



MULTIVARIATE STATISTICAL MODELING

Lecture-1

Applied → Refers to using base to make it useful

Normal Distribution :

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



Variable :

→ Takes different values

→ $x = 1D$

i	x_1	x_2	x_3	\vdots	x_n
1	x_{11}	x_{12}	x_{13}	\vdots	x_{1p}
2	x_{21}	x_{22}	x_{23}	\vdots	x_{2p}
n	x_{n1}	x_{n2}	x_{n3}	\vdots	x_{np}

- ① Fixed / Deterministic
- ② Random / Probabilistic

→ Different features

→

i	$x_1 = ID$	$x_2 = DD$	$x_3 = TB$	-	x_p
1	x_{11}	x_{12}	x_{13}	\vdots	x_{1p}
2	x_{21}	x_{22}	x_{23}	\vdots	x_{2p}
n	x_{n1}	x_{n2}	x_{n3}	\vdots	x_{np}

→ Observations on multiple variables

$$x_i^o = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \rightarrow \text{Each row can be taken out separately}$$

ith Observation on 'p' variables

When each observation has multiple values, we can it a multivariate

Linear combination of variables x_1, x_2, x_3 like $\beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

with empirically determined weights $\beta_1, \beta_2, \beta_3$ is known as variable

$p = 1$, Univariate

$p = 2$, Bivariate

$p \geq 2$, Multivariate

→ Field Based Observation

→ Expt

→ Naturalistic

Lecture 2

$$x_i^o = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \rightarrow \text{Check if features are correlated or Covariance}$$

Correlation Matrix : $\begin{bmatrix} & & \\ & & \\ p \times p & & \end{bmatrix}_{p \times p}$

$$\text{Mean Matrix } M^o = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix}_{p \times 1}$$

Behaviour should be analysed for design, improvement
 → Will lose substantial information if covariance/correlation is not considered

Random Variable :

① Discrete

② Continuous

Data Types :

1) Nominal :

→ Provide identity to some things or items

→ Examples : Month, dept name, brand name

→ Computational Limitation :
Mathematical operations can't be done

3) Interval :

→ Provides continuous data in a range

→ Example : Temperature

→ Computational Limitation :
• Division can't be done

2) Ordinal :

→ Provide order or rank to some item or things

→ Examples : Low, Medium, Good

→ Computational Limitations :
Mathematical operations can't be done

→ Better than nominal cuz we get a rank for organisation and comparison

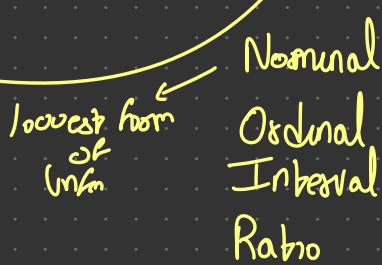
4) Ratio :

→ Contains absolute zero

→ Highest form of data

→ Examples : Hours, % etc

→ No computational limitations



Data Sources:

Primary: Collected from source

Secondary: Stored in repo or collected by someone else

Tertiary: Common knowledge through Wikipedia etc

Model:

- Mimic reality
- Explain regularity of phenomenon
- Can be extended to law or theory

Modelling:

- Process of building a model
- Physical, Mathematical, Statistical

$$\text{Data} = \text{Pattern} + \text{Error}$$

Regular/Systematic
Component

Principles of modelling

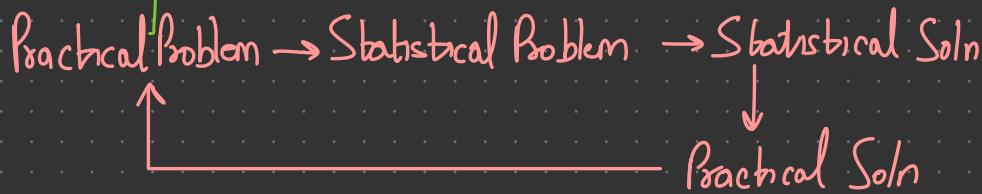
- Build simple if possible
- ^{Build} Modelling prob to fit technique
- Design phase must be rigorous
- Verify model before validation
- Don't take model liberally

- Prob specific
- ~~Don't oversell model~~
- Primary benefits associated with process of developing model
- Can't replace decision makers
- Limited by info provided

Lecture-3

Data = Pattern + Errors

Extract thus using statistical approaches to problem solving



Example :-

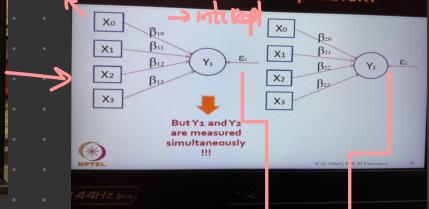
Threat to profit
↳
"Sales volume"



Response = Dependent Variables

cond.

