

# Week 1

- o) Probability = How likely something is to happen  
= The number of events(A) that might happen  
Total number of event that might happen
- o) Statistics = Branch of mathematics dealing with the way of collecting, processing, analyzing, and presenting data conclusions.

Statistics

- Descriptive = Using data gathered on a group to describe or reach conclusions about that same group.
- Data collection, data presentation, characterize data
- Inferential = Gathers data from a sample and uses the statistics generated to reach conclusions about the population from which the sample was taken.

o) Population = All items of interest  
Characteristic-nya = Parameters

o) Sample = Subset of population  
Representative of the population  
Randomly selected

Sampling method

- Random sampling : S By random
- Stratified sampling : S By grouping data into several level
- Systematic sampling : P By particular order, S By random
- Cluster sampling : P By several parts /clusters, S By random

o) Nominal ⊂ Ordinal ⊂ Interval ⊂ Ratio

o) Datatypes

- Categorical (Quali)
- Numerical (Quant)

o) Measurement level

- Nominal
  - Categorical, no order
    - ↳ address, name
- Ordinal
  - Categorical, some order
    - ↳ dislike, neutral, like
- Interval
  - Numerical, no true zero
    - ↳ temperature
- Ratio
  - Numerical, has true zero
    - ↳ exam score

o) Data collection techniques

- Interview
- Questionnaire
- Observation
- Test and objective scale
- Projective method

## Week 2

- c) Quantitative data graphs
- o Histogram
  - o Frequency polygon
  - o Ogive
  - o Dot plot
  - o Stem-and-leaf plot
  - o Time series plot
  - o Scatter plot

### A. Histogram

- o Series of contiguous bars/ rectangles that represent the frequency of data in given class intervals.

### B. Frequency Polygon

- o Each class frequency is plotted as a dot at the class midpoint, and the dots are connected by a series of line segments.

### C. Ogive / Cumulative Frequency Polygon

- o Most useful when the decision maker wants to see running totals.

### D. Dot Plot

- o Useful for observing the overall shape of the distribution of data points along with identifying data values or intervals for which there are groupings and gaps in the data.

### E. Stem-and-leaf Plot

- o Constructed by separating the digits for each number of the data into 2 groups.

o) Left most digits = Stem = Higher valued digits

o) Rightmost digits = Leaf = Lower values

## F. Time Series Plot

o) Used to visualize data sets against time

## G. Scatter Plot

o) Type of plot / mathematical diagram using Cartesian coordinates to display values for typically 2 variables for a set of data.

o) Qualitative Data Graphs :

- o Pie charts
- o Bar charts
- o Pareto charts

## A. Pie Charts

o) A circle = 100% data, and slices of pie = Percentage breakdown of the sublevels

o) Shows the relative magnitudes of the parts to the whole.

## B. Bar Charts

o) Contains 2 / more categories along 1 axis and a series of bars, one for each category, along the other axis.

o) Length of bar = Magnitude of the measure for each category

## C. Pareto Charts

o) Produce vertical bar chart that display the most common types of defects, ranked in order of occurrence from left to right.

## Week 3

o) Random variables = Formed by assigning a numerical value to each outcome in the sample space of a particular experiment

o) Continuous random variable = Take any value within a continuous intervals, take uncountably many values

o) Discrete random variable = Used to express variables in discrete values finite values

↳ e.g. o Byk siswa dpt A

- o Byk mobil msk parking lot from 8-9am
- o Byk siswa bawa laptop to class

o p.m.f = Probability mass function  
 = A set of probability values  $p_i$  assigned to each of the values  $x_i$  taken by the discrete random variable.

o Syarat p.m.f o  $P(X=x_i) = p_i$

- o  $0 \leq p_i \leq 1$
- o  $\sum_i p_i = 1$

o c.d.f = Cumulative distribution function  
 = Alternative way of specifying the probabilistic properties of a random variable  $X$

o Syarat c.d.f o  $F(x) = P(X \leq x) = \int_{-\infty}^x f(y) dy$

- o  $F(x) = \sum_{y:y \leq x} P(X=y)$
- o  $P(X=x) = F(x) - F(x^-)$
- o  $P(a \leq x \leq b) = P(x \leq b) - P(x \leq a) = F(b) - F(a)$
- o  $\lim_{x \rightarrow -\infty} F(x) = 0$  &  $\lim_{x \rightarrow \infty} F(x) = 1$
- o  $f(x) = \frac{dF(x)}{dx}$

