangement having a probability of $p^3(1-p)^2$.

Prob. Distribution of Discrete r.v.

Let X be the r.v. equal to the total number of successes in n trials. To calculate the probability of obtaining x successes, we have:

$$P(X = x) = \binom{n}{x} \cdot p^x \cdot (1 - p)^{n-x}$$

↑
prob. of
 x S's

arrangements of
 x S's and $(n-x)$ F's prob. of
 $(n-x)$ F's

The distribution of the count of successes is called the binomial distribution with two parameters, n and p . We say $X \sim B(n, p)$.

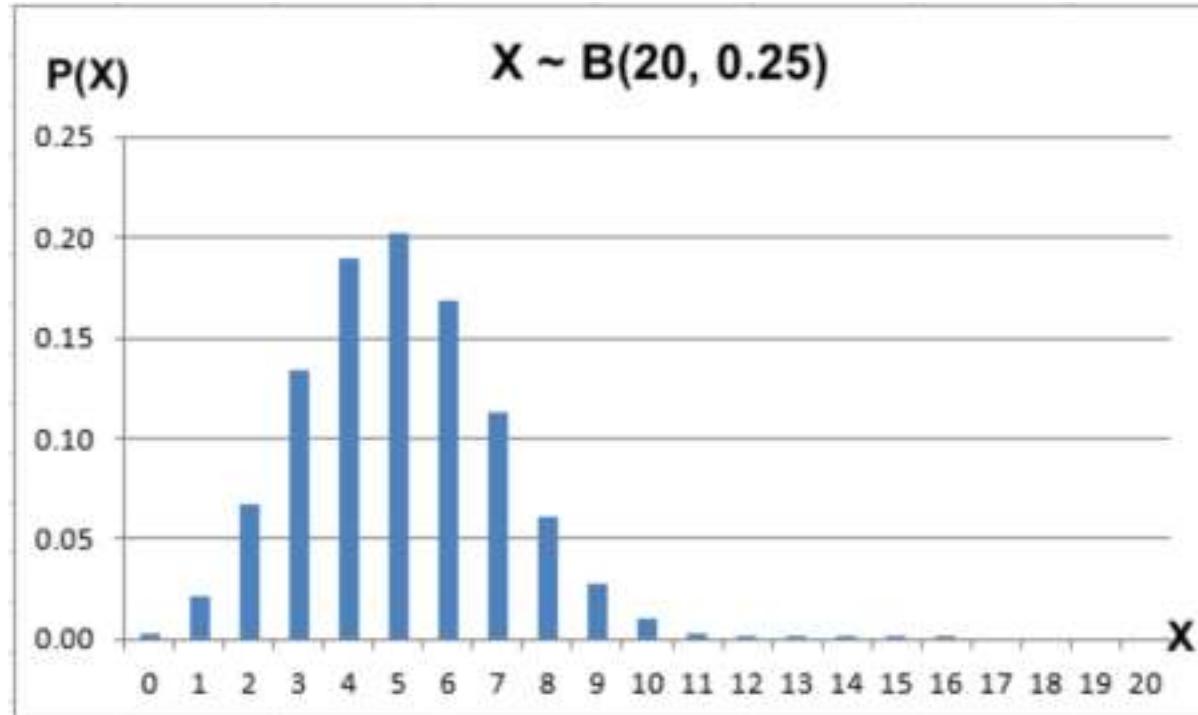
X	0	1	...	n
$P(X=x)$	${}^nC_0 p^0 (1-p)^n$	${}^nC_1 p^1 (1-p)^{n-1}$...	${}^nC_n p^n (1-p)^{n-n}$

Prob. Distribution of Discrete r.v.

Binomial Distribution – $X \sim B(n,p)$

$$P(X = x) = \binom{n}{x} \cdot p^x \cdot (1 - p)^{n-x}$$

Example:



Prob. Distribution of Discrete r.v.

Mean and variance of $X \sim B(n,p)$

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

The mean of X is given by:

$$E[X] = \sum_{x=0}^n x \binom{n}{x} p^x (1 - p)^{n-x} = np$$

And the variance is:

$$Var[X] = E[X^2] - E[X]^2 = np(1 - p) = npq$$

where $q = (1 - p)$

Prob. Distribution of Discrete r.v.

Example:

A football team plays 3 games. Assume each game is a Bernoulli trial with $\text{prob}(\text{win}) = 0.5$. Let the r.v. X be the number of wins. What is the probability that the team will win exactly 2 games?

$$\binom{3}{2} 0.5^2 (0.5)^1 = \frac{3!}{2!1!} \cdot 0.25 \cdot 0.5 \\ = 0.375$$

Solution

X has binomial distribution with $n = 3$ and $p = 0.5$, outcome Win(W) or Lose(L) on each trial. i.e.

$$X \sim B(3, 0.5)$$

Prob. Distribution of Discrete r.v.

Using the formula for Binomial probabilities:

$$P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$$
$$\approx \binom{3}{2} 0.5^2 (1-0.5)^{3-2}$$

$$\approx 0.375$$

There are 10 blue and 30 green balls in a bag. A ball is random picked, the colour is noted and it is put back into the bag. This process is repeated 10 times. Let X be the number of times a blue ball was picked, Calculate $E[X]$ and $\text{Var}[X]$.

Go to www.wooclap.com and use the code **CX2100**

There are 10 blue and 30 green balls in a bag. A ball is random picked, the...

1

$$P(x) = \frac{10}{40} = 0.25$$

$$E[X] = nP = 10 \cdot 0.25 = 2.5$$

The expected value of X is

The Variance of X is

$$\text{Var}[X] = nPq = 10 \cdot 0.25 \cdot 0.75 = 1.875$$

Prob. Distribution of Discrete r.v.

Example:

A quality control system selects a sample of 10 items from each batch of products for testing. If 2 or more of the items are defective the whole batch is rejected.

If the probability of an item being defective is 0.05, what is the probability of

- (i) having 2 defectives in the sample?
- (ii) the batch being rejected?

Prob. Distribution of Discrete r.v.

Soln: Let X be the number of defectives in the sample of $n = 10$ items. $X \sim B(10, 0.05)$

$$(i) P(X=2) = \binom{10}{2} \cdot 0.05^2 \cdot (0.95)^8 = 0.0746$$

$$\begin{aligned}(ii) P(X \geq 2) &= 1 - P(X=1) - P(X=0) \\&= 1 - \binom{10}{1} \cdot 0.05 \cdot 0.95^9 - \binom{10}{0} 0.05^0 0.95^{10} \\&= 1 - 0.3151 - 1 \cdot 1 \cdot 0.5987 \\&= 0.0861\end{aligned}$$

Prob. Distribution of Discrete r.v.

Poisson Distribution – Notation: $X \sim \text{Pois}(\mu)$

- Let r.v. X be the number of “successes” in a given time interval where X is a non-negative integer
- The r.v. X is a Poisson r.v. with probability distribution given by

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

where μ is the average number of successes and e is a constant = 2.71828...

