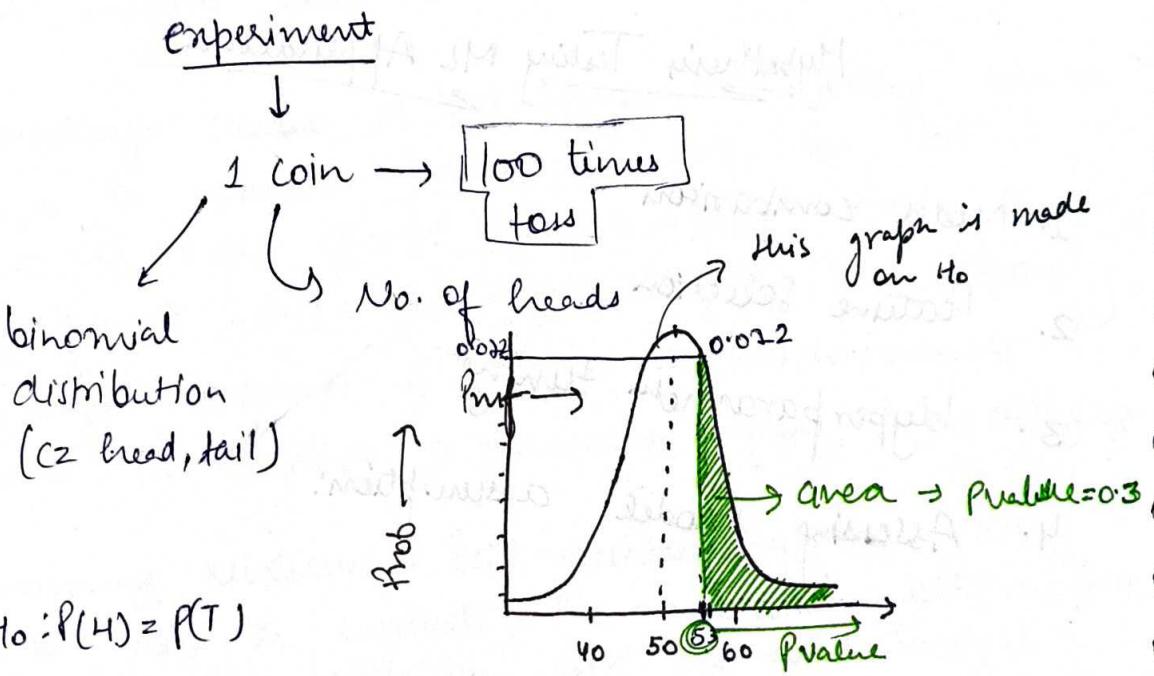


P-value

P-value is the prob of getting a sample as or more extreme (having more evidence against H_0) than our own sample given the Null hypo (H_0) is true.



again did exp → 100 time toss

(53) heads

No. of heads →

if exp → 100 coins → 100 time toss
 $P = 0.5$ then
 30 times → 53 head came.

Interpreting p-value

with significance value

$$\alpha = 0.05 / 0.01$$

$$P\text{-Value} \leq \alpha$$

reject Null Hypo

Without significance value

1. Very small p-value (e.g. $p < 0.01$) indicate strong evidence against the null hypo, suggesting that the observed effect or diff is unlikely to have by chance alone.
2. Small p-value ($0.01 < p \leq 0.05$) indicate moderate evidence against the null hypo, suggesting that the observed effect or diff is less likely to have occurred by chance alone.
3. Large p-value (e.g. $0.05 \leq p \leq 0.1$) indicate weak evidence against the null hypo, suggesting that the observed effect or difference might have occurred by chance alone, but there is still some level of uncertainty.
4. Very large p-value ($p \geq 0.1$) indicate weak or no evidence against the null hypothesis, suggesting that the observed effect or diff is likely to have occurred by chance alone.