o) Bottom step =  $F(x^-)$   
 Top step =  $F(x)$   
 Height of step =  $P(X=x)$

o) Continuous random variables = Express variables whose values  
 are continuous.  
 ↳ Intinya ada interval

o) The probabilistic properties of a continuous random variable  
 defined through probability density function

o) p.d.f = Probability density function

o) Syarat p.d.f :  
 $P(a \leq X \leq b) = \int_a^b f(x)dx$   
 •  $\int_{\text{state space}} f(x) dx = 1$   
 •  $f(x) \geq 0$

o) Expectation = Mean = For summary = Average value of random variable

o)  $E(x) = \sum_i p_i x_i$  = Expected value of a discrete random variable  
 with p.m.f  $P(X=x_i) = p_i$

o)  $E(x) = \int_{\text{state space}} x f(x) dx$  = Expected value of a continuous  
 random variable with p.d.f

o) If a continuous random variable  $X$  has a p.d.f  $f(x)$  that is  
 symmetric about a point  $\mu$ , so that  $f(\mu+x) = f(\mu-x)$ ,  $\forall x \in \mathbb{R}$ , then  
 $E(x) = \mu$ , so that  $E(x) = \mu$  Point of symmetry.

o) If a continuous random variable  $X$  has a p.d.f  $f(x)$  that is  
 symmetric about a point  $\mu$ , then both the  $Q_2 \wedge E(x) = \mu$

o) Median / 2<sup>nd</sup> quartile :  $F(x) = 0.5$   
 1<sup>st</sup> quartile :  $F(x) = 0.25$   
 3<sup>rd</sup> quartile :  $F(x) = 0.75$

$$\textcircled{1} \quad \text{Var}(X) = E((X - E(X))^2) = E(X^2) - (E(X))^2$$

\textcircled{2} Variance always positive

\textcircled{3} Larger values of variance = Greater spread in distribution of random variable about the  $\mu$  value

\textcircled{4} Variance = Spread / Variability in the values taken by random variable.

\textcircled{5} The  $p$ th quantile of a random variable  $X$  with a c.d.f  $F(x)$  is defined :  $F(x) = p$

↳ First quartile  $\Rightarrow p = 0.25 \rightarrow F(x) = 0.25$   
 ↳ Third quartile  $\Rightarrow p = 0.75 \rightarrow F(x) = 0.75$  } Find  $x$

\textcircled{6} A linear function of random variables  $X$  is another random variable  $Y = aX + b$  for some numbers  $a, b \in \text{Real}$ .

$$\begin{aligned} \text{↳ } E(Y) &= a E(X) + b \\ \text{↳ } \text{Var}(Y) &= a^2 \text{Var}(X) \end{aligned}$$

### \textcircled{7} Linear Combinations of Random Variables

- Consider now a sequence of random variables  $X_1, \dots, X_n$  together with some constants  $a_1, \dots, a_n$  and  $b$ , and define a new random variable  $Y$  to be the *linear combination*

$$Y = a_1 X_1 + \dots + a_n X_n + b$$

- The **expectation and variance of the linear combination** is given as

$$E(a_1 X_1 + \dots + a_n X_n + b) = a_1 E(X_1) + \dots + a_n E(X_n) + b$$

$$\text{Var}(a_1 X_1 + \dots + a_n X_n + b) = a_1^2 \text{Var}(X_1) + \dots + a_n^2 \text{Var}(X_n)$$

- Note : assumed that they are *independent* random variables.



## Week 4

### o) Binomial / Bernoulli Distributions

- o A Bernoulli random variable with parameter  $p$ ,  $0 \leq p \leq 1$  with  $P(X=0) = 1-p$  and  $P(X=1) = p$ , then  $E(X) = p$   
 $\text{Var}(X) = p(1-p)$
- o Consider an experiment consisting of  $n$  Bernoulli trials
  - o independent
  - o each have constant probability of success
- o , then total number of success  $X$  is a random variable that has a binomial distribution with parameter  $n$  &  $p$ , written  $X \sim B(n,p)$