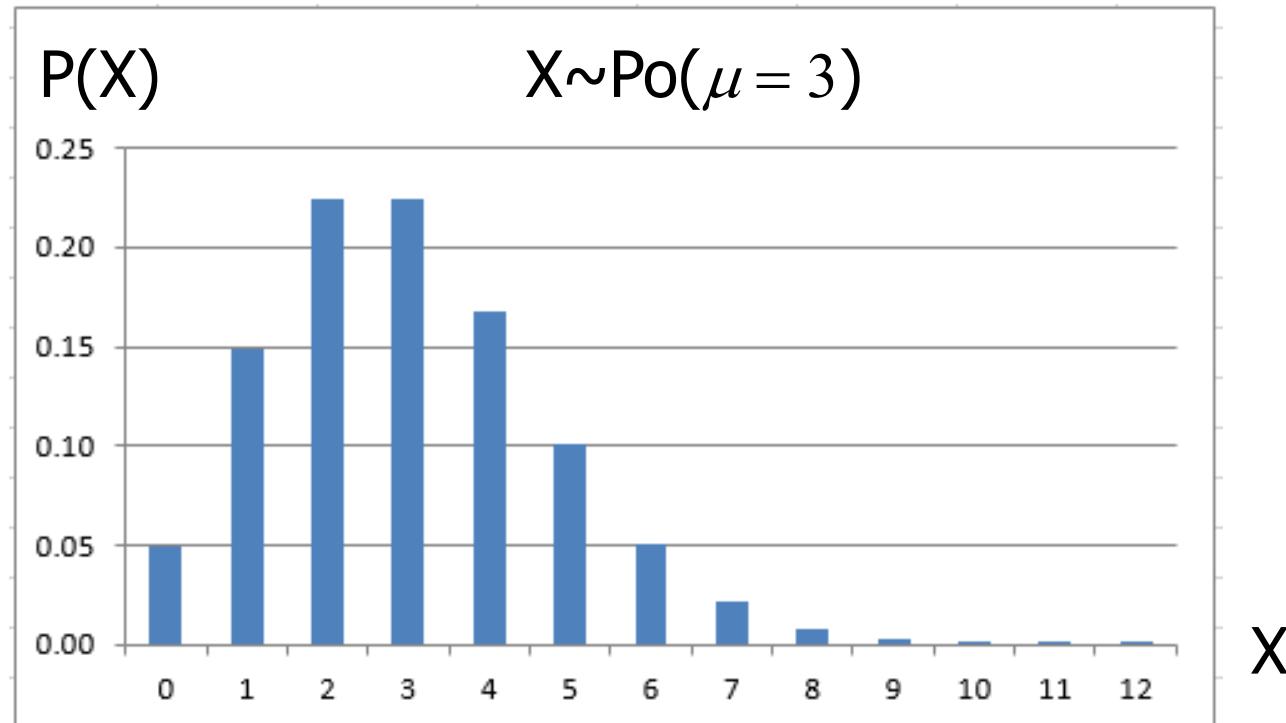
- The expected value and the variance of X are both equal to μ .

Prob. Distribution of Discrete r.v.

Poisson Distribution – $X \sim \text{Pois}(\mu)$

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

Example:



Prob. Distribution of Discrete r.v.

Eg:

1 day A salesman sells on the average 3 iphones per day. The number of iphones sold per day is a Poisson r.v.

- (i) Calculate the probability that in a given day he will sell some iphones.
- (ii) Given that there are 8 working hours per day, what is the probability that in a given hour he will sell one iphone?

Prob. Distribution of Discrete r.v.

Soln:

(i) $P(x=1) = 3$ $P(x=0)$
 $P(\text{iphone}) = 1 - P(\text{no iphone})$
 $P(x=0) = 1 - \frac{e^{-3} 3^0}{0!} = 1 - 0.05 = 0.95$

(ii) Avg iphone sold/hour = $3/8 = 0.375$

$$P(x=1) = \frac{e^{-0.375} (0.375)^1}{1!} = 0.258;$$

Eg: A shop selling a certain product of which the weekly demand is a Poisson variable with mean 3. On 1st day of each month, the stocks are replenished. Obtain the minimum number of the product in stock on 1st day of a month so that the shop is at least 50% sure of being able to meet the demands for the month.

$X \leftarrow$ no of product

$X \sim \text{Pois}(3 \times 4)$
monthly

Find largest X that $P(X \leq x) \geq 0.5$

choose $x = 12 + \text{AVG}(3 \times 4)$

$$P(X \leq 12) = 0.576$$

$$X = 12$$

$$\text{check } P(X \leq 11) = 0.4616 \quad X$$

Prob. Distribution of Discrete r.v.

Poisson Distribution – Approx to Binomial Dist

For Binomial Distribution $X \sim B(n, p)$ with large n :
i.e. $n \rightarrow \infty$, $p \rightarrow 0$ and $np \rightarrow$ a constant μ

Then it can be shown that :

$$\frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \rightarrow \frac{e^{-\mu} \mu^x}{x!}$$

$\therefore X \sim B(n,p) \approx \text{Pois}(\mu)$, where $\mu = np$, we have:

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

Poisson Distribution - Approximation to Binomial Distribution $X \sim B(n, p)$, $\mu = np$

For large n :

$$\begin{aligned} & \lim_{n \rightarrow \infty} \binom{n}{x} p^x (1-p)^{n-x} \quad \text{and } \mu = np \Rightarrow p = \frac{\mu}{n} \\ &= \lim_{n \rightarrow \infty} \frac{n!}{(n-x)! x!} \left(\frac{\mu}{n}\right)^x \left(1-\frac{\mu}{n}\right)^{n-x} \\ &= \lim_{n \rightarrow \infty} \underbrace{\frac{n!}{(n-x)!}}_{\substack{n(n-1)(n-2)\dots(n-x+1) \\ n \cdot n \cdot n \dots n \quad x \text{ terms}}} \underbrace{\frac{1}{n^x}}_{\substack{\frac{\mu^x}{x!}}} \underbrace{\left(1-\frac{\mu}{n}\right)^n}_{\substack{\lim_{n \rightarrow \infty} \left(1-\frac{\mu}{n}\right)^n \\ = e^{-\mu}}} \underbrace{\left(1-\frac{\mu}{n}\right)}_{\substack{\rightarrow 1 \text{ as } n \rightarrow \infty}} \rightarrow 1 \text{ as } n \rightarrow \infty \\ &= \frac{\mu^x}{x!} e^{-\mu} \quad (\text{if } X \sim \text{Pois}(\mu)) \\ & \quad \mu = np \end{aligned}$$

Prob. Distribution of Discrete r.v.

Poisson Distribution – Approx to Binomial Dist

Conditions for which the probability of $X \sim B(n, p)$ can be approximated by the probability $X \sim \text{Pois}(\mu)$:

- large n ($n \geq 100$)
- small p ($p \leq 0.01$)
- constant $\mu = np$ ($\mu \leq 20$)

Note that when $p \rightarrow 0$, then $(1 - p) \rightarrow 1$.

\therefore the expected value and the variance of $X \sim B(n, p)$ are approximately equal.

Intuitive Explanation:

To obtain the mean and variance of the Poisson distribution, consider the binomial distribution under the following conditions:

$$n \rightarrow \infty, p \rightarrow 0 \text{ and } np \rightarrow \text{a constant } \mu$$

The mean and variance of the binomial distribution are given by:

$$E[X] = np$$

$$Var[X] = npq$$

Since $q = 1 - p$, and as $p \rightarrow 0, q \rightarrow 1$ we have:

$$E[X] = \mu \text{ and } Var[X] = npq = \mu$$

Example: Suppose in a production line, 1 in 200 LED produced is defective. A random sample of 1000 LEDs are selected. What is the probability that at most 2 of them are defective?

let x = no of defective LED

$$P(1) = \frac{1}{200} = 0.005$$

$$NP = 1000 \cdot 0.005 = 5$$

$n > 100$, $P < 0.01$ $NP < 70$
use Poisson dist to approx. Binomial dist

$$\begin{aligned} P(x \leq 2) &= P(x=0) + P(x=1) + P(x=2) \\ &= e^{-5} + 5e^{-5} + \frac{5^2 e^{-5}}{2!} \\ &= 0.725, \end{aligned}$$

Prob. Distribution of Discrete r.v.

Geometric Distribution – Notation: $X \sim G(p)$

- Consider a sequence of independent Bernoulli trials, where each trial is an “experiment” with exactly 2 possible outcomes, "success" and "failure".
- Assume that the probability of success (S) is the same for all trials, $P(S) = p$, and the probability of failure $P(F) = 1 - p$.
- Each experiment is performed consecutively until the first success is obtained.
- Let $X = \text{no. of trials}$. Hence there are $X - 1$ failures before the 1st success is obtained.

Prob. Distribution of Discrete r.v.

