

$$y = mx + c \quad (\text{general})$$

$$y = b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_0$$

$y \approx x$ variables

function $\{b_0, b_1, b_2, b_3, \dots\}$

objective:- find b_0, b_1, b_2, \dots in such a way that $\sum \text{error (residual)}$ is minimum.

$$m = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

① statistical method.
(closed form of solution)
linear data
small size data.

② open form soln.
Neural nets
deep learning \rightarrow ml method.
(gradient based method)

$$y = b_0 + b_1 x_1$$

ensure error is minimum.

x_1	y (salu)
1000	10
5000	15
3000	12
\vdots	\vdots

statistical approach $\rightarrow b_0, b_1$

$\checkmark \checkmark$
 $\rightarrow b_0 + b_1 x_1$

x	y	y (pred)	residual	residual ²
1000	10	9	1	1
5000	15	6	9	81
3000	12	8	4	16
\vdots	\vdots	\vdots	\vdots	\vdots
14	11	11	8	64
23	18	18	5	25
\vdots	\vdots	\vdots	\vdots	\vdots

$$\sum \text{error} = \frac{1}{n} \sum = \text{SSE}$$

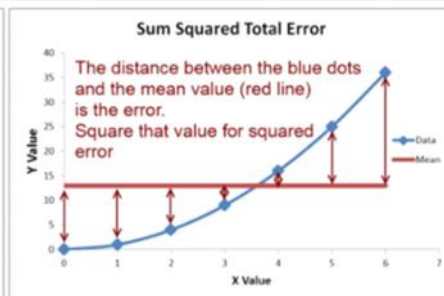
70

R^2 (how good is the model)

$\rightarrow b_0, b_1, \dots$
 $\rightarrow \text{error (min)}$

$\rightarrow 1.1$ 10. model

- b_0, b_1, \dots
- error (min)
- quantify the model performance.

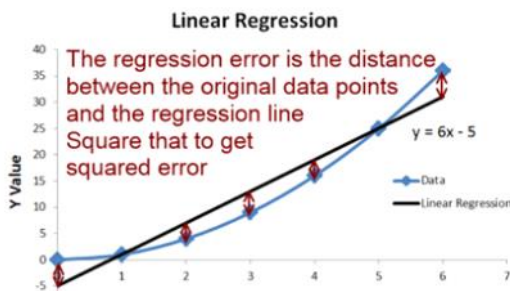


no regression model

w. a ~~reg. model~~

X	y	\hat{y}	residual
10			10
15			5
25	20		5
30			10

$$SSE = 100 \text{ (SST)}$$



← with reg model.
 b_0, b_1, \dots

X	y	\hat{y}	residual
10		12	-
15		13	-
25	20	22	-
30		18	-

$$\sum = SSE = 80$$

$$\frac{SSE}{SST} = \frac{80}{100} = 0.8$$

80% of error which existed with the avg model
→ still exist.

$$1 - \frac{SSE}{SST} = 1 - 0.8 = 0.2$$

↑ R^2

$$R^2 = 0.2$$

11 remained 20% of

$R^2 = 0.2$
 reg model remain 20% of errors with the avg model.
 → (0 — 1) ideally.

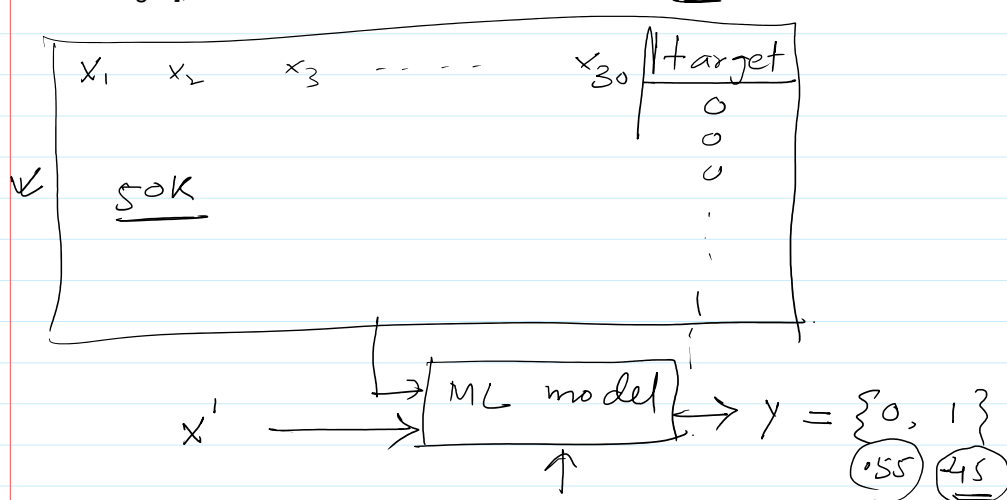
Probability

Medical: Disease Diagnosis Probability

Example: Calculating the probability of a patient having a certain disease given their symptoms and test results.

