

The Chi-Squared Distribution

Every hypothesis test we have met so far has been about a single *parameter*: a mean, a proportion, a correlation coefficient. We now ask a more ambitious question: does a data set fit a whole *distribution*? To answer it we need a way of measuring the discrepancy between data and model, and we need to know how that measure behaves when the model is actually true. The distribution that does this job is the chi-squared distribution.

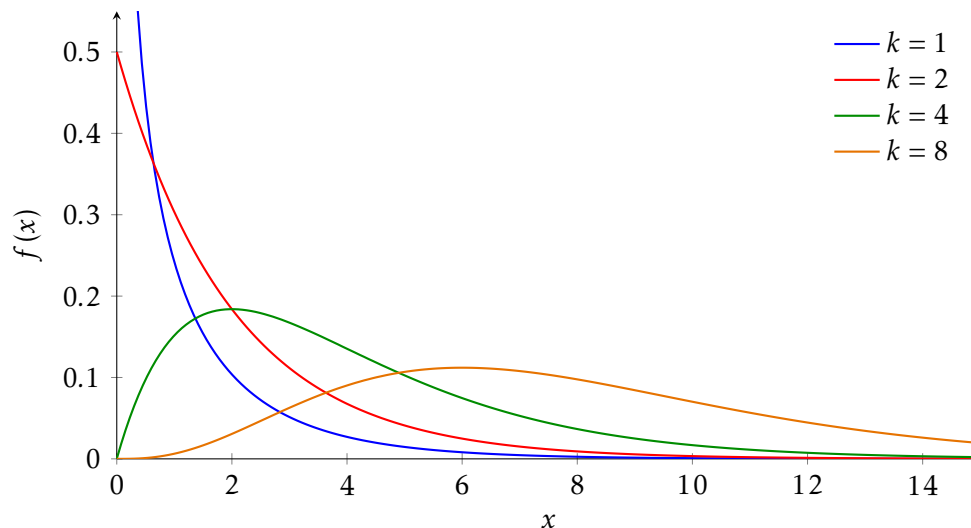
Definition. Let Z_1, Z_2, \dots, Z_k be *independent* standard normal random variables, so $Z_i \sim N(0, 1)$ for each i . Then

$$X = Z_1^2 + Z_2^2 + \dots + Z_k^2$$

has the **chi-squared distribution** with k **degrees of freedom**. We write $X \sim \chi_k^2$.

Remark. χ is the Greek letter chi, and χ^2 should be read as a single symbol — there is no random variable called χ . Since X is a sum of squares, $X \geq 0$ always.

The shape of the distribution



Some observations, which you should check by simulation (e.g. square and sum columns of random normal values in a spreadsheet, or play with sliders in Geogebra):

- For $k = 1$, $X = Z^2$ and most of the mass is squashed up against 0 (squaring a number between -1 and 1 makes it smaller); the density is unbounded as $x \rightarrow 0^+$.
- For $k = 2$ the density is a decreasing exponential, $f(x) = \frac{1}{2}e^{-x/2}$.
- For $k \geq 3$ the curve is humped, with its peak at $x = k - 2$, and a long right tail.
- As k grows the shape becomes increasingly symmetric and bell-like — unsurprising, since χ_k^2 is a sum of k independent identically distributed random variables (Z_i^2), so the Central Limit Theorem applies.

Fact — If $X \sim \chi_k^2$ then

$$\mathbb{E}[X] = k \quad \text{and} \quad \text{Var}[X] = 2k.$$

Moreover, if $X \sim \chi_m^2$ and $Y \sim \chi_n^2$ are independent, then $X + Y \sim \chi_{m+n}^2$ (a sum of m squares plus a sum of

n squares is a sum of $m + n$ squares).

Example

Given that $\mathbb{E}[Z^4] = 3$ for $Z \sim N(0, 1)$, prove the fact above.

For each i ,

$$\mathbb{E}[Z_i^2] = \text{Var}[Z_i] + (\mathbb{E}[Z_i])^2 = 1 + 0 = 1,$$

so by linearity $\mathbb{E}[X] = \mathbb{E}[Z_1^2] + \dots + \mathbb{E}[Z_k^2] = k$.

