

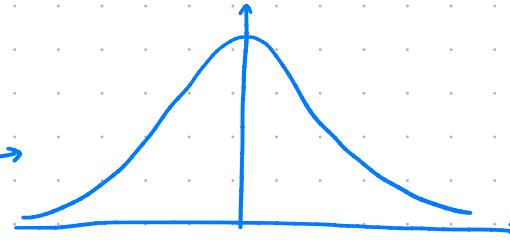
# $\chi^2$ CHI-SQUARED Distribution / Tests

(Not "CHEE"!)

Suppose  $X_1, X_2, X_3, \dots, X_n$  are independent and identically distributed (i.i.d.)

$$X_i \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2)$$

$$Z_i = \frac{X_i - \mu}{\sigma} \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$$

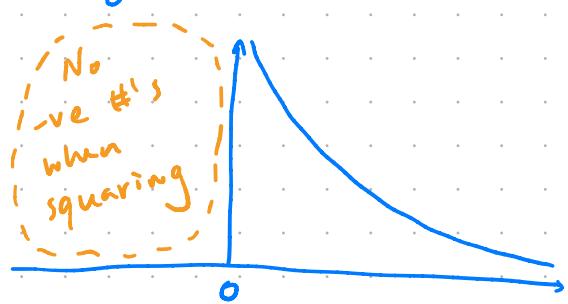


↓ square each member

$$Z_i^2 \sim \chi_i^2 \leftarrow \begin{array}{l} \text{chi-squared} \\ \text{distn with} \\ \text{"1 degree of freedom"} \end{array}$$

$\chi^2$  "chi"  
(pronounced kye)

Just some parameter for now.



$$\sum_{i=1}^n Z_i^2 \sim \chi_n^2 \quad \begin{array}{l} \text{Sum of } n \text{ (normal distributions)}^2 \\ \downarrow \\ \chi^2 \text{ with } n \text{ degrees of freedom} \end{array}$$

As  $n \rightarrow \infty$ ,  $\chi_n^2$  approaches a **normal distribution**

$\chi_n^2$  distn is used to perform a Goodness-of-Fit Test

Is it reasonable to assume that the observed data has been generated by some statistical method?

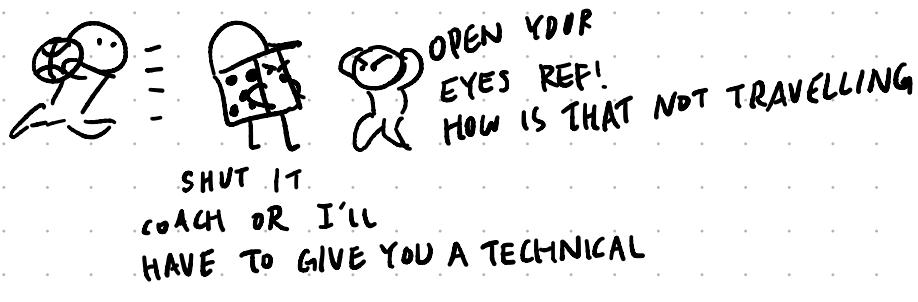
# GOODNESS-OF-FIT

How do we decide (objectively) whether a 6-sided die is fair?

1. Roll die many times
2. Compare set of observed outcomes against the set of expected outcomes
3. Develop some measure of similarity to assess the difference

Roll a die 60 times

$$\begin{cases} H_0: \text{die is fair} \\ H_1: \text{die isn't fair} \end{cases}$$



Category	i	1	2	3	4	5	6
observed	$O_i$	12	7	14	10	8	9
expected	$E_i$	10	10	10	10	10	10

(experiment)  
(assuming  $H_0$ )

How similar are  $O_i$  and  $E_i$ ?

$$\sum_{i=1}^6 (O_i - E_i)$$

Kind of makes sense. (Sum the differences)

Sometimes +ve, sometimes -ve, so differences would cancel.

$$\sum_{i=1}^6 |O_i - E_i|$$

Now always +ve! (I guess...)

$$\sum_{i=1}^6 (O_i - E_i)^2$$

Also would work, but not what we do :)

$$\sum_{i=1}^6 \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_r$$

This is what we use!  
(Because we know it is distributed some way)

$\sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$  is the goodness-of-fit statistic

What the f\*\*k are degrees of freedom? (degrees of freedom =  $v$ )

