Cyanide and Happiness © Explosm.net

**Statistics** is a branch of mathematics that deals with the collection, analysis, interpretation, presentation, and organization of data. It involves methods for gathering, summarizing, and analyzing data to draw conclusions or make decisions in the presence of uncertainty.

Basic statistics concepts are the backbone of data analysis, enabling us to make sense of the vast amounts of data generated daily. It's like conversing with data, where statistics helps us ask the right questions and understand the stories data tries to tell.

From predicting future trends and making decisions based on data to testing hypotheses and measuring performance, statistics is the tool that powers the insights behind data-driven decisions. It's the bridge between raw data and actionable insights, making it an indispensable part of data science.

**Exploratory data analysis** helped you understand how to discover patterns in data using various techniques and approaches. As you've learnt, EDA is one of the most important parts of the data analysis process. It is also the part on which data analyst spend most of their time. However, sometimes, you may require a very large amount of data for your analysis which may need too much time and resources to acquire. In such situations, you are forced to work with a smaller sample of the data, instead of having the entire data to work with.

Also, understanding where our data comes from and how we gather it is essential before diving deep into the ocean of data. This is where populations, samples, and various sampling techniques come into play.

Imagine we want to know the average height of people in a city. It's practical to measure everyone, so we take a smaller group (sample) representing the larger population. The trick lies in how we select this sample. Techniques such as random, stratified, or cluster sampling ensure our sample is represented well, minimizing bias and making our findings more reliable.

By understanding populations and samples, we can confidently extend our insights from the sample to the whole population, making informed decisions without the need to survey everyone.

Situations like these arise all the time at big companies like Amazon. For example, say the Amazon QC department wants to know what proportion of the products in its warehouses are defective. Instead of going through all of its products (which would be a lot!), the Amazon QC team can just check a small sample of 1,000 products and then find, for this sample, the defect rate (i.e. the proportion of defective products). Then, based on this sample's defect rate, the team can **"infer"** what the defect rate is for all the products in the warehouses. This process of **"inferring"** insights from sample data is called **"Inferential Statistics".**

## Types of Data and Measurement Scales

**Quantitative Data:** This type of data is all about numbers. It's measurable and can be used for mathematical calculations. Quantitative data tells us "how much" or "how many," like the number of users visiting a website or the temperature in a city. It's straightforward and objective, providing a clear picture through numerical values.

**Qualitative Data:** Conversely, qualitative data deals with characteristics and descriptions. It's about "what type" or "which category." Think of it as the data that describes qualities or attributes, such as the color of a car or the genre of a book. This data is subjective, based on observations rather than measurements.

## Scales of Measurement

1. **Nominal Scale:** This is the simplest form of measurement used for categorizing data without a specific order. Examples include types of cuisine, blood groups, or nationality. It's about labeling without any quantitative value.

2. **Ordinal Scale:** Data can be ordered or ranked here, but the intervals between values aren't defined. Think of a satisfaction survey with options like satisfied, neutral, and unsatisfied. It tells us the order but not the distance between the rankings.

3. **Interval Scale:** Interval scales order data and quantify the difference between entries. However, there's no actual zero point. A good example is temperature in Celsius; the difference between 10°C and 20°C is the same as between 20°C and 30°C, but 0°C doesn't mean the absence of temperature.

4. **Ratio Scale:** The most informative scale has all the properties of an interval scale plus a meaningful zero point, allowing for an accurate comparison of magnitudes. Examples include weight, height, and income. Here, we can say something is twice as much as another.

# Descriptive Statistics

Descriptive statistics has two main types: central tendency and variability measures.

**Measures of Central Tendency:** These are like the data's center of gravity. They give us a single value typical or representative of our data set.

**Mean:** The average is calculated by adding up all the values and dividing by the number of values. It's like the overall rating of a restaurant based on all reviews. The mathematical formula for the average is given below:

$$\bar{x} = \frac{\sum_{i=1}^{n} x_i}{n}$$

where $x_i$ represents each value in the data set, and $n$ is the number of values.

**Median**: The middle value when the data is ordered from smallest to largest. If the number of observations is even, it's the average of the two middle numbers. It's used to find the middle point of a bridge.

If n is even, the median is the average of the two central numbers.

- For odd $n$: Median is the $\frac{(n+1)}{2}$th value.
- For even $n$: Median is the average of the $\frac{n}{2}$th and $\left(\frac{n}{2} + 1\right)$th values.

**Mode**: It is the most frequently occurring value in a data set. Think of it as the most popular dish at a restaurant.

**Measures of Variability:** While measures of central tendency bring us to the center, measures of variability tell us about the spread or dispersion.

**Range**: The difference between the highest and lowest values. It gives a basic idea of the spread.

$$\text{Range} = \text{Maximum value} - \text{Minimum value}$$

**Variance**: Measures how far each number in the set is from the mean and thus from every other number in the set. It's calculated as:

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

where $x_i$ is each individual value, $\bar{x}$ is the mean of the data set, and $n$ is the number of values in the sample. For a population, divide by $n$ instead of $n-1$.

**Standard Deviation**: The square root of the variance provides a measure of the average distance from the mean. It's like assessing the consistency of a baker's cake sizes. It is represented as :

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

It gives us an idea of how spread out the data points are from the mean. Like variance, for a population, you would divide by $n$ instead of $n-1$ in the formula.

# Probability Basics

**Random Variables:** Before performing any kind of statistical analysis on a problem, it is advisable to quantify its outcomes by using random variables. So, the random variable X basically converts outcomes of experiments to something measurable. A random variable is a variable whose possible values are outcomes of a random phenomenon. It assigns numerical values to the outcomes of a random experiment.

A **discrete variable** is a variable that can only take on a "countable" number of values. If you can count a set of items, then it's a discrete variable. An example of a discrete variable is the outcome of a dice. It can only have 1 of 6 different possible outcomes and is therefore discrete. A discrete random variable can have an infinite number of values. For example, the whole set of natural numbers (1,2,3,etc.) is countable and therefore discrete.

A **continuous variable** takes on an "uncountable" number of values. An example of a continuous variable is length. Length can be measured to an arbitrary degree and is therefore continuous.

**Definition of Probability:** Probability is a measure of the likelihood of an event occurring. It is a number between 0 and 1, where 0 indicates impossibility and 1 indicates certainty.

It's about the chance or likelihood of events happening. Understanding concepts in probability is essential for interpreting statistical results and making predictions.