For the variance,

$$\text{Var}[Z_i^2] = \mathbb{E}[Z_i^4] - (\mathbb{E}[Z_i^2])^2 = 3 - 1 = 2,$$

and since the Z_i (hence the Z_i^2) are independent, variances add:

$$\text{Var}[X] = 2 + 2 + \dots + 2 = 2k.$$

Remark (Where $\mathbb{E}[Z^4] = 3$ comes from). The value $\mathbb{E}[Z^4] = 3$ comes from integration by parts on $\int z^4 \cdot \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz$.

More slickly: χ_k^2 is the Gamma distribution with parameters $\alpha = \frac{k}{2}$, $\lambda = \frac{1}{2}$, and the substitution $u = \frac{z^2}{2}$ in the integral $\mathbb{E}[e^{tZ^2}]$ produces the Gamma function (see the Gamma Function notes). The MGF route gives $M_{\chi_k^2}(t) = (1 - 2t)^{-k/2}$, from which both moments drop out.

Measuring Goodness of Fit

Example

A die is rolled 120 times, with these results:

| | | | | | | |
|-----------|----|----|----|----|----|----|
| Score | 1 | 2 | 3 | 4 | 5 | 6 |
| Frequency | 25 | 17 | 15 | 23 | 24 | 16 |

If the die is fair we “expect” 20 of each score. The data is clearly not exactly that — but data never is. Is it *far enough* from 20, 20, ..., 20 to convince us the die is biased?

We need a single number measuring how far the **observed frequencies** O_i are from the **expected frequencies** E_i .

Definition. The **goodness-of-fit statistic** is

$$X^2 = \sum_i \frac{(O_i - E_i)^2}{E_i},$$

summed over all the cells. The individual terms $\frac{(O_i - E_i)^2}{E_i}$ are called the **contributions** to the test statistic (exam questions often ask for these).

Squaring stops positive and negative discrepancies cancelling; dividing by E_i makes the measure relative (being 10 out matters far more when you expected 15 than when you expected 1500).

Why chi-squared? The two-cell case

Theorem

For two cells, with expected frequencies large enough, X^2 has approximately the χ_1^2 distribution.

The key is that with two cells there is really only *one* free quantity, and it is approximately normal.

Suppose each of n independent trials lands in cell 1 (probability p_1) or cell 2 (probability $p_2 = 1 - p_1$). Then $O_1 \sim B(n, p_1)$, with $E_1 = np_1$ and $E_2 = np_2$. Since the total is fixed, $O_2 = n - O_1$, so

$$O_2 - E_2 = (n - O_1) - (n - E_1) = -(O_1 - E_1),$$

and hence

$$X^2 = (O_1 - E_1)^2 \left(\frac{1}{E_1} + \frac{1}{E_2} \right) = (O_1 - E_1)^2 \cdot \frac{E_1 + E_2}{E_1 E_2} = \frac{(O_1 - np_1)^2}{np_1 p_2}.$$

By the normal approximation to the binomial, $O_1 \approx N(np_1, np_1 p_2)$, so

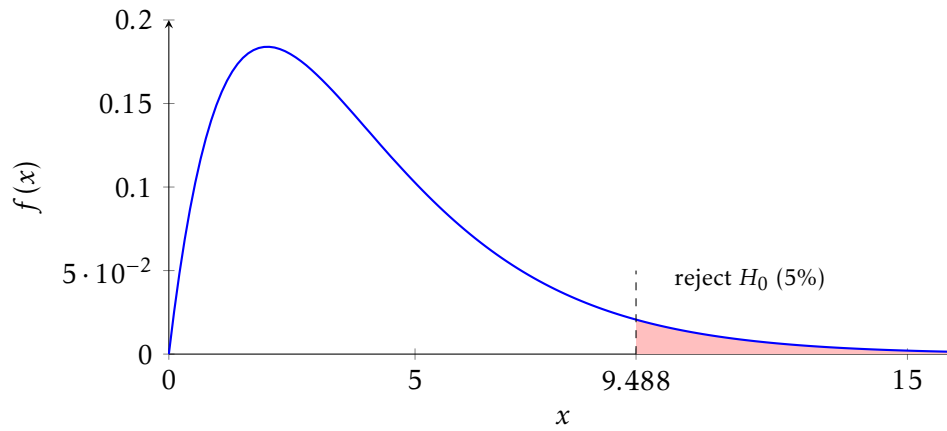
$$X^2 \approx \left(\frac{O_1 - np_1}{\sqrt{np_1 p_2}} \right)^2 \approx Z^2 \sim \chi_1^2.$$

Remark (More cells). The same idea with the multinomial distribution (and much more work) shows that for c cells, $X^2 \approx \chi_{c-1}^2$ when the cell probabilities are fully specified. Notice the degrees of freedom: c cells minus 1 constraint, because the frequencies must total n — exactly as O_2 was determined by O_1 above.

