

## Distribution

1. Simulate 30 rolls with =RANDBETWEEN(1,6). What is the probability of rolling a 3 exactly 5 times? (Hint: Use BINOM.DIST)

**Answer:** To find the probability of rolling a 3 exactly 5 times in 30 rolls, where each roll has a probability

p=1, Use the binomial distribution:

6

=BINOM.DIST(5, 30, 1/6, FALSE)

### **Result:**

The probability is approximately:

0.1921 (above 19.21%)

This means there is about a 19% chance of getting exactly five 3s in 30 rolls.

**2. Generate 100 values in Excel using the continuous uniform distribution RAND() and plot a histogram. Describe the shape of the distribution.**

**Answer:** If we generate 100 values in Excel using the continuous uniform distribution with:

=RAND()

And then plot a histogram, here's what you should expect:

Shape of the Distribution

A histogram of 100 RAND() values will generally show:

- Roughly equal frequency across all bins
- A flat, rectangular shape (because the continuous uniform distribution gives every value between 0 and 1 the same probability)
- Some random variation—with only 100 samples, it won't be perfectly flat, but it should still look fairly even overall.

**In words:**

The distribution is approximately uniform—every interval between 0 and 1 is equally likely, so the histogram looks flat rather than peaked.

**3. A dataset has a mean of 50 and a standard deviation of 5. What percentage of values lie between 45 and 55 if the data follows a normal distribution?**

**Answer:** Since the data follows a normal distribution:

- Mean ( $\mu$ ) = 50
- Standard deviation ( $\sigma$ ) = 5
- The range 45 to 55 is  $\mu \pm 1\sigma$

Using the Empirical Rule (68–95–99.7 rule):

- About 68% of the data lies within 1 standard deviation of the mean.

**Result:**

68%

**4. What is the concept of standardization (z-score), and why is it important in data analysis? Explain the formula and how standardization transforms a dataset.**

**Answer: Standardization (Z-Score): Concept & Importance**

**Standardization**, also called z-score normalization, is a statistical technique used to convert raw data values into a common scale by expressing each value in terms of how many standard deviations it is away from the mean of the dataset.

### **Z-Score Formula**

$$z = \frac{x - \mu}{\sigma}$$

**Where:**

- $x$  = original data value
- $\mu$  = mean of the dataset
- $\sigma$  = standard deviation of the dataset
- $z$  = z-score (standardized value)

When a dataset is standardized:

- The **mean becomes 0**

- The standard deviation becomes 1
- The shape of the distribution remains the same
- Units are removed, making values dimensionless

## Example

If:

- Mean = 50
- Standard deviation = 5
- Value = 60

$$z = \frac{60 - 50}{5} = 2$$

This means **60 is 2 standard deviations above the mean.**

### 1. Comparison Across Different Scales

Standardization allows comparison between variables measured in different units (e.g., salary vs. age).

### 2. Identifying Outliers

Values with very high or low z-scores (e.g.,  $z>3$  or  $z<-3$ ) may be considered outliers.

### **3. Required for Many Algorithms**

Many machine learning algorithms perform better or require standardized data:

- K-means clustering
  
  
  
  
  
  
- Linear & logistic regression
  
  
  
  
  
  
- Principal Component Analysis (PCA)
  
  
  
  
  
  
- Support Vector Machines (SVM)

### **4. Probability & Normal Distribution Analysis**

Z-scores help calculate probabilities and percentiles using standard normal distribution tables.

### **5. Improves Model Performance**

Prevents features with large scales from dominating model outcomes.

## **Summary**

- Standardization converts data to a common scale

- Mean = 0, Standard Deviation = 1
- Helps in comparison, modeling, and outlier detection
- Essential for statistical and machine learning analysis

## 5. What is Kurtosis and their type?

**Answer:** **Kurtosis** is a statistical measure that describes the shape of a data distribution, specifically how peaked or flat the distribution is and how heavy or light the tails are compared to a normal distribution.

In simple terms, kurtosis tells us how extreme the values (outliers) in a dataset are.

### Types of Kurtosis

#### 1. Mesokurtic

- **Description:** Similar to a normal distribution
- **Kurtosis value:**  $\approx 3$  (or **Excess Kurtosis = 0**)
- **Characteristics:**

- Moderate peak
  - Moderate tails
- 
- **Example:** Normal distribution
- ## 2. Leptokurtic
- **Description:** More peaked than normal
  - **Kurtosis value:**  $> 3$  (Excess Kurtosis  $> 0$ )
  - **Characteristics:**
    - Sharp peak
    - **Heavy tails**
    - Higher chance of extreme values (outliers)

- **Example:** Stock market returns

### 3. Platykurtic

- **Description:** Flatter than normal

- **Kurtosis value:**  $< 3$  (Excess Kurtosis  $< 0$ )

- **Characteristics:**

- Flat peak
- **Light tails**
- Fewer extreme values

- **Example:** Uniform distribution

#### Kurtosis Formula

$$\text{Kurtosis} = \frac{\sum(x-\mu)^4}{n\sigma^4}$$

**Excess Kurtosis** is commonly used:

$\text{Excess Kurtosis} = \text{Kurtosis} - 3$

### **Why Kurtosis Is Important**

- Helps identify outliers
- Useful in risk analysis and finance
- Describes data distribution shape
- Assists in choosing appropriate statistical models

### **Quick Summary Table-**

Type	Shape	Tail Weight	Excess Kurtosis
Mesokurtic	Normal	Medium	0
Leptokurtic	Highly peaked	Heavy	Positive
Platykurtic	Flat	Light	Negative

## **6. Explain why the uniform distribution is a good model for the outcome of rolling a fair die.**

**Answer:** A uniform distribution is a good model for the outcome of rolling a fair die because each possible outcome has the same probability of occurring.

