# The Application of Statistics in Daily Life

## Delun Xu

**Abstract:**

Mathematics is a subject of study with major relevance in the scientific method, decision-making methods, and countless other fields. Statistics is a complex and diversified field of study located within the mathematics domain. In the field of data science, various techniques and strategies are used to collect, analyze, interpret, present, and organize data of various kinds. I am a textual user. When it comes to navigating the complexities of uncertainty inherent in real-world events, the discipline of statistics is a powerful weapon that researchers, analysts, and decision-makers may use to their advantage.

**Keywords:** mathematics, decision-making methods, scientific method, domain of mathematics

Mathematics is a subject of study with major relevance in the scientific method, decision-making methods, and countless other fields. Statistics is a complex and diversified field of study located within the mathematics domain. In the field of data science, various techniques and strategies are used to collect, analyze, interpret, present, and organize data of various kinds. I am a textual user. When it comes to navigating the complexities of uncertainty inherent in real-world events, the discipline of statistics is a powerful weapon that researchers, analysts, and decision-makers may use to their advantage.

My first perception of statistics was a visual representation consisting of a bar graph displaying varying heights of bars. Upon further enhancing my comprehension of statistics, I have realized that this particular graphical representation is formally referred to as a histogram. A histogram is a visual depiction of the frequency distribution of a given dataset. A histogram is a graphical representation that displays a given dataset's overall frequency distribution, whether continuous or discrete. Histograms are extensively used in many domains, such as statistics, data analysis, and quality control, to visually examine and comprehend the attributes of a dataset. An intuitive method is offered to comprehend the distributional features and patterns present in the data.

At its core, the primary objective of statistics is to derive useful information from data to enable practitioners and researchers to arrive at well-informed conclusions about whole populations or phenomena by using a sample that is both representative and well-selected. The process at hand requires the utilization of several different statistical methodologies and approaches. These range from the most fundamental descriptive statistics, which enumerate and depict the primary characteristics of a dataset, to the more advanced inferential statistics, which enable projections and generalizations at the population level.

On the other hand, even though we are now doing in-depth research on statistics, our goal is not only to acquire knowledge. More than anything else, it is about better understanding how to apply it in real life.

Descriptive statistics, such as the mean, median, mode, and standard deviation, may provide a concise overview of a dataset's fundamental patterns and variations. Summary statistics provide a foundation for understanding the essential elements of the data. In contrast, inferential statistics use probability theory to produce educated approximations regarding population parameters. Inferential statistics recognize the inherent uncertainty associated with working with samples instead of whole populations.

One may confront interconnected issues in life. Suppose you are researching to determine the average height of a group of students attending a certain school. In that case, the following should be considered: A representative sample of one hundred students is gathered, and their heights are measured. A measurement of the spread of heights in your sample is 10 centimeters, and the average height in your sample is 160 centimeters. Using inferential statistics, we can make a prediction with a degree of confidence of 95% about the range in which the mean height of the whole student population is anticipated to fall.

By using statistical analysis, we can effectively tackle this challenge. We are now using inferential statistics to estimate a confidence range for the population's mean height. This estimation is based on the data obtained from the sample. The formula for calculating a confidence interval for the population mean ($\mu$) is as follows: The formula for calculating the confidence interval is $x \pm Z(s/\sqrt{n})$. Within the context of this equation, the variable x stands for the sample mean(which is equal to 160cm, the variable s stands for the sample standard deviation(which

is equal to 10cm), the variable n stands for the sample size(which is equal to 100 students), and the variable Z stands for the Z-score that is connected with the desired confidence level(which is equal to 95%). For a 95% confidence interval, the Z-score is approximately 1.96 (this value is commonly used for large samples). To establish a range, it is necessary to add or remove 1.96 from the average height of the class when we get this value. We can get the confidence interval (160 -1.96, 160 +1.96). After that, we can figure this out and get the answer ( 158.04, 161.96 ). With a confidence level of 95%, we may conclude that the average height of the total student population is expected to be between 157.04 cm and 162.96 cm.