**Sample Space:** The sample space, denoted by S, is the set of all possible outcomes of a random experiment. Each element of the sample space is called a sample point or an outcome.

**Event:** An event is any subset of the sample space. It consists of one or more sample points.

**Probability of an Event:** The probability of an event A, denoted by P(A), is the sum of the probabilities of the sample points in A. It can be calculated using the formula

$$P(A) = \frac{\text{Number of favorable outcomes}}{\text{Total number of possible outcomes}}$$

- **Independent and Dependent Events:**
  - **Independent Events:** One event's outcome does not affect another's outcome. Like flipping a coin, getting heads on one flip doesn't change the odds for the next flip.
  - **Dependent Events:** The outcome of one event affects the result of another. For example, if you draw a card from a deck and don't replace it, your chances of drawing another specific card change.

**Probability Rules:**

•**Complement Rule:** The probability of the complement of event A, denoted by P(A') or P(A⁻), is equal to 1−P(A).

•**Union Rule:** The probability of the union of two events $A$ and B, denoted by P(A∪B), is equal to the sum of the probabilities of $A$ and B minus the probability of their intersection:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

(for mutually exclusive events, P(A∩B)=0).

•**Intersection Rule:** The probability of the intersection of two events A and B, denoted by P(A∩B), is the probability that both events occur simultaneously.

$$P(A \cap B) = P(A) \cdot P(B)$$

**Conditional Probability:** Conditional probability measures the probability of an event B occurring given that another event A has already occurred. It is denoted by P(B|A), which reads as "the probability of B given A". The formula for conditional probability is:

$$P(B \mid A) = \frac{P(A \cap B)}{P(A)}$$

where:

- $P(B \mid A)$ is the conditional probability of event $B$ given $A$.

- $P(A \cap B)$ is the probability of the intersection of events $A$ and $B$.

- $P(A)$ is the probability of event $A$.

**Independence:** Two events A and B are independent if the occurrence of one event does not affect the occurrence of the other. Mathematically, two events are independent if $P(A \cap B) = P(A) \times P(B)$.

**Multiplication Rule for Independent Events:** If events A and B are independent, then the probability of both events occurring is the product of their individual probabilities: $P(A \cap B) = P(A) \times P(B)$.

**Bayes' Theorem:** Bayes' theorem is a fundamental theorem in probability theory that describes the probability of an event, based on prior knowledge of conditions that might be related to the event. It is expressed as:

$$P(A \mid B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- $P(A \mid B)$ is the posterior probability of $A$ given $B$, which we want to find.

- $P(B \mid A)$ is the conditional probability of $B$ given $A$, often easier to determine or estimate than $P(A \mid B)$.

- $P(A)$ is the prior probability of $A$, which is our initial belief about the probability of $A$ before considering $B$.

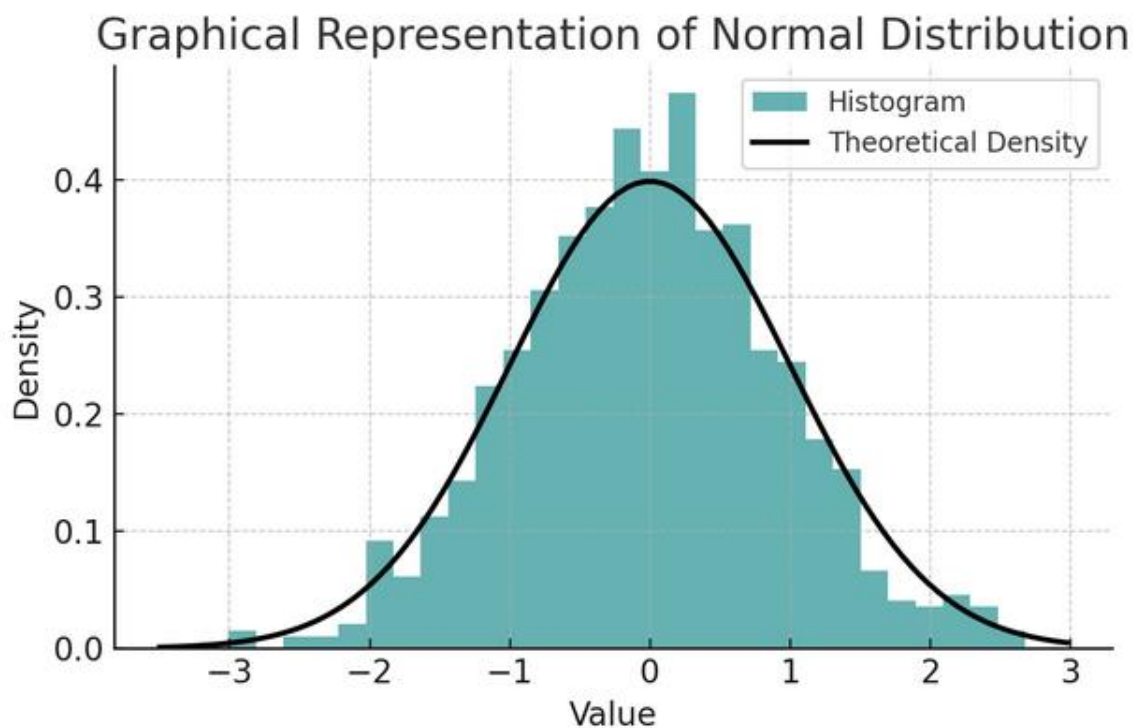- $P(B)$ is the total probability of $B$, also known as the marginal likelihood of $B$.

Probability provides the foundation for making inferences about data and is critical to understanding statistical significance and hypothesis testing.

In statistics,

- We represent a **distribution of discrete variables with PMF's (Probability Mass Functions) and CDF's (Cumulative Distribution Functions).**
- We represent distributions of **continuous variables with PDF's (Probability Density Functions) and CDF's.**
- The PMF defines the probability of all possible values x of the random variable.
- A PDF is the same but for continuous values.
- The CDF represents the probability that the random variable X will have an outcome less or equal to the value x.
- The name CDF is used for both discrete and continuous distributions.

