Chi-Squared Tests Cheat Sheet

Goodness of fit tests are used to determine how well an observed frequency distribution fits a known distribution. For example, you might have a set of data that you suspect follows a Poisson distribution. To figure out whether the data is suitable to be modelled by a Poisson distribution, you can conduct a goodness of fit test. We will learn to carry out goodness of fit tests for discrete uniform, Poisson, binomial, and geometric distributions in addition to contingency tables.

The idea behind goodness of fit tests

Let's say that you have a dice, and you are wondering whether it is fair or not. You might decide to throw the dice 240 times and record the results. If the dice is fair, then in theory, the results should look like this:

Number on dice, n	1	2	3	4	5	6	The expected results follow a discrete uniform
Expected frequency	40	40	40	40	40	40 🔵	distribution.

However, since you have taken a sample, you likely won't roll each side of the dice exactly 40 times. In practice, your results might look something like this

Number on dice, n	1	2	3	4	5	6
Observed frequency	36	41	47	35	38	44

So, even if the dice was in fact unbiased, the scores won't necessarily exactly match the theoretical expected values. But how can we decide whether our results are 'good enough' to conclude that the dice is fair (i.e. the observed frequencies do indeed follow a discrete uniform distribution)? Well, we can set up a hypothesis test.

We take the null hypothesis to be that there is no difference between the theoretical frequency distribution and the observed frequency distribution, and that any differences can be put down to natural variations. This would mean that the dice is unbiased (i.e. the score on the dice follows a discrete uniform distribution)

We take the alternative hypothesis to be that there is in fact a difference between the theoretical frequency distribution and the observed frequency distribution that cannot be explained by natural variations. This would mean that the dice is in fact biased. (i.e. the score on the dice does not follow a discrete uniform distribution)

In order to test these hypotheses, we need a measure of how closely correlated the expected and theoretical frequencies are. In other words, we need a measure for the goodness of fit between the expected and observed frequencies. This is where the X^2 test statistic comes in:



If we were to calculate X^2 for the above data, we would have:



You can see that when there is a small difference between observed and expected frequencies, the value of $\frac{(O-E)^2}{n}$ is small, but when there is a bigger difference between observed and expected frequencies then the value is larger. This should convince you that the statistic X^2 is a suitable measure for the goodness of fit between the observed and expected frequencies.

Chi-squared (χ^2) family of distributions and degrees of freedom

The χ^2 family of distributions are a set of distributions that can be used as approximations for the X^2 test statistic. In other words, X^2 follows a chi-squared distribution. We use the chi-squared distributions to find critical values for our hypothesis test. You will need to calculate the number of degrees of freedom for your data before you can find a critical value. The number of degrees of freedom is what distinguishes the different members of the χ^2 family.

Degrees of freedom are a measure of how much information from the sample data has not been used up. Every time a statistic is calculated from a sample, a degree of freedom is used up. For the above data

- v = 6 1 = 5Number of degrees of freedom = number of columns - number of constraints (after any combining) -• (v)
 - \Rightarrow Any parameter that you calculate relating to the distribution counts as a constraint
 - \Rightarrow We always subtract.at least 1 because the requirement that the total expected and observed frequencies must match counts as a constraint.
 - \Rightarrow If there are any expected values less than 5, you will need to combine columns so that there are no expected values less than 5. This is because X^2 is only approximated well by χ^2 when the expected frequencies are all above 5.
- To find critical values for your goodness of fit test, you will need to use the table of values for the percentage points of the χ^2 . distribution, which is given in the formula booklet (or p192 in the textbook). Some calculators (e.g. Casio CG-50) will allow you to directly find these critical values, eliminating the need to refer to the table.
- To make clear which member of the χ^2 distribution you are talking about, you write χ^2_{ν} . For example, χ^2_3 is used to denote the χ^2 distribution with 3 degrees of freedom.
- The critical value of χ^2 that is exceeded with probability 5% is written as $\chi^2_{\nu}(5\%)$ or $\chi^2_{\nu}(0.05)$.

If we were carrying out a test at 5% significance level for the above data, our critical value would be $\chi_5^2(5\%) = 11.070$ (from tables).