Hence X is a random variable with a Geometric distribution $X \sim G(p)$:

$$P(X = x) = (1 - p)^{x-1} p$$

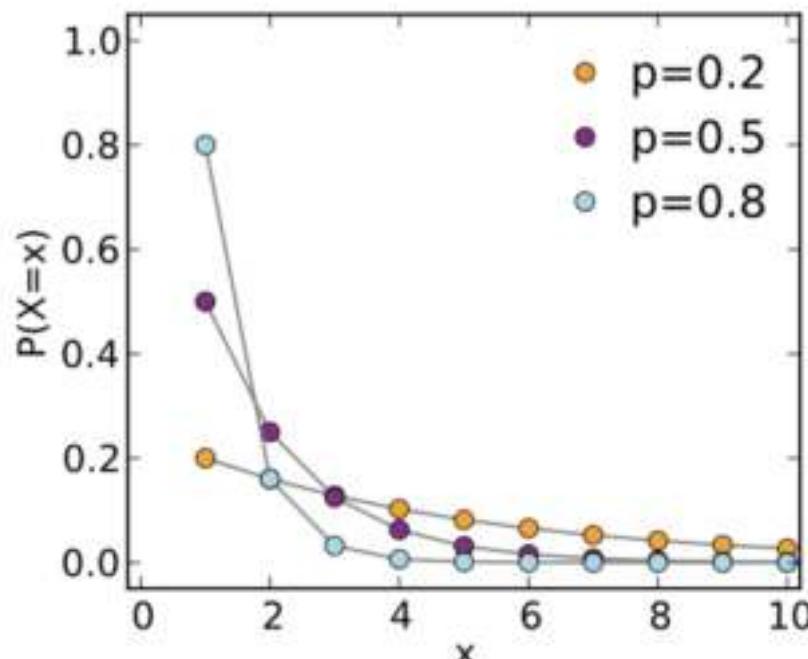
Expected value:

$$E[X] = \frac{1}{p}$$

Variance:

$$\text{Var}[X] = \frac{1-p}{p^2}$$

Examples:



Mean and Variance of Geometric Distribution $X \sim G(p)$

$$E[X] = \sum_{k=1}^{\infty} x \cdot (1-p)^{k-1} p$$

let $q = 1-p$, $\therefore E[X] = p \sum_{k=1}^{\infty} k q^{k-1}$ the term is $\frac{d}{dq} = k q^{k-1}$

($1^{\text{st}} \text{ term } k=0 \Rightarrow 0$)

$$= p \frac{d}{dq} \sum_{k=0}^{\infty} q^k$$

is infinite geometric series $\approx \frac{1}{1-q}$

$$= p \frac{d}{dq} (1-q)^{-1} = \frac{p}{(1-q)^2}$$

$$= \frac{1}{p} \quad \text{since } 1-q=p$$

$$\text{Var}[X] = E[X^2] - E[X]^2 = E[X(X-1)] + E[X] - E[X]^2$$

$$E[X(X-1)] = \sum_{k=1}^{\infty} x(x-1) p (1-p)^{x-1}$$

Similar to the above equation we have $\sum_{k=1}^{\infty} k(k-1) q^{k-1}$

$$\text{Var}[X] = E[X^2] - E[X]^2 = E[X(X-1)] + E[X] - E[X]^2$$

$$E[X(X-1)] = \sum_{x=1}^{\infty} x(x-1)p(1-p)^{x-1}$$

Similar to the above approach, using $\frac{d}{dq}(X-1)^q$,
we will get $\frac{2(1-p)}{p^2}$.

$$\therefore \text{Var}[X] = \frac{2(1-p)}{p^2} + \frac{1}{p} + \left(\frac{1}{p}\right)^2 = \frac{1-p}{p^2}$$

In a data transmission, the probability that a data frame is received in error is 0.1. If a data frame is

Example:

In a data transmission, the probability that a data frame is received in error is 0.1. If a data frame is received in error, a retransmission will be requested.

- What is probability that a data frame will be successfully received in no more than 3 transmissions?
- Determine the average number of transmission per frame.

$$a) P = 1 - 0.1 = 0.9$$

$$P(X \leq 3) = 0.1^{3-1} \cdot 0.9 = 0.01 \cdot 0.9 = 0.009$$

$$\begin{aligned} P(X \leq 3) &= P(X=1) + P(X=2) + P(X=3) \\ &= 1 \cdot 0.9 + 0.1 \cdot 0.9 + 0.009 \\ &\approx 0.999 \end{aligned}$$

$$b) E[X] = 1/p = 1/0.1 = 10$$

Probability Distributions for Continuous Variables

- Introduction
- Special Probability Distributions
 - Uniform Distribution
 - Exponential Distribution
 - Normal Distribution



Prob. Distribution of Continuous r.v.

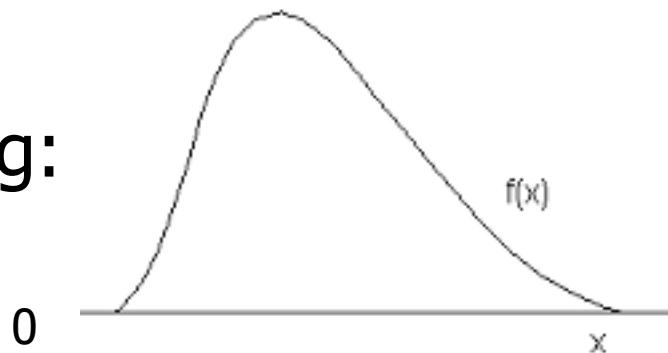
A continuous random variable is a variable used for continuous measurements such as distance, weight, and time.

Probability Density Function (pdf)

For a continuous random variable X , we can describe the probability distribution by some function $f(X)$, such that

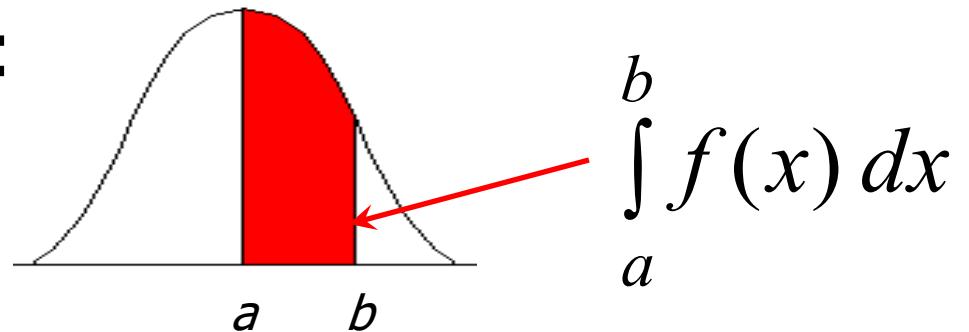
- (i) $f(X) \geq 0$ for all X

eg:



Prob. Distribution of Continuous r.v.

- (ii) $P(a < X < b)$ is the area under the curve between a and b :



- (iii) Total area under curve = 1.

The function $f(x)$ is called the probability density function (or pdf) of X .

The probability $P(X < x)$ is known as the Cumulative distribution function (or CDF) of X .

Woodlap Ex.

Given that the probability density function of X is:

$$f(x) = \frac{x^3}{4} \text{ where } 0 < x < c$$

Find the value of the constant c.

Given that the probability density function of X is:

$$f(x) = \frac{x^3}{4} \text{ where } 0 < x < c$$

Find the value of the constant c.

$$\int_0^c \frac{x^3}{4} dx = 1 \quad \text{Solve for } c$$

$$= \frac{x^4}{16} \Big|_0^c = 1 \quad c=2$$

Determine the CDF of X.

$$\begin{aligned}\text{CDF of } X &= \int_{-\infty}^x f(x) dx \\ &= \int_{-\infty}^x \frac{x^3}{4} dx = \frac{x^4}{16} \Big|_0^x = \frac{x^4}{16}\end{aligned}$$

Prob. Distribution of Continuous r.v.

For a continuous random variable X the probability that X equals to an exact value is zero.

Eg: for the length of a bolt

$$P(X = 1.999965700 \text{ cm}) = 0$$

i.e., $P(X=x) = 0$. Hence, for continuous random variables, probabilities are always associated with a range of values.

Eg: $P(X \leq 2.00 \text{ cm})$ or $P(X > 1.95 \text{ cm})$.

Prob. Distribution of Continuous r.v.

The expected value of a continuous r.v. is defined

$$\text{as: } E[X] = \int_{x=-\infty}^{x=\infty} x f(x) dx = \mu$$

and its variance is:

$$\text{var}[X] = \int_{x=-\infty}^{x=\infty} (x - \mu)^2 f(x) dx$$

$$= E[X^2] - E[X]^2$$

$$\text{where } E[X^2] = \int_{x=-\infty}^{x=\infty} x^2 f(x) dx$$

For discrete r.v.

