

"Check al-e-gs from
slides"

Chi-Squared :-

What if we have categorical data and we want to make predictions about that?

↗ / tests

e.g (i) Whether gender is related to voting preference?

(ii) Whether smoking is related to lungs cancer?

To answer this we use chi-squared distribution.

~~Imp!~~

Date _____

if Z_1, \dots, Z_k are independant std normal random variables then.
 $Z_1^2 + \dots + Z_k^2$, they have χ^2 dist with k degrees of freedom.

the picture in the slides shows the probability density function (PDF) of chi-squared distribution of various degrees of freedom. \downarrow
 (χ^2)

Chi squared dist is a continuous probability distribution which is often used in confidence intervals and hypothesis testing for variance.

(i) It is only defined for $x \geq 0$, since it represents squared values (of Z perhaps).

(ii) The shape of the dist depends upon degree of freedom

x -axis \rightarrow random variable $f_k(x)$

y -axis \rightarrow Probability density of how likely values of x are given for k .

e.g. $k=1$.

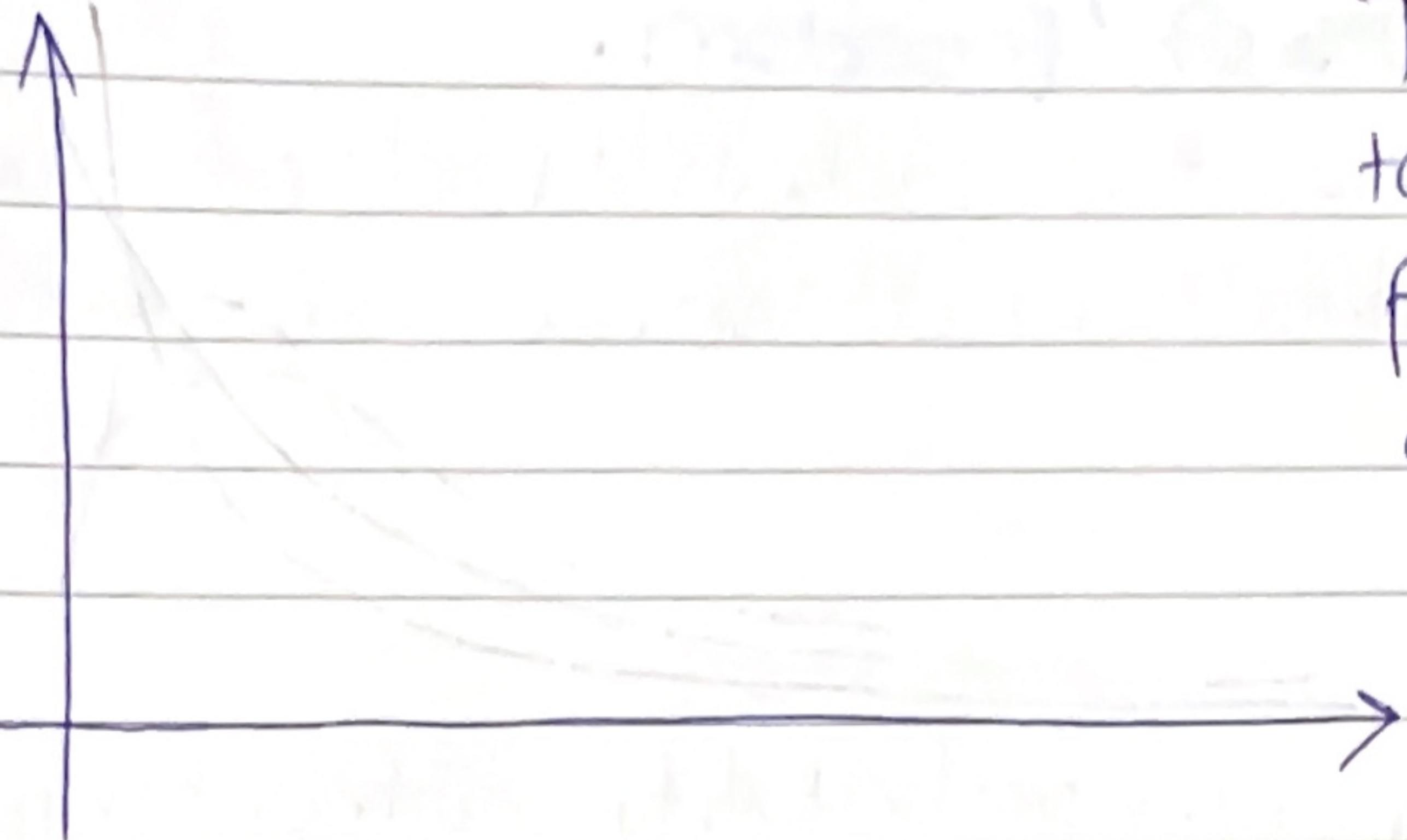
When $k=1$ the distribution is heavily skewed to the right, with most of the probability density concentrated near zero.