- o p.m.f of  $B(n,p)$  random variable:

$$P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}, \text{ for } x = 0, 1, 2, \dots, n$$

$$E(X) = np$$

$$\text{Var}(X) = np(1-p)$$

### o) Multinomial Distribution

- o Mirip binomial distribution tapi bisa banyak outcome, tidak seperti binomial yang hanya 2 outcomes.
- o p.m.f of multinomial distribution:

$$P(X=x_1, \dots, X_k=k) = \frac{n!}{x_1! \dots x_k!} \cdot p_1^{x_1} \cdot \dots \cdot p_k^{x_k}, \text{ where } x_1 + x_2 + \dots + x_k = n$$

$$p_1 + p_2 + \dots + p_k = 1$$

$$E(x_i) = np_i$$

$$\text{Var}(X_i) = np_i(1-p_i)$$

## ▷ Poisson Distribution

- A random variable  $X$  distributed as Poisson random variable with parameter  $\lambda$ , written ◦

$$X \sim P(\lambda)$$

- p.m.f of Poisson distribution ◦

$$P(X=x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}, e = 2.71828\dots, \text{ for } x = 0, 1, 2, 3, \dots$$

- Useful to model the number of times that a certain event occurs per unit of time, distance, or volume, and has mean, var equal to parameter  $\lambda$ .

# Week 5

## o) Uniform Distributions

- Random variable  $X$  with a flat p.d.f between 2 points  $a$  and  $b$ , so that :

$$f(x) = \frac{1}{b-a}$$

,  $a \leq x \leq b$ , written  $X \sim U(a, b)$

- The c.d.f is :  $F(x) = \frac{x-a}{b-a}$
- $E(X) = \frac{a+b}{2}$
- $\text{Var}(X) = \frac{(b-a)^2}{12}$
- $p$ th quantile :  $(1-p)a + pb$
- interquartile range :  $(b-a)/2$

## o) Exponentials Distributions

- An exponential distribution with parameter  $\lambda > 0$  has p.d.f :

$$f(x) = \lambda \cdot e^{-\lambda x}, \text{ for } x \geq 0, \text{ and } f(x) = 0 \text{ for } x < 0$$

- The c.d.f is :  $F(x) = 1 - e^{-\lambda x}$

- Useful for modelling failure times & waiting times

$$E(X) = \frac{1}{\lambda}$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

## o) Gaussian / Normal Distribution

- The p.d.f is :  $f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$ , for  $-\infty \leq x \leq \infty$

$$E(X) = \mu$$

$$\text{Var}(X) = \sigma^2$$

- Bell-shaped curve, written  $X \sim N(\mu, \sigma^2)$

## o) Standard Normal Distributions

- A normal distribution with mean  $\mu = 0$  and  $\sigma^2 = 1$ . Its p.d.f has notation  $\Phi(x)$ , is  $\Phi$

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

- If  $X \sim N(\mu, \sigma^2)$ , then  $Z = \frac{X-\mu}{\sigma} \sim N(0,1)$

$$P(a \leq x \leq b) = \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right)$$

Mean = Median



- o) Correlation  $\Phi$ :  $r = r_{xy} =$

$$\frac{n \cdot \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \cdot \sum_{i=1}^n y_i}{\sqrt{n \cdot \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \cdot \sqrt{n \cdot \sum_{i=1}^n y_i^2 - (\sum_{i=1}^n y_i)^2}}$$

- o) Covariance  $\Phi$ :  $\text{Cov}(x,y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N}$

$$\text{o) Variance } \Phi: \sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N}$$

$$\text{o) Kurtosis } \Phi: \alpha_4 = \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{n \cdot s^4}$$

- o) Standard deviation  $\Phi: \sigma = \sqrt{\text{Varians'}}$

$$\text{o) Skewness } \Phi: \alpha = \frac{\bar{x} - \text{modus}}{\sigma}$$

$$\text{o) Ganjil } \left\{ \begin{array}{l} Q_1 = \frac{1}{4}(n+1) \\ Q_2 = \frac{1}{2}(n+1) \\ Q_3 = \frac{3}{4}(n+1) \end{array} \right.$$

$$= \frac{3(\bar{x} - \text{median})}{\sigma}$$

$$\text{Genap } \left\{ \begin{array}{l} Q_1 = \frac{1}{4}(n+2) \\ Q_2 = \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1}) \\ Q_3 = \frac{1}{4}(3n+2) \end{array} \right.$$

$$\text{o) Interquartile Range = Hamongan} \\ = H = IQR = X_{Q_3} - X_{Q_1}$$