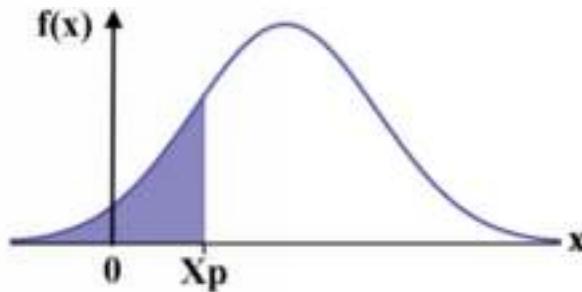
$$E(X) = \sum_{i=1}^n p_i x_i$$

For discrete r.v.

$$\text{Var}[X] = \sum_{i=1}^n (x_i - \mu)^2 p_i$$

Prob. Distribution of Continuous r.v.

How do we determine the p^{th} percentile from the probability distribution of a continuous r.v.?



The p^{th} percentile is the value X_p such that the probability that $x < X_p$ is p percent.

$$\int_{-\infty}^{X_p} f(x)dx = \frac{p}{100}$$

p^{th} percentile = X_p

Wooclap Ex

Given that the probability density function of x is

$$f(x) = e^{-x}, \text{ where } x > 0$$

determine the 25th percentile.

$$\begin{aligned} \int_0^{x_p} e^{-x} dx &= \frac{25}{100} \\ [-e^{-x}]_0^{x_p} &= 0.25 \\ (-e^{x_p}) - (-e^0) &= 0.25 \\ 1 - e^{-x_p} &= 0.25 \\ e^{-x_p} &= 0.75 \\ \ln e^{-x_p} &= \ln 0.75 \\ -x_p &= \frac{\ln 0.75}{\ln e} = \frac{\ln(0.75)}{1} \\ x_p &= 0.288 \end{aligned}$$

Prob. Distribution of Continuous r.v.

Uniform Distribution

One example of a pdf for a continuous r.v. is the uniform continuous distribution.

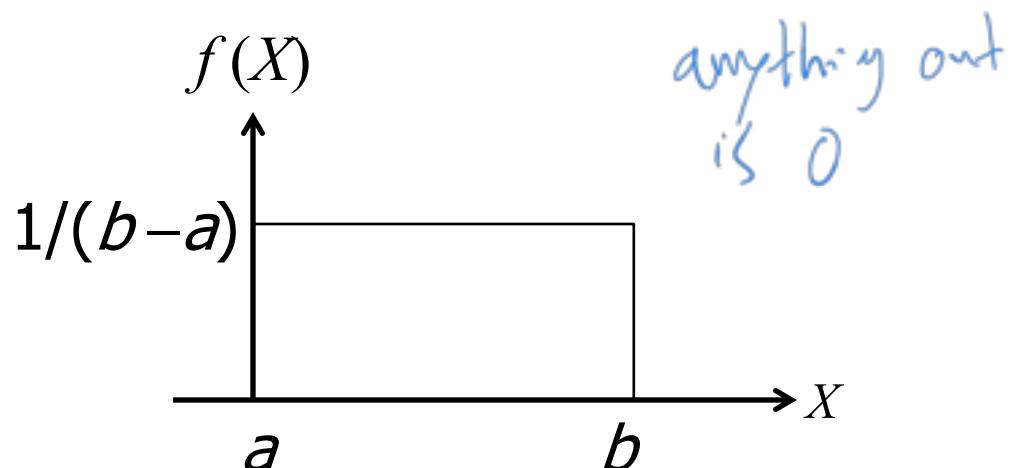
X can take any real value between a and b with uniform probability over this interval.

Notation: $X \sim U(a, b)$.

Total area = 1

Length = $b - a$

Height = $1/(b - a)$

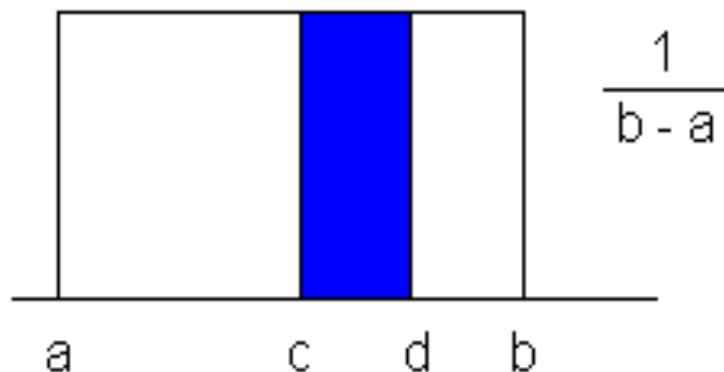


Prob. Distribution of Continuous r.v.

Thus the pdf is $f(X) = \frac{1}{b-a}$ for $a \leq X \leq b$
 $= 0$, otherwise

For any values c and d between a and b

$$P(c < X < d) = \frac{1}{b-a} (d - c)$$



Exercise: Obtain the CDF.

Obtain the expressions for $E[X]$ and $\text{Var}[X]$ of r.v. $X \sim U(a,b)$.

$$\therefore \int_a^b x \cdot f(x) dx \quad (\text{by definition})$$

$$E[X] = \int_a^b x \frac{1}{b-a} dx = \frac{1}{b-a} \left[\frac{x^2}{2} \right]_a^b = \frac{ab}{2}$$

$$E[X^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \frac{1}{b-a} \int_a^b x^2 dx = \frac{a^2 + ab + b^2}{3}$$

$$\begin{aligned} \text{Var}[X] &= E[X^2] - E[X]^2 = \frac{a^2 + ab + b^2}{3} - \frac{a^2 + ab + b^2}{4} \\ &= \frac{a^2 - 2ab + b^2}{12} = \frac{(a-b)^2}{12} \end{aligned}$$

Prob. Distribution of Continuous r.v.

Exponential Distribution

If the number of arrivals during an interval is Poisson distributed, then the interarrival times x are exponentially distributed. The pdf of x is:

$$f(x) = \lambda e^{-\lambda x}$$

where λ is the mean arrival rate and the mean interarrival time is $1/\lambda$.

Notation: $X \sim \text{Exp}(\lambda)$

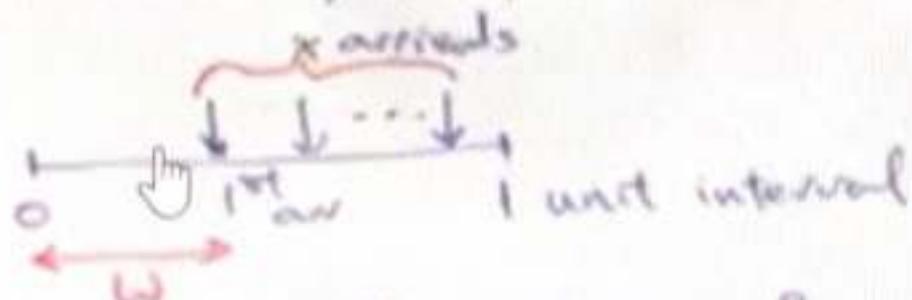
$$\mu = \frac{1}{\lambda}$$

If the number of arrivals in an interval is Poisson distributed, we show that the interarrival times are exponentially distributed.

If the number of arrivals in an interval is Poisson distributed, show that the interarrival times are exponentially distributed.

Let $X = \text{no. of arrivals in 1 unit interval}$

$$X \sim \text{Pois}(\mu), \text{ i.e. } P(X=x) = \frac{e^{-\mu} \mu^x}{x!}$$



w is the arrival time from $t=0$
we want to find
W = interarrival time and
the p.d.f. of W

CDF $P(W \leq w) = \frac{1 - P(W > w)}{1 - P(\text{no arrival in } [0, w])}$

$= e^{-\mu} \cdot P(X=0)$

$$\begin{aligned}
 \text{CDF } P(W \leq w) &= 1 - P(W > w) \\
 &= 1 - P(\text{No arrival in } [0, w]) \\
 &= 1 - \underbrace{P(X=0 \text{ in } [0, w])}_{\text{Poisson with mean } \lambda w} \\
 &= 1 - e^{-\lambda w}
 \end{aligned}$$

$$\therefore \text{pdf} = \frac{d}{dw} \text{CDF} = \boxed{\lambda e^{-\lambda w}}$$

exponential dist

where $\lambda = \text{mean arrival rate}$

OR $\frac{1}{\lambda}$ is average interarrival time

Prob. Distribution of Continuous r.v.