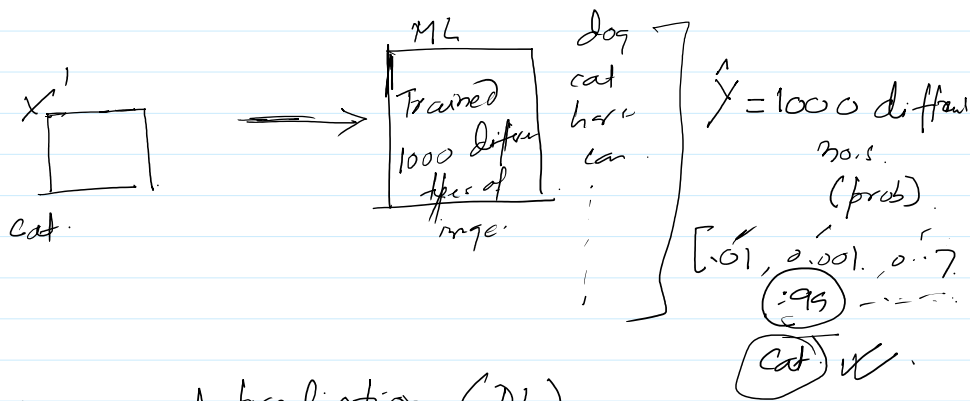
Method: Bayes' theorem is commonly used in medical diagnosis to calculate probabilities. This theorem incorporates prior probabilities (such as the prevalence of the disease in the population) and conditional probabilities (such as the likelihood of certain symptoms given the presence or absence of the disease) to determine the probability of the disease given observed evidence.

Index(['mean radius', 'mean texture', 'mean perimeter', 'mean area', 'mean smoothness', 'mean compactness', 'mean concavity', 'mean concave points', 'mean symmetry', 'mean fractal dimension', 'radius error', 'texture error', 'perimeter error', 'area error', 'smoothness error', 'compactness error', 'concavity error', 'concave points error', 'symmetry error', 'fractal dimension error', 'worst radius', 'worst texture', 'worst perimeter', 'worst area', 'worst smoothness', 'worst compactness', 'worst concavity', 'worst concave points', 'worst symmetry', 'worst fractal dimension', 'target'],

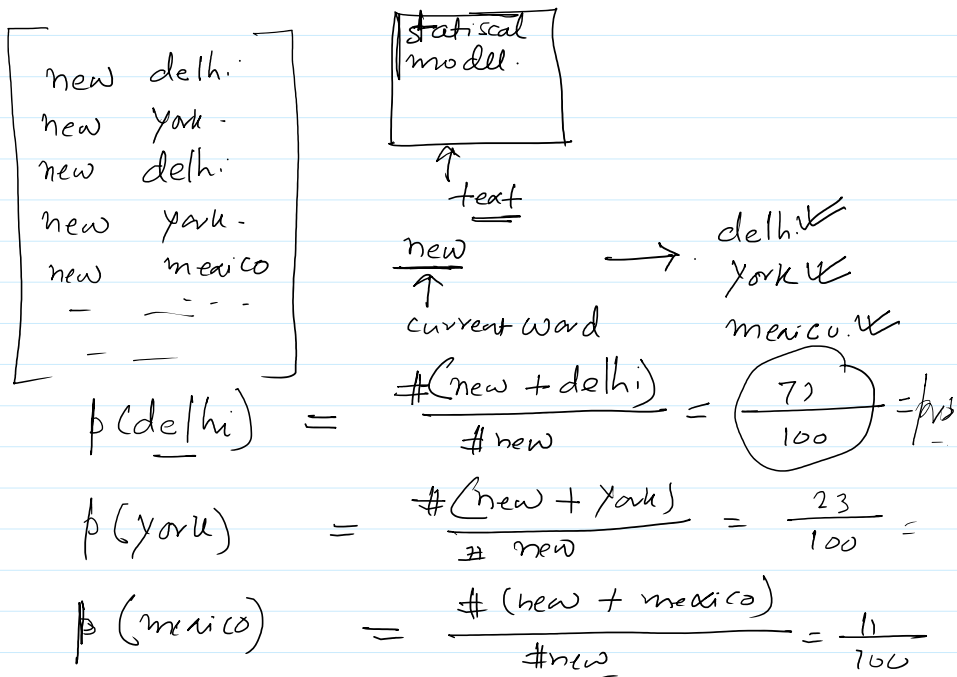


Example: Predicting the probability of a stock's price increasing or decreasing within a certain time frame.