Example: statistical problem



Purpose of multivariate modelling :-

Description → Explanation → Prediction

System details
Type of data
Problem faced

Dependence of relationships
of features

Predict future outcomes

2 types of Multivariate Models

- ① Dependence model
- ② Interdependence model

① Response Variables ← Explanatory Variable
Explains

② [All variables under one bracket]

Dependence Models:

- MANOVA
- Regressions
- Path Models
- Structural Equation Modeling
- Discriminant Analysis

Interdependence Models:

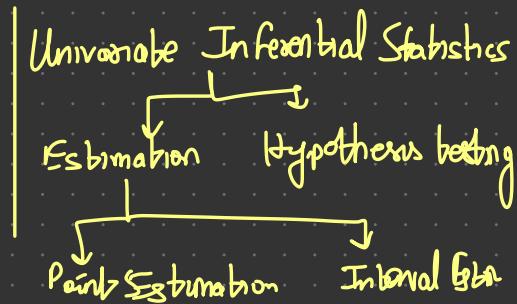
- Principal Component Analysis
- FACTOR ANALYSIS
- Cluster Analysis

MANOVA : Multivariate Analysis of Variance

Lecture 5

Univariate Descriptive Statistics

↳ Central tendency
↳ Dispersion



Population:

- Entity / Totality of the whole
- Example: Population etc
- Generally a large no of items

① Finite:

Size is known

Sampling with replacement

② Infinite:

Size is infinite

Sampling with replacement

Population is characterised by different variables applicable to the population

Population

Describe parameters: $[\mu, \sigma]$

$$\text{Mean} = \mu = E(x) = \sum_{\text{all } x} x f(x)$$

↓

Discrete RV

$$\text{Variance} = \sigma^2 = E(x - \mu)^2 = \sum_{\text{all } x} (x - \mu)^2 f(x)$$

↓

Probability Mass Function

μ, σ^2 are constants for a population (generally)

With population parameters, we get a mathematical interpretation

Continuous Population Parameters: $[\mu, \sigma]$

$$\text{Mean} = E(x) = \mu = \int_{-\infty}^{\infty} x f(x) dx$$

↓

Probability Density Function

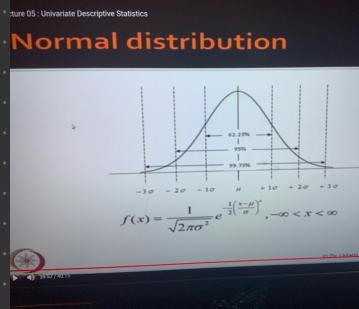
$$\text{Variance} = \sigma^2 = E(x - \mu)^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

Discrete

- Binomial
- Poisson
- Negative Binomial
- Geometric
- Hypergeometric

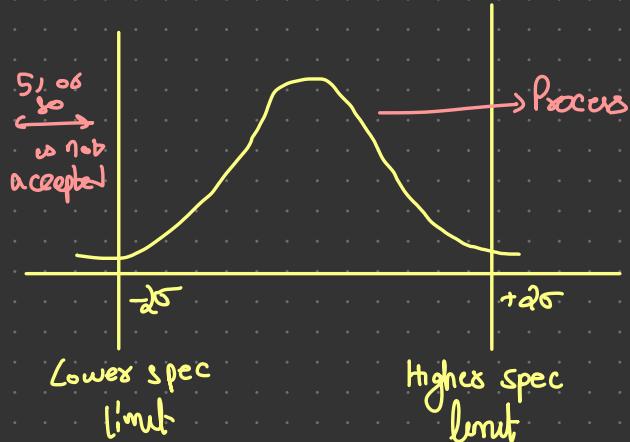
Continuous

- Normal
- Lognormal
- Exponential
- Weibull
- Gamma



μ

Lecture 6



Example :

Lecture 06 : Univariate Descriptive Statistics (Contd)

Example-1

The quality of service provided, measured on a 100 point scale, at three service centres A, B and C is normally distributed as $N(80, 9)$, $N(80, 16)$ and $N(90, 9)$, respectively. Comment on the performance of the three service centres.

Service Centers	Mean (μ)	Variability (σ^2)
A	80	9
B	80	16
C	90	9

C's performance is better than A and B.

144Hz 3ms

→ A is better than B

Higher mean and lesser variability

Variability is very important



Same but same as it can be shifted

Population parameters \rightarrow Prob distribution of variable of interest (x)

Before data collection

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix} \rightarrow \text{Random & Unknown}$$

After data collection

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \text{Sample of size } n \text{ which is known and fixed}$$

\rightarrow Most common

Mean, Median, Mode

Average

$$\left[\frac{N+1}{2} \right]^{\text{th}} \text{ position}$$

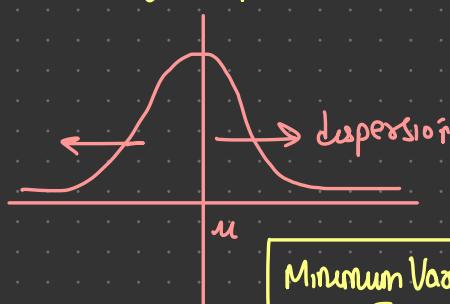
If ~~NP~~ decimal,
take average of floor
ceiling

Dot Plot.



No. of values helps to
find mode

Measure of dispersion



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

$$\rightarrow \text{Range} = \text{Max} - \text{Min}$$

$$\rightarrow \text{Interquartile Range} = Q^{\text{3rd}} - Q^{\text{1st}}$$

Minimum Variance Unbiased
Estimators

$$3 \left[\frac{N+1}{2p} \right]^{\text{th}} \quad \left[\frac{N+1}{4} \right]^{\text{th}}$$