k increases, the distribution becomes more symmetric and spreads out, resembling a normal distribution.

e.g. $k=9$

* The mean of chi squared dist $\Rightarrow k$ \rightarrow degree of freedom.
The std deviation of $\dots \Rightarrow \sqrt{2k}$

how value of k - affects the chi-squared distribution.



The Chi-squared dist is often skewed to the right for small degrees of freedom because, the sum of a very few normal random variables often results in small values near to zero but occasional larger values.

As the k increases, the graph becomes; that is why a long tail more spread out and starts resembling a symmetric, bell shaped curve.

(i) This happens because adding more squared normal variables averages out the variability and reduces the skewness.

Chi-squared distribution uses:-

(i) goodness of fit;

to observe a frequency distribution fits a specific expected distribution, also known as the goodness of fit.

(ii) to test whether the categorical variables are independant or if there is some association b/w them.

(testing independence in contingency tables)

(iii) estimating variance.

Date

Example :- test at a significance level of 0.05 or 5%.

H_0 : The die was fair
thus the probability of each toss from the 90 tosses is equal to $\frac{1}{6}$.

$P(X=x)$ where $x = 1, 2, 3, 4, 5, 6$.

$\Rightarrow H_A$: The die is not fair.

$\therefore P(X=x) \neq \frac{1}{6}$.

The expected frequency thus $\Rightarrow (90) \frac{1}{6} \Rightarrow 15$.

Observed frequencies are $\Rightarrow 20, 15, 12, 17, 9$ and 17 .

$$\frac{(20-15)^2}{15} + \frac{(15-15)^2}{15} + \frac{(12-15)^2}{15} + \dots + \frac{(17-15)^2}{15}$$

$$\Rightarrow 5.2$$

number of categories.

$$df = 6 - 1 = 5 \quad \left. \right\} \quad \chi^2 \Rightarrow 11.1$$

$$11.1 > 5.2$$

therefore we accept the NULL hypothesis.

Example :-

	Boy	Girls	Total
observed values	117	130	247
Grades	50	91	141
Popular	60	30	90
Sports	227	251	478
Total			

Test at a significance level of 5%.

Date

$H_0 \Rightarrow$ The gender and categories are independent.

$H_A \Rightarrow$ " " are dependent.

Solve for specifics separately.

for boys

$$\frac{(117 - 117.299)^2}{117.299} + \frac{(50 - 66.960)^2}{66.960} + \frac{(60 - 42.741)^2}{42.741} \Rightarrow 11.3032.$$

for girls:

$$\frac{(130 - 129.701)^2}{129.701} + \frac{(91 - 74.040)^2}{74.040} + \frac{(30 - 47.259)^2}{47.259}.$$

$$\chi^2_e \Rightarrow 11.1999.$$

total Chi-squared.

$$\text{rows } \chi^2 = 11.3032 + 11.1999 \Rightarrow 22.5031.$$

$$df \Rightarrow (3-1)(2-1) \Rightarrow (3-1)(2-1) \Rightarrow 2.$$

$$\alpha = 0.05. \quad \chi^2_c = 5.99 \quad \text{columns.} \quad 5.99 < 22.5031.$$

as the critical Chi-squared value is ^{not} greater than the Chi-squared statistic we reject the NULL hypothesis.