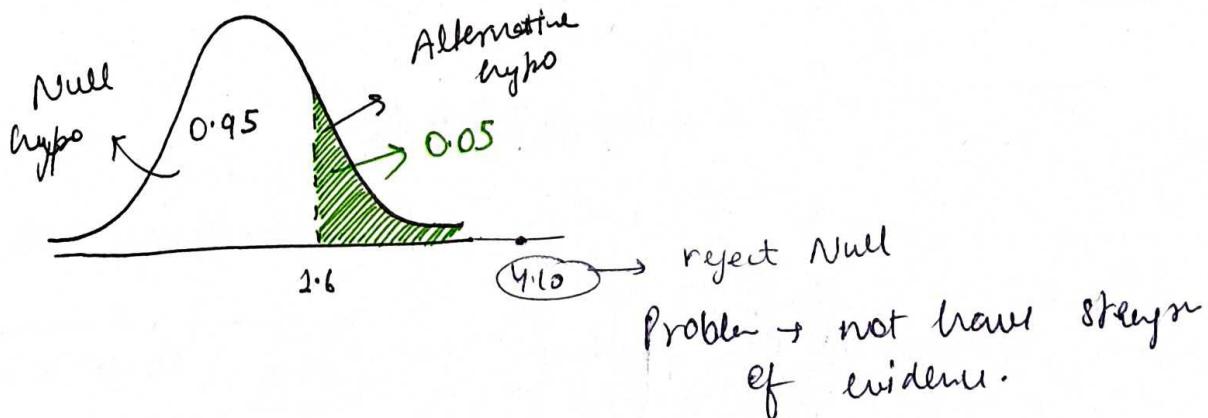
P-value in context of Z-test

Suppose a company is evaluating the impact of a new training program on the productivity of its employees. The company has data on the avg productivity of its employees before implementing the training program. The avg productivity was 50 units per day. After implementing the training program the company measures the productivity of a random sample of 30 employees. The sample has an avg productivity of 53 units per day and the pop std is 4. The company wants to know if the new training program has significantly increased productivity.

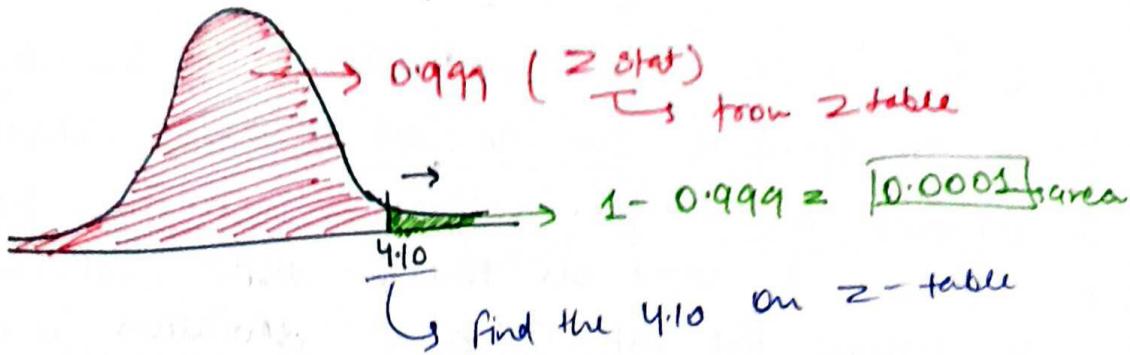
$$\mu = 50 \quad n = 30 \quad \bar{x} = 53 \\ \sigma = 4 \quad \alpha = 0.05$$

$$H_0: \mu = 50 \\ H_a: \mu > 50$$

$$Z\text{-stat} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{53 - 50}{4 / \sqrt{30}} = \frac{3}{4} \times \sqrt{30} = 4.10$$



P-value \rightarrow critical point



P-value < 0.05

$0.0001 < 0.05 \rightarrow$ reject Null hypo

Q Suppose a snack food company claims that their Lays wafer packets contains an avg weight of 50gm per packet. To verify this claim, a consumer watchdog organisation decides to test a random sample of Lays wafer packets. The organisation wants to determine whether the actual avg weight differs significantly from the claimed 50 gms. The organisation collects a random sample of 40 Lays wafer packet and measure their weights. They find that sample has an average weight of 49 grams with a standard deviation of 5 grams.

$$\mu = 50 \quad n = 40 \quad \bar{x} = 49 \quad \sigma = 5 \quad \alpha = 0.05$$

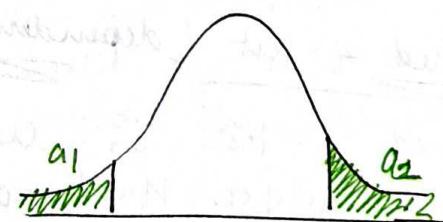
$$z = \frac{49 - 50}{\frac{5}{\sqrt{40}}} = \frac{-5}{5} = -1.26$$

$$H_0: \mu = 50 \\ H_a: \mu \neq 50$$

both side bec we don't know direction ($H_1: \mu > 50$ or $\mu < 50$)

$$P\text{-value} = \alpha_1 + \alpha_2$$

$$\alpha_1 = 0.103$$



$$P\text{-value} = 0.103 + 0.103 = 0.206$$

P-value (0.206) $> 0.05 \rightarrow$ Not reject Null hypo

T-test

A t-test is a statistical test used in hypothesis testing to compare the means of two samples or to compare a sample mean to a known pop mean. The t-test is based on the t-distribution, which is used when the pop standard deviation is unknown and the sample size is small.

There are three main types of t-tests:

One sample t-test: The one-sample t-test is used to compare the mean of a single sample to a known pop mean. The null hypothesis states that there is no significant diff bet'n the sample mean and the pop mean, while the alternative hypo states that there is a significant difference. $\textcircled{O} \rightarrow X$ pop std \times

Independent two-sample t-test: The independent two sample t-test is used to compare the means of two independent samples. The null hypo states that there is no significant diff bet'n the means of the two samples, while the alternate hypo states that there is a significant diff.