Fact (The $E \geq 5$ rule) — The argument above relies on the binomial-to-normal approximation, which fails when np is small. To keep the approximation honest, **every expected frequency must be at least 5**: where necessary, adjacent classes are combined until this holds. Combining classes is called **pooling**.

A right-tail test only

If the model is true, X^2 behaves like a χ^2_ν random variable. A *large* value of X^2 means observed and expected frequencies disagree badly — evidence against the model. A *small* value just means the fit is good. So the chi-squared test is always a test on the **right-hand tail only**: we reject H_0 when X^2 exceeds the critical value.



The χ^2_4 density with the 5% rejection region shaded: the critical value is 9.488.

Remark (Too good to be true?). There is no significance test on the left-hand tail at A Level. But a value of X^2 deep in the left tail *should* raise an eyebrow: real data is noisy, and a suspiciously perfect fit suggests the data may have been tidied up. R. A. Fisher famously analysed Gregor Mendel's pea-breeding data and found the fit to Mendel's 3:1 ratios was far better than chance should allow — the combined chi-squared statistic was so small that data this good would arise only a few times in 100 000. Whether Mendel (or an enthusiastic assistant) trimmed the data remains controversial. For your exams: tests are right-tailed, always.

Remark. This test is in some sense backwards compared to other hypothesis tests. Rejecting H_0 gives evidence the data was *not* drawn from the specified distribution; but failing to reject does *not* give evidence that it *was* — it just means we couldn't tell the difference.

The Goodness-of-Fit Test

- Tip (The procedure)**
1. State H_0 and H_1 (see below for the phrasing).
 2. Compute the expected frequencies E_i under H_0 . **Pool** adjacent cells until every $E_i \geq 5$.
 3. Compute $X^2 = \sum \frac{(O_i - E_i)^2}{E_i}$.
 4. Degrees of freedom: $\nu = (\text{number of cells after pooling}) - 1$.
 5. Compare with the critical value of χ^2_ν from the formula booklet at the given significance level (right tail).
 6. Conclude *in context*, without over-claiming.

Remark (Phrasing the hypotheses). Mark schemes write the hypotheses as, for example,

$$H_0 : \text{the data are consistent with the distribution } B(4, 0.5)$$

$$H_1 : \text{the data are not consistent with the distribution } B(4, 0.5)$$

Strictly this phrasing is poor — the real null assumption is that the data were drawn from a *population* with the specified distribution, and “consistency” is a property of the sample, not a hypothesis about the population. But this is how the mark schemes phrase it, so it is best to do the same.

Example (Given ratio)

A genetics model predicts that four types of flower should occur in the ratio 9 : 3 : 3 : 1. In a sample of 160 flowers the observed counts are 86, 35, 26 and 13. Test at the 5% significance level whether the data are consistent with the model.

H_0 : the four types occur in the ratio 9 : 3 : 3 : 1. H_1 : the types do not occur in the ratio 9 : 3 : 3 : 1.
 Expected frequencies: $160 \times \frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16}$.

| | Type 1 | Type 2 | Type 3 | Type 4 |
|-----------------------------|--------|--------|--------|--------|
| O_i | 86 | 35 | 26 | 13 |
| E_i | 90 | 30 | 30 | 10 |
| $\frac{(O_i - E_i)^2}{E_i}$ | 0.1778 | 0.8333 | 0.5333 | 0.9000 |

All $E_i \geq 5$, so no pooling needed.

$$X^2 = 0.1778 + 0.8333 + 0.5333 + 0.9000 = 2.444$$

Degrees of freedom: $\nu = 4 - 1 = 3$. Critical value at 5%: $\chi^2_3(5\%) = 7.815$.

Since $2.444 < 7.815$, we do not reject H_0 . There is insufficient evidence to suggest that the flower types do not occur in the ratio 9 : 3 : 3 : 1.

Example (Discrete uniform)

Using the die data from earlier ($O_i = 25, 17, 15, 23, 24, 16$ in 120 rolls), test at the 5% level whether the die is fair.

