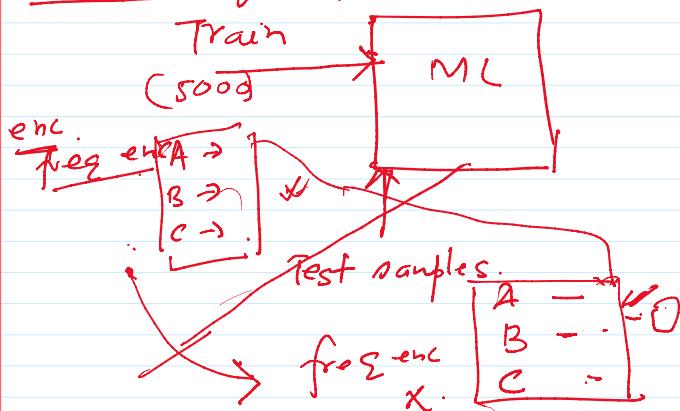


- Encoding - best practices
- Scaling
- Imbalanced data
- Data splitters
- Adv statistics

Data leakage



Next topic : Scaling

Scaling

Definition

- Scaling involves transforming the values of a variable to a specific range, making it easier to compare and interpret.
- standardize the range of independent variables or features of the data.

Methods:

- Min-Max Scaling:** Scales the values between 0 and 1.
- Standardization (Z-score Scaling):** Scales the values to have a mean of 0 and a standard deviation of 1.
- Robust Scaling:** Scales the values based on the median and interquartile range.

Min-Max Scaling

- Min-Max Scaling, also known as Min-Max Normalization, is a technique used to transform the values of a numerical variable to a specific range, typically between 0 and 1.

$$X_{\text{scaled}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- X is the original value of the variable.
- X_{\min} is the minimum value of the variable in the dataset.
- X_{\max} is the maximum value of the variable in the dataset.

Robust scaler

Formula

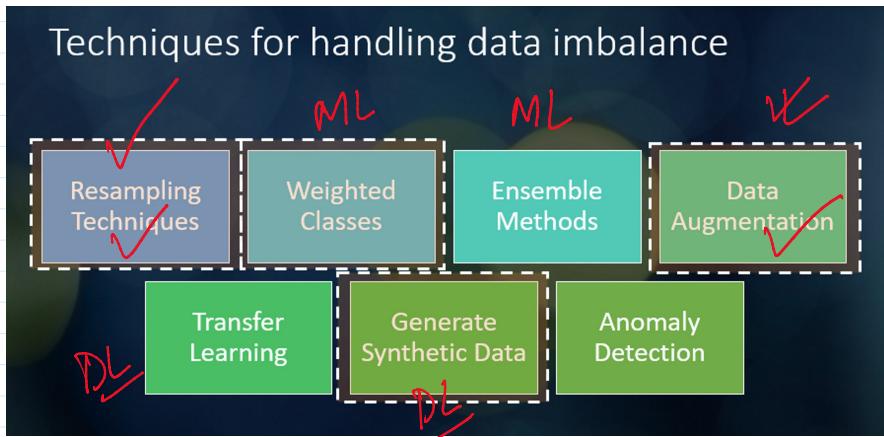
$$X_{\text{robust-scaled}} = \frac{X - \text{median}(X)}{\text{IQR}(X)}$$

where:

- X is the original value of the variable.
- $\text{median}(X)$ is the median of the variable.
- $\text{IQR}(X)$ is the interquartile range of the variable.

Next topic : Imbalanced datasets

Define data imbalance



why image aug??
SMOTE ??
↑ 6 mins!
(self-walkthru).

Step 1 - Identify Minority Class Instances

This step involves determining which class in the dataset has fewer instances.

In the context of SMOTE, this is typically the class you want to oversample, referred to as the minority class.

Step 2 and 3

Select a Minority Instance:

A single instance is randomly chosen from the minority class.

Find k Nearest Neighbors:

The k -nearest neighbors of the selected instance within the minority class are identified. The choice of k is a parameter set by the user, defining the number of neighbors to consider.

Step 4

Generate Synthetic Instances:

- For each neighbor, a synthetic instance is created by calculating the vector between the selected instance and the neighbor.
- This vector is scaled by a random value between 0 and 1.
- The scaled vector is added to the selected instance to generate a new synthetic instance.