# of degrees of freedom = # of categories - # of restrictions

The restriction that's always there:  $\sum_{i=1}^n O_i = \sum_{i=1}^n E_i = N$

# of trials is the same for both observed and expected

Ex 6B

$$\textcircled{7} \quad P(X > x) = 0.95 \quad X \sim \chi^2_5$$

Check the table of values ( $v=5$ ,  $p=0.95$ )

$$P(\chi^2_5 > 1.145) = 0.95 \quad \therefore x = 1.145$$

$$\textcircled{8} \quad Y \sim \chi^2_{12} \quad \text{a)} \quad P(Y < y) = 0.5 \quad P(Y > y) = 0.95 \\ P(\chi^2_{12} > 5.226) = 0.95$$

$$\therefore y = 5.226$$



These come from statistical tables

$v$  0.995 0.990 0.975 0.950 0.900 ..

1  
2  
3  
4  
5  
:

BLAH BLAH BLAH

## Goodness-of-Fit Testing

Expected count in category  $i = \text{total count} \times P(\text{fall in category } i)$

assuming  $H_0$  is true

Ex 6c

① die rolled 72 times	#	1	2	3	4	5	6
	observed $O_i$	16	11	13	15	8	9
	expected $E_i$	12	12	12	12	12	12

$N=72$   $H_0$ : die is fair (follows uniform dist $\cong$ )

$H_1$ : die is unfair (follows non-uniform dist $\cong$ )

$$E_i = \frac{1}{6} \times 72 = 12 \quad \sum_{i=1}^6 \frac{(O_i - E_i)^2}{E_i} = \frac{4^2}{12} + \frac{1^2}{12} + \frac{1^2}{12} + \frac{3^2}{12} + \frac{4^2}{12} + \frac{3^2}{12}$$

$$= 4.33 \leftarrow \chi^2$$

$$\begin{matrix} N=12 \text{ is one restriction} \\ \uparrow \\ v = 6 - 1 = 5 \end{matrix}$$

6 categories

$$\alpha = 0.05 \text{ (5%)}$$

from table

$$P(\chi^2_5 > 11.070) = 0.05$$

$$\chi^2_5(0.05) = 11.070 \neq 4.33$$

$\therefore$  We fail to rej.  $H_0$  at  $\alpha = 0.05$  or 5%

$\therefore$  Insufficient evidence to suggest biased die

② LOTTERY! Tickets fall out of spinning thing (tombola)  $N=120$

If ends in 0 or 5, win else, lose

i Win Lose

In 10 numbers, 2 are winning

$O_i$  15 105

$\frac{2}{10} = \frac{1}{5}$  chance to win if

$E_i$  24 96

tombola is fair  $\alpha = 0.05$

$$\chi^2 = \frac{(15-24)^2}{24} + \frac{(105-96)^2}{96} = 4.21875$$

$\{ H_0$ : tombola fair ( $\frac{1}{5}$  chance to win)

$H_1$ : tombola biased (not  $\frac{1}{5}$ )

$v = 2 - 1 = 1$  degree of freedom

$$\chi^2_{1(0.05)} = 3.84 < 4.21875$$

$\therefore$  Sufficient evidence to rej.  $H_0$   
Tombola is BIASED!

④ N=100 households, # of dogs

# of dogs	0	1	2	3	4	5	>5
$O_i$	45	19	11	8	7	6	4
$E_i$	55	20	10	7	4	3	1
						Combine	
$O_i$	0	1	2	3	24		
$E_i$	55	20	10	7	8		

Combine categories when

Any  $E_i$  are less than 5

$$4+3+1 = 8$$

$$\sum \frac{(O_i - E_i)^2}{E_i} \sim \chi^2$$

only if  $E_i \geq 5$

$$v = 5 - 1 = 4$$

↑ 5 categories

$$\alpha = 5\%$$

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$$= 12.236$$

$$\chi^2_{(0.05)} = 9.488 < 12.236$$

∴ Sufficient disagreement

∴ # of dogs not distributed as such.

Why do we have to combine categories when  $E_i < 5$ ?