# Common Probability Distributions

- **Normal Distribution:** Often called the bell curve because of its shape, this distribution is characterized by its mean and standard deviation. It is a common assumption in many statistical tests because many variables are naturally distributed this way in the real world.
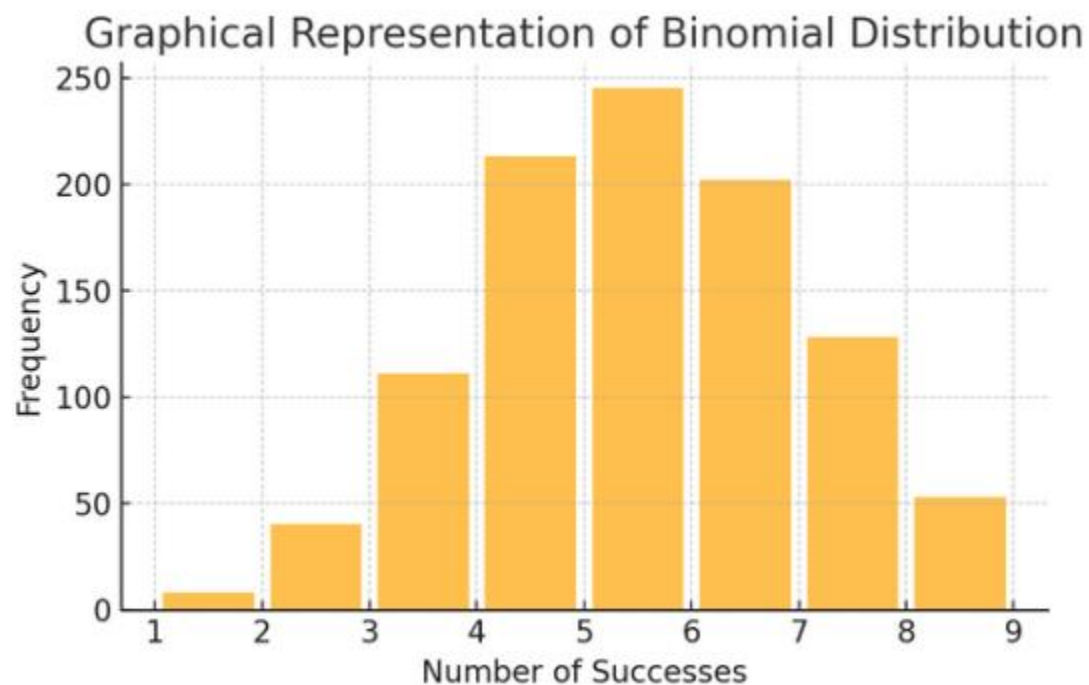


A set of rules known as the empirical rule or the 68-95-99.7 rule summarizes the characteristics of a normal distribution, which describes how data is spread around the mean.

**68-95-99.7 Rule (Empirical Rule)**

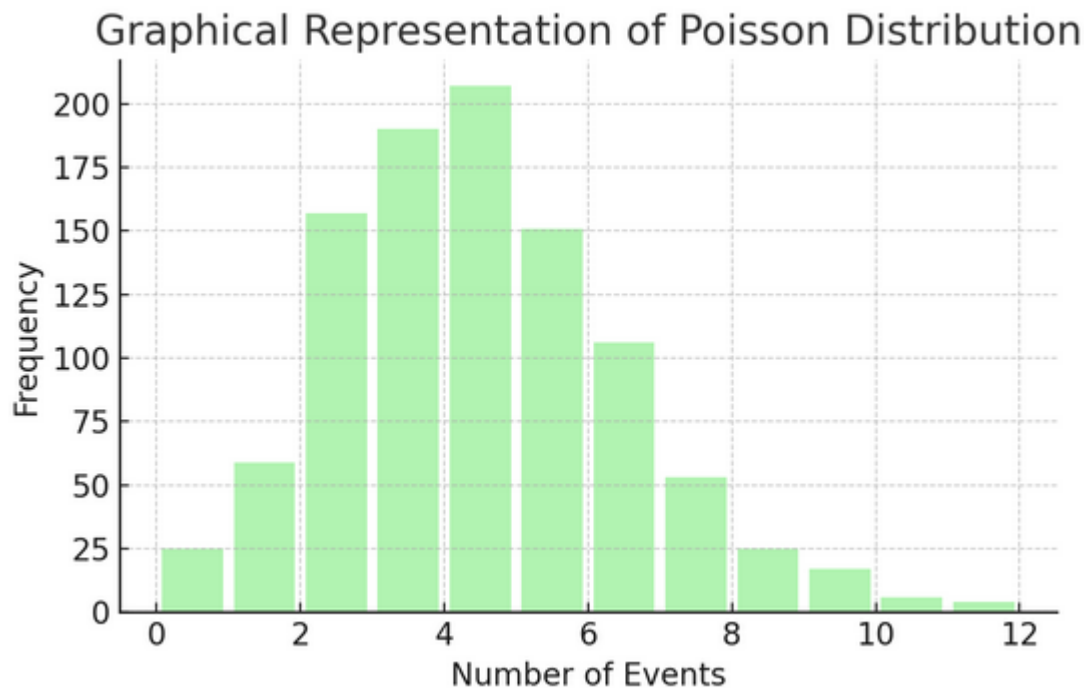This rule applies to a perfectly normal distribution and outlines the following:

- 68% of the data falls within one standard deviation ($\sigma$) of the mean ($\mu$).

- 95% of the data falls within two standard deviations of the mean.

- Approximately 99.7% of the data falls within three standard deviations of the mean.

**Binomial Distribution:** This distribution applies to situations with two outcomes (like success or failure) repeated several times. It helps model events like flipping a coin or taking a true/false test.

## Graphical Representation of Binomial Distribution

**Poisson Distribution** counts the number of times something happens over a specific interval or space. It's ideal for situations where events happen independently and constantly, like the daily emails you receive.