Illustrating step 4

- Suppose we have a minority class instance A with two features:
 - Instance A: $(x_A, y_A) = (3, 5)$
- Let's say we choose $k=2$, and the two nearest neighbors of instance A in the minority class are B and C:
 - Neighbor B: $(x_B, y_B) = (4, 4)$
 - Neighbor C: $(x_C, y_C) = (2, 6)$
- Now, we'll generate synthetic instances for instance A based on its neighbors.

Generate synthetic data

1. Calculate Vectors:

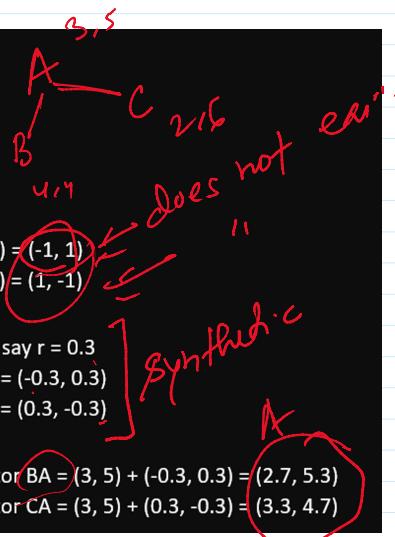
- Vector BA = $(x_A - x_B, y_A - y_B) = (3 - 4, 5 - 4) = (-1, 1)$
- Vector CA = $(x_A - x_C, y_A - y_C) = (3 - 2, 5 - 6) = (1, -1)$

2. Scale Vectors:

- Choose a random value between 0 and 1, let's say $r = 0.3$
- Scaled Vector BA = $r * \text{Vector BA} = 0.3 * (-1, 1) = (-0.3, 0.3)$
- Scaled Vector CA = $r * \text{Vector CA} = 0.3 * (1, -1) = (0.3, -0.3)$

3. Generate Synthetic Instances:

- Synthetic Instance 1 = Instance A + Scaled Vector BA = $(3, 5) + (-0.3, 0.3) = (2.7, 5.3)$
- Synthetic Instance 2 = Instance A + Scaled Vector CA = $(3, 5) + (0.3, -0.3) = (3.3, 4.7)$



Data splitters

Various types of datasets

1. Training Set:

- Purpose:** Used to train the machine learning model.
- Size:** Typically the largest subset (60-80% of the data).
- Usage:** The model learns patterns, relationships, and features from this set.

?? dev

2. Validation Set:

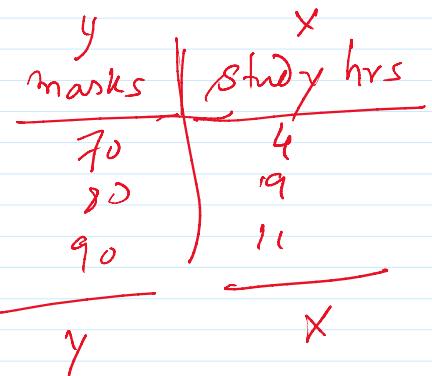
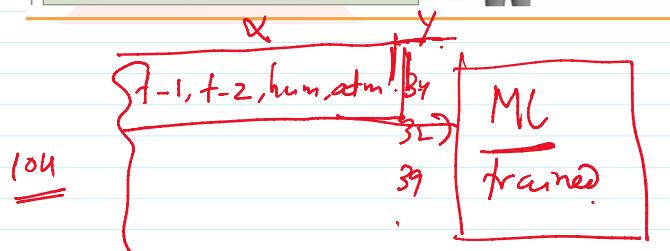
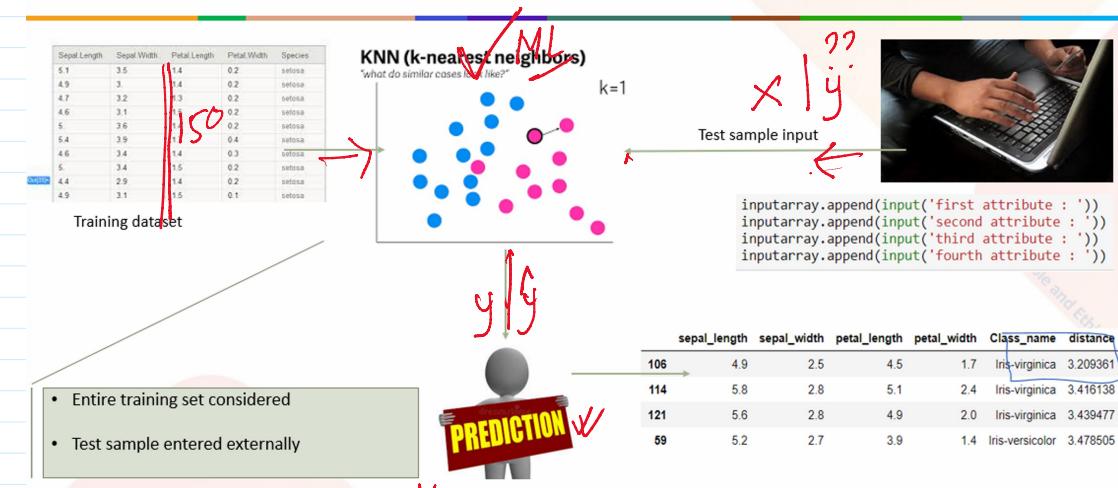
- Purpose:** Used for hyperparameter tuning and model selection during training.
- Size:** Smaller than the training set (usually 10-20% of the data).
- Usage:** Helps prevent overfitting by fine-tuning model parameters without contaminating the test set.