Method: Statistical models such as time series analysis, Monte Carlo simulation, or option pricing models (e.g., Black-Scholes model) can be used to calculate probabilities in finance.



next word prediction (DL).



Conditional Probability

- Fundamental from DS/ML angle also
- Prob of some event occurring GIVEN that another event has taken place $P(A|B)$

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

Color	Shape	Ripe
Red	Round	Yes
Green	Round	No
Yellow	Oval	Yes
Red	Oval	No
Green	Round	Yes
Yellow	Oval	No

- **Event A:** Fruit is red.
- **Event B:** Fruit is round.
- **Event C:** Fruit is ripe.

- calculate $P(C \cap A \cap B)$, the joint probability of a fruit being ripe, red, and round.

$$P(C \cap A \cap B) = \frac{\text{Number of occurrences of ripe, red, and round fruits}}{\text{Total number of fruits}}$$

$$P(C \cap A \cap B) = \frac{1}{6}$$

Sample	Class	Long	Sweet
1	Orange	Yes	Yes
2	Orange	Yes	Yes
3	Banana	Yes	No
4	Orange	No	Yes
5	Orange	Yes	Yes
6	Banana	Yes	Yes
7	Banana	No	Yes
8	Banana	No	No
9	Orange	No	Yes
10	Banana	No	Yes
11	Banana	Yes	No
12	Orange	Yes	Yes
13	Banana	Yes	No
14	Banana	No	No
15	Orange	No	No
16	Orange	Yes	Yes
17	Orange	Yes	Yes
18	Orange	No	Yes
19	Banana	No	No
20	Orange	No	Yes



$P(\text{Banana}) =$
 $(\text{Count}(\text{Banana}) + 1) / (\text{Total Samples} + \text{Number of Classes})$
 $= (10 + 1) / (20 + 2) = 11/22$



$P(\text{Orange}) =$
 $(\text{Count}(\text{Orange}) + 1) / (\text{Total Samples} + \text{Number of Classes})$
 $= (10 + 1) / (20 + 2) = 11/22$

Prior probabilities of features - For 'Long'


$$P(\text{Long=Yes} \mid \text{Banana}) = (\text{Count}(\text{Long=Yes, Banana}) + 1) / (\text{Count}(\text{Banana}) + 2) = (6 + 1) / (10 + 2) = 7/12$$

$$P(\text{Long=No} \mid \text{Banana}) = (\text{Count}(\text{Long=No, Banana}) + 1) / (\text{Count}(\text{Banana}) + 2) = (4 + 1) / (10 + 2) = 5/12 \quad (**)$$


$$P(\text{Long=Yes} \mid \text{Orange}) = (\text{Count}(\text{Long=Yes, Orange}) + 1) / (\text{Count}(\text{Orange}) + 2) = (4 + 1) / (10 + 2) = 5/12$$

$$P(\text{Long=No} \mid \text{Orange}) = (\text{Count}(\text{Long=No, Orange}) + 1) / (\text{Count}(\text{Orange}) + 2) = (6 + 1) / (10 + 2) = 7/12 \quad (**)$$


Prior probabilities of features - For 'Sweet'




$$P(\text{Sweet=Yes} \mid \text{Banana}) = (\text{Count}(\text{Sweet=Yes, Banana}) + 1) / (\text{Count}(\text{Banana}) + 2) = (8 + 1) / (10 + 2) = 9/12 \quad (**)$$



$$P(\text{Sweet=Yes} \mid \text{Orange}) = (\text{Count}(\text{Sweet=Yes, Orange}) + 1) / (\text{Count}(\text{Orange}) + 2) = (7 + 1) / (10 + 2) = 8/12 \quad (**)$$



$$P(\text{Sweet=No} \mid \text{Banana}) = (\text{Count}(\text{Sweet=No, Banana}) + 1) / (\text{Count}(\text{Banana}) + 2) = (2 + 1) / (10 + 2) = 3/12$$



$$P(\text{Sweet=No} \mid \text{Orange}) = (\text{Count}(\text{Sweet=No, Orange}) + 1) / (\text{Count}(\text{Orange}) + 2) = (3 + 1) / (10 + 2) = 4/12$$