Paired t-test (dependent two sample t-test): The paired t-test is used to compare the means of two samples that are dependent or paired such as pre-test and post-test scores for the same group of subjects or measurements taken on

The same subjects under two diff condition.
The null hypo states that there is ~~is~~ no significant difference betn the means of the paired diff, while the alternative hypo states that there is a significant diff.

Single Sample t-test

A one-sample t-test checks whether a sample mean differ from the pop mean.

Assumption for a single sample t-test

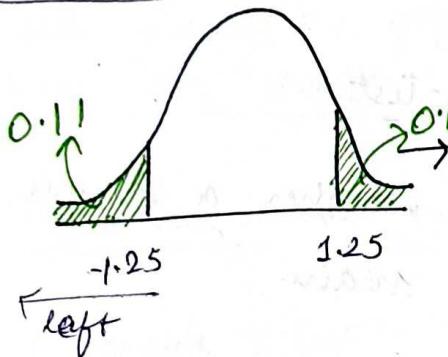
1. Normality: Population from which the sample is drawn is normally distributed
2. Independence: The observation in the sample must be independent
3. Random Sampling
4. Unknown pop std.

Suppose a manufacture claims that the average weight of their new chocolate bars is 50 grams, we highly doubt that and want to check this so we drew out a sample of 25 chocolate bars and measured their weight, the sample mean came out to be 49.7 grams and the sample std dev was 1.2 grams. consider the significance level to be 0.05.

$$\mu = 50 \quad n = 25 \quad \bar{x} = 49.7 \quad \alpha = 0.05$$

$$S = 1.2$$

$H_0: \mu = 50$
 $H_{ac}: \mu \neq 50$



assuming it is normal

$$t = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{49.7 - 50}{1.2/\sqrt{25}} = -1.25$$

Student 't' distn^b

for t-value require

we don't have this

info so we not use
t-table

(n-1)^a
degree of freedom

→ cumulative distn
or one tail or two tail

from scipy.stats import t

cdf-value = t.cdf(t-value, df) → area from left

$$P\text{-value} = 0.11 + 0.11 = 0.22$$

P-value = 0.22 Not reject Null Hypo

P-value > 0.05

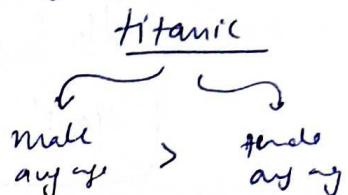
Shapiro-Wilk → use for check whether sample data is Normal distn or not

→ P-value < 0.05 not normal

Independent 2 sample t-test

An Independent two-sample t-test, also known as an unpaired t-test, is a statistical method used to compare the means of two independent groups to determine if there is a significant diff betn them.

Assumption for the test:



1. Independence of observations!
2. Normality
3. Equal variance (Homoscedasticity): $(\sigma_A^2 = \sigma_B^2)$

test equal var \rightarrow F-test or levene

* If this assumption is not true then use Welch's t-test which does not require equal param.

Suppose a website owner claims that there is no diff in the avg time spent on their website betn desktop and mobile users. To test this claim, we collect data from 30 mobile users regarding the time spent on the website in minutes. The sample statistics are as follows:

$$\text{desktop users} = [12, 15, 18, 16, 20, \dots]$$
$$\text{mobile users} = [10, 12, 14, 13, 16, \dots]$$

Desktop user stats

- o Sample size (n_1): 30
- o Sample mean (mean₁): 18.5 minutes
- o Sample standard dev (Std-dev) = 3.5 minutes

Mobile users: Using desktop & laptop

- Sample size (n_2): 30
 - Sample mean (mean₂): 14.3 minutes
 - Sample std (std-dev₂) = 2.7 minutes
- we will use a significance (α) of 0.05 for the hypothesis test.

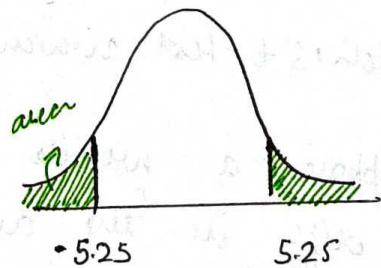
$$\boxed{H_0: \mu_{\text{desktop}} = \mu_{\text{mobile}}}$$
$$H_a: \mu_d \neq \mu_m$$

- Shapiro test → To check Normality
- Breniere test → To check equal variance ($p\text{-value} > 0.05$)
(then equal varian)

$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$ → this formula for 2 sample test

$$\sqrt{\frac{(3.5)^2}{30} + \frac{(2.1)^2}{30}}$$

$$t = \frac{18.5 - 14.3}{\sqrt{\frac{(3.5)^2}{30} + \frac{(2.1)^2}{30}}} = [5.25]$$



$$\text{degree of freedom} = n_1 + n_2 - 2$$

from `scipy.stats import t`

$$\text{cdf-value} = t.cdf(t_value, df)$$

$\rightarrow P\text{-value} < 0.05$
(reject Null Hypo)

Paired 2 Sample t-test

A paired two-sample t-test, also known as a dependent or paired-samples t-test is a statistical test used to compare the means of two related or dependent groups.