Identifying the components of the formula, establishing the right Z-score for the desired degree of confidence, plugging in the values, simplifying the expression, and interpreting the result in the context of the situation are the many phases involved in the process. Using the data from the sample, this method generates an estimate supported by statistical evidence of the probable range for the population's mean.

Before talking about probability, we must mention Bayes's Rule. Before discussing probability, it is essential to acknowledge Bayes's Rule. The reason for this is because it is the fundamental principle of probability. The formula gives the expression of Bayes' Rule: $P(A \mid B) = P(B \mid A) \times P(A) / P(B)$. The key aspect of this theorem is in the recognition of the correlation between A and B and the computation of the transformed correlation between A and B.

Here's an example of what I thought: Suppose there is a rare disease and a diagnostic test has been developed to detect it. The prevalence of the disease in the population is known to be 1 in 1000 (P(Disease)=P(D)=0.001). The sensitivity of the test (probability of a positive result given the presence of the disease) is 95% (P(Positive $\mid$ Disease)=P(Pos $\mid$ D)=0.95). The specificity (probability of a negative result given the absence of the disease) is 90% (P(Negative $\mid$ NoDisease)=P(Neg $\mid$ D)=0.90). If a person tests positive, what is the probability of having the disease?

Bayes' Rule may be immediately applied to this question using the formula. We put the name into the formula: P(Disease|Positive) = (P(Positive|Disease)*P(Disease))/P(Positive). We can calculate the marginal probability of a positive test(P(Positive): P (Positive) = P(Positive|Disease) * P(Disease) + P(Positive|NoDisease) * P(NoDisease). The probability of a positive test given no disease(P(Positive|NoDisease)) is 1-P(Negative|NoDisease) since specificity is the probability of a true negative.Also, P(NoDisease)=

1-P(Disease).We can put the data in each function, and we can get: P(Positive) = ( 0.95 * 0.001)+ ((1-0.90) * (1- 0.001)). We also can get the function of P(Disease | Positive), which is equal to (0.95 * 0.001)/(Calculated P(Positive). And then, we can find out the answer to this question.

Based on a positive test result, this computation will determine the chance of a person having the illness. Please be aware that the precise numerical outcome depends on the exact values assigned to sensitivity, specificity, and prevalence in the given problem.

Probability, a basic statistical notion, is the foundation for several statistical procedures by measuring the probability of different outcomes. This probabilistic framework is crucial for generating accurate statistical inferences and making informed decisions. Hypothesis testing is a crucial component of statistical analysis. It entails creating hypotheses about population characteristics and using statistical tests to assess the accuracy of these assumptions using sample data.

An illustrative instance that may be seen in real-life situations is: Imagine a scenario where a pharmaceutical corporation creates a novel medication specifically formulated to reduce hypertension. The mean systolic blood pressure in the general population is 120 mm Hg. The corporation asserts that its novel medication effectively lowers blood pressure and conducts a study. A sample of 50 persons who consumed the drug was gathered, and the average systolic blood pressure of the sample was found to be 115 mm Hg, with a standard variation of 8 mm Hg. By doing hypothesis testing at a significance level of 5%, we can ascertain if there is enough data to substantiate the company's assertion that the medicine effectively reduces blood pressure.

The first stage in statistical analysis involves formulating hypotheses. The null hypothesis is that the new medication does not impact blood pressure, with a mean value of 120. The alternative hypothesis is that the new medicine significantly reduces blood pressure, with a mean value of less than 120. We choose a significance level (α) of 0.05, corresponding to a 5% significance level. We must also choose the suitable statistical test: A one-sample t-test compares a sample mean to a population mean, especially when the sample size is small (n < 30). We must gather and examine data and then input it into the formula to compute the test statistic. In this equation, the sample mean (x) is 115 mm Hg, the sample standard deviation (s) is 8 mm Hg, and the sample size (n) is 50 persons who took the drug. The test statistic may be calculated using the formula t = (s / √n) * (x̄ - μ₀), where s represents the standard deviation, √n represents the square root of the sample size, x̄ represents the sample

