

## **ASSIGNMENT-2**

**Task: Take any one domain and draw the graph (normal distribution using empirical formula).**

### **1. Selected Domain: Employee Salary Data**

#### **Domain Description**

Consider a real-time dataset of **monthly salaries of employees in a mid-sized IT company**.

Such salary data often follows a normal distribution because:

- Most employees earn around an average salary
- Very few employees earn extremely low or extremely high salaries

This type of data is commonly used in:

- Salary prediction models
- HR analytics
- Workforce planning systems

### **2. Assumptions for the Dataset**

To apply the empirical rule, assume:

- Mean salary ( $\mu$ ) = **₹50,000**
- Standard deviation ( $\sigma$ ) = **₹10,000**

These values help illustrate how employee salaries are distributed.

### **3. Normal Distribution in Machine Learning**

A **normal distribution** is a symmetric, bell-shaped curve where:

- The centre represents the **mean ( $\mu$ )**
- The spread is controlled by **standard deviation ( $\sigma$ )**

#### **Importance in Machine Learning**

- Used in **feature normalization**
- Helps in **outlier detection**
- Applied in **probabilistic models**
- Used for **confidence interval estimation**

#### 4. Empirical Rule (68–95–97 Rule)

For a normally distributed dataset:

- **68%** of values lie within  $\mu \pm 1\sigma$
- **95%** of values lie within  $\mu \pm 2\sigma$
- **99.7%** of values lie within  $\mu \pm 3\sigma$

#### 5. Application of Empirical Rule to Salary Data

##### Calculations

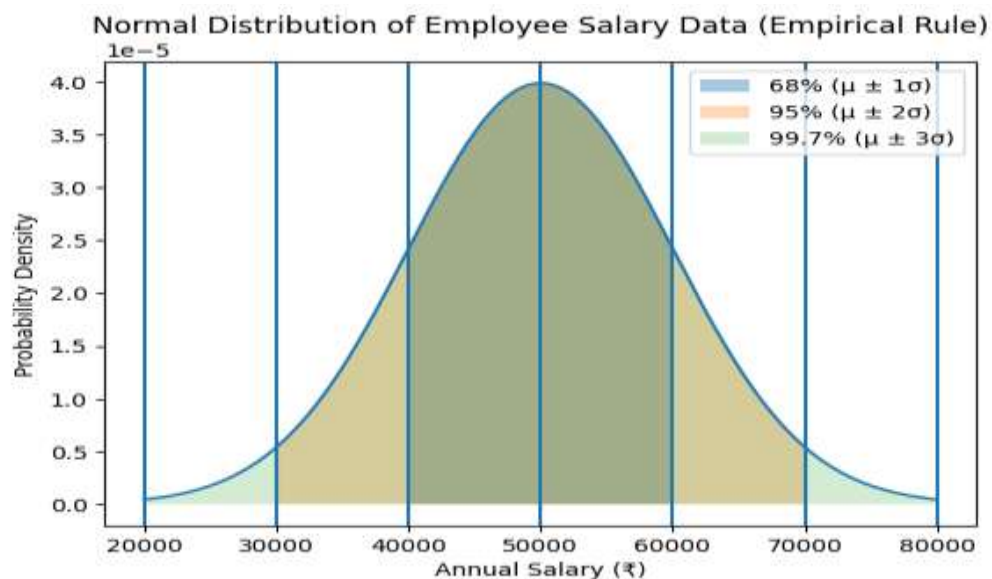
- $\mu \pm 1\sigma \rightarrow ₹50,000 \pm ₹10,000 \rightarrow ₹40,000 \text{ to } ₹60,000$
- $\mu \pm 2\sigma \rightarrow ₹50,000 \pm ₹20,000 \rightarrow ₹30,000 \text{ to } ₹70,000$
- $\mu \pm 3\sigma \rightarrow ₹50,000 \pm ₹30,000 \rightarrow ₹20,000 \text{ to } ₹80,000$

##### Interpretation

- About **68%** of employees earn between **₹40,000 and ₹60,000**
- About **95%** of employees earn between **₹30,000 and ₹70,000**
- About **99.7%** of employees earn between **₹20,000 and ₹80,000**

Employees outside this range can be treated as **outliers** in machine learning models.

#### 8. Normal Distribution Graph (Empirical Rule)



## 9. Significance in Machine Learning

Using the empirical rule in this domain helps to:

- Identify **salary outliers**
- Improve **data preprocessing**
- Design better **regression models**
- Avoid biased predictions due to extreme values