Common scenarios where a paired two sample t-test is used include:

1. Before-and after studies: Comparing the performance of a group before and after an intervention or treatment.
2. Matched or correlated groups: Comparing the performance of two groups that are matched or correlated in some way, such as siblings or pairs of individuals with similar characteristics.

Assumption

1. Paired observations: The two sets of observations must be related or ~~possibly~~ paired in some way, such as before-and-after measurements on the same subjects or observations from matched or correlation groups.
2. Normality: histogram, Q-Q plots or statistical test (Shapiro-Wilk test).
3. Independence of pairs: Each pair of observations should be independent of other pairs. In other

words, the outcome of one pair should not affect the outcome of another pair. This assumption is generally satisfied by appropriate study design and random sampling.

	I	II	d (difference (I - II))
A	50	55	-5
B	60	60	0
C	20	60	10
D	40	60	-20
E	25	100	25

not relate A to B

must be normally

Q Let's assume that a fitness center is evaluating the effectiveness of a new 8-week weight loss program. They enroll 15 participants in the program and measure their weights before and after the program. The goal is to test whether the new weight loss program leads to a significant reduction in the participants' weights.

Before the program:

[80, 92, 95, - - - - -]

After the program

[28, 93, - - - - -]

Significance level (α) = 0.05

$H_0: \bar{M}_{\text{before any weight}} = \bar{M}_{\text{after any weight}}$

$H_1: \bar{M}_{\text{before}} > \bar{M}_{\text{after}}$

name	wt. before	wt. after	diff
A	-	-	-
B	-	-	-
C	-	-	-
⋮	-	-	-
k	-	-	-

① Normal check
② $\bar{x}_{\text{diff}} \rightarrow \text{Find}$
③ $S_{\text{diff}} \rightarrow \text{Find}$

$$t = \frac{\bar{x}_{\text{diff}} - (\bar{M}_{\text{diff}})}{S_{\text{diff}} / \sqrt{n}} \rightarrow \text{assumption is there's no diff}$$

so we take as a 0.

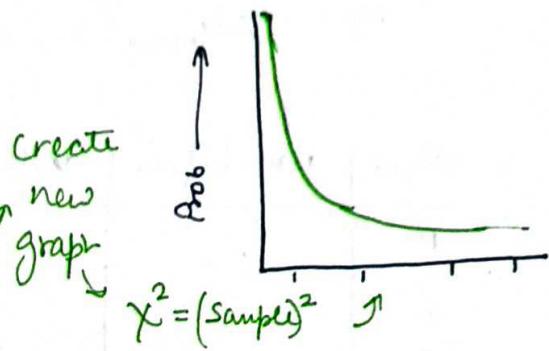
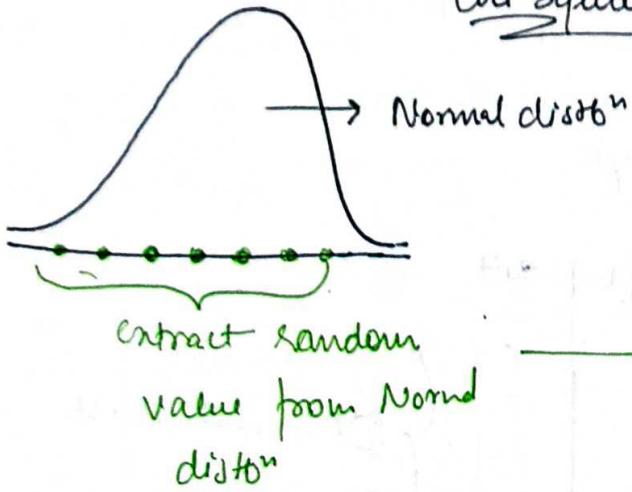
$$S_{\text{diff}} / \sqrt{n}$$

Check code

$$t = \frac{\bar{x}_{\text{diff}} - 0}{S_{\text{diff}} / \sqrt{n}}$$

Chi Square Distribution

→ continuous
function

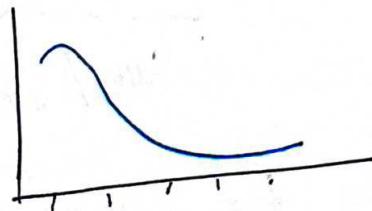


Now take multiple sample value from Normal distn

Sample value (1)

Sample value (2)

$$\chi^2 = (\text{Samp 1})^2 + (\text{Samp 2})^2$$



~~Sample~~ degree of freedom (extract data from pop) = df

df ↑↑ = graph → Normal distn

(no. of
sample var)

$$\chi^2 = z_1^2 \Rightarrow df = 1$$

$$\chi^2 = z_1^2 + z_2^2 \Rightarrow df = 2$$

$$\chi^2 = z_1^2 + z_2^2 + z_3^2 \Rightarrow df = 3$$

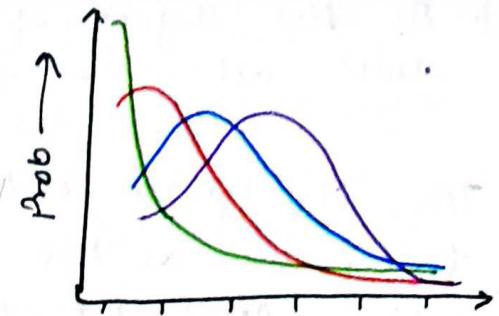
$$\chi^2 = \sum_{i=1}^k z_i^2 \quad \boxed{df = k}$$