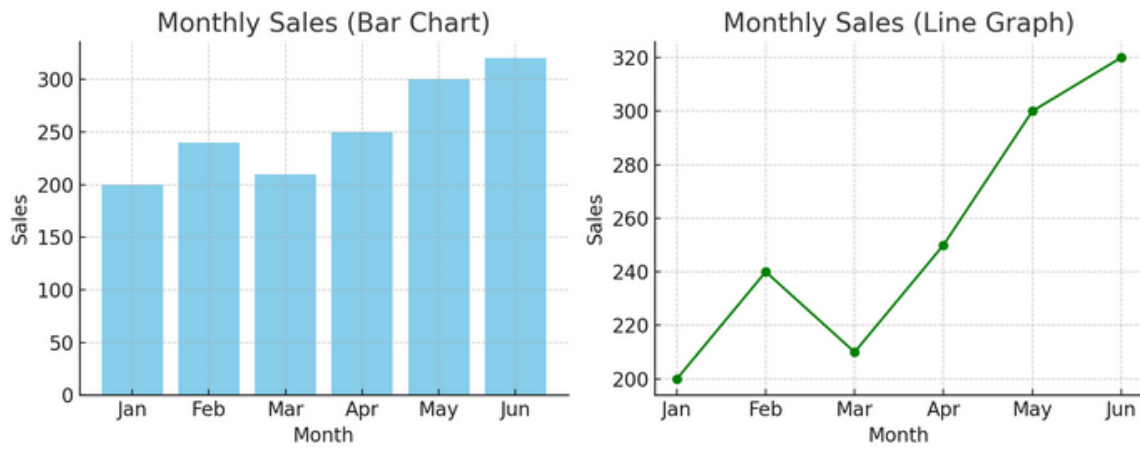## Graphical Representation of Poisson Distribution



Each distribution has its own set of formulas and characteristics, and choosing the right one depends on the nature of your data and what you're trying to find out. Understanding these distributions allows statisticians and data scientists to model real-world phenomena and predict future events accurately.

# Data visualization

Data visualization,is the art and science of telling stories with data. It turns complex results from our analysis into something tangible and understandable. It's crucial for exploratory data analysis, where the goal is to uncover patterns, correlations, and insights from data without yet making formal conclusions.
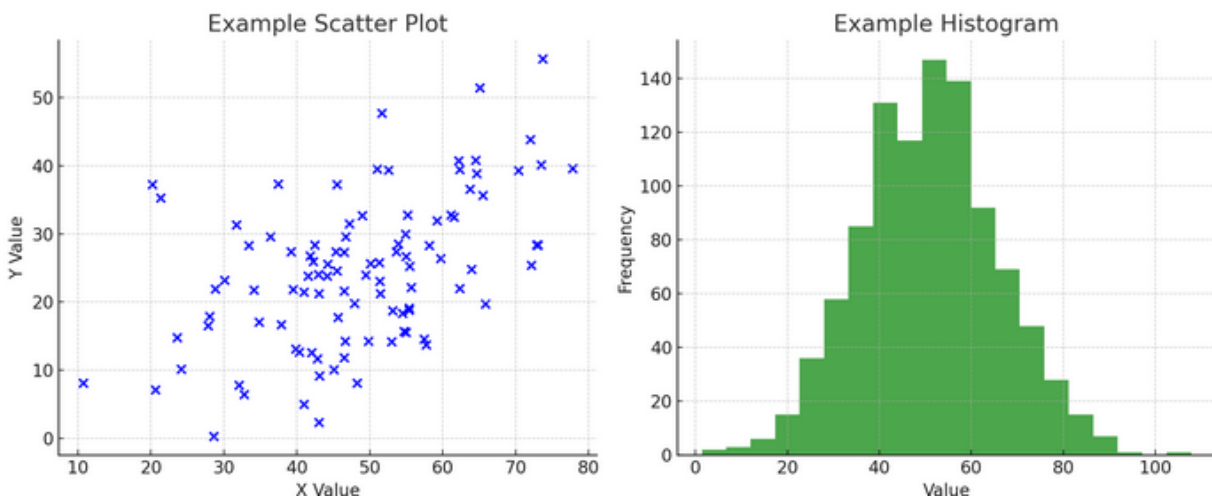
- **Charts and Graphs:** Starting with the basics, bar charts, line graphs, and pie charts provide foundational insights into the data. They are the ABCs of data visualization, essential for any data storyteller.

We have an example of a bar chart (left) and a line chart (right) below.



- **Advanced Visualizations:** As we dive deeper, heat maps, scatter plots, and histograms allow for more nuanced analysis. These tools help identify trends, distributions, and outliers.

Below is an example of a scatter plot and a histogram



Visualizations bridge raw data and human cognition, enabling us to interpret and make sense of complex datasets quickly.

**Use Case Scenario:**

**Insurance Underwriting:** Insurance companies use probability models to assess risk and set premiums for various insurance policies, such as life insurance, health insurance, and property insurance. Probability calculations help determine the likelihood of an event (e.g., illness, accident, natural disaster) occurring and the potential financial impact, allowing insurers to price policies accurately and manage their risk exposure.

Example: An insurance company uses logistic regression, a probability-based algorithm, to assess the likelihood of an individual filing a claim based on their demographic and lifestyle factors (e.g., age, gender, occupation, smoking status). The algorithm calculates the probability of a claim occurrence and determines the appropriate premium to charge for a life insurance policy.

You can use libraries such as scikit-learn to implement logistic regression for insurance underwriting. Here's a simplified example:

```python
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
# Assuming X and y are your feature matrix and labels
X = np.array([[1, 2], [3, 4], [5, 6]])
y = np.array([0, 1, 0])
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Create and train the logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)
```

**Financial Risk Management:**

Financial institutions use probability models to quantify and manage various types of risk, including credit risk, market risk, and operational risk. Probability calculations help estimate the likelihood of default on loans, fluctuations in asset prices, and the occurrence of adverse events, enabling banks and investment firms to make informed decisions about lending, trading, and investment strategies.

Example: A bank employs Monte Carlo simulation, a probabilistic modeling technique, to assess the risk of default on a portfolio of loans. The simulation generates thousands of possible future scenarios for economic conditions and borrower behavior, estimating the probability distribution of potential losses and informing the bank's capital allocation and risk mitigation strategies.

Example: Monte Carlo Simulation

Python Implementation: You can use libraries such as NumPy to implement Monte Carlo simulation. Here's a simplified example:

```python
import numpy as np
# Generate random samples for the distribution of risk factors
risk_factor_samples = np.random.normal(mean, std_dev, num_samples)
# Run simulations to estimate potential losses
simulated_losses = np.maximum(0, portfolio_value - risk_factor_samples)
```

**Marketing and Customer Segmentation:**

Companies use probability models to analyze customer data and segment their customer base into different groups based on demographic, behavioral, and transactional attributes. Probability calculations help identify patterns and trends in customer behavior, predict future purchasing behavior, and target marketing campaigns to specific customer segments more effectively, increasing customer engagement and sales.

Example: A retail company utilizes k-means clustering, a probabilistic clustering algorithm, to segment its customer base into distinct groups based on their purchasing behavior and preferences. By analyzing the probability of customers belonging to different segments, the company can tailor marketing campaigns and promotions to target specific customer segments effectively.