3. Testing Set (or Test Set):

- Purpose:** Reserved for evaluating the model's performance on unseen data.
- Size:** Independent subset not used during training or validation (10-20% of the data).
- Usage:** Provides an unbiased assessment of how well the model generalizes to new, unseen instances.

✓... often used at end of dev
... more test sets.

EVALUATION – 1 (RANDOM INPUT TEST SAMPLE)



given $x \rightarrow y$
 $y = f(x)$.

$y = f(\text{adver} + \text{emp skill})$ ML

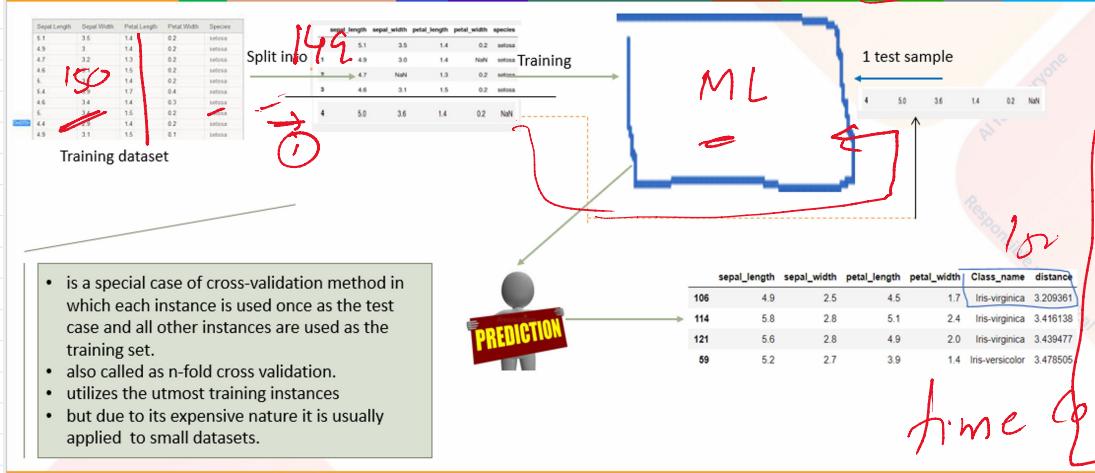
$y = f(x) \rightarrow$ learn

$y = b_1 x_1$

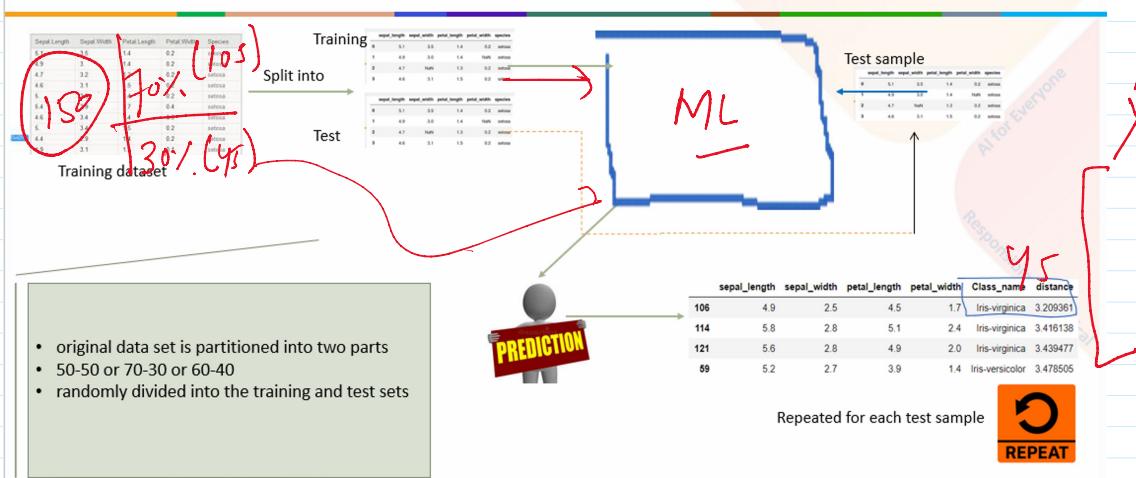
$b = (x^T x)^{-1} (x^T y)$

linear fn.

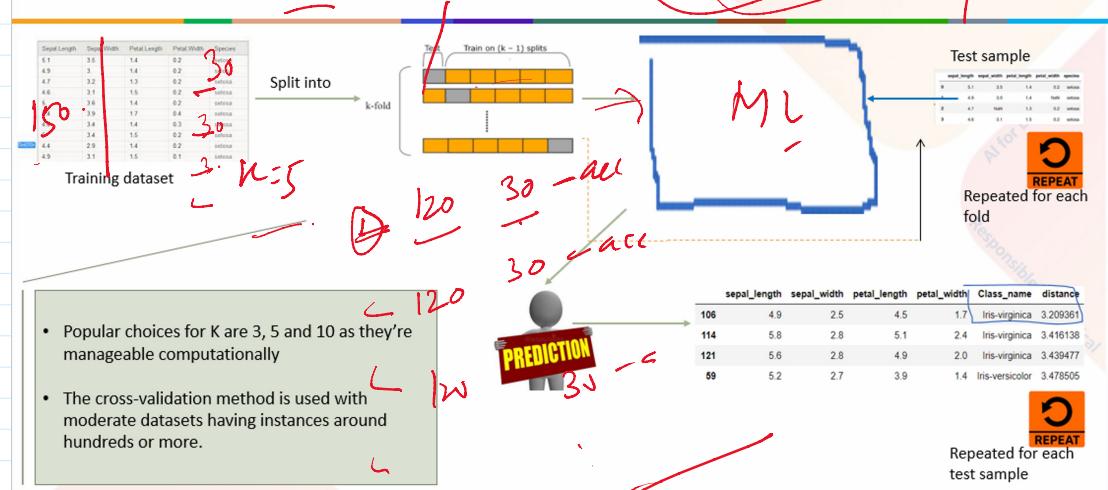
EVALUATION – 2 (LEAVE-ONE-OUT CROSS-VALIDATION - LOOCV)



EVALUATION – 3 (HOLDOUT METHOD)



EVALUATION – 4 (K-FOLD CROSS-VALIDATION METHOD)



Comprehensive? ?

LOOCV → very sparingly; cluster of machines.
 holdout method → very often; dev cycle.
 K-fold → at end of major milestones.
 ↘ (Variations).
 ↘ 3/4

Next topic : Inferential statistics

Define inferential statistics



is a branch of statistics that involves making inferences or predictions about a population based on a sample of data.



allows business leaders and managers to draw conclusions and make decisions about a larger group based on a smaller subset of data.

Why learn inferential statistics?

Generalization from Samples:

Example: A pharmaceutical company tests a new drug on a sample of 1,000 patients to determine its effectiveness. Using inferential statistics, they can generalize the findings to the entire population of potential patients, estimating how effective the drug would be if given to millions of people.