Take a test fruit, Long = 'No', Sweet = 'Yes'
Calculate the **posterior probabilities** for each class:

- For Banana:
 - $P(\text{Banana} \mid \text{Long=No, Sweet=Yes})$
 - $= P(\text{Long=No} \mid \text{Banana}) * P(\text{Sweet=Yes} \mid \text{Banana}) * P(\text{Banana})$
 - $= (5/12) * (9/12) * (11/22) = 0.142$
- For Orange:
 - $P(\text{Orange} \mid \text{Long=No, Sweet=Yes})$
 - $= P(\text{Long=No} \mid \text{Orange}) * P(\text{Sweet=Yes} \mid \text{Orange}) * P(\text{Orange})$
 - $= (7/12) * (8/12) * (11/22) = 0.212$

Conditioning ??

So, the test fruit is more likely to be classified as an **Orange**.

$$P(F|E) = \frac{P(F \cap E)}{P(E)}$$

Where:

- $P(F \cap E)$ is the joint probability of both events F and E .
- $P(E)$ is the probability of event E .

Bayes theorem

(TBD)

Inferencing

- Making some statement or predictions about a population, based on sample data

Methods for inferring

- Hypothesis testing
- Regression analysis

Null Hypothesis (H_0):	Alternative Hypothesis (H_a):
<ul style="list-style-type: none"> • statement that there is <u>no significant difference</u> or <u>effect</u>. It often includes an equal sign (=). • Example: $H_0: \mu=50$ (population mean is equal to 50). 	<ul style="list-style-type: none"> • statement that contradicts the null hypothesis, indicating a significant difference or effect. • Example: $H_1: \mu \text{ not equal } 50$ (population mean is not equal to 50).

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]`.

Example - Training Effectiveness

- 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

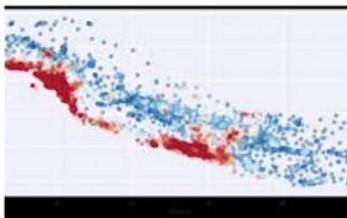
Example: Medical Treatment

- **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



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

Example – test

01

Hypothesis Test -

Age of Houses and
House Value:

02

Null Hypothesis (H_0):

There is no significant correlation between the age of houses and median house values in California districts.

03

Alternative Hypothesis (H_1): Older houses have significantly lower median values compared to newer houses.

Process of Hypothesis testing

1. Objective setting
 - a. Feature
 - i. Independent
 - ii. Dependent
2. Choose the statistical method
 - a. t-test
 - b. ANOVA
 - c. Chi2
 - d. ...
 - e. ...

Test statistic

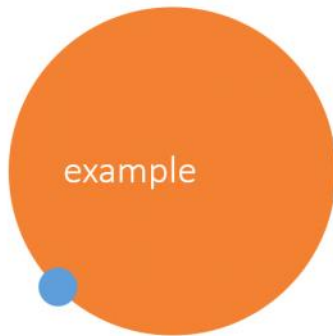
- Number/
- Output from the statistical fn (chosen)
- Used for concluding

p-value

- Number, prob
- Output
- Used for concluding

Critical value

- Threshold
 - Based on the level of significance you have chosen for your study



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

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"> $\mu \text{ not } = 65$ (population mean height is not equal to 65 inches).

e-commerce conversion rate

- H_0 - no change due to changes in website design

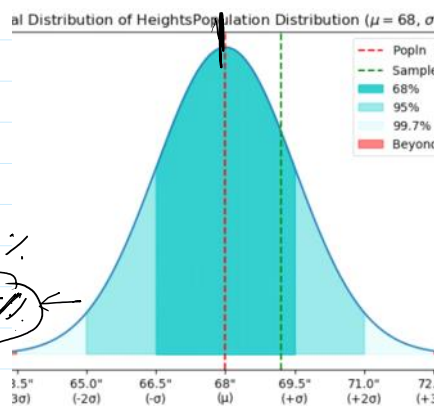
Churn analysis

- Introduced some loyalty program
- H_0 : ...