mean, and $\mu_0$ represents the hypothesized population mean. The calculation can be applied as t = (115 - 120)/(8 √50). We must now ascertain the crucial region. To get the crucial t-value for a one-sample t-test with a significance level of 0.05 and degrees of freedom (df = 49) using a t-table or statistical software, we need to discover the value corresponding to the 0.05 significance level and the given degrees of freedom. We should choose our course of action by comparing the computed t-value with the important t-value. The null hypothesis should be rejected if the computed t-value is inside the crucial zone. If the calculated t-value falls inside the critical region, it may be inferred that there is sufficient evidence to reject the null hypothesis. In this context, it suggests that the new medicine effectively lowers blood pressure levels.

This step-by-step solution illustrates the use of hypothesis testing to evaluate the pharmaceutical company's assertion on the efficacy of a novel medication for reducing blood pressure. The process involves developing hypotheses, setting a significance level, selecting the suitable statistical test, gathering and analyzing data, computing the test statistic, identifying the crucial area, making a choice, and drawing a conclusion based on the evidence derived from the sample.

At the same time, unpredictability presents a second challenge in life. Imagine that you are employed by a manufacturing business that is responsible for the production of light bulbs. According to the manufacturer, its light bulbs' typical lifetime is one thousand to one thousand hours. On the other hand, you are the kind of person who is dubious and believes that the real average lifetime is different. After selecting fifty light bulbs randomly from the manufacturing line, you discover that the sample mean lifetime is 980 hours, with a standard variation of fifty hours. This information is based on the production line. Can you present evidence to back up your skepticism using a significance threshold of 0.05?

The problem-solving concept is the same as the one discussed in the last pharmaceutical problem. The first step in our process is to formulate hypotheses based on the parameters of the investigation. The Null Hypothesis posits that the average lifespan of light bulbs is 1000 hours ($\mu$=1000), based on the assumption. (Ha)The alternative hypothesis may be stated as follows: The average lifespan of light bulbs deviates from 1000 hours ($\mu\neq$1000), which is worth noting. Currently, we are obligated to choose a degree of relevance. The significance level ($\alpha$) represents the probability of erroneously rejecting the null hypothesis when it is true.

The term "significance level" refers to the degree of importance applied in the hypothesis testing process. One often used value is 0.05, which represents a 5% probability of committing a Type I error when a legitimate null hypothesis is rejected incorrectly. The subsequent task is selecting the statistical test that is most appropriate for our specific circumstances. A Z-test is the appropriate statistical approach to apply when comparing the mean of a sample to the known mean of a population, especially when the population's standard deviation is also known. Next, we will continue to calculate the measurement statistic. The answer obtained using the t formula mentioned in the preceding question, $Z = (x - \mu) / (\sigma / \sqrt{n})$, will have the same look. To address this problem, compiling a comprehensive inventory of all the data provided to us is necessary. The sample's mean is 980 hours, the population's standard deviation is 50 hours, the sample size consists of fifty randomly selected bulbs, and the significance level is 0.05, which is comparable to five percent. By substituting the data back into the Z formula, we can calculate that Z equals the result of (980 - 1000) divided by (50 divided by the square root of 50).

Upon the conclusion of this computation, we will have the result of the Z test statistic. We consult a table representing the standard normal distribution with a significance level of 0.05 (using a two-tailed test) to get the critical Z-values. The cumulative probability of 0.025 is associated with each tail, and these values represent the chance of occurrence. Compare the magnitude of the Z-test statistic produced to the critical Z-values. If the magnitude of the value exceeds the critical Z-value, it may be inferred that the null hypothesis is invalid and should be rejected. Another decision must be made is whether the test statistic falls inside the critical region. If such is the case, we will refute the null hypothesis. The null hypothesis can only be rejected if this condition is not met. Finally, we need to analyze and interpret the result. Arrive to a conclusion that is based on the decision made. We have gathered information that supports the alternative hypothesis, indicating that the average lifespan of the light bulbs deviates from 1000 hours. If we choose not to accept the null hypothesis, we have gathered evidence that supports the alternative hypothesis.