H_0 : the die is fair (scores follow a discrete uniform distribution); H_1 : the die is not fair. $E_i = 20$ for each score. Contributions: $\frac{25}{20}, \frac{9}{20}, \frac{25}{20}, \frac{9}{20}, \frac{16}{20}, \frac{16}{20}$, so $X^2 = \frac{100}{20} = 5$. With $\nu = 6 - 1 = 5$, the 5% critical value is 11.07. Since $5 < 11.07$, insufficient evidence to suggest the die is not fair.

Example (OCR MEI Further Statistics, June 2023 (part))

An eight-sided dice has its faces numbered 1, 2, ..., 8. A student thinks that the dice may be biased. To investigate this, the student decides to roll the dice 80 times and then carry out a χ^2 goodness of fit test of a uniform distribution. The spreadsheet below shows the data for the test, where some of the values have been deliberately omitted.

| | A | B | C | D |
|---|-------|--------------------|--------------------|--------------------------|
| 1 | Score | Observed frequency | Expected frequency | Chi-squared contribution |
| 2 | 1 | 14 | 10 | 1.6 |
| 3 | 2 | 4 | 10 | 3.6 |
| 4 | 3 | 10 | 10 | 0 |
| 5 | 4 | 15 | 10 | |
| 6 | 5 | 6 | 10 | 1.6 |
| 7 | 6 | 11 | 10 | 0.1 |
| 8 | 7 | 7 | 10 | 0.9 |
| 9 | 8 | | 10 | 0.9 |

- Explain why all of the expected frequencies are equal to 10.
- Determine the missing values in cells B9 and D5.
- Carry out the χ^2 test at the 5% significance level.

(i) Under H_0 each score has probability $\frac{1}{8}$, so every expected frequency is $80 \times \frac{1}{8} = 10$.

(ii) The observed frequencies must total 80, so $B9 = 80 - (14 + 4 + 10 + 15 + 6 + 11 + 7) = 13$. The missing contribution is $D5 = \frac{(15-10)^2}{10} = 2.5$. (Note B9 is consistent with the given contribution: $\frac{(13-10)^2}{10} = 0.9$. ✓)

(iii) H_0 : the scores follow a discrete uniform distribution; H_1 : they do not.

$$X^2 = 1.6 + 3.6 + 0 + 2.5 + 1.6 + 0.1 + 0.9 + 0.9 = 11.2$$

All $E_i = 10 \geq 5$, no pooling. Degrees of freedom: $\nu = 8 - 1 = 7$; critical value $\chi_7^2(5\%) = 14.07$. Since $11.2 < 14.07$, we do not reject H_0 : there is insufficient evidence to suggest that the dice is biased.

Textbook Exercises: [CUP.S] Ch 6 §3; [S3&4] S3 Ch 5

Estimating Parameters from the Data

Sometimes the hypothesis specifies only the *family* of distribution — “the data follow a Poisson distribution” — without giving the parameter. We then estimate the parameter from the data itself (e.g. $\hat{\lambda} = \bar{x}$) before computing expected frequencies. This convenience has a price.

Fact (Degrees of freedom with estimated parameters) — Each parameter estimated from the data imposes one extra constraint on the expected frequencies (we have forced the model to match the data in one more way), so

$$\nu = (\text{cells after pooling}) - 1 - (\text{number of parameters estimated}).$$

Tip

Your conclusion must **not** mention the value of the estimated parameter, because it was never part of the null hypothesis. Write “insufficient evidence to suggest the data do not follow a Poisson distribution”, *not* “. . .do not follow Po(2)”.

Example (Poisson with estimated mean)

The number of calls arriving at a helpdesk was recorded for each of 80 one-minute intervals:

| | | | | | | |
|------------------|----|----|----|----|---|---|
| Calls per minute | 0 | 1 | 2 | 3 | 4 | 5 |
| Frequency | 11 | 21 | 22 | 13 | 9 | 4 |

Test, at the 5% significance level, whether a Poisson distribution is a suitable model.

H_0 : the data can be modelled by a Poisson distribution. H_1 : the data cannot be modelled by a Poisson distribution.