$$\chi_1^2 = z_1^2 \Rightarrow df = 1$$

$$\chi_2^2 = z_1^2 + z_2^2 \Rightarrow df = 2$$

$$\chi_3^2 = z_1^2 + z_2^2 + z_3^2 \Rightarrow df = 3$$

$$\chi_4^2 = z_1^2 + z_2^2 + z_3^2 + z_4^2 \Rightarrow df = 4$$



The chi-square distⁿ, also written as χ^2 distⁿ is a continuous prob distⁿ that is widely used in statistical hypothesis testing, particularly in the context of goodness-of-fit tests and ^{test} for independence in contingency tables. It arises when the sum of the squares of independent standard normal random variable follows this distribution.

The chi-square distⁿ has a single parameter, the degree of freedom (df), which influences the shape and spread of the distⁿ. The degree of freedom are typically associated with the num. of independent var or constraints in a statistical problem.

Some key properties of the chi-square distⁿ are:

- a. It is continuous distⁿ, defined for non-negative values.
- b. It is positively skewed, with the degree of skewness decreasing as the degrees of freedom increase.
- c. The mean of the chi-square distⁿ is equal to its degrees of freedom and its variance is equal to twice the degree of freedom.

c. As the degree of freedom increase, the chi-square distⁿ approaches the normal distribution in shape.

The chi-square distⁿ is used in various statistical tests, such as the chi-square goodness-of-fit test, which evaluates whether an observed frequency distⁿ fits an expected theoretical distⁿ, and the chi-square test for independence, which checks the association between categorical variables in a contingency table.

Chi Square Test

The Chi-Square test is a statistical hypo test used to determine if there is a significant association between categorical variables or if an observed distⁿ of categorical data differs from an expected theoretical distⁿ. It is based on the chi-square (χ^2) distⁿ and it is commonly applied in two main scenarios:

~~test~~
 χ^2 -square

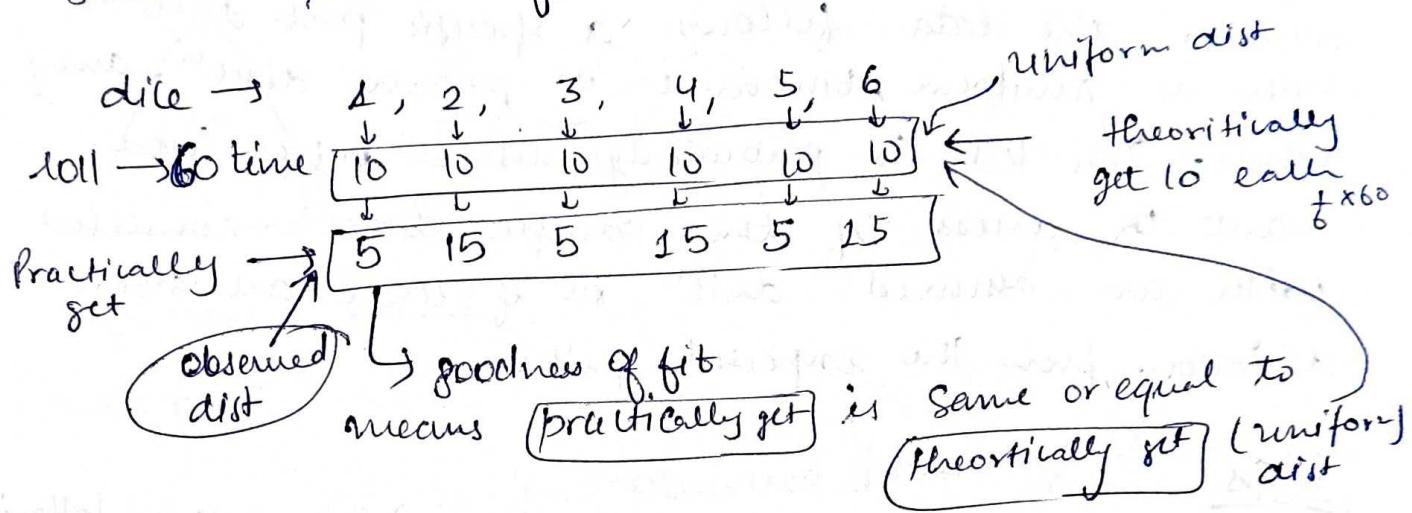
Goodness of fit test

1 categorical col

test for independence

2 categorical diff col

1. Chi-Square Goodness-of-Fit Test: This test is used to determine if the observed distⁿ of a single categorical variable matches an expected theoretical distribution. It is often applied to check if the data follows a specific prob distⁿ, such as the uniform or binomial distⁿ.