Date

Variance Estimation:-

Chi-squared dist is often used for testing hypothesis about the variance of a population but the estimating variance in it, refers to constructing an interval of confidence where the true population variance can lie.

$$\text{sample variance.} \quad \text{unbiased population variance estimator.}$$
$$(S^2) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

The sample variance S^2 can be transformed to represent the chi-squared distribution (to account for the smaller sample sizes).

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

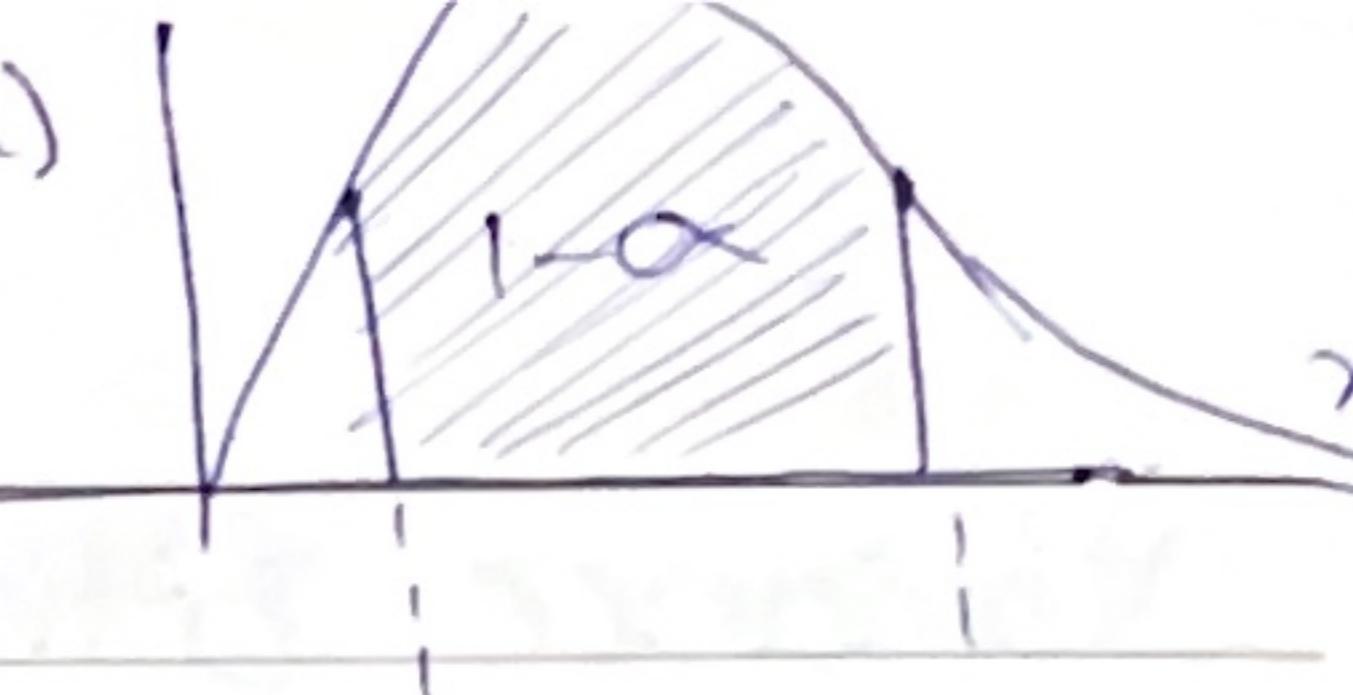
or.