Exponential Distribution – examples:

Pdf:

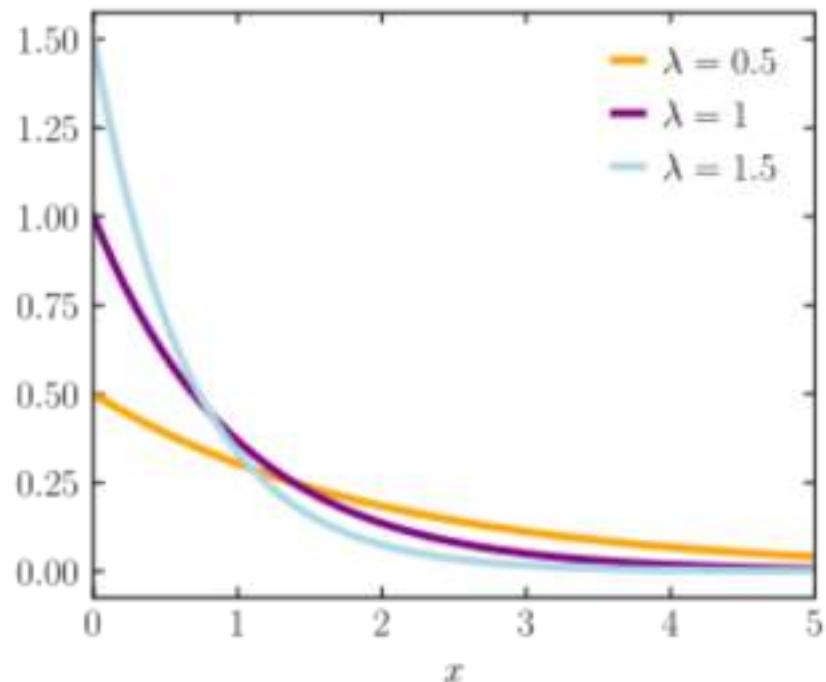
$$f(x) = \lambda e^{-\lambda x}$$

Expected value:

$$E[X] = \frac{1}{\lambda}$$

Variance:

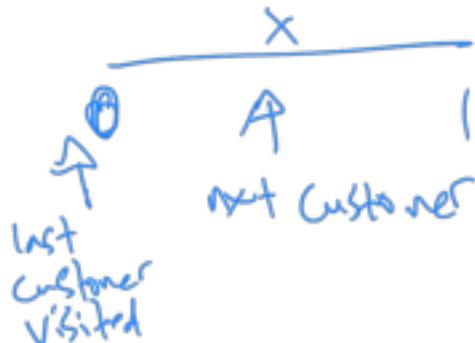
$$\text{Var}[X] = \frac{1}{\lambda^2}$$



Eg: Given that an average of number of customers walk into a shop per hour is 20 and the time between arrivals is exponentially distributed. After a customer arrives, what is the probability that the next customer will arrive within 1 minute?

$$\text{avg arrived rate} \Rightarrow = 20/\text{hr} \text{ or } \frac{20}{60} / \text{min}$$

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



$$x \sim \text{Exp}(\lambda)$$

$$P(x < 1) = \int_0^1 e^{-\lambda x} dx$$

$$= [-e^{-\lambda x}]_0^1$$

$$= 0.283$$

Prob. Distribution of Continuous r.v.

Normal Distribution

A common function for a continuous distribution is the normal (bell-shaped) curve. Its pdf depends on μ and σ , and is given by:

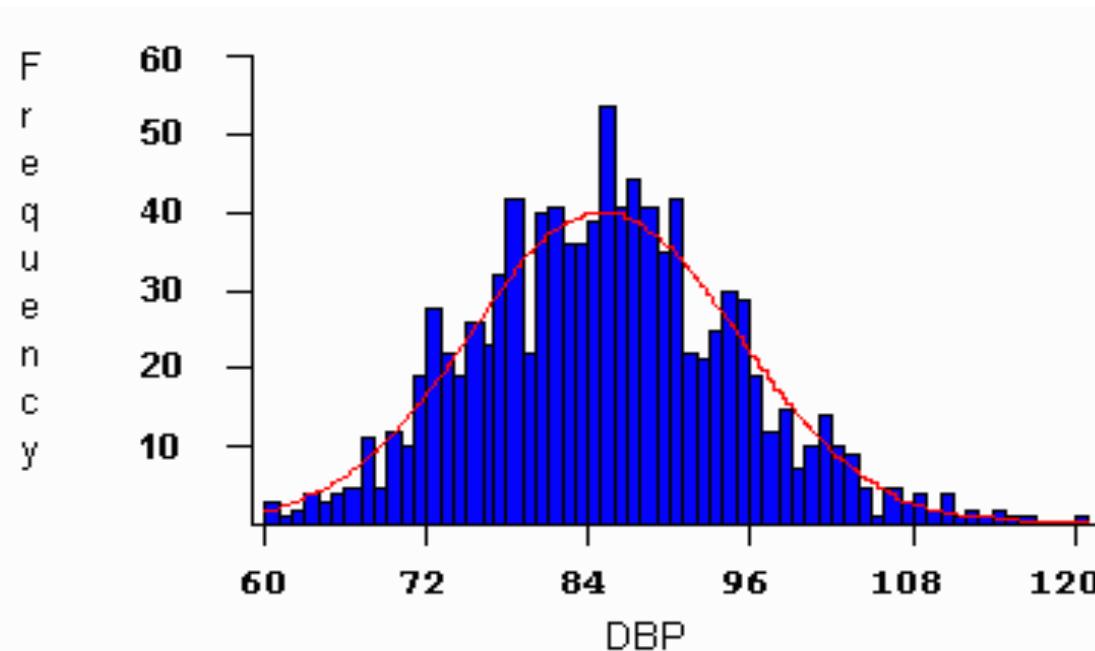
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \text{where } -\infty < x < \infty$$

do not remember this formula

The normal distribution is by convention written as $X \sim N(\mu, \sigma^2)$. Eg: if $X \sim N(10, 16)$ then this implies that $\mu = 10$ and $\sigma^2 = 16$, and (hence $\sigma = 4$).

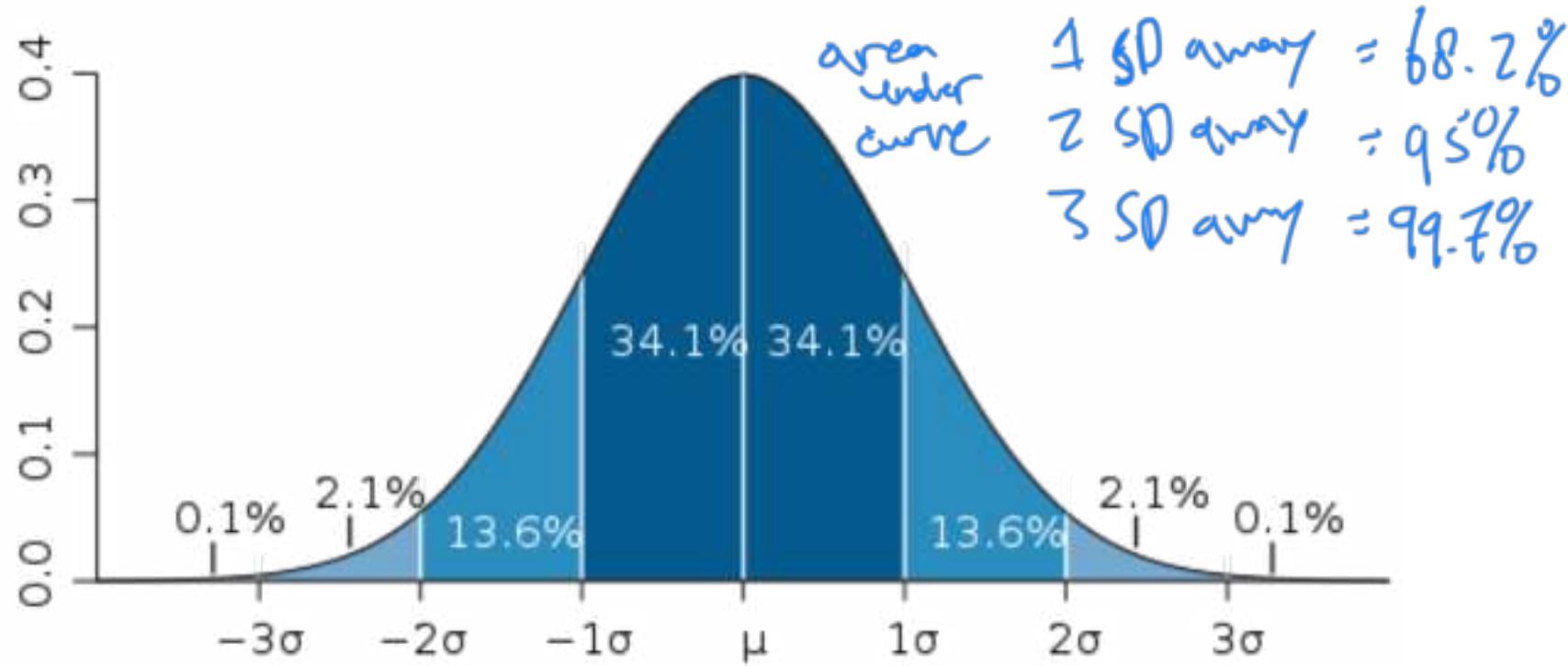
Prob. Distribution of Continuous r.v.