Estimate the mean from the data:

$$\hat{\lambda} = \bar{x} = \frac{0(11) + 1(21) + 2(22) + 3(13) + 4(9) + 5(4)}{80} = \frac{160}{80} = 2.$$

Expected frequencies $E_i = 80 \mathbb{P}(X = i)$ with $X \sim \text{Po}(2)$, the final cell being $\mathbb{P}(X \geq 5)$:

| | | | | | | |
|-------|-------|-------|-------|-------|------|----------|
| Calls | 0 | 1 | 2 | 3 | 4 | ≥ 5 |
| O_i | 11 | 21 | 22 | 13 | 9 | 4 |
| E_i | 10.83 | 21.65 | 21.65 | 14.44 | 7.22 | 4.21 |

The final cell has $E = 4.21 < 5$, so pool it with the previous cell:

| | | | | | |
|-----------------------------|--------|--------|--------|--------|----------|
| Calls | 0 | 1 | 2 | 3 | ≥ 4 |
| O_i | 11 | 21 | 22 | 13 | 13 |
| E_i | 10.83 | 21.65 | 21.65 | 14.44 | 11.43 |
| $\frac{(O_i - E_i)^2}{E_i}$ | 0.0028 | 0.0198 | 0.0055 | 0.1428 | 0.2160 |

$$X^2 = 0.387 \text{ (3 s.f.)}$$

Degrees of freedom: 5 cells, minus 1 for the total, minus 1 for the estimated parameter λ : $\nu = 5 - 1 - 1 = 3$. Critical value: $\chi_3^2(5\%) = 7.815$.

Since $0.387 < 7.815$, we do not reject H_0 . There is insufficient evidence to suggest that the number of calls per minute cannot be modelled by a Poisson distribution.

(Note the conclusion does not mention $\lambda = 2$. Note also how very small X^2 is here — a fit this good is in the far left tail; no test is performed there, but it might prompt a sceptical second look at the data.)

We met the next context in the Poisson notes: bird-watchers recording the number, N , of separate bursts of chaffinch song in 5 minute periods, with sample mean 3.55 and variance 5.6475 over 60 periods. There we judged the Poisson model informally by comparing mean and variance; the chi-squared test makes the comparison of model and data formal.

Example (OCR Further Stats, June 2024 (parts))

The complete results for the 60 periods are shown in the table.

| | | | | | | | | | | |
|-----------|----|---|---|---|----|---|---|---|---|----------|
| n | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ≥ 9 |
| Frequency | 10 | 3 | 7 | 8 | 13 | 6 | 6 | 2 | 5 | 0 |

The bird-watchers carry out a χ^2 goodness of fit test at the 5% significance level.

- State suitable hypotheses for the test.
- Determine the contribution to the test statistic for $n = 3$.
- The total value of the test statistic, obtained by combining the cells for $n \leq 1$ and also for $n \geq 6$, is 9.202, correct to 4 significant figures. Complete the goodness of fit test.

(a) H_0 : N has a Poisson distribution; H_1 : N does not have a Poisson distribution. (Not “Po(3.55)” — the parameter is estimated from the data, so it is not part of the hypothesis.)

(b) Estimate $\hat{\lambda} = \bar{x} = 3.55$. Then

$$E_3 = 60\mathbb{P}(N = 3) = 60 \cdot \frac{e^{-3.55} 3.55^3}{3!} = 12.85, \quad \frac{(O_3 - E_3)^2}{E_3} = \frac{(8 - 12.85)^2}{12.85} = 1.83.$$

(c) After pooling there are 6 cells ($\leq 1, 2, 3, 4, 5, \geq 6$), so

$$\nu = 6 - 1 - 1 = 4$$

(one constraint for the total, one for the estimated parameter λ). Critical value: $\chi_4^2(5\%) = 9.488$. Since $9.202 < 9.488$, we do not reject H_0 : there is insufficient evidence that the number of bursts of song in 5-minute periods does not have a Poisson distribution.

(Agonisingly close — and remember from the Poisson notes that the variance 5.6475 was already suspiciously large compared with the mean 3.55. Forgetting the lost degree of freedom and using $\nu = 5$ would give critical value 11.07 and a much more comfortable-looking, but wrong, margin.)

Textbook Exercises: [CUP.S] Ch 6 §3, Ch 7 §9; [S3&4] S3 Ch 5

Contingency Tables

Definition. A **contingency table** splits a sample according to two attributes simultaneously, tabulating the frequency of each combination. A table with r rows and c columns is called an $r \times c$ contingency table (same convention as matrices).