Understanding Uncertainty:

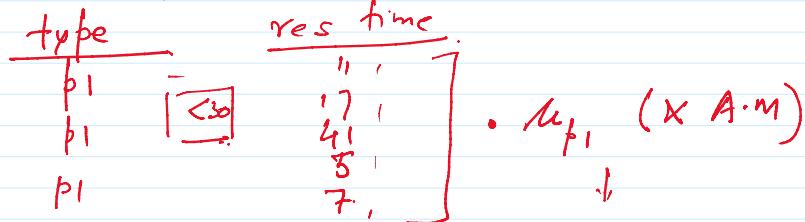
Example: A polling organization surveys 1,200 voters to predict the outcome of an election. They use inferential statistics to calculate a confidence interval, saying that Candidate A is likely to receive 48% to 52% of the vote with 95% confidence. This helps them understand the potential range of the true voter percentage.

Hypothesis testing

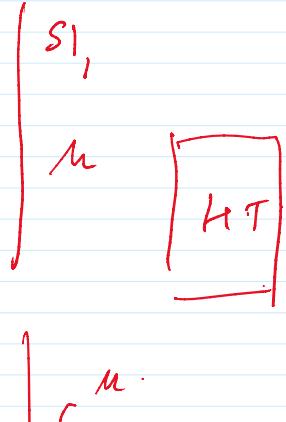
is a statistical method used to make inferences about a population based on a sample of data.

involves formulating a hypothesis about a population parameter, collecting data, and using statistical techniques to determine whether the data provide enough evidence to reject or not reject the null hypothesis.

Ticket/complaint mgmt.



similair? ?



$$\begin{bmatrix} p_2 \\ \vdots \\ p_r \end{bmatrix} = \boxed{\text{exp}} \left[\begin{bmatrix} 77 \\ 41 \\ 11 \\ 2 \\ 81 \end{bmatrix} \right] \cdot \mu_{12} (\times A \cdot m) \cdot \begin{bmatrix} \mu \\ \zeta_1 = \end{bmatrix}$$

CA. housing.

$$\frac{\# \text{ rooms}}{\mu_1 = \boxed{\square}} \sim \text{simla}$$

$$\mu_2 = \boxed{\square}$$

$$\vdots$$

$$\mu_3 = \boxed{\square}$$

ML (discriminatory nature).

- \rightarrow ML (large amt; wise; -----)
- \rightarrow statistical (inferential stats).

(limited to sample). -

ocean-proximity

- { island }
- { near shore }
- { < 1H }

TYPES OF INFERRENTIAL STATISTICS

HYPOTHESIS TESTING

- serves to scrutinize assumptions and make inferences about the entire population using the data at hand.
- entails establishing both a **null hypothesis** and an alternative hypothesis,
- Test using statistical analysis for significance.
- evaluating the test statistic against **critical values** and confidence intervals, ultimately leading to a conclusion.

T-test | F-test | Z-test | ANOVA |

REGRESSION ANALYSIS

- Regression analysis is used to quantify how one variable will change with respect to another variable.
- There are many types of regressions available such as
- simple linear, multiple linear, nominal, logistic, and ordinal regression.

(ML).

Hypothesis (notion, beliefs, accepted norm).

Null Hypothesis (H_0):

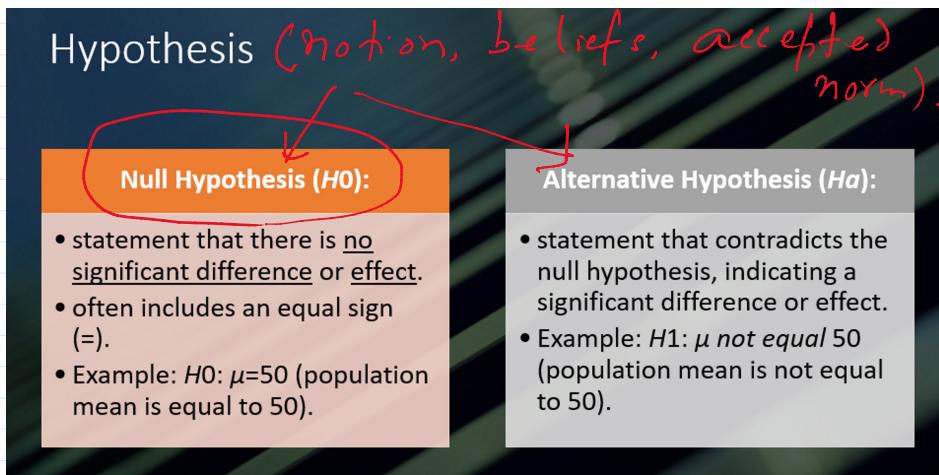
- statement that there is no significant difference or effect.
- often includes an equal sign