Eg: The distribution of blood pressure in patients can be approximated as a normal distribution with mean 85 mmHg and standard deviation 20 mmHg. A histogram of 1,000 observations and the normal curve is shown below.



Prob. Distribution of Continuous r.v.

The normal distribution has about 68% of the observations lying within one standard deviation of the mean, 95% within two standard deviations and 99.7% within 3 standard deviations.



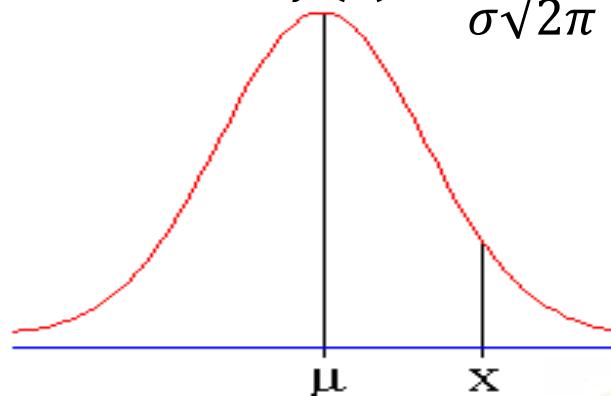
Prob. Distribution of Continuous r.v.

Standard Normal Distribution

To obtain the area under the density curve for a normal distribution, it is necessary to express any value of X in terms of the number of standard deviation units away from μ , i.e.

$$X = \mu + Z\sigma.$$

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



This can be rearranged to give $Z = \frac{X - \mu}{\sigma}$

Prob. Distribution of Continuous r.v.

Since X is a r.v., so is Z . Using the formulas for functions of random variables we can obtain

$$E[Z] = \frac{1}{\sigma} E[X] - \frac{1}{\sigma} \mu = \frac{1}{\sigma} \mu - \frac{1}{\sigma} \mu = 0$$

$$\text{var}[Z] = \frac{1}{\sigma^2} \text{var}[X] = \frac{\sigma^2}{\sigma^2} = 1$$

Also it can be shown that Z has a normal distribution. Thus $Z \sim N(0, 1)$.

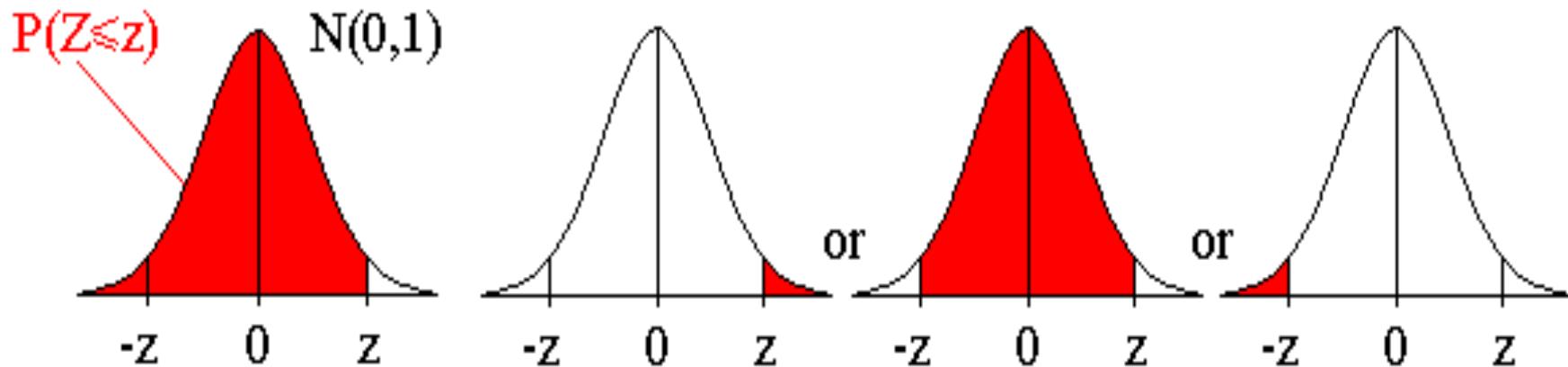
This is called the standard normal distribution.

↳ mean = 0
var = 1

Prob. Distribution of Continuous r.v.

Most statistics textbooks include a table of the standard Normal distribution. It tabulates the area under the curve of the $N(0,1)$ density function.

Other normal tables may give other areas, e.g.



$P(Z < 0.21) = 0.5832$ Table of Normal Distribution, $P(Z \leq z)$

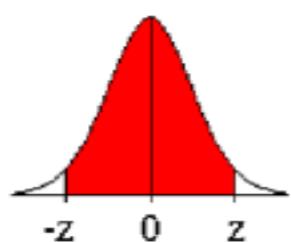
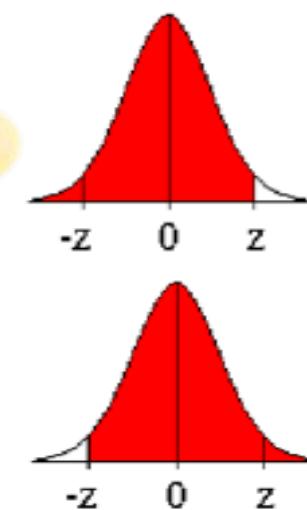
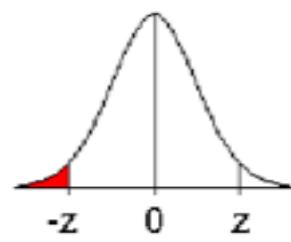
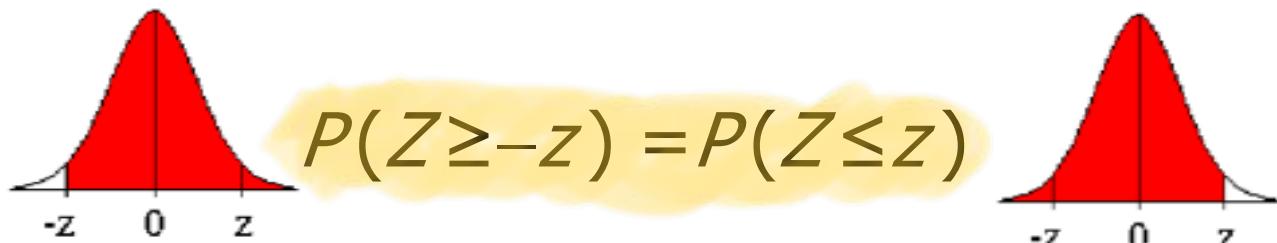
z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7703	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990

Prob. Distribution of Continuous r.v.

The $P(Z \leq z)$ table gives area under the curve from $-\infty$ to z , where $0 < z < 3.09$

For other values, we perform some manipulations.

Eg:



Prob. Distribution of Continuous r.v.