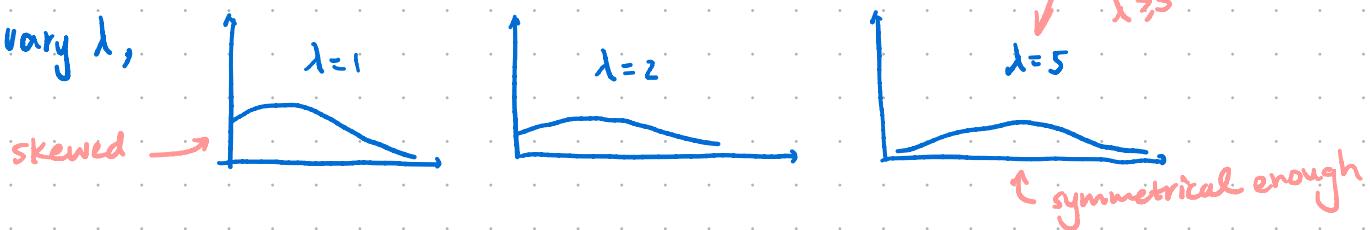
1 Recall that  $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$   $\rightarrow Z_i \stackrel{iid}{\sim} \chi^2_1$

2  $Z_i = \frac{X_i - \mu}{\sigma} \stackrel{iid}{\sim} N(0, 1)$   $\rightarrow \sum_{i=1}^n Z_i^2 \stackrel{iid}{\sim} \chi^2_n$

Suppose observed counts:  $O_i = 1, 2, 3, \dots, n$

Assume  $O_i \sim Po(\lambda) \rightarrow \lambda = E_i$  ( $\lambda$  is the expected occurrences)  
 $= E[O_i] = \text{Var}(O_i)$

When we vary  $\lambda$ ,



If  $E_i \geq 5$ ,  $O_i \sim N(E_i, E_i)$

$$Z_i = \frac{O_i - E_i}{\sqrt{E_i}} \sim N(0, 1) \quad Z_i^2 = \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_1$$

$$\sum_{i=1}^n Z_i^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{n-1}$$

If  $E_i < 5$ ,  $O_i$  no longer normal,  $Z_i^2$  no longer  $\chi^2 \sim$

1 degree of freedom lost

$$\because \sum O_i = \sum E_i$$

Ex (6D)

① a)	x	0	1	2	3	4	5	>5
O <sub>i</sub>	12	23	24	24	12	5	0	
E <sub>i</sub>	13.53	27.07	27.07	18.04	9.02	3.61	1.66	

From  $X \sim Po(2)$   $H_0$ : Data can be modelled by  $Po(2)$   $\leftarrow$  Just a list of  $P(X=x) \cdot N$   
 $H_1$ : Data cannot be modelled by  $Po(2)$  How well does this chosen dist<sup>n</sup> fit the data?

$$\chi^2 = \frac{(12-13.53)^2}{13.53} + \dots + \frac{(5-3.61)^2}{3.61} = 4.101 \quad v = 6-1=5$$

$$\chi^2_{5(0.05)} = 11.070 < 4.101$$

$\therefore$  Fail to reject  $H_0$  at significance level of 0.05

b) Estimate  $\lambda$  (for a better model than  $Po(2)$ )

$$\lambda = E[X] \quad X \sim Po(\lambda)$$

$$\hat{\lambda} \text{ (estimate of } \lambda) = \bar{x} = \frac{12+0+\dots+5+5}{12+23+\dots+5} = \frac{216}{100} = 2.16$$

new estimate of  $\lambda$

$$df = v = 6 - 1 = 5 \quad \sum O_i = \sum E_i \quad \text{another restriction!}$$

estimates reduce degrees of freedom

③ a)	# of accidents	0	1	2	3	4	>4
	O <sub>i</sub> (weeks)	15	13	9	13	0	0

$$\bar{x} = \frac{13+18+39}{15+13+9+13} = \frac{7}{5} = 1.4$$

$H_0$ : Data follows Po dist<sup>n</sup>

$H_1$ : Data doesn't follow Po dist<sup>n</sup>

b) Approximating using  $Po(1.4)$ , ( $\hat{\lambda} = \bar{x}$ )

# of accidents	0	1	2	3	4	>4
p	0.2465	0.3452	0.2416	0.1127	0.0394	0.015

E <sub>i</sub>	12.325	17.26	12.08	5.635	1.97	0.73
----------------	--------	-------	-------	-------	------	------

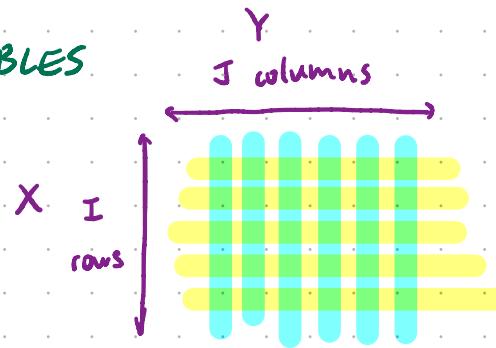
$$\chi^2 = \frac{(15-12.325)^2}{12.325} + \dots = 5.0282 \quad df = 4-2=2$$

combine

$$\chi^2_{2(0.1)} = 4.605 < 5.0282 \quad \therefore \text{sufficient evidence to rej. } H_0$$

