Dr. Upendra Pratap Singh

LNMIIT

September 8, 2023

Introduction to Inferential Statistics

Motivation

Introduction

- Inferential statistics involves drawing conclusions/making inferences about a population based on a sample of data from that population.
- ② Useful when it is not feasible to get the population data.

Motivation

- Cost and Resource Constraints: expensive data collection, logistic chain complexity.
- 2 Time Efficiency: conducting a study on an entire population might take an impractically long time.
- Destructive Testing: in fields such as medical research, destructive testing might be involved, where samples are consumed or destroyed during the testing process.
- Population Variability: populations are often diverse and heterogeneous. By studying a sample, researchers can capture a range of variation that is representative of the broader population.



Motivation

- Infeasibility: inferential statistics offers a way to work with manageable subsets while still making meaningful inferences.
- Ethical Considerations: it might not be ethical to gather data from an entire population.
- Testing Hypotheses: researchers often have specific hypotheses or questions they want to address. Inferential statistics provides a framework for testing these hypotheses using sample data and making conclusions about the population.

Introduction

What it involves?

- Sampling: This is the process of selecting a subset of individuals or items from a larger population to represent it.
- 4 Hypothesis Testing: This involves making educated guesses about a population parameter and using sample data to test whether the hypotheses is likely to be true or not.
- Onfidence Intervals: finding a range around a sample statistic (like the sample mean) that is likely to contain the corresponding population parameter.

Introduction

What it involves?

- Regression Analysis: model the relationship between dependent and independent variables and use that model to make predictions.
- Estimation: involves estimating population parameters based on sample statistics.
- Significance Testing: process of determining whether an observed effect in a sample is statistically significant or if it could have occurred by chance.
- Probability: to quantify the likelihood of different outcomes or events occurring, given a certain set of assumptions.

Statistic Vs Parameter

Statistic

- A statistic is a numerical measure calculated from a sample.
- May differ from sample to sample reflecting the inherent variability in the sampling process.
- **3** Examples include: sample mean (\bar{X}) , sample standard deviation and sample median, etc.

Statistic Vs Parameter

Parameter

- A parameter is a numerical measure that describes a specific characteristic of an entire population.
- Fixed value that may not be observable in its entirety.
- Serve as targets of inference in statistical analysis.
- **1** Examples include: population mean (μ) , population standard deviation (σ) and variance.

Estimator

- An estimator is a statistic or a method used to calculate an estimate of an unknown parameter.
- In other words, an estimator is a formula, a calculation, or a procedure that takes observed sample data as input and produces an estimated value for a population parameter as output.

Estimator Types

- **1** Point Estimator: A point estimator is a specific value calculated from the sample data that serves as an estimate for the true parameter value; example: sample mean (\bar{x}) .
- Interval Estimator: An interval estimator provides a range of values within which the true parameter value is likely to fall; example: confidence intervals.
- Maximum Likelihood Estimator (MLE): finds parameter values that maximize the likelihood of observing the given sample data.

Estimator Types

• Bayesian Estimator: an estimator incorporates both prior beliefs and the likelihood of the data to calculate a posterior distribution.

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$
 (1)

- 2 Robust Estimator: is less sensitive to outliers or deviations.
- An ideal estimator would have low bias and low variability.



Properties of a good estimator

- **Unbiased:** The center of the sampling distribution for the estimate is the same as that of the population.
- The estimate has the smallest standard error when compared to other estimators.
- Standard error: measures the variability of sample statistics across multiple samples.

$$Standard Error = \frac{Standard Deviation}{\sqrt{N}}$$
 (2)



Point Estimate: Mean

- Estimating the unknown population parameter is straightforward.
- ② Calculate the sample mean (\bar{x}) and use it as an estimate of true population mean μ .
- **3** In other words, $\hat{\mu} = \bar{x}$.



Point Estimate: Standard Deviation

- Sample needs to be reasonable sized. Why?. For example, in a sample sized 1, standard deviation will always be 0.
- Minimum of two samples needed to observe any variability.
- **3** Sample deviation s is usually small compared to the population standard deviation σ . Moreover, sample deviation s increases with the size of the sample.

$$s^{2} = \frac{1}{N} \sum_{i=1}^{i=N} (X_{i} - \bar{X})^{2}$$
 (3)

Point Estimate: Standard Deviation

- In other words, sample variance s^2 is the biased estimator of true population variance σ^2
- Solution: Bessels's correction

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^{N} (X_i - \bar{X})^2$$
 (4)

Homework: Find the point estimate for a median.