Probabilities for normal distributions other than the standard normal distribution $N(0, 1)$ are obtained by using the formula:

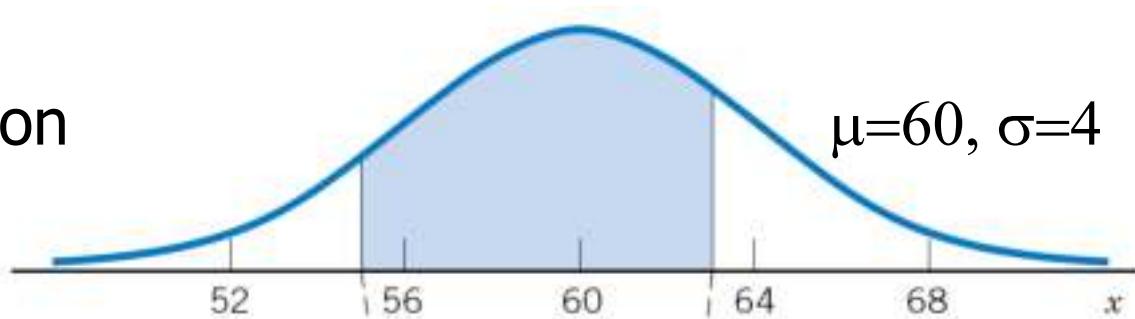
$$Z = \frac{X - \mu}{\sigma}$$

to convert from $X \sim N(\mu, \sigma^2)$ to $Z \sim N(0, 1)$ and then using the table of probabilities for the standard normal distribution $N(0, 1)$.

Prob. Distribution of Continuous r.v.

Normal distribution

$$X \sim N(\mu, \sigma^2)$$



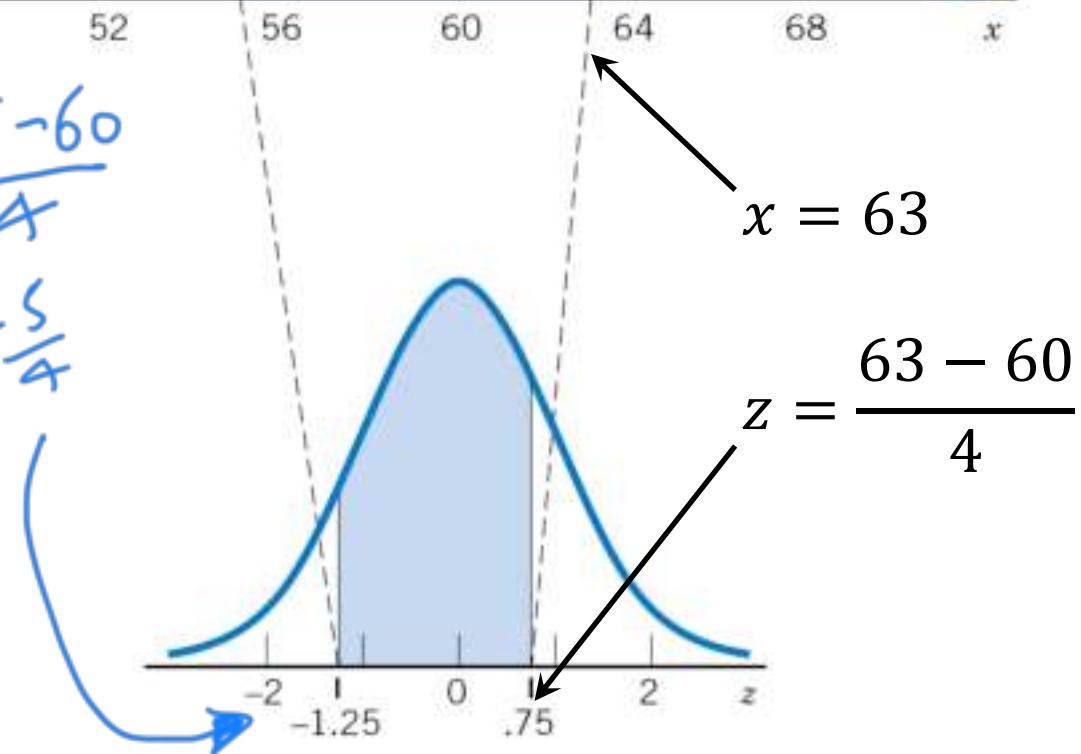
Standard Normal

$$Z \sim N(0, 1)$$

$$\text{where } Z = (X - \mu) / \sigma$$

$$z = \frac{55 - 60}{4}$$

$$= -\frac{5}{4}$$



Prob. Distribution of Continuous r.v.

Eg: A machine fills bottles with 300 ml of soft drinks. However, the actual quantity filled varies according to the normal distribution with $\mu = 298$ ml and $\sigma = 3$ ml.

What is the probability that an individual bottle contains less than 295 ml?

$$Z = \frac{295 - 298}{3} = -1.00$$
$$X \sim N(298, 3^2)$$
$$P(X < 295) = P(Z < -1)$$

See graph

$$= 1 - P(Z < 1)$$
$$\approx 1 - 0.8413$$
$$= 0.1587 \approx 1.6\%$$

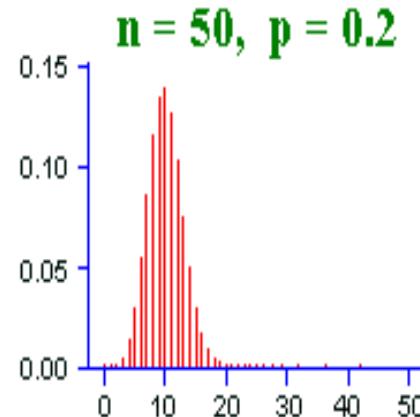
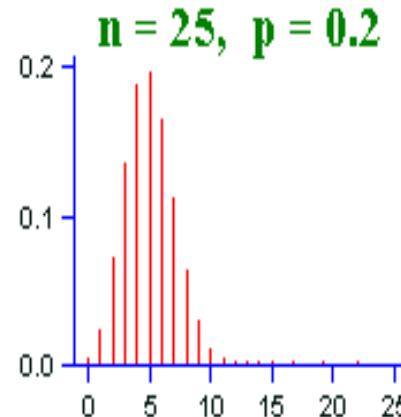
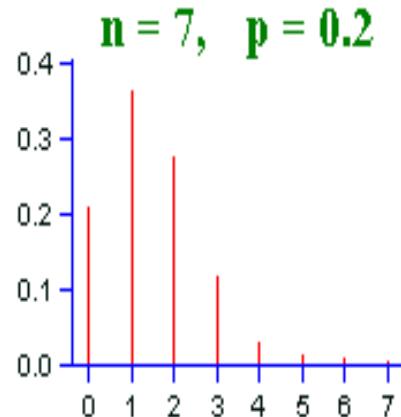
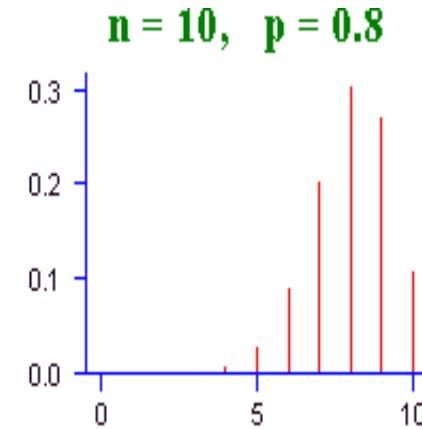
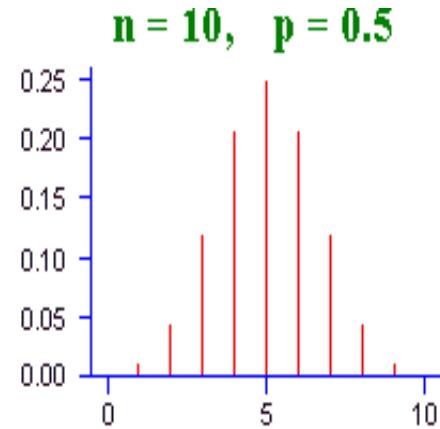
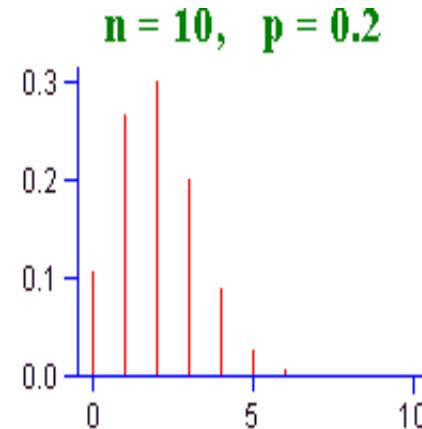
Prob. Distribution of Continuous r.v.