### **1. Equal likelihood of outcomes**

A fair die has six possible outcomes:

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

Since the die is fair, each number has an equal chance:

$$P(1) = P(2) = \dots = P(6) = \frac{1}{6}$$

This matches the key property of a uniform distribution, where all outcomes are equally likely.

### **2. Discrete uniform distribution**

Rolling a die is a **discrete random experiment** because outcomes are whole numbers.

The discrete uniform distribution assigns the same probability to each discrete value in a finite set, which fits the die outcomes perfectly.

### **3. No bias or preference**

A fair die is designed so that:

- All faces have equal size and weight
- No face is favoured over another

Because there is **no bias**, no outcome is more frequent than another over many rolls—another defining feature of a uniform distribution.

## 4. Long-run behaviour (Law of Large Numbers)

When a die is rolled many times:

- Each number appears approximately the same number of times
- Relative frequencies approach  $\frac{1}{6}$

This long-term behavior supports modeling the outcomes using a uniform distribution.

## Conclusion

The uniform distribution is an appropriate model for rolling a fair die because:

- All outcomes are **equally likely**
- The experiment is **discrete**
- There is **no bias**
- Long-term frequencies are equal

Hence, the outcomes of rolling a fair die follow a **discrete uniform distribution**.

## 7. Use Excel to compute the probability of getting at least 8 successes in 15 trials with success probability 0.5.

**Answer:** To compute this in Excel, we use the binomial distribution.

## Given

- Number of trials ( $n$ ) = 15
- Probability of success ( $p$ ) = 0.5
- We want:  $P(X \geq 8)$

## Concept

“At least 8 successes” means:

$$P(X \geq 8) = 1 - P(X \leq 7)$$

Excel can directly calculate cumulative binomial probabilities.

## Excel Formula (Recommended Method)

Step 1: Use the BINOM.DIST function

In any cell, enter:

$$=1 - \text{BINOM.DIST}(7, 15, 0.5, \text{TRUE})$$

### Explanation

- 7 → number of successes (up to 7)
- 15 → total trials
- 0.5 → probability of success
- TRUE → cumulative probability
- 1 - converts  $P(X \leq 7)$  to  $P(X \geq 8)$

### Result:

$$P(X \geq 8) \approx 0.5000$$

(Exactly 0.5 due to symmetry when  $p = 0.5$  and  $n = 15$ )

The probability of getting at least 8 successes in 15 trials is 0.5 (50%)

## **8. How does log transformation help in stabilising variance and making data more normally distributed?**

**Answer:** A log transformation is a common data-preprocessing technique used to handle skewed data, non-constant variance, and non-normal distributions. Here's how it helps, step by step.

### **1. Stabilising Variance (Reducing Heteroscedasticity)**

#### **Problem:**

In many real-world datasets (e.g., income, sales, population), the **variance increases with the mean**.

- Small values vary a little
- Large values vary a lot

This violates the assumption of **constant variance (homoscedasticity)** required by many statistical methods (e.g., regression, ANOVA).

#### **How log helps:**

The log function **compresses large values more than small values**.

#### **Example:**

- Raw scale:  
 $10 \rightarrow 100 \rightarrow 1000$
- Log scale:  
 $\log(10)=1, \log(100)=2, \log(1000)=3$

→ Differences at high values shrink.

#### **Result:**

- Large observations are pulled closer together
- Spread becomes more uniform across the range
- Variance becomes approximately constant

## 2. Making Data More Normally Distributed

### Problem:

Many datasets are **right-skewed** (long tail to the right):

- Income
- Sales
- Time-to-complete tasks
- Biological measurements

Such skewness violates the **normality assumption** used in many statistical tests.

### How log helps:

The log transformation:

- Pulls in extreme high values
- Reduces right skewness
- Makes the distribution more symmetric

### Visual intuition:

- Right-skewed distribution → log transform → bell-shaped distribution

→ The data often becomes **closer to a normal distribution**.

---

## 3. Converting Multiplicative Relationships into Additive Ones

### Problem:

Some variables grow **multiplicatively**, not additively.

### Example:

- Revenue = Price × Quantity
- Population growth
- Compound interest

## How log helps:

Taking logs converts multiplication into addition:

$$\log(ab) = \log(a) + \log(b)$$

→ Linear models work better after log transformation.

## 4. Handling Outliers

- Large outliers have **less influence** after log transformation
- This makes statistical estimates more stable and robust

## 5. When to Use Log Transformation

Use log transformation when:

- Data is **positively skewed**
- Variance increases with mean
- Values span several orders of magnitude

Avoid when:

- Data contains **zero or negative values**  
(unless you use  $\log(x+1)$  or another variant)

## In short:

- Log transformation stabilises variance by compressing large values and makes data more normally distributed by reducing right skewness, which helps meet the assumptions of many statistical models and tests.
- If you want, I can also show a before-and-after example with a dataset or Excel steps, since you often work with practical data analysis.