Example: K-means Clustering

Python Implementation: You can use libraries such as scikit-learn to implement k-means clustering. Here's a simplified example:

```python
from sklearn.cluster import KMeans
# Assuming X is your feature matrix
kmeans = KMeans(n_clusters=3)
clusters = kmeans.fit_predict(X)
```

**Supply Chain Management:**

Supply chain managers use probability models to forecast demand, plan inventory levels, and optimize production and distribution processes. Probability calculations help estimate the likelihood of demand fluctuations, stockouts, and delivery delays, enabling companies to minimize inventory costs, improve customer service levels, and reduce supply chain risks.

Example: A manufacturing company employs a demand forecasting model based on autoregressive integrated moving average (ARIMA), a time series analysis technique, to predict future demand for its products. By estimating the probability distribution of future demand, the

company can optimize inventory levels, production schedules, and procurement decisions to minimize stockouts and inventory holding costs.

Example: Autoregressive Integrated Moving Average (ARIMA)

Python Implementation: You can use libraries such as statsmodels to implement ARIMA for demand forecasting. Here's a simplified example:

```python
from statsmodels.tsa.arima.model import ARIMA
# Assuming data is your time series data
model = ARIMA(data, order=(5,1,0))
results = model.fit()
```

### Quality Control and Manufacturing:

Manufacturers use probability models to monitor and control product quality, detect defects, and ensure compliance with quality standards. Probability calculations help assess the likelihood of defects occurring during the production process, identify root causes of quality issues, and implement corrective actions to improve product quality and reduce manufacturing waste and rework.

Example: An automotive manufacturer uses statistical process control (SPC), a probability-based method, to monitor the quality of its production processes and detect deviations from quality standards. By analyzing the probability distribution of key process parameters (e.g., dimensions, tolerances), the company can identify and address potential sources of variation, ensuring consistent product quality and minimizing defects.

Example: Statistical Process Control (SPC)
Python Implementation: You can use libraries such as scipy and pandas to implement statistical process control charts. Here's a simplified example:

```python
import pandas as pd
from scipy.stats import norm
# Assuming data is your process measurements
mean, std_dev = np.mean(data), np.std(data)
control_limits = norm.interval(0.997, loc=mean, scale=std_dev)
```

### Healthcare and Clinical Trials:

Healthcare providers and pharmaceutical companies use probability models to analyze patient data, diagnose diseases, and evaluate the effectiveness of medical treatments. Probability calculations help estimate the likelihood of disease occurrence, predict patient outcomes, and design clinical trials to test new drugs and therapies, leading to improved patient care and medical advancements.

Example: A pharmaceutical company conducts a clinical trial to evaluate the efficacy of a new drug for treating a specific medical condition. The trial uses survival analysis, a probability-based

modeling technique, to estimate the probability of patients experiencing certain health outcomes (e.g., disease progression, survival) over time. By analyzing the probability of treatment success, the company can make data-driven decisions about the drug's effectiveness and regulatory approval.

Example: Survival Analysis

Python Implementation: You can use libraries such as lifelines to implement survival analysis. Here's a simplified example:

```
from lifelines import KaplanMeierFitter
# Assuming T is your time-to-event data and E is your event indicator
kmf = KaplanMeierFitter()
kmf.fit(T, event_observed=E)
```

In each example, probability-based algorithms play a critical role in analyzing data, making predictions, and informing decision-making processes in various business contexts. These algorithms leverage probability theory to quantify uncertainty, assess risk, and derive actionable insights from data, ultimately helping businesses optimize their operations and achieve their strategic objectives.

In each example, Python libraries such as scikit-learn, statsmodels, scipy, pandas, and lifelines provide implementations of probability-based algorithms that can be readily applied to solve real-world business problems.

# Hypothesis Testing

Think of hypothesis testing as detective work in statistics. It's a method to test if a particular theory about our data could be true. This process starts with two opposing hypotheses:

- Null Hypothesis (H0): This is the default assumption, suggesting that there's 'No' effect or difference. It's saying, "Not" ing new here."

- Al "alternative Hypothesis (H1 or Ha): This challenges the status quo, proposing an effect or a difference. It claims, "Something is interesting going on."

Example: Testing if a new diet program leads to weight loss compared to not following any diet.

**Null Hypothesis (H0):** The new diet program does not lead to weight loss (no difference in weight loss between those who follow the new diet program and those who do not).

**Alternative Hypothesis (H1):** The new diet program leads to weight loss (a difference in weight loss between those who follow it and those who do not).

Hypothesis testing involves choosing between these two based on the evidence (our data).

**Type I and II Error and Significance Levels:**

**Type I Error:** This happens when we incorrectly reject the null hypothesis. It convicts an innocent person.
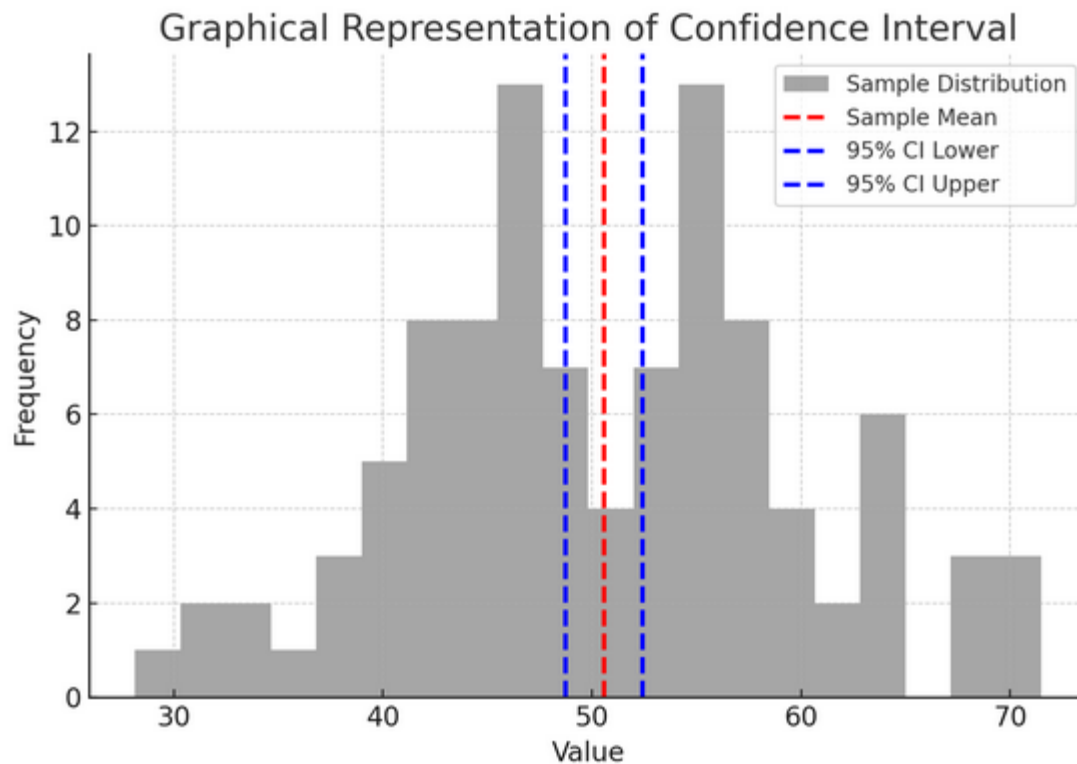
**Type II Error:** This occurs when we fail to reject a false null hypothesis. It lets a guilty person go free.