# CONTINGENCY TABLES

## GENERALISING!



$X$  has  $I$  categories

$Y$  has  $J$  categories

$O_{ij} \rightarrow$  observed count in row  $i$  column  $j$

$$\text{Sum of row } 1: \sum_{j=1}^J O_{1,j} = O_{1,1} + O_{1,2} + O_{1,3} + \dots + O_{1,J-1} + O_{1,J} = O_{1\cdot}$$

$$\text{Sum of column } 1: \sum_{i=1}^I O_{i,1} = O_{1,1} + O_{2,1} + O_{3,1} + \dots + O_{I-1,1} + O_{I,1} = O_{\cdot 1}$$

↓  
we have sums!

		Y						sum of row $s$
		$O_{1,1}$	$O_{1,2}$	$O_{1,3}$	...	$O_{1,J-1}$	$O_{1,J}$	
X	$O_{2,1}$	$O_{2,2}$	$O_{2,3}$	...	$O_{2,J-1}$	$O_{2,J}$	$O_{2\cdot}$	
	$O_{3,1}$	$O_{3,2}$	$O_{3,3}$	...	$O_{3,J-1}$	$O_{3,J}$	$O_{3\cdot}$	
	:	:	:	⋮	⋮	⋮	⋮	
	$O_{I-1,1}$	$O_{I-1,2}$	$O_{I-1,3}$	...	$O_{I-1,J-1}$	$O_{I-1,J}$	$O_{I-1\cdot}$	
	$O_{I,1}$	$O_{I,2}$	$O_{I,3}$	...	$O_{I,J-1}$	$O_{I,J}$	$O_{I\cdot}$	
								grand total
Sum of columns →		$O_{\cdot 1}$	$O_{\cdot 2}$	$O_{\cdot 3}$	...	$O_{\cdot J-1}$	$O_{\cdot J}$	$N$

$$\hat{P}(X=1) = \frac{O_{1\cdot}}{N}$$

$$\hat{P}(Y=1) = \frac{O_{\cdot 1}}{N}$$

$$\hat{P}(X=i) = \frac{O_{ii}}{N}$$

$$\hat{P}(Y=j) = \frac{O_{\cdot j}}{N}$$

$$\hat{P}(X=1 \cap Y=1) = P(X=1) \cdot P(Y=1) = \frac{O_{1\cdot} \times O_{\cdot 1}}{N^2}$$

$$\hat{E}_{ij} = N \times P(X=1 \cap Y=1) = \frac{O_{1\cdot} \times O_{\cdot j}}{N}$$

use the data to estimate the probability of each cell, that's all.

# of degrees of freedom = # of cells that are free to vary

= total # of cells - # of cells that aren't free to vary

In every row/column, the last cell cannot vary since they have to sum to a number

Total # of cells =  $IJ$     cells that can't vary = margins of the table  $= (I-1) + (J-1) + 1$

# of degrees of freedom =  $IJ - ((I-1) + (J-1) + 1) = IJ - I - J + 1 = (I-1)(J-1)$

Criterion 1 : Grade  
Criterion 2 : Teacher

OBSERVED	A*	A	≤B	
Dr Green	0	2	13	15
Mr ?	4	6	5	15
Monkey	13	2	0	15
	17	10	18	45

45 pupils cross-classified to 2 categorical variables.

$H_0$ : grade independent of teacher

$H_1$ : grade dependent of teacher

To test  $H_0$ , compare table of observed counts to table of expected counts

calculated assuming  $H_0$  is true

EXPECTED	A*	A	≤B	
Dr Green	5.66	3.33	6	15
Mr ?	5.66	3.33	6	15
Monkey	5.66	3.33	6	15
	17	10	18	45

all rows same since independent

Now we test the observed table against expected table (Goodness of fit)

15 people distributed according to the column sums

Goodness of fit:  $\sum \frac{(O_i - E_i)^2}{E_i}$  for every cell and do the test!

$$G = 33.18 \quad df = (3-1)(3-1) = 4$$

$$\chi^2_{4(0.05)} = 9.488 < 33.18 = G \quad \therefore \text{sufficient evidence to rej. } H_0!$$

$\therefore$  Grades dependent of teacher