Interval Estimate: Confidence Intervals

- Onfidence interval: an interval of values computed from the sample data that is likely to cover the true parameter.
- For a good CI, the center of the interval should be the point estimate for the parameter of interest.

Interval Estimate: Confidence Intervals

General form of confidence interval:

Confidence interval = sample statistic
$$\pm$$
 margin of error (5)

② General form of margin of error:

margin of error =
$$M \times \hat{SE}$$
 (6)

where \hat{SE} is the standard error of sample statistic and M is a multiplier based on how confident we want to be in.

Interpretation: We are 95% confident that the true population parameter lies within the calculated confidence interval.



Confidence Intervals: Advantages and Limitations

Advantages

- An easy interpretation to quantifies uncertainty.
- Incorporates the sample size.
- Maybe adjusted for different confidence intervals (M needs to change accordingly).

Limitations

- Doesn't ensure true parameter inclusion.
- No information about distribution shape, bias.
- Sensitive to outliers.



Hypothesis Testing: P-Value

- Hypothesis Testing: a more formal method for testing whether a given value (hypothesized value) is a reasonable value of a population parameter.
- ② Fundamental steps:
 - Compare data from a sample to a hypothesized parameter.
 - Compute the probability that a population with the specified parameter would produce a sample statistic as extreme or more extreme to the one we observed in our sample.
 - The probability value computed in the above step is the p-value.
- Goal: to determine whether there is enough evidence to support a particular claim.

Hypothesis Testing: Null and Alternate Hypothesis

- The null hypothesis (\mathbf{H}_0) and the alternative hypothesis (\mathbf{H}_1) are two competing statements used to make inferences about a population parameter based on sample data.
- The above statements play a central role in the process of hypothesis testing to determine whether there is enough evidence to support a particular claim.

Hypothesis Testing: Null Hypothesis

- Represents the status quo, a default assumption, or a lack of an expected effect.
- In other words, the null hypothesis is a statement of no effect
- Mathematically, null hypothesis often includes equality statements; for example, population mean = 25.5



Hypothesis Testing: Alternate Hypothesis

- Contradicts the null hypothesis.
- It represents what we are trying to find evidence for a claim of an effect, a difference, or a change in the population parameter.
- It can be one-sided (greater than or less than) or two-sided (not equal to) depending on the research question we are trying to answer.

Hypothesis Testing: Type I Error

Type I Error (False Positive)

- Occurs when we reject the null hypothesis when it is actually true; in essence, we conclude that there is an effect or difference when there isn't one in the population.
- ② The probability of making a Type I error is denoted by the significance level (α) .
- Example: convicting an innocent person in a court trial.

Hypothesis Testing: Type II Error

Type II Error (False Negative)

- Occurs when we fail to reject the null hypothesis when it is actually false; in essence, we conclude that there is no effect or difference when there is one in the population.
- ② The probability of making a Type II error is denoted by β ; power given by $(1-\beta)$
- Second Example: Not diagnosing a disease in a person who actually has the disease.

Hypothesis Testing: Type II Error

Type II Error (False Negative)

- Occurs when we fail to reject the null hypothesis when it is actually false; in essence, we conclude that there is no effect or difference when there is one in the population.
- ② The probability of making a Type II error is denoted by β ; power given by $(1-\beta)$
- Example: Not diagnosing a disease in a person who actually has the disease.

In practice, there's often a trade-off between Type I and Type II errors. Adjusting the significance level (α) can impact the balance between these errors.



Hypothesis Testing: Significance level

- **Significance level** (α): a predetermined threshold used in hypothesis testing to determine the level of evidence required to reject the \mathbf{H}_0 .
 - A lower significance level (e.g., 0.01) requires stronger evidence to reject \mathbf{H}_0 ; reduces the risk of falsely rejecting a true null hypothesis.
 - A higher significance level (e.g., 0.10) requires less evidence to reject the null hypothesis but it also increases the risk of making a Type I error.
- 2 Interpretation: If the result is statistically significant, it means that the observed data is unlikely to have occurred under the assumption that the null hypothesis is true.

Hypothesis Testing: Test Statistic

- Test statistic: is a numerical value calculated from sample data that is used to assess the strength of evidence against the null hypothesis.
- ② In other words, it quantifies the difference between the sample data and the null hypothesis's expected values.
- **3** The test statistic's behavior helps determine whether the observed data provides enough evidence to either reject or fail to reject \mathbf{H}_0 .
- **①** Choice of test statistic depends on the data itself; common examples include **Z** score, **T** score and χ^2 -square test.