**Significance Level (α):** This is the threshold for deciding how much evidence is enough to reject the null hypothesis. It's often set at 5% (0.05), indicating a 5% risk of a Type I error.

# Confidence Intervals

Confidence intervals give us a range of values within which we expect the valid population parameter (like a mean or proportion) to fall with a certain confidence level (commonly 95%). It's like predicting a sports team's final score with a margin of error; we're saying, "We're 95% confident the true score will be within this range."

Constructing and interpreting confidence intervals helps us understand the precision of our estimates. The wider the interval, our estimate is less precise, and vice versa.



The above figure illustrates the concept of a confidence interval (CI) in statistics, using a sample distribution and its 95% confidence interval around the sample mean.

Here's a breakdown of the critical components in the figure:

**Sample Distribution (Gray Histogram):** This represents the distribution of 100 data points randomly generated from a normal distribution with a mean of 50 and a standard deviation of 10. The histogram visually depicts how the data points are spread around the mean.

**Sample Mean (Red Dashed Line):** This line indicates the sample data's mean (average) value. It serves as the point estimate around which we construct the confidence interval. In this case, it represents the average of all the sample values.

**95% Confidence Interval (Blue Dashed Lines):** These two lines mark the lower and upper bounds of the 95% confidence interval around the sample mean. The interval is calculated using the standard error of the mean (SEM) and a Z-score corresponding to the desired confidence level (1.96 for 95% confidence). The confidence interval suggests we are 95% confident that the population mean lies within this range.

## Correlation and Causation

Correlation and causation often get mixed up, but they are different:

**Correlation:** Correlation refers to a statistical measure that describes the extent to which two variables are related or move together in a systematic way. When two variables are correlated, changes in one tend to be associated with changes in the other. Correlation is measured by a correlation coefficient ranging from -1 to 1. A value closer to 1 or -1 indicates a strong relationship, while 0 suggests no ties.

**Causation:** Causation on the other hand, implies that one variable directly causes a change in another variable. Establishing causation typically requires more than just observing a correlation. It involves rigorous experimentation, control of confounding factors, and often temporal precedence (i.e., the cause must precede the effect).

Here's why it's crucial to distinguish between the two:

Correlation: It's valuable for identifying relationships and making predictions. For instance, knowing that there is a strong positive correlation between smoking and lung cancer risk helps in understanding health patterns. Just because two variables are correlated does not mean one causes the other. This is a classic case of not confusing "correlation" with "causation."

Causation: This is essential for making informed decisions and interventions. For example, concluding that smoking causes lung cancer enables public health campaigns to target smoking cessation.

In research and decision-making, mistaking correlation for causation can lead to erroneous conclusions and inappropriate actions. Therefore, it's important to use statistical methods and research designs that can help differentiate between correlation and causation when exploring relationships between variables.

**The purpose of Hypothesis Testing** is to see if the Null Hypothesis can be rejected or not. In most cases, rejecting the Null Hypothesis does not imply that the alternative hypothesis is true. Rejecting the Null Hypothesis, on the other hand, can lead to acceptance of the Alternative Hypothesis in some cases.

There are two types of errors that can occur when executing a Hypothesis Test:

• <u>Type-I Error:</u> Rejecting the Null Hypothesis when it is actually true is a Type-I Error.