The goodness of fit tests that you will be expected to carry out

You need to be able to test whether a given set of data can be modelled by a:

- Discrete uniform distribution
- Binomial distribution
- Poisson distribution
- Geometric distribution

Additionally, you need to be able to carry out goodness of fit tests to figure out whether two events are independent or not using data given in a contingency table. The process for such questions is slightly different to the others but is covered in this cheat sheet

We have already briefly covered how a goodness of fit test is done with a discrete uniform distribution. With the help of examples, we will now look at how it is done with other distributions.



	a resources a tuition a courses		
		c) Stating your hypotheses clearly, use a 1	0% lev
The step-by-step procedure Regardless of which distribution you are testing for	or, the procedure for carrying out a goodness of fit test is largely the same:	Start by stating your hypotheses.	
1) Form your null and alternative hypotheses. Th the alternative hypothesis should be that it is not	e null hypothesis is always that the given distribution is a good model for your data while a good model.	values and use it to find X ² . Remember that the had an expected frequency less than 5 so we con the final two cells together before calculating X ²	ed last cell nbine !
2) Work out the expected frequencies. To do this, for a binomial distribution then you will need to u	you will need to consider the distribution you are testing for. For example, if you are testing se the appropriate binomial distribution to figure out the expected frequencies.	Find degrees of freedom.	
 Create a table that contains the expected frequ If so, you will need to pool adjacent cells together 	encies and observed frequencies. Check to see if any expected frequencies are less than 5. so that there are no expected frequencies less than 5.	State the critical value.	
4) Calculate $X^2 = \sum_{l=1}^{N} \frac{(O_l - E_l)^2}{E}$. Remember that the	is should be done after any columns are pooled together.	Compare X^2 with the critical value and write conclusion in context.	e a
5) Calculate the number of degrees of freedom to critical value using the Chi-squared percentage po significance level, which will be given to you in the	for your data using the formula $v = no. of cells - no. of constraints - 1$ and find the ints table in the formula booklet. You will need the number of degrees of freedom and the equestion.	Using contingency tables Another type of question you can be asked to conduct a goodness of fit test to figure	d is one out wh
6) If X^2 exceeds the critical value, then you shoul conclusion in the context of the question. Refer to	d reject the null hypothesis. Otherwise accept the null hypothesis. Remember to give your o the next examples for more detail.	1) Form your null and alternative hypoth alternative hypothesis is always that the transmission of the second seco	eses. T wo var
Testing the binomial distribution as a mode		2) Work out the expected frequency for early a second seco	ach cel
Example 1: A research station is doing some work containing 7 seeds in each row. The number of se	on the germination of a new variety of genetically modified wheat. They planted 120 rows eds germinating in each row was recorded. The results are as follows:	3) Find $X^2 = \sum_{l=1}^{N} \frac{(O_l - E_l)^2}{E}$.	
No. seeds germinating in each row Observed number of rows	0 1 2 3 4 5 6 7 2 6 11 19 25 32 16 9	4) Work out the degrees of freedom using	the fo
a) Write down two reasons why a binomial distrib	ution may be a suitable model.	5) Find the critical value using v and the sine to us the sine to be set of the sinet	gnifica
required for the binomial distribution to be used, that	imploits -cach seed either germinates of obesit (jeach that is either a success of railine), tare satisfied. -The germination of each seed is independent of the others (independency).	Example 3: A researcher investigates the	result
b) Show that the probability of a randomly selected	ed seed from this sample germinating is 0.6 .	sample of 620 candidates gave the followi	ng resi
probability= total germinations total seeds planted	$p = \frac{0(2) + 1(6) + 2(11) + 3(19) + 4(25) + 5(32) + 6(16) + 7(9)}{120 \times 7} = \frac{504}{840} = 0.6$	Results	
The research station used a binomial distribution decimal places. The results are as follows:	with probability 0.6 of a seed germinating. The expected frequencies were calculated to 2		To
No. seeds germinating in each row Expected number of rows	0 1 2 3 4 5 6 7 0.20 2.06 s 23.22 t 31.35 15.68 3.36	Test, at the 5% level of significance, wheth centre. State your hypotheses and show y	ner the our wo