Set the significance level

- A probability
- Threshold
- Beyond the threshold
 - o Reject the H_0

$$\alpha = 1$$



$$\begin{aligned} \mu &= 65 \\ \pm 1\sigma &= [64 - 66] \rightarrow 68\% \\ \pm 2\sigma &= [63 - 67] \rightarrow 95\% \\ \pm 3\sigma &= [62 - 68] \rightarrow 99.7\% \end{aligned}$$

known to you:
68 - 95 - 99.7

height = 70 inches?

outlier

Common (something wrong with H_0)

$$\alpha = 95\% \mid 99.5\%$$

setting aside some room for wrong prediction.
how much?

Carry out the test

- We collect a sample of 100 people
- Measure heights
- Choose/execute statistical function
- obtain the test statistic (e.g., t-statistic), and
- obtain a p-value of 0.03.

obtain a p-value of 0.03.

- **P-value Interpretation:**

- **Test Statistic and Critical Value:**

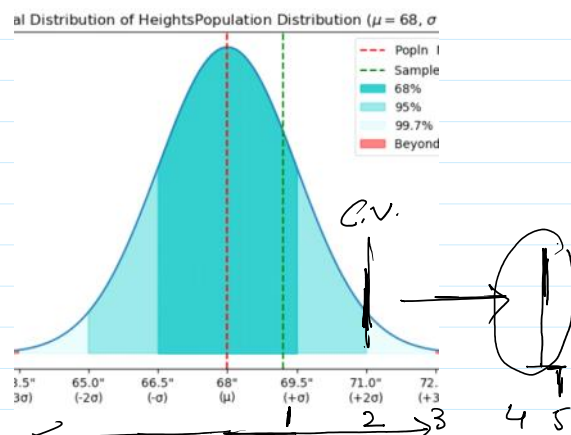
-
- al Distribution of Heights Population Distribution ($\mu = 68, \sigma$)
- Legend:
- Popln
 - Sample
 - 68%
 - 95%
 - 99.7%
 - Beyond
- Handwritten notes:
- reject H_0
 - H₀s, fail to reject H_0
 - rejection region
 - fail to reject H_0
 - fn statistic

Conclude:

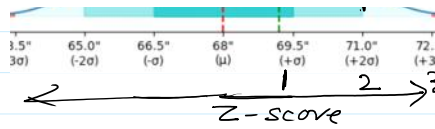
conclude:

- ① if $t\text{-statistic} > \frac{\text{statistic for sis level}}{\text{Critical value}}$
- ② if $p\text{-value} < \frac{0.05}{0.05}$

$\alpha = 95\%$ (5%)
↓
Critical value.



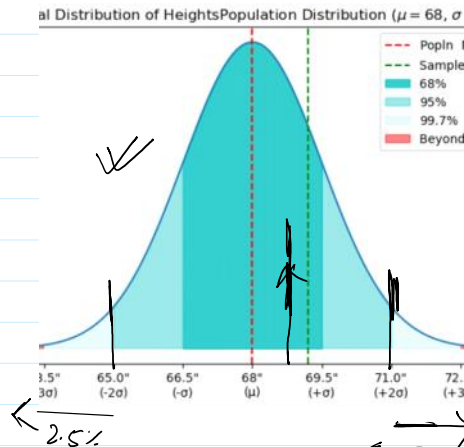
* Critical value.