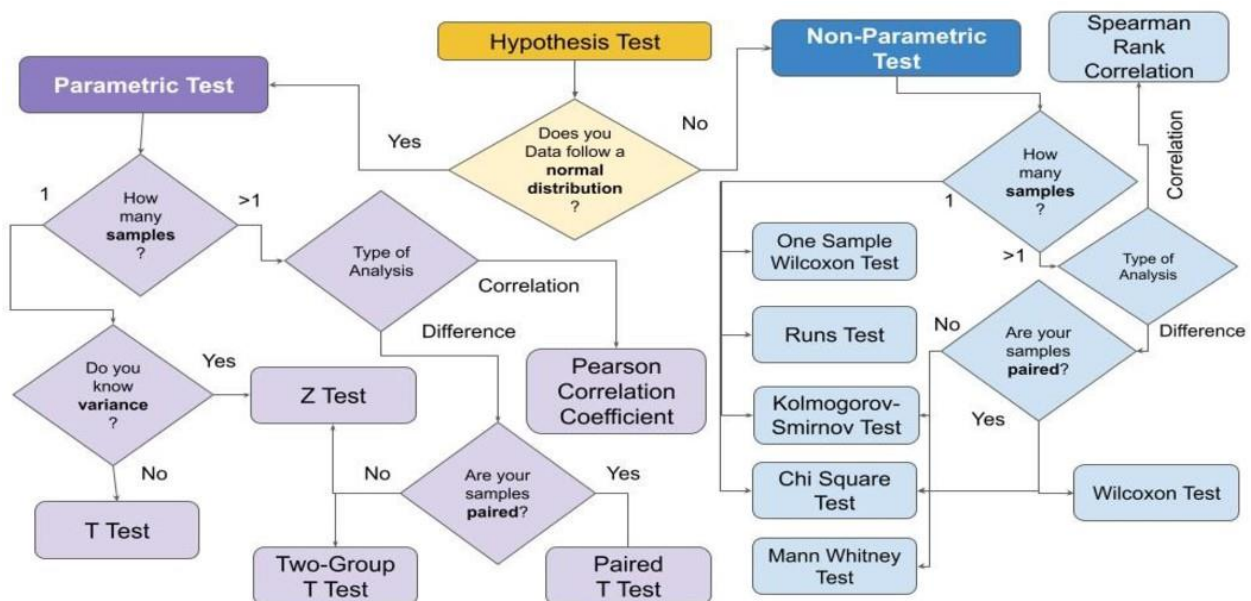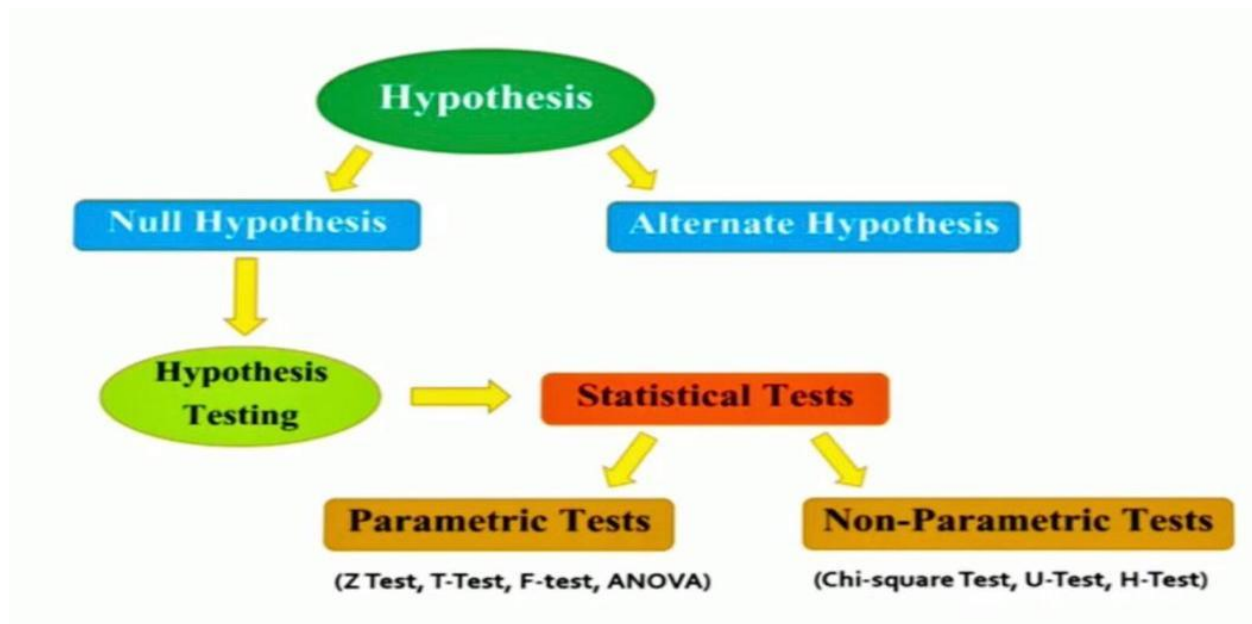• <u>Type-II Error:</u> Accepting the Null Hypothesis when it is false is a Type-II Error.

|  | Truth | |
|---|---|---|
| **Hypothesis Testing** | **The Null Hypothesis Is True** | **The Alternative Hypothesis Is True** |
| **The Null Hypothesis Is True** | Accurate | Type II Error |
| **The Alternative Hypothesis Is True** | Type I Error | Accurate |

_(Row header label: **Research**)_

# Types of Hypothesis Testing

- **Hypothesis tests are divided into two categories:**

1) **Parametric tests –** are used when the samples have a normal distribution. In general, samples with a mean of 0 and a variance of 1 follow a normal distribution.

2) **Non-Parametric tests –** If the samples do not follow a normal distribution, non-parametric tests are used.

- **Two types of Hypothesis Testing can be created depending on the number of samples to be compared:**

• **One Sample –** If there is only one sample that must be compared to a specific value, it is called a single sample.

• **Two Samples –** if you're comparing two or more samples. Correlation and sample difference are two tests that could be used in this situation. Samples can be paired or not in both circumstances. Dependent samples are sometimes known as paired samples, while independent samples are known as unpaired samples. Natural or matched couplings occur in paired samples.

# Parametric and Non-Parametric Tests

A parametric test is one in which the parameters are predetermined and the population distribution is always known. A mean value is used to calculate the central tendency. These tests are common, which makes conducting research relatively simple and time-consuming. The Non-parametric test makes no assumptions and measures using the median value. Kruskal-Wallis, Mann-Whitney, and other non-parametric tests are examples.

**Industry demonstration of hypothesis testing:**

1.  This demonstration illustrates how hypothesis testing is applied in the pharmaceutical industry to assess the efficacy of new drugs before they are brought to market.

**Scenario:** A pharmaceutical company has developed a new drug that they believe can lower blood pressure in patients more effectively than the current leading medication in the market. Before they can start marketing the new drug, they need to prove its effectiveness through hypothesis testing.

**Hypothesis:**

**Null Hypothesis (H0):** The new drug is not more effective at lowering blood pressure compared to the current leading medication.

**Alternative Hypothesis (H1):** The new drug is more effective at lowering blood pressure compared to the current leading medication.

**Sample Selection:** The company selects a sample of patients with high blood pressure who are willing to participate in a clinical trial. They randomly assign half of the participants to receive the new drug (experimental group) and the other half to receive the current leading medication (control group).

**Data Collection:** Blood pressure readings are taken from both groups before and after the treatment. The data includes systolic and diastolic blood pressure measurements.

**Statistical Analysis:**

a. Calculate Descriptive Statistics: Calculate the mean and standard deviation of blood pressure for both groups before and after treatment.

b. Perform Hypothesis Test: Use a paired t-test to compare the mean reduction in blood pressure between the two groups. This test is appropriate because it compares the means of two related groups (before and after treatment) and determines whether there is a statistically significant difference

**Interpretation of Results:**

If the p-value obtained from the t-test is less than the significance level (e.g., $\alpha = 0.05$), then we reject the null hypothesis.

If the null hypothesis is rejected, it means there is sufficient evidence to conclude that the new drug is more effective at lowering blood pressure compared to the current leading medication.

If the p-value is greater than the significance level, we fail to reject the null hypothesis, indicating that there is not enough evidence to conclude that the new drug is more effective.

**Conclusion:** Based on the results of the hypothesis test, the pharmaceutical company can make an informed decision about whether to proceed with marketing the new drug. If the null hypothesis is rejected, they can move forward with confidence, knowing that they have statistical evidence supporting the effectiveness of their product.