$$\frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2}$$

Now to define the confidence interval for true population mean

$$P \left\{ \chi_{1-\alpha/2}^2 \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi_{\alpha/2}^2 \right\} = 1 - \alpha.$$

taking reciprocal.

$$\frac{1}{\chi^2_{1-\alpha/2}} \leq \frac{\sigma^2}{(n-1)s^2} \leq \frac{1}{\chi^2_{\alpha/2}}$$


Multiplying $(n-1)s^2$.

$$\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi^2_{\alpha/2}}$$

$$\therefore \left[\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}}, \frac{(n-1)s^2}{\chi^2_{\alpha/2}} \right] \checkmark$$

Date

Chi-Squared Test.

i) Goodness of fit:-

In such type of test we want to analyze how well our observed data is being mapped to the Expected data.

Consider the example by Joshua Emmanuel.

Expected probabilities of different outcomes:-

$$p_1 = 0.28 \text{ (being white)}$$

H_0 = the given expected values are true.

$$p_2 = 0.25 \text{ (" black)}$$

H_A = at least one p_i is not the same as specified in the H_0 .

$$\text{No. of categories} = 4 \quad \checkmark$$

$$\text{No. of degrees of freedom} = 4 - 1 \quad \checkmark \\ \Rightarrow 3.$$

$$\text{Chi test} = \sum \frac{f_{\text{obs}} - f_{\text{Exp}}}{f_{\text{Exp}}}^2$$

$$\alpha = 0.05$$

Using the chisquared table

$$\text{Chi Critical Value} = 7.815$$

Now rejection region $\Rightarrow \chi^2 > 7.815$.

for that just divide categories there could be a scenario where you with 4 → (more) below may be given observed frequency and asked if they are equally distributed.

Now

using this Chi-squared test Chi-squared formula.

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

samplespace.

→ 140

total no. of objects

O	E	O-E	$(O-E)^2$	$(O-E)^2/E$	$E_i = n p_i$	expected frequency
39	140(0.28)	-0.2	0.04	0.001		
29	140(0.25)	-6.0	36.00	1.028		or probability
24	140(0.16)	1.6	2.56	0.114		
28	140(0.31)	4.6	21.16	0.4876		
140	140					

$$\chi^2 = 1.6315$$

hence we accept the NULL hypothesis.

X

After dividing you will know how much percentage each category will have, form a expected table from that and solve as the same.

EXAMPLE FROM SLIDES :-

	Observed	Expected	$O-E$	$(O-E)^2$	$(O-E)^2/E$
DICE1	20	16(1/6)			
DICE2	15	16(1/6)			
DICE3	12	"			
DICE4	17	"			
DICE5	9	"			
DICE6	17	"			

$$n = 6 + 84 \Rightarrow 90 \text{ tosses}$$

Remaining steps are similar.

(ii) Chi-Squared test of Independence

USES

↳ (i) Cross Tabulation

↳ (ii) Contingency

↳ (iii) 2-way table

↳ (iv) 2-categorical table.

can be
ranked

Nominal/ordinal

We aim to describe if the two categories are independent or not using the Chi-squared test.

Consider the example by Joshua Emmanuel.

$H_0 \Rightarrow$ No relationship b/w the categories.

$H_A \Rightarrow$ There is a relationship.

Now lets make table to draw χ^2 for each category.

Yes	O	E	$\frac{\text{rowtotal} \times \text{column total}}{\text{sample size}} \rightarrow \frac{490 \times 899}{1440}$
under 30	312	275.28	$(O-E)^2 / E$
30 -> 50	277	256.18	
Over 50.	224	277.53	

Similarly for No and sum

$$df = (r-1)(c-1) \Rightarrow 2$$

$\alpha = 0.05$

$$\chi^2_{\text{Yes}} + \chi^2_{\text{No}} = 57.26$$

Compare.