Hypothesis Testing: Steps Involved

- f 0 State $f H_0$ and $f H_1$
- 2 Choose a significance level α
- Collect the data and then compute the test statistic
- Calculate the p-value
 - A p-value smaller than the significance level (α) suggests stronger evidence against the null hypothesis; the null hypothesis is rejected in favor of the alternative hypothesis.
 - Alternatively, a p-value greater than α means that the null hypothesis is not rejected; there is not enough evidence to support the claim made by the alternative hypothesis.

Hypothesis Testing: One Sample Test

One Sample T Test

- 1 AKA single sample t-test.
- ② Is a statistical hypothesis test used to determine whether the mean calculated from sample data collected from a single group is different from a designated value.
- The designated value does not come from the sample itself.
- The different hypothesis are:
 - ullet $oldsymbol{H}_0$: The population mean equals the specified mean value
 - ullet \mathbf{H}_1 : The population mean is different from the specified mean value
- **5** Test determines if there is enough evidence to reject \mathbf{H}_0 in favor of \mathbf{H}_1 .

Hypothesis Testing: One Sample Test

One Sample T Test: Use Case

- Assumptions:
 - Sample size must be large $N \ge 30$.
 - Population must be known to be normally distributed.
 - \bullet α is predetermined.
- 4 Hypothesis:
 - \mathbf{H}_0 : $\mu = \mu_0$
 - $\mathbf{H}_1 : \mu \neq \mu_0$
- Ompute Test Statistic: one group mean

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

(7)

where s is the sample standard deviation n is the sample size.

Hypothesis Testing: One Sample Test

One Sample T Test: Use Case

- **Outpute** p value: use a t-distribution to find a p value.
- Oecision Making:
 - If $p \leq \alpha$, reject the null hypothesis \mathbf{H}_0 .
 - If $p > \alpha$, fail to reject the null hypothesis.
- 3 State the real world research finding based on the above step.

Hypothesis Testing: Two Sample Test for Mean

Two Sample Test:

- Determines if two samples correspond to the same population mean.
- Requires some extra steps:
 - If samples are independent.
 - Inherent variability in each sample.



Hypothesis Testing: Two Sample Test for Mean

Independent and Dependent Samples

- Independent samples: if the samples selected from one of the populations have no relationship with the samples selected from the other population.
- ② Dependent samples: if each measurement in one sample is matched or paired with a particular measurement in the other sample

Hypothesis Testing: Two Sample Test for Mean

- Assumption: samples are reasonable sized and population follows a normal distribution.
 - If $\mu_1 \mu_2 = 0$: then there is no difference between two population means.
 - If $\mu_1 \mu_2 \neq 0$: then the two population means are different.



Hypothesis Testing: Two Sample Test for Mean

Computation of t-statistics

● The t—statistic is computed as:

$$t^* = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_{x_1}^2}{n_{x_1}} + \frac{s_{x_2}^2}{n_{x_2}}}}$$
 (8)

where $s_{x_1}^2$ and $s_{x_2}^2$ are the sample variances and n_{x_1} and n_{x_2} are the corresponding sample sizes.

Regression

Simple Linear Regression

- Statistical method that allows us to *summarize* and *study relationships* between predictor *x* and response **y**.
- Simple: only one predictor variable is considered; multiple linear regression also exists.
- The relationship we study is statistical and not deterministic.
- Visualization tool: scatter plot is suited.



Simple Linear Regression

Best Fit Curve

- 1 To determine the best fit curve, the following notations are useful:
 - y_i: observed response
 - \hat{y}_i : predicted response
 - x_i : predictor value
- 2 With the above notation in place, the best fit line may be written as:

$$\left[\hat{y}_i = \beta_0 + \beta_1 x_i\right] \tag{9}$$

Olearly, the predictions made may not be correct, so we define the residual error as:

$$e_i = y_i - \hat{y}_i \tag{10}$$



Simple Linear Regression

Best Fit Curve

- **1** The best-fit line should have minimum residual error $e_i \forall i = 1, 2, ...N$.
- In other words, the best-fit line minimizes the sum of total squared residuals given by:

$$Q(\beta_0, \beta_1) = \sum_{i=1}^{i=N} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{i=N} (y_i - (\beta_0 + \beta_1 x_i))^2$$
(11)

Method of Least Squares: Squared differences are considered so that they don't cancel the effect of one another.