Alternative Hypothesis (H_a):

- statement that contradicts the null hypothesis, indicating a significant difference or effect

$\mu_{p_1} = 20 \text{ min.}$

null hyp: H_0 .



$\mu_{p1} = 20 \text{ min}$
 null hyp; H_0

$\mu_{p1} \neq 20$
 > 20 < 20 H_a

Examples – Market research

- Suppose we have a sample of 100 consumers who participated in the market research study.
- categorize their age groups into three categories: Young (18-30), Middle-aged (31-50), and Older (51 and above).
- For screen size preference, we'll consider two categories: Small and Large.
- TASK** : the relationship between screen size preference and age group

Consumer ID	Age Group	Screen Size Preference
1	Young	Large
2	Middle-aged	Small
3	Older	Large
4	Young	Large
5	Middle-aged	Small
6	Middle-aged	Large
7	Older	Small
8	Young	Large
9	Young	Large
10	Older	Small
...
100	Middle-aged	Large

Example - Employee Productivity

- A business leader wants to determine if there is a significant difference in productivity between two teams within the organization.
- use a hypothesis test to compare the average productivity scores of Team A and Team B based on specific metrics (e.g., sales volume, project completion time).
- test would help the leader infer whether the observed difference in productivity between the two teams is statistically significant.

Team A Sales Volume: [100, 120, 110, 90, 105, 115, 95, 105, 115, 100]

Team B Sales Volume: [110, 105, 115, 100, 125, 115, 120, 110, 115, 105]

The customer satisfaction scores **before training** are represented by the list [80, 75, 85, 70, 75, 78, 82, 79, 80, 77].

The customer satisfaction scores **after training** are represented by the list [85, 82, 88, 75, 80, 84, 87, 86, 85, 82].

- A business leader invests in a training program for customer service representatives
- To assess the effectiveness of the training program, the leader could conduct a hypothesis test by comparing the customer satisfaction scores before and after the training.
- test would help infer whether there is a significant improvement in customer satisfaction as a result of the training program

Example - Training Effectiveness

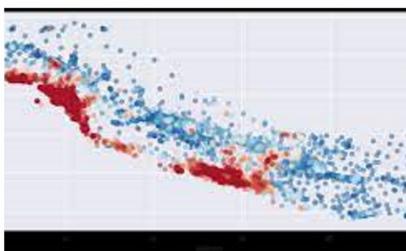
Example: Medical Treatment

- **Null Hypothesis (H_0):** The new drug has no effect on patient recovery.
- **Alternative Hypothesis (H_1):** The new drug improves patient recovery.
- **Interpretation:** Evaluating whether a medical intervention has a significant impact on patient outcomes.

Example: Educational Intervention

- **Null Hypothesis (H_0):** The teaching method has no impact on student performance.
- **Alternative Hypothesis (H_1):** The teaching method improves student performance.
- **Interpretation:** Investigating the effectiveness of a particular teaching approach on student learning.

Example

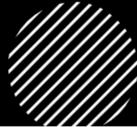


18-08-2024

© Copyright - GROKKERS - Prepared by Nitin

- + • **MedInc (Median Income):** Median income of households within a district.
- o • **HouseAge (Housing Age):** Median age of houses within a district.
- o • **AveRooms (Average Rooms):** Average number of rooms per household within a district.
- o • **AveBedrms (Average Bedrooms):** Average number of bedrooms per household within a district.
- o • **Population:** Total population of the district.
- o • **AveOccup (Average Occupancy):** Average household occupancy within a district.
- o • **Latitude:** Latitude coordinate of the district's location.
- o • **Longitude:** Longitude coordinate of the district's location.
- l • **MedHouseVal (Median House Value):** **target/y**
l Median house value for households within a district.

Example – test



01

Hypothesis Test -
Age of Houses and
House Value:

02

Null Hypothesis (H0):
There is no significant correlation between the age of houses and median house values in California districts.