c) Find the value of s and the value of t .		Start by stating your hypotheses.	H ₀ : H ₁ :
To find the expected values, we need to use the suggester binomial distribution. The probabilities can be found usin calculator.	Let $X =$ the number of seeds germinating in one row. Then $X \sim B(7, 0.6)$ P(X = 2) = 0.0774. This is the probability that there are 2 germinated seeds in a given row. So, for 120 rows: the expected number of rows for which there are 2 germinated seeds is $120 \times 0.0774 = 9.29$ $\therefore s = 9.29$.	Find the expected frequencies for each cell of the table using the formula expected frequency = row total × column total tota	For
To find t , use the fact that the total expected frequency mequal to the total observed frequency (120).	ust be $t = 120 - (0.20 + 2.06 + 9.29 + 23.22 + 31.35 + 15.68 + 3.36) = 34.84$	grana totat	
d) Stating your hypotheses clearly, test at the 1%	level of significance, whether or not the data can be modelled by a binomial distribution.	Now calculate $X^2 = \sum_{l=1}^{N} \frac{(O_l - E_l)^2}{E}$.	X ² :
Start by stating your hypotheses.	H_0 : A binomial distribution is a suitable model for these data. H_0 : A binomial distribution is not a suitable model for these data.	Find degrees of freedom.	v =
Create a table showing the expected and observed values and use it to find X ² . You will need to pool the first three cells and the last two cells, so that all	No. seeds germinated 0-2 3 4 5 6-7 Observed number of rows 19 19 25 32 25 Expected number of rows 11.55 23.22 34.84 31.35 19.04 0 - 5/* 4.995 0.777 0.770 0.013 1.966	Find critical value, compare it to our X^2 value and write a conclusion.	7.5. cen
expected frequencies are above 5.	$\therefore X^{2} = 4.805 + 0.767 + \dots + 1.866 = 10.23$	Testing the geometric distribution as	a mo
Find degrees of freedom.	There are 5 cells after pooling. We also calculated p in part (a) so that counts as a constraint. $\therefore p = 5 - 1 - 1 = 3$	stops, over the course of 100 days.	
State the critical value.	v = 3, significance level = 1% critical value = $\chi_3^2(1\%) = 11.345$	Attempts Frequency	
Compare X^2 with the critical value and write a conclusion in context.	$10.23 < 11.345$ \div result is insignifcant, accept H_0 . Evidence suggests that a binomial distribution is a suitable model for the given data.	Katie thinks she can model the number of a) Using the observed frequencies, find an	attem estim
Testing the Poisson distribution as a model		$p = \frac{\text{total successes}}{\text{total attempts}} \qquad p = \frac{1}{1(76)}$	+ 2(17)
Example 2: The number of accidents on a particul summarized in the following table:	ar stretch of motorway was recorded each day for 200 consecutive days. The results are	b) Conduct a goodness of fit test at the 2.5	% sign
No. accident Frequency	s 0 1 2 3 4 5 47 57 46 35 9 6	Start by stating your hypotheses.	
a) Show that the mean number of accidents per d	ay for these data is 1.6. 47) + 1(57) + 2(46) + 3(35) + 4(9) + 5(6) - 320 - 1.6	Next, we work out the expected frequencies. We with $p = 0.741$ to do so. The formula $P(X = i) =$	use the = p(1 -
A motorway supervisor believes that the number distribution. She uses the mean found in part (a) t	200 -200 -200 $-1.0of accidents per day on this stretch of motorway can be modelled by a Poisson$	Remember that the last expected value will be gi (sum of the other expected values) since you expected frequency matches up with the total of	ven by 1 i need to oserved f
following table. No. accidents Expected frequency 4	0 1 2 3 4 5 or more 10.38 64.61 r s 11.03 t	Now we put the expected and observed values in were some expected values less than 5. We need columns to make sure that all expected frequence	i a table. I to pool ties are a
b) Complete the table, giving your answers to 2 de	ecimal places.		
To find the expected values, we need to use the suggested Poisson distribution. The number of accidents follow the distribution Po(1.6). The probabilities can be found with some view.	Let $x = the number of accidents in one day. Then X-Po(1.6)P(X = 2) = 0.2584$. This is the probability that there are 2 accidents in a given day. P(X = 3) = 0.1378. This is the probability that there are 2 accidents in a given day. So, for 200 days: the expected number of days for which there are 2 accidents is $200 \times 0.2584 = 51.69$	We use our table to find X ² .	tract on 1
iound using a calculator.	the expected number of days for which there are 3 accidents is $200 \times 0.1378 = 27.57$ \therefore r = 51.69 and s = 27.57.	estimated p in part a.	aut dh 6