In this example, we conducted a hypothesis test to investigate the assertion made by a manufacturing company on the average lifespan of their light bulbs. The analyst harbored skepticism and believed the mean lifespan diverged from the firm's stated value. The company said that the mean duration of existence is one thousand hours.

Regression analysis is an advanced statistical approach exploring many variables' relationships. This approach enables the identification and measurement of patterns and trends via data analysis. This method is particularly valuable in topics like economics, sociology, and

environmental science, where it is crucial to understand the interrelationships among variables.

Multiple regression is a statistical technique that expands on basic linear regression by allowing the prediction of a dependent variable using two or more independent variables. It facilitates the investigation of the connections between a variable influenced by other variables and numerous predictors, enabling a more thorough comprehension of the elements that impact the final result. The standard format of a multiple regression model including two independent variables (X1 and X2) may be expressed as: $Y = \beta_0 + \beta_1 * X_1 + \beta_2 * X_2 + \epsilon$. In the above formula, Y stands for the dependent variable, $\beta_0$ stands for the y-intercept, $\beta_1$, and $\beta_2$ stand for the coefficients for X and X2 correspondingly, and $\epsilon$ stands for the error term. An instance that comes to mind is: Suppose you are studying the factors influencing employee performance, considering the number of hours spent on training (X1) and years of work experience (X2). Can you build a multiple regression model to predict employee performance based on these two factors? Now I can give the number of training hours equal to 24 hours, and the years of work experience is equal to 2 years.

The info supplied enables us to resolve the problem at this juncture. The value of X1 is equivalent to twenty-four hours, whereas the value of X2 is equivalent to two years. To assume the estimated coefficients, we may consider $\beta_0$ as 75 (the y-intercept), $\beta_1$ as 3 (the coefficient of X1), and $\beta_2$ as 5 (the coefficient of X2). We must input the data into the equation to do multiple regression and compute the result. The equation is $y = 75 + 3(24) + 5(2)$. Based on this, we can ascertain that the result equals 157. Here is the current explanation: The multiple regression model predicts that the employee's performance (Y) will be 157 units, given that the employee has 24 hours of training (X1=24) and two years of work experience (X2=2).

Furthermore, we may ascertain the significance of the coefficients by The y-intercept of 75, which indicates that the expected performance is 75 units when both the number of training hours and the amount of work experience are zero. Assuming that the amount of work experience remains constant, it is expected that the anticipated performance will rise by three units for each extra hour of training (X1). The $\beta_1$ idea is based on this assumption. The value of $\beta_2$ is contingent upon the assumption that the total number of training hours remains constant. The forecasted performance is expected to grow by five units for each extra year of work experience (X2).

Multiple regression is a valuable technique in real-world scenarios for comprehending and predicting the relationships between various variables and the desired outcome. This provides valuable perspectives for decision-making in several fields, including economics, social sciences, and human resources.

The function of statistics will certainly continue to develop as the world becomes more data-driven. The extraction of meaningful insights from complicated datasets is facilitated by combining cutting-edge statistical approaches and developing technological developments. Statistical analysis is vital for making sense of the huge and complex landscapes of data, which eventually contributes to improved decision-making and problem-solving. This is true in various fields, including the social sciences, business, healthcare, and other areas. Statistics is widespread throughout various disciplines, including but not limited to economics, biology, psychology, sociology, engineering, and many more. By offering tools for analyzing and interpreting data, it contributes to discovering patterns, trends, and correlations, which ultimately play a significant role in decision-making, research, and the scientific process. The significance of statistics in gleaning significant insights from complicated datasets is becoming more crucial as technological advancements and data collection techniques evolve.