03

Alternative Hypothesis (H1): Older houses have significantly lower median values compared to newer houses.



Example

- Credit Scoring in Finance:
 - **Hypothesis:** A new credit scoring model is more accurate in predicting loan default.
 - **Null Hypothesis (H0):** The new credit scoring model is not more accurate.
 - **Alternative Hypothesis (H1):** The new credit scoring model is more accurate.
 - **Test:** Collect data on loans, apply both the old and new credit scoring models, and compare their accuracy in predicting loan defaults.

Process



Objective:

Evaluate the importance of each feature **independently** in relation to the target variable.



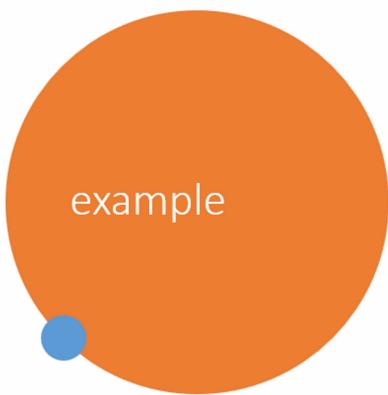
Statistical Measures:

commonly involves statistical tests or scoring methods to quantify the significance or relevance of each feature. Examples include **t-tests**, **ANOVA** F-tests, **mutual information**, and **chi-square tests**

Limitations

Does Not Consider **Interactions** or Relationships Between Features:

- **Explanation:** Univariate methods evaluate features in isolation, neglecting potential synergies or dependencies between features that may contribute jointly to predictive power.
- **Mitigation:** Consider employing multivariate feature selection methods that capture relationships between features.



- Let's say we want to test whether the average height of a certain population is **different** from 65 inches.

hypothesis	
Null Hypothesis (H_0)	Alternative Hypothesis (H_1):
<ul style="list-style-type: none">• $\mu = 65$ (population mean height is equal to 65 inches).	<ul style="list-style-type: none">• $\mu \neq 65$ (population mean height is not equal to 65 inches).

$\neq 65$ ①
 > 65 ②
 < 65 ③

Null vs alternate



Null Hypothesis:

Think of the null hypothesis as a statement that nothing special is happening. It's like saying "there's no change" or "there's no effect." It's a way of being cautious and assuming that any observed differences are just due to random chance.



Alternative Hypothesis:

The alternative hypothesis is like the opposite of the null hypothesis. It's where you suggest that something interesting or different is happening. It's like saying "there is a change" or "there is an effect."

gender

pol. part

my

F

I

A

B

C

A

I

$H_0 = \text{No relation}$

$H_a = \text{gender impact for } \rightarrow (M \rightarrow A, F \rightarrow B)$

Simply stating



In layman words

Null Hypothesis (H_0): "Nothing unusual is going on."

Alternative Hypothesis (H_1): "Something interesting or different is happening."



If a new medicine helps people sleep better

Null Hypothesis (H_0): "Taking this medicine doesn't make any difference in sleep quality."

Alternative Hypothesis (H_1): "Taking this medicine does make a difference in sleep quality."

$$\rightarrow M = 5.6 \\ \text{S: } H_0.$$

$$(D) M > 5.6$$

Null Hypothesis (H_0) - examples



E-commerce Conversion Rates: "The recent website design change did not lead to a significant change in the conversion rate."

Ad Click-Through Rates: "The new ad campaign did not result in a higher click-through rate compared to the old campaign."

Customer Churn Analysis: "The introduction of a new loyalty program did not result in a significant reduction in customer churn."

Product Recommendation Engine: "The personalized product recommendations do not lead to a higher average purchase value compared to non-personalized recommendations."

- ① Significance level]
 ② Test statistic .]
 ③ p-value]
 ④ Critical value]
 info related

4 Critical value ↴

#	Topic	Number of Use cases/ Examples discussed	Explained Y/N, comments if any
1	Scaling	Min max scaler Std scaler Robust scaler Partial scaling Best practices	
2	Data imb	-SMOTE (std data) -Image aug (PIL package)	
3	Data splitters	-LOOCV, Leave-Pout -Hold out -K-fold	
4.	Inferential stats	-What they are, -Examples HT	