4 5

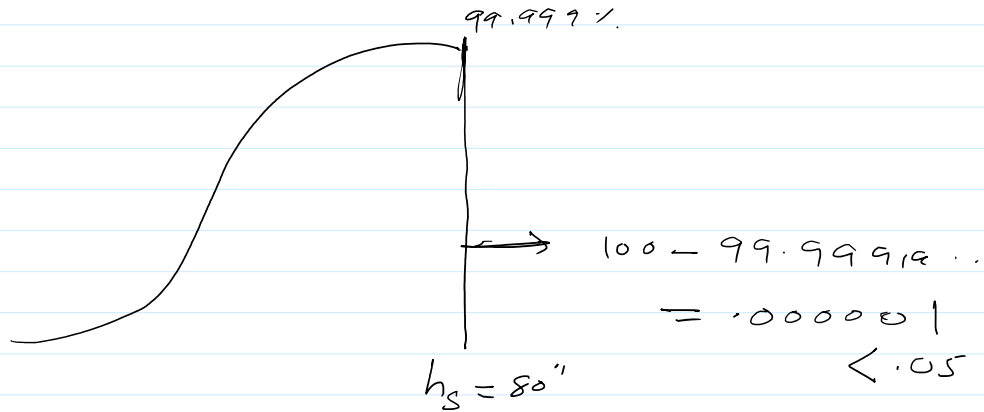
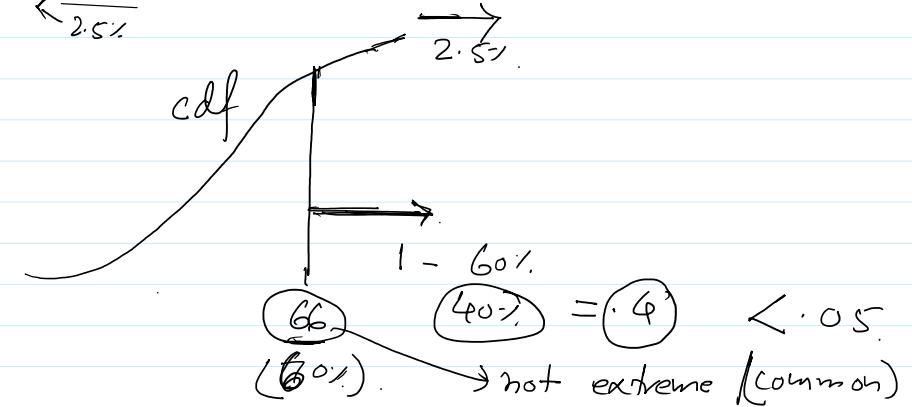
$h_s = 66$
likely to happen

pdf
cdf



$$\boxed{h_s = 80''}$$

$$\boxed{z = 76}$$



t-test

- Check if there are differences in the data

o 2 columns

support team CITES (----)

Complaint types

p_1
 p_1
 p_1
 p_2
 p_2
 p_2

time to resolve

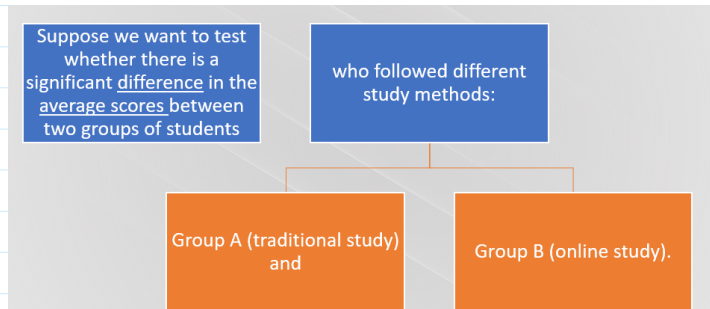
4
6
11
9
100
70
60
2

μ_{p_1}

μ_{p_2}

Same/similar.
 $\mu_1 = \mu_2$

$\left\{ \begin{matrix} p_1 \\ p_2 \\ \vdots \end{matrix} \right\}$
 $\left\{ \begin{matrix} 60 \\ 2 \\ 53 \end{matrix} \right\}$
 μ_{p_2}



Hypothesis

$H_0: \mu_A = \mu_B$
(No significant difference)

$H_1: \mu_A \text{ Not } = \mu_B$
(Significant difference)

test

Data:

- Group A: [78, 85, 92, 88, 76]
- Group B: [82, 90, 88, 79, 84]

Calculate T-Statistic:

- Use the statistical method to calculate the **t-statistic**.

Degrees of Freedom:

- For an independent samples t-test, degrees of freedom = $n_1 + n_2 - 2$.

P-Value:

- Use the **t-statistic** and **degrees of freedom** to find the **p-value**.

Calculate T-Statistic

the formula for calculating the t-statistic in an independent samples t-test is:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

where:

- \bar{X}_1 and \bar{X}_2 are the sample means of Group A and Group B, respectively.
- s_1^2 and s_2^2 are the sample variances of Group A and Group B, respectively.
- n_1 and n_2 are the sample sizes of Group A and Group B, respectively.

$$f(t; df) = \frac{\Gamma\left(\frac{df+1}{2}\right)}{\sqrt{df\pi}\Gamma\left(\frac{df}{2}\right)} \left(1 + \frac{t^2}{df}\right)^{-\frac{df+1}{2}}$$

where:

- t is the value of the random variable (t-statistic).
- df is the degrees of freedom.
- Γ denotes the gamma function.