		1.1					-
Find	critical	value,	compare	it to	our X^2	value	ar



To find *t*, use the fact that the total expected frequency t = 200 - (40.38 + 64.61 + 51.69 + 27.57 + 11.03) = 4.72.

 $\textcircled{\begin{time}{0.5ex}}$

Edexcel FS1

el of significance to test the motorway supervisor's belie

	H_0 : Po(1.6) is a suitable mo H_0 : Po(1.6) is not a suitable	odel for these e model for tl	edata. hese data.				
ved	No. accidents	0	1	2	3	4 or more	1
e last cell	Observed frequency Expected frequency	47	57 64.61	46	35	9+6=15 11.03+4.72=15.75	
ombine ⁽² .	$\frac{(O_i - E_i)^2}{E}$	0.164	0.896	0.626	2.002	0.035	1
	$\therefore X^2 = 0.164 + 0.896 + \cdots$	+ 0.035 = 3.72	2				
	There are 5 cells. We also $v = 5 - 1 - 1 = 3$	calculated the	e mean in p	oart (a) so i	hat counts	s as a constraint.	
	v = 3, significance level critical value = $\chi_3^2(10\%)$	l = 10% %) = 6.251					
te a	3.72 < 6.251 ∴ result is ir with mean 1.6 is a suitable	nsignifcant, a	ccept H ₀ . E e number	vidence su	ggests tha	t a poisson distribution on the motorway.	on

e where you are given a contingency table (sometimes called a two-way table) and you need nether two variables are independent. The process for such questions is slightly different.

The null hypothesis is always that the two variables are independent (no association). The iables are not independent (there is an association)

Il using the formula expected frequency $= \frac{row total \times column total}{grand total}$.

v = (rows - 1)(columns - 1).

ance level. Compare this to your value for X^2 and decide whether to accept or reject the null

s of candidates who took their driving test at one of three driving test centres. A random ults

		A	В	С	Total
	Pass	99	110	68	277
	Fail	108	116	119	343
To	tal	207	226	187	620

ere is an association between the results of candidates' driving tests and the driving test orking clearly. You should state your expected frequencies correct to 2 decimal places.

	pected frequencies	for Pass/A will	be $\frac{2778207}{620} =$	92.48. Doin	g this for a	l cells gives us:
EXPECTED	FREQUENCIES	Α	В	С	Total	
Results	Pass	92.48	100.97	83.55	277	
	Fail	114.52	125.03	103.45	343	
Total		207	226	187	620	

del

a taxi to take her home. She records how many taxis she needs to flag down before one

1	2	3	4	5	Total
76	17	4	2	1	100

pts each day using a geometric random variable $X \sim Geo(p)$.

ate for p (to 3 d.p.).

 $\frac{1}{1} + 3(4) + 4(2) + 5(1) = 0.741$

ificance level, and say whether a geometric random variable is a good model for the data. H_0 : A geometric distribution is a suitable model

	H ₁ : A geometric distributi	H ₁ : A geometric distribution is not a suitable model.								
Ve use the geometric distribution $p = p(1 - p)^{l-1}$ is used. given by 100 – bu need to make sure the total observed frequency (100).	$ \begin{array}{l} Expected \ frequency = \\ \therefore \ E_1 = 100 \times (0.741)(1) \\ E_2 = 100 \times (0.741)(1) \\ E_1 = 100 \times (0.741)(1) \\ E_1 = 100 \times (0.741)(1) \\ E_{25} = 100 - (74.1 + 19), \end{array} $	$100 \times P$ - 0.741) ⁶ 0.741) ¹ = 0.741) ² = 0.741) ³ = 2 + 4.98	(X = i) = 74.1 = 19.2 = 4.98 = 1.29 + 1.29)	, since ti 1 = 0.42	here are 100) days.				
in a table. Notice that there	Attempts	1	2	≥3	Total	1				
ed to pool the final three	Observed frequency	76	17	7	100					
ncies are above 5.	Expected frequency	Expected frequency 74.1 19.2 6.71 100		100]					
	Observed frequency	76	17		7					
	Expected frequency	74.1	19.2	19.2 6.7						
	$\frac{(O_i - E_i)^2}{E}$	0.049	0.25	2 0	.013					
	$\therefore X^2 = 0.049 + 0.252 + $	0.013 = 0).313							
btract an extra 1 because we	v = 3 - 1 - 1 = 1									
nd write a conclusion.	critical value = χ_1^2 (2.5° 0.313 < 5.024 \therefore result is association between the c	%) = 5.0 significan triving tes	v = 3 - 1 - 1 = 1 critical value = $\chi_1^2(2.5\%) = 5.024$ 0.313 < 5.024 · result is significant. Evidence suggests that there is in fact an							

