

BT1101 Introduction to Business Analytics

Hands-on Tutorial 04

Probability & Statistics Concepts

Probability : ① Random Variable

Discrete
p.m.f

Continuous
p.d.f

② Distribution { mean, variance
parameters



Statistics

Descriptive Stats : ① descriptive statistics ② graphing

Mathematical Stats : ① parameter estimation

② hypothesis testing e.g. Shapiro test

* Random Variable X :

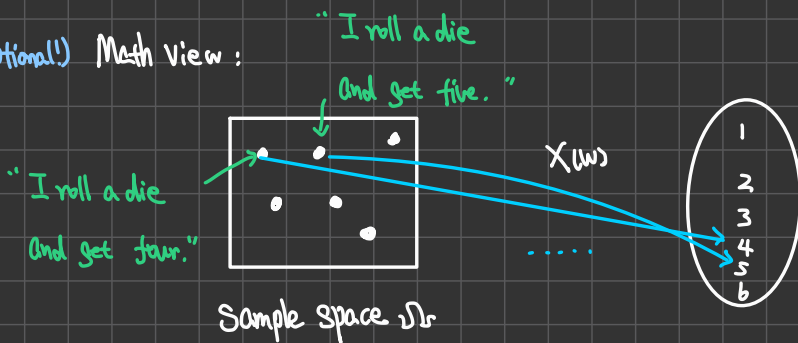
Intuition : X has some possible values . unknown before doing the experiment.

Each possible value has a chance (probability)

e.g. Rolling a die. Let X denote the number appears.

X has possible values : 1, 2, 3, 4, 5, 6

(optional!) Math View :



Random Variable is a function $X: \Omega \rightarrow \mathbb{R}$ (or \mathbb{N})

Sample space : the set of all possible outcomes

Discrete Random Variable & Continuous Random Variable

e.g. Die, Coin, books I have... e.g. height, Waiting time.

Probability Distribution : ① How to describe? ② How to calculate $P(X \in A)$?

Discrete R.V : Die Number

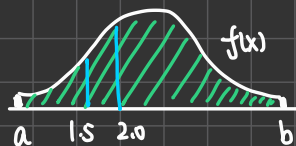
Die Number	1	2	3	4	5	6
we define the probability	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{2}{6}$

p.m.f : $P(X=w_i) = p_i$

$\sum_{all i} p_i = 1$ assign!

Continuous R.V : Height . Hard to assign the probability to each value.

p.d.f :



$$\int_a^b f(x) dx = 1. \quad f(x) \geq 0$$

$$P(X \in [1.5, 2.0]) = \int_{1.5}^{2.0} f(x) dx$$

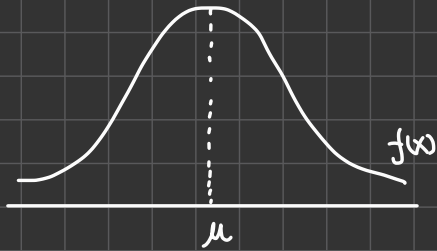
Rmk: Using p.m.f or p.d.f to describe the probability distribution of a random variable.

*important!

Special Distributions: $\left\{ \begin{array}{l} \text{Binomial, Poisson, Geometric} \dots \\ \text{Normal, Gamma} \dots \end{array} \right.$

Fixed the "framework" of the pmf/pdf but leave some parameters flexible!

e.g. Normal Distribution $X \sim N(\underline{\mu}, \underline{\sigma^2})$ $f(x) = \frac{1}{\sqrt{2\pi}\underline{\sigma}} \exp\left\{-\frac{1}{2\underline{\sigma}^2}(x-\underline{\mu})^2\right\}$



Question: Assume the heights of all Singaporean Men are normally distributed.

Want to explore the distribution!

⇒ The nature of this question is to find μ and σ^2 .

Mathematical Statistics:

Samples $\xrightarrow{\text{inference}}$ population X known: it has a normal distribution
 before observed: $(X_1, X_2, \dots, X_n) \leftarrow n \text{ random variables}$ $N(\underline{\mu}, \underline{\sigma^2})$ unknown: parameters μ, σ^2
 after observed: $(x_1, x_2, \dots, x_n) \leftarrow n \text{ fixed constants}$ Height of a person in the world.

Descriptive Statistics:

* NOTE: Sample Variance: $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ ↖ Sample mean population variance = $\frac{1}{N} \sum_{i=1}^N (X_i - \mu)^2$

Outlier Analysis:

Method 1 IQR:

Method 2 Z-Score:

Chebyshev's Thm: For any dataset, no matter what distribution, we can calculate the mean value and s.d. σ .

Then, the interval $[\text{mean} - z \cdot \sigma, \text{mean} + z \cdot \sigma]$ will include $(1 - \frac{1}{z^2}) \times 100\%$ data of this dataset.

e.g. we set $z=3$. Then $[\text{mean} - 3\sigma, \text{mean} + 3\sigma]$ will include around 89% data of this dataset.

(the rest 11% data we can regard as outliers)

Empirical Rule: For a dataset, normal distribution. Then:

$[\text{mean} - 2\sigma, \text{mean} + 2\sigma]$ will include around 95.44%

$[\text{mean} - 3\sigma, \text{mean} + 3\sigma]$ will include around 99.72%

Z-score of the samples of a random Variable:

$$\boxed{\text{Z-score Amount}} = \frac{\text{Amount} - E[\text{Amount}]}{\Delta(\text{Amount})}$$

Sample mean = 0

Sample sd = 1

Rmk: After transforming the sample data into z-score,

we can easily identify who is in the interval

$$[\text{mean} - 3\sigma, \text{mean} + 3\sigma]$$

$$= [0 - 3 \times 1, 0 + 3 \times 1] = [-3, 3]$$