2. Chi-Square Test for Independence: (Chi-Square Test for Association) This test is used to determine whether there is a significant association betn. two categorical variables in a sample.

(titanic)	
Pclass	Survived
1	0
2	1
3	0

* Check using chi square for independence → Pclass and survived are dependent or each other or not?

Goodness of Fit Test

The Chi-square goodness - of - fit test is a statistical hypo test used to determine if the observed distⁿ of a single categorical variable matches an expected theoretical distⁿ. It helps to evaluate whether the data follows a specific prob distⁿ, such as uniform, binomial or poisson distⁿ, among others. This test is particularly useful when you want to assess if the sample data is consistent with an assumed distⁿ or if there are significant deviations from the expected pattern.

Steps

The Chi- square Goodness -of- Fit test involves the - following steps:

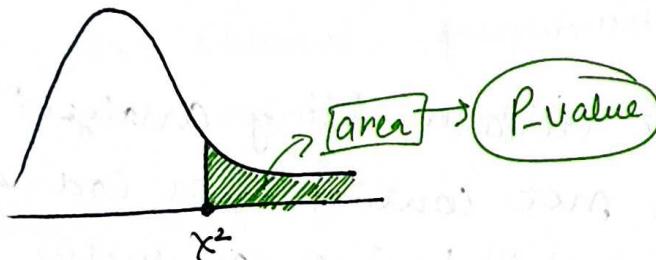
- Define the null hypothesis (H_0) and the alternative hypo (H_1):
 - H_0 : The observed data follows the expected theoretical distⁿ.
 - H_1 : The observed data does not follow the expected theoretical distⁿ.
- Calculate the expected frequencies for each category based on the theoretical distⁿ and the sample size.
- Compute the chi-square test statistic (χ^2) by comparing the observed and expected frequencies.
The test statistic is calculated as:

$$\chi^2 = \sum \frac{(\text{Observed}_i - \text{Expected}_i)^2}{\text{Expected}_i}$$

RHS	Dice \rightarrow	1	2	3	4	5	6
Observed \rightarrow	5	15	5	5	15	15	
Theoretical \rightarrow	10	10	10	10	10	10	

(expected frequency)
like uniform distb

$$\chi^2 = \frac{(5-10)^2}{10} + \left(\frac{15-10}{10}\right)^2 + \dots + \left(\frac{15-10}{10}\right)^2$$



$$df = n - 1 = 6 - 5 = 5$$

$$\alpha = 0.05$$

→ where O_i is the observed frequency in category i , and E_i is the expected frequency in category i , and the summation is taken over all categories.

- Determine the degree of freedom (df), which is typically the number of categories minus one ($df = k - 1$), where k is the number of categories.
- calculate the p-value for the test statistic using the chi-square distribution with the calculated degree of freedom.
- Compare the test statistic to the critical value or the p-value.

Assumption

1. Independence: The observation in the sample must be independent of each other. This means that the outcome of one observation should not influence the outcome of another observation.
2. Categorical data: The variable being analysed must be categorical, not continuous or ordinal. The data should be divided into mutually exclusive and exhaustive categories.
3. Expected frequency: Each category should have an expected frequency of at least 5. This guideline helps ensure that the chi-square distribution is a reasonable approximation for the distribution of the test statistic. Having small expected frequencies can lead to an inaccurate estimation of the chi-square distribution potentially increasing the likelihood of a Type I error (incorrectly rejecting the null hypothesis) or a Type II error (incorrectly failing to reject the null hypothesis).
4. Fixed distribution: The theoretical distribution being compared to the observed data should be specified before the test is conducted. It is essential to avoid choosing a distribution based on the observed data, as doing so can lead to biased results.

Example 1

Suppose we have a six-sided fair die, and we want to test if the die is indeed fair, we roll the die 60 times and record the num. of times each side comes up. We'll use the chi-square Goodness-of-fit test to determine if the observed frequencies are consistent with a fair die (i.e. uniform dist'n of the sides).

Observed freq.

→ Side 1: 12 times

→ Side 2: 8 times

→ Side 3: 11 times

→ Side 4: 9 times

→ Side 5: 10 times

→ Side 6: 10 times

H_0 : die is fair \rightarrow uniform

H_1 : die is not fair.