Simple Linear Regression

Best Fit Curve

- **1** In the determination of the best-fit line, the parameters β_0 and β_1 need to be estimated.
- ② Using differential calculus, the estimates may be obtained by setting:

$$\left| \frac{\partial Q(\beta_0, \beta_1)}{\partial \beta_0} = 0 \right| \tag{12}$$

and

$$\frac{\partial Q(\beta_0, \beta_1)}{\partial \beta_1} = 0 \tag{13}$$

Simple Linear Regression

Best Fit Curve

1 By solving the above partial differential equations, the optimal estimates β_0^* and β_1^* are obtained as:

$$\beta_0^* = \bar{y} - \beta_1^* \bar{x} \tag{14}$$

where

$$\beta_1^* = \frac{\sum_{i=1}^{i=N} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{i=N} (x_i - \bar{x})^2}$$
(15)

② The line represented by the coefficients, β_0^* and β_1^* is **least squares** regression line.



Analysis of Variance

- Also known as ANOVA.
- It helps understand whether the differences observed in sample means are likely due to actual differences in the populations they represent or if they could have occurred by random chance.
- Categorization:
 - One-way: for comparing means of groups.
 - Two-way: for examining the effect of two independent variables on a dependent variable.

Analysis of Variance

Hypothesis

- **1** $\mathbf{H}_0: \mu_1 = \mu_2 = ... = \mu_K$
- **2** $\mathbf{H}_1: \mu_i \neq \mu_j$ for some i and j where $i \neq j$
- Interpretation:
 - Null hypothesis suggests that the factor did not have any significant impact on the results obtained; in other words, the different groups may be considered to be the part of the same population.
 - Alternate hypothesis states that the means of at least one of the pairs is not equal.



Analysis of Variance

Test Statistics: F-Ratio

• For more than 2 groups, we use test statistic F- ratio defined as the ratio of between-group sample variance and the within-group-sample variance

$$F = \frac{\text{between sample group variance}}{\text{within sample group variance}}$$
 (16)

- Interpretation:
 - Measures whether the means of different samples are significantly different.
 - *F*—ratio higher than the *F*—critical value suggests evidence against the null hypothesis.



Analysis of Variance

Between-Group Sample Variability

- Computation similar to standard deviation.
- Sample deviations are weighted by sample sizes.
- Between sample variability is computed as:

$$MS_{between} = \left(\frac{n_1(\bar{x}_1 - \bar{x}_G)^2 + n_2(\bar{x}_2 - \bar{x}_G)^2 + \dots + n_K(\bar{x}_K - \bar{x}_G)^2}{K - 1}\right) \tag{17}$$

where $\bar{x_G}$ is the grand mean.



Analysis of Variance

Within-Group Sample Variability

- No interactions between the samples.
- Measured by looking at how much each value in the sample differs from its respective sample mean
- Within-group sample variability is computed as:

$$MS_{within} = \left(\frac{\sum (x_{ij} - \bar{x}_j)^2}{N - K}\right) \tag{18}$$

where j corresponds to a group.

• Conclude given F - ratio computed and F-critical value.



Analysis of Variance

Limitation

- Tells us that at least two groups are different.
- Output
 However, it does not tell us which groups are different.

Relation with t - test

- With only two samples, t-test and ANOVA give same result;

Introduction

- AKA Chi-Square Test of independence.
- Provides statistical evidence of an association or relationship between the two categorical variables.
- On not confuse it with correlation. Why?



Introduction

Assumptions:

- The sample must be been obtained *randomly*.
- Variables must be *categorical*: nominal or ordinal.

② Data Representation:

- Contingency table (also known as cross-tab).
- Tabulates the relationship between two categorical variables.

Contingency Table

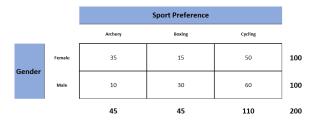


Figure: Contingency Table

Contingency Table

Steps:

- Hypothesis:
 - H_0 : the two categorical variables are independent.
 - **H**₁: the two categorical variables are not independent.
- 2 Fix α and hence, χ^2_{α} .
- Now, given the observed counts in the contingency table, compute the expected counts under H₀.

$$E = \frac{(row\ total)(column\ total)}{total\ sample\ size}$$
(19)



Contingency Table

Steps:

- What question is being answered?
 - Are the observed counts so different from the expected counts that we can conclude a relationship exists between the two variables?
- ② Compute χ^2 statistic: in a summary table of $r \times c$ entries:

$$\chi^{2*} = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} + \dots + \frac{(O_{rc} - E_{rc})^2}{E_{rc}}$$
(20)

3 Conclude by comparing χ^{2*} and χ^2_{α} .