Solution

Let X denote the quantity in an individual bottle.
Given $X \sim N(298, 3^2)$, we want to find $P(X < 295)$

Approx. to Binomial Distribution

For large n or p near 0.5 the binomial distribution approximately follows the normal distribution.



Approx. to Binomial Distribution

If $X \sim B(n, p)$ where n is large and p not too near 0 or 1, then X can be approximated by a normal distribution with $E(X) = np$ and $\text{var}(X) = npq$, where $q = 1 - p$.

So $Z = \frac{X - np}{\sqrt{npq}}$ is approximately $N(0, 1)$.

This approximation is reasonably good when $np > 5$ and $n(1 - p) > 5$.

Approx. to Binomial Distribution

Continuity Correction and Accuracy

For more accurate values of binomial probabilities, the approximation is improved using the continuity correction. This method considers that each whole number occupies the interval from 0.5 below to 0.5 above it.

When an outcome X is to be included in the probability calculation, the normal approximation uses either $(X-0.5)$ or $(X+0.5)$.

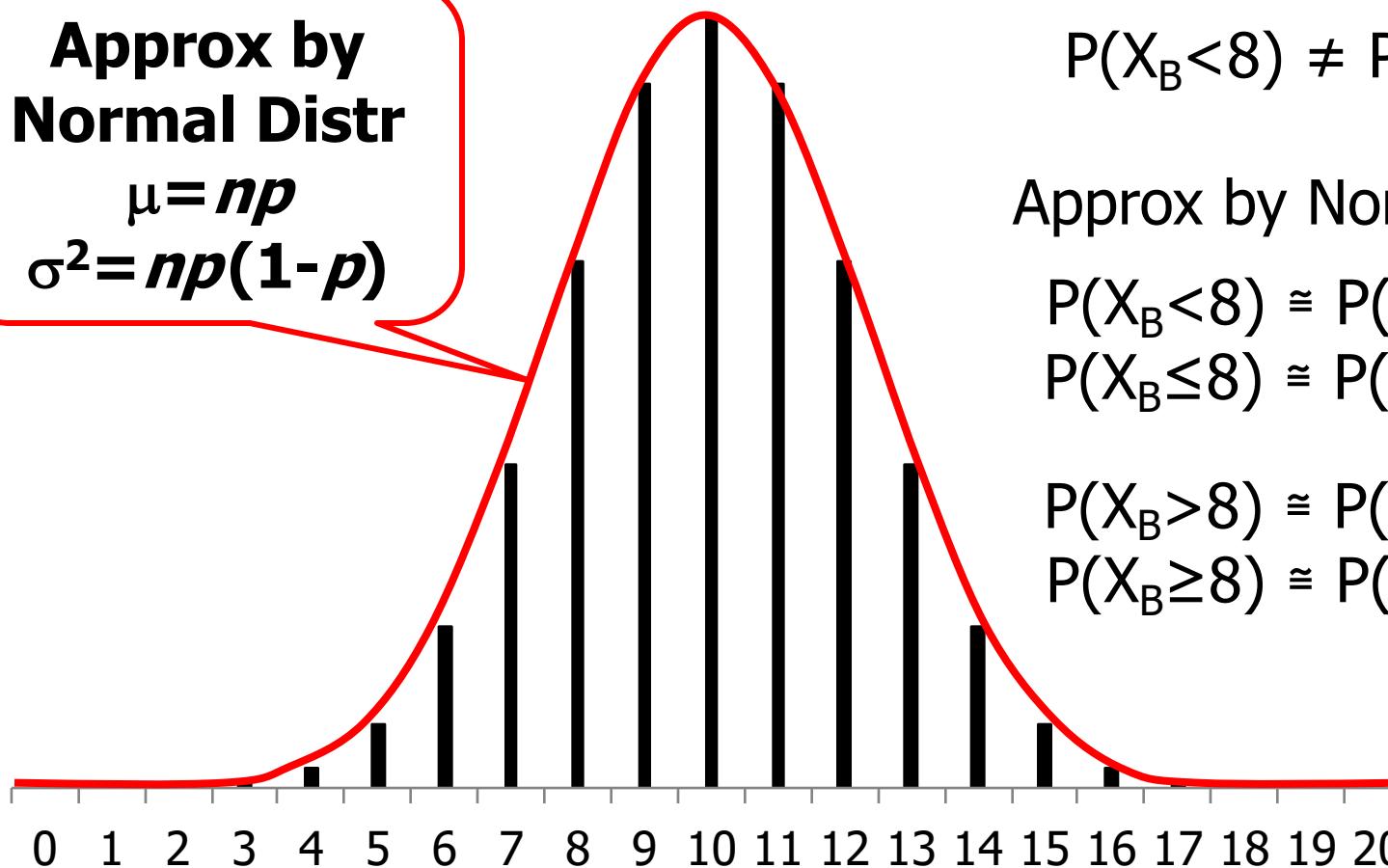
Approx. to Binomial Distribution

Binomial Distribution $X_B \sim B(n, p)$

**Approx by
Normal Distr**

$$\mu = np$$

$$\sigma^2 = np(1-p)$$



Binomial r.v.

$$P(X_B < 8) \neq P(X_B \leq 8)$$

Approx by Normal distr

$$P(X_B < 8) \approx P(X_N < 7.5)$$

$$P(X_B \leq 8) \approx P(X_N < 8.5)$$

$$P(X_B > 8) \approx P(X_N > 8.5)$$

$$P(X_B \geq 8) \approx P(X_N > 7.5)$$

Approx. to Binomial Distribution

Eg: In a particular faculty 60% of students are men and 40% are women. In a sample of 50 students what is the probability that more than half are women?

(let X = number of women in the sample n assume $X \sim B(50, 0.4)$)

$$E(X) = np = 50 \cdot 0.4 = 20$$
$$npq = 50 \cdot 0.4 \cdot 0.6 = 30$$
$$\text{Var}(X) = npq = 50 \cdot 0.4 \cdot 0.6 = 12$$

we approximate X by $X_N \sim N(20, 12)$

$$\begin{aligned} P(X > 25) &= P(Z > \frac{25-20}{\sqrt{12}} = 1.44) \\ &= 1 - P(Z \leq 1.44) \\ &= 1 - 0.9251 = 0.075 \end{aligned}$$

$$\begin{aligned} P(X > 25) &= P(X=26) + \dots + P(X=50) \\ &\rightarrow 0.0573 \end{aligned}$$

Approx. to Binomial Distribution

We need to find $P(X > 25)$. Note: not $P(X \geq 25)$

The exact ans. calculated from binomial probabilities
is:

Approx. to Binomial Distribution

The approximate probability, using the continuity correction, is

$$\begin{aligned} P(X > 25) &= P\left(Z > \frac{25.5 - 20}{\sqrt{12}}\right) \\ &= P(Z > 1.59) \\ &= 1 - P(Z \leq 1.59) \\ &\approx 0.0559 \end{aligned}$$

more accurate approximation than the original
approximation (without correction) of 0.075

Approx. to Binomial Distribution

The value 25.5 was chosen as the outcome 25 was not to be included but the outcomes 26, 27, ..., 50 were to be included in the calculation.

Similarly, if the example required the probability that less than 18 students were women, the continuity correction would require the calculation

$$P(X < 18) = P\left(Z < \frac{17.5 - 20}{\sqrt{12}}\right)$$

Summary

- Random Variables
- Expected Values and Variance
- Discrete Distributions
 - 1) Binomial
 - 2) Poisson
 - 3) Geometric
- Continuous Distributions
 - 1) Uniform
 - 2) Exponential
 - 3) Normal and Standard Normal
- Normal Approximation to Binomial

Alternate approach to obtain $\binom{n}{x}$

R1

1 1

R2

1 2 1

R3

1 3 3 1

.

1 4 6 4 1

.

1 5 10 10 5 1

.

1 6 15 20 15 6 1

1 7 21 35 35 21 7 1

1 8 28 56 70 56 28 8 1

⋮
⋮
⋮