1 2 3 4 5

Observed =

12	8	11	10	10
----	---	----	----	----

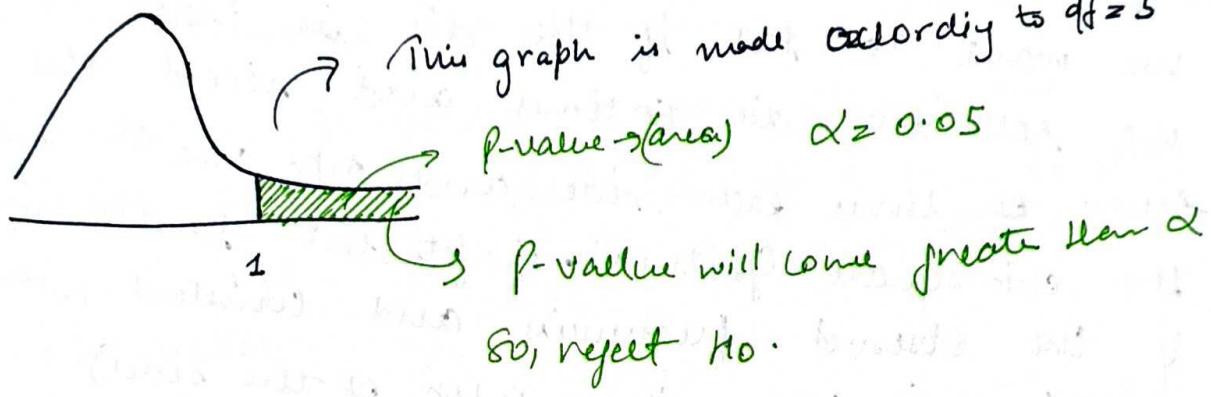
Expected =

10	10	10	10	10
----	----	----	----	----

 \rightarrow uniform

$$\chi^2 = \frac{(12-10)^2}{10} + \frac{(8-10)^2}{10} + \frac{(11-10)^2}{10} + \frac{(9-10)^2}{10} = \frac{10}{10} = 1$$

$$\chi^2 = 1 \rightarrow \text{Chi-square} \quad [df = n-1 = 5]$$



Example 2

Suppose a marketing team at a retail company wants to understand the distⁿ of visits to their website by day of the week. They have a hypo that visits are uniformly distributed across all days of the week, meaning they expect an equal num. of visits on each day. They collected data on website visits for four weeks and want to test if the observed distⁿ matches the expected uniform distⁿ.

Observed frequencies (no. of website visits per day of the week for four weeks):

- Monday: 420
- Tuesday: 380
- Wednesday: 410
- Thursday: 400
- Friday: 410
- Saturday: 430
- Sunday: 390

	Mon	Tue	Wed	Thur	Fri	Sat	Sun
Obs.	420	380	410	450	410	430	390
Expected	405	405	405	405	405	405	405

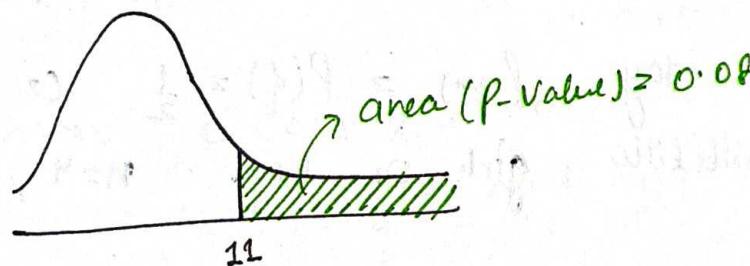
Expected = $\frac{420 + 380 + 410 + 450 + 410 + 430 + 390}{7} = 405$

H_0 : Uniform distribution

H_1 : Not Uniform distribution

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} = \frac{(420 - 405)^2 + (380 - 405)^2 + \dots}{405}$$

$$\chi^2 = 11, \text{ d.f.} = n - 1 = 7 - 1 = 6$$



Example 3

A survey of 800 families in a village with 4 children each revealed the following distn:

# boys :	0	1	2	3	4
----------	---	---	---	---	---

# families :	32	178	290	236	64
--------------	----	-----	-----	-----	----

Is that data consistent with the result that male and female births are equally probable?

Survey \rightarrow village \rightarrow 800 families

\downarrow
4 children

800

$$P(\text{son}) = P(\text{daughter}) = \frac{1}{2}$$

H₀: $P(\text{m}) = P(\text{f})$

H_a: $P(\text{m}) \neq P(\text{f})$

* In this example we use binomial distⁿ not uniform distⁿ.

In binomial we use p and q where $q = 1 - p$.

We can say $P(\text{m}) = P(\text{f}) = \frac{1}{2}$ (i.e. male)

two possibilities girl or boy. $n=4$, $p=\frac{1}{2}$

	0	1	2	3	4	$\sum n \times p^x (1-p)^{n-x}$
Obs	32	178	290	286	64	640

no of boys $\Rightarrow n=1, 2, 3, 4 = 4$
 Boys 0 is not possible

For 0 $\Rightarrow P(0) = {}^4C_0 \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^4 = \frac{1}{16}$ \rightarrow this prob. we want number 10.

$$\frac{1}{16} \times 800 = 50 \rightarrow \text{Number}$$

$$\text{For } 1 \Rightarrow P(1) = {}^4C_1 \left(\frac{1}{2}\right)^1 \left(\frac{1}{2}\right)^3 = \frac{1}{4}$$

$$\frac{1}{4} \times 800 = 200$$

$$\text{For } 2 \Rightarrow 4C_2 \left(\frac{1}{2}\right)^4 = \frac{6}{16} \times 800 = 300$$

$$\text{For } 3 \Rightarrow 4C_3 \left(\frac{1}{2}\right)^4 = 1200$$

For 4 = $\boxed{50}$ → after solving it will be

$$\chi^2 = \frac{(32-50)^2}{50} + \frac{(128-200)^2}{200} + \frac{(290-300)^2}{300} + \\ \frac{(236-200)^2}{200} + \frac{(64-50)^2}{50} = 18.93$$

$$df = 5-1=4$$

~~p-value~~ p-value using chi-square is
 $= 0.00081$

which means. $0.00081 < \alpha(0.05)$

reject Null Hypo

Test for Independence

The chi-square test for independence, also known as the chi-square test for association, is a statistical test used to determine whether there is a significant association betⁿ two categorical variables in a sample. It helps to identify if the occurrence of one variable is dependent on the occurrence of the other variable, or if they are independent of each other.