2. This demonstration showcases how hypothesis testing can be applied in the automotive industry to assess the performance of new products, such as engine oils, before they are introduced to consumers.

**Scenario:** An automotive manufacturer is developing a new type of engine oil that they claim can improve fuel efficiency in vehicles compared to the currently popular engine oils in the market.

Hypothesis:

Null Hypothesis (H0): The new engine oil does not improve fuel efficiency compared to the current leading engine oils.

Alternative Hypothesis (H1): The new engine oil improves fuel efficiency compared to the current leading engine oils.

**Sample Selection:** The automotive manufacturer selects a sample of vehicles from different models and makes. They ensure the vehicles are in similar conditions and have similar mileage. The sample includes vehicles that will be using the new engine oil (experimental group) and vehicles using the current leading engine oils (control group).

**Data Collection:** The vehicles are driven for a specified distance under controlled conditions, and their fuel consumption is measured. Data includes fuel efficiency metrics such as miles per gallon (MPG) or liters per 100 kilometers (L/100 km).

**Statistical Analysis:**

a. Calculate Descriptive Statistics: Calculate the mean and standard deviation of fuel efficiency for both groups.
b. Perform Hypothesis Test: Use a two-sample t-test to compare the mean fuel efficiency between the group using the new engine oil and the group using the current leading engine oils. This test determines whether there is a statistically significant difference in fuel efficiency between the two groups.

**Interpretation of Results:**

> If the p-value obtained from the t-test is less than the significance level (e.g., $\alpha = 0.05$), then we reject the null hypothesis.
> Rejecting the null hypothesis indicates that there is sufficient evidence to conclude that the new engine oil improves fuel efficiency compared to the current leading engine oils.
> If the p-value is greater than the significance level, we fail to reject the null hypothesis, indicating that there is not enough evidence to conclude that the new engine oil improves fuel efficiency.

> **Conclusion:** Based on the results of the hypothesis test, the automotive manufacturer can decide whether to proceed with marketing the new engine oil. If the null hypothesis is rejected, they can confidently market the new product, knowing that they have statistical evidence supporting its improved fuel efficiency.

Packages:

PIP INSTALL SCIKIT-LEARN IN YOUR CMD PROMPT or Anaconda prompt or in Jupyter notebook

The **specific packages** you'll need for machine learning can depend on what you're trying to accomplish, but some of the most commonly used ones include:

1. **NumPy**: NumPy is a fundamental package for numerical computing in Python. It provides support for arrays, matrices, and mathematical functions.

2. **Pandas**: Pandas is a powerful library for data manipulation and analysis. It provides data structures like DataFrame, which are essential for handling structured data.

3. **scikit-learn**: Scikit-learn is one of the most popular machine learning libraries in Python. It provides a wide range of tools for classification, regression, clustering, dimensionality reduction, and more.

4. **TensorFlow** or **PyTorch**: These are deep learning frameworks that allow you to build and train neural networks. TensorFlow is developed by Google, while PyTorch is developed by Facebook. Both are widely used and have extensive documentation and community support.

5. **Matplotlib** or **Seaborn**: These libraries are used for data visualization in Python. Matplotlib is more low-level and provides fine-grained control over plots, while Seaborn is built on top of Matplotlib and offers more high-level functions for statistical visualization.

6. **Keras**: Keras is a high-level neural networks API that runs on top of TensorFlow or other deep learning frameworks. It provides a simpler interface for building and training neural networks, especially for beginners.

7. **SciPy**: SciPy is a collection of mathematical algorithms and convenience functions built on top of NumPy. It provides additional functionality for optimization, integration, interpolation, and more.

8. **NLTK** or **spaCy**: If your machine learning tasks involve natural language processing (NLP), you'll likely need one of these libraries. NLTK is more traditional and provides a wide range of tools for text processing and analysis, while spaCy is more modern and optimized for performance.

9. **XGBoost** or **LightGBM**: These are gradient boosting libraries that are highly efficient and widely used for structured/tabular data problems, especially in competitions like Kaggle.

10. **OpenCV**: If your machine learning tasks involve computer vision, OpenCV is an essential library. It provides tools and algorithms for image and video analysis, including object detection, face recognition, and more.

These are just some of the most commonly used packages in machine learning, but there are many others out there depending on your specific needs and preferences.

**TensorFlow, Keras, and scikit-learn (often abbreviated as sklearn)** are all popular libraries in the field of machine learning and deep learning, but they serve different purposes and are used in different contexts:

TensorFlow:

TensorFlow is an open-source deep learning library developed by Google.

It is widely used for building and training deep learning models, particularly neural networks, including both traditional feedforward networks and more advanced architectures like convolutional neural networks (CNNs) and recurrent neural networks (RNNs).

TensorFlow provides a flexible framework for building and deploying machine learning and deep learning models at scale, with support for both CPU and GPU computing.

It is often used for research, production-level deployment, and projects requiring low-level control over model architecture and training process.

Keras:

Keras is a high-level neural networks API, originally developed as an independent project but later integrated into TensorFlow as its official high-level API.

Keras provides a user-friendly interface for building and training neural networks, allowing developers to quickly prototype and experiment with different model architectures.

It offers a simple and intuitive API for defining layers, configuring model architecture, compiling models with different optimizers and loss functions, and training models on training data.

Keras is often preferred for its ease of use, readability, and flexibility, making it suitable for both beginners and experienced deep learning practitioners.

scikit-learn (sklearn):

scikit-learn is a popular machine learning library in Python, providing a wide range of tools and algorithms for supervised and unsupervised learning, as well as data preprocessing, model evaluation, and model selection.

It offers a consistent and simple API for implementing a variety of machine learning algorithms, including linear and logistic regression, decision trees, random forests, support vector machines (SVM), k-nearest neighbors (k-NN), clustering algorithms, and more.

scikit-learn is designed to be user-friendly, efficient, and accessible to users of all skill levels, making it suitable for both beginners and experienced practitioners.

It is often used for data analysis, building predictive models, and solving classification, regression, and clustering tasks on structured data.

In summary:

Use TensorFlow for deep learning projects requiring low-level control over model architecture and training process.

Use Keras for quickly prototyping and experimenting with different neural network architectures, especially when working with TensorFlow.

Use scikit-learn for implementing a wide range of machine learning algorithms and tasks, including data preprocessing, model training, and evaluation on structured data.