The test is based on comparing the observed frequencies in a contingency table (a table that displays the frequency distribution of the variables with the frequency that would be expected under the assumption of independence between the two variables).

Pass \leftrightarrow Survived

↳ check whether these ~~are~~ are related to each other or not.
dependent

Steps

1. State the null hypo (H_0) and alternative hypo (H_1):
 - H_0 : There is no association betⁿ the two categorical variables (they are independent).
 - H_1 : There is an association betⁿ the two

Categorical Variables (they are depended).

2. Create a contingency table with the observed frequencies for each combination of the categories of the two variables.

		Pclass		
		1	2	3
Sex	0	23	42	111
	1	100	200	300

Observed

3. Calculate the expected frequencies for each cell in the contingency table assuming that the null hypothesis is true (i.e. the variables are independent)

		Pclass		
		1	2	3
Survived	0	57	50	100
	1	113	150	250

Expected

4. Compute the Chi-square test statistic!

$$\chi^2 = \sum [(O_i - E_i) / E_i]$$

5. Determine the degree of freedom \rightarrow $d.f. = (\text{no. of rows} - 1) * (\text{no. of cols} - 1)$

6. Obtain the critical value or p-value using the Chi-square distribution table or a statistical software/calculator with the given degree of freedom and significance level ($\alpha = 0.05$).

Assumptions

1. Independence of observation: The observations in the sample should be independent of each other. This means that the occurrence of one observation should not affect the occurrence of another observation. In practice, this usually implies that the data should be collected using a simple random sampling method.

2. Categorical Variable: both column categorical

3. Adequate Sample Size!

	1	2	3
0	1	1	1
1	1	1	1

value > 5 (every cell value greater than 5)

4. Fixed Margin total

	1	2	3
0	1	1	1
1	1	1	1

Total Total Total

Example 1

A researcher wants to investigate if there is an association bet'n the level of education (categorical variable) and the preference for a particular type of exercise (categorical variable) among a group of 150 individuals. The researcher collects data and create the following contingency table.

		Exercise Type			Total
		Yoga	Running	Swimming	
Education	High school	15	20	20	45
	Bachelor's	20	30	15	65
	Master's or PhD	5	15	20	40
	Total	40	65	45	150

Education \leftrightarrow Exercise (independent)

H₀: They are independent

H₁: They are associated

	Yoga	Run	Swim
High	$\frac{45 \times 40}{150}$	$\frac{65 \times 45}{150}$	$\frac{45 \times 45}{150}$
Bach	-	-	-
PhD	-	-	-

expected ↗

$$\left(\frac{45}{150} \right) \times \left(\frac{40}{150} \right) = \frac{180}{150^2}$$

prob

multiply 150 bcz same
number chahiye prob
nhi.

expected

	yoga	run	swim
High	12	19	13.5
Base	17	28	20
Low	20	17	12

$$\sum \left(\frac{x_i - E_i}{E_i} \right)^2 \rightarrow \left(\frac{15 - 12}{15} \right)^2 + \left(\frac{20 - 19}{20} \right)^2 + \dots$$

$$\chi^2 = 9.95$$

$$df = (3-1) \times (3-1) = 4$$

$$P\text{Value} = 0.04$$

$P\text{Value} < 0.05$ (reject Null Hypo)

Applications in Machine Learning

1. Feature Selection: chi-square test can be used as a filter-based feature selection method to rank and select the most relevant categorical features in a dataset. By measuring the association betn each categorical feature and the target variable, you can eliminate irrelevant or redundant features, which can help improve the performance and efficiency of machine learning models.
2. Evaluation of classification models: For multiclass classification problems, the chi-square test can be used to compare the observed and expected class frequencies in the confusion matrix. This can help assess the goodness of fit of the classification model, indicating how well the model's prediction align with the actual class distnb.
3. Analysing relationships betn categorical features:
4. Discretization of continuous variables: When converting continuous variable into categorical variables (binning), the chi-square test can be used to determine the optimal number of bins or intervals that best represent the relationship betn the continuous variable and the target variable.

5. Variable selection in decision trees: Some decision tree algo., such as the CHAID (Chi-squared Automatic Interaction Detection) algo., use the Chi-square test to determine the most significant splitting variables at each node in the tree. This helps construct more effective and interpretable decision trees.