Think of a contingency table as an observed sample from a *bivariate joint distribution*: each individual carries a pair of values (X, Y) . Recall that X and Y are **independent** iff

$$\mathbb{P}(X = x \cap Y = y) = \mathbb{P}(X = x)\mathbb{P}(Y = y) \quad \text{for all } x, y.$$

We do not know the population distribution of X or Y — but, just as in a goodness-of-fit test, if the observed cell counts differ significantly from what independence would predict, we have reason to suspect that the attributes are associated.

Expected frequencies under independence

Fact — With row totals R_i , column totals C_j and grand total n ,

$$E_{ij} = \frac{(\text{row total}) \times (\text{column total})}{\text{grand total}} = \frac{R_i C_j}{n}.$$

Where does this formula come from? Derive it from the definition of independence.

Estimate the marginal probabilities from the data:

$$\hat{\mathbb{P}}(\text{row } i) = \frac{R_i}{n}, \quad \hat{\mathbb{P}}(\text{column } j) = \frac{C_j}{n}.$$

Under H_0 (independence) the probability of landing in cell (i, j) is the product of the marginals, so the expected frequency is

$$E_{ij} = n \cdot \frac{R_i}{n} \cdot \frac{C_j}{n} = \frac{R_i C_j}{n}.$$

Degrees of freedom

Fact — For a test of independence in an $r \times c$ contingency table,

$$\nu = (r - 1)(c - 1).$$

As with goodness of fit, rows or columns, as appropriate, should be combined so that each expected frequency is at least 5.

Why $(r - 1)(c - 1)$? Count how many cells you are genuinely free to choose.

The expected frequencies are built from the row and column totals, so the row and column totals of the E_{ij} automatically match those of the table. In an $r \times c$ table, once you fill in an $(r - 1) \times (c - 1)$ block of cells, every remaining cell is forced by the totals — so there are $(r - 1)(c - 1)$ free cells. (Equivalently: rc cells, minus 1 for the grand total, minus $(r - 1) + (c - 1)$ for the marginal probabilities estimated from the data.)

Example (Test for independence)

150 people were asked their opinion on a proposed bypass:

| | For | Against | Undecided | Total |
|-------------|-----|---------|-----------|-------|
| Under 40 | 32 | 28 | 15 | 75 |
| 40 and over | 18 | 42 | 15 | 75 |
| Total | 50 | 70 | 30 | 150 |

Test, at the 5% significance level, whether opinion is independent of age group.

H_0 : opinion and age group are independent (not associated). H_1 : opinion and age group are not independent (associated).

Expected frequencies $E_{ij} = \frac{R_i C_j}{150}$, e.g. $E_{11} = \frac{75 \times 50}{150} = 25$:

| E_{ij} | For | Against | Undecided |
|-------------|-----|---------|-----------|
| Under 40 | 25 | 35 | 15 |
| 40 and over | 25 | 35 | 15 |

All $E_{ij} \geq 5$, so no pooling is required. Contributions:

| $\frac{(O-E)^2}{E}$ | For | Against | Undecided |
|---------------------|------------------------|------------------------|-----------|
| Under 40 | $\frac{49}{25} = 1.96$ | $\frac{49}{35} = 1.40$ | 0 |
| 40 and over | 1.96 | 1.40 | 0 |

$$X^2 = 2(1.96 + 1.40) = 6.72$$

Degrees of freedom: $\nu = (2 - 1)(3 - 1) = 2$. Critical value: $\chi^2_2(5\%) = 5.991$.

Since $6.72 > 5.991$, we reject H_0 . There is evidence at the 5% level to suggest that opinion on the bypass is associated with age group.

Example (OCR S3, June 2012)

A study was carried out into whether patients suffering from a certain respiratory disorder would benefit from particular treatments. Each of 90 patients who agreed to take part was given one of three treatments A, B or C as shown in the table.

| Treatment | A | B | C |
|-----------------|----|----|----|
| Number in group | 31 | 25 | 34 |

- (i) It is claimed that each patient was equally likely to have been given any of the treatments. Test at the 5% significance level whether the numbers given each treatment are consistent with this claim.
- (ii) After 3 months the numbers of patients showing improvement for treatments A, B and C were 14, 18 and 25 respectively. By setting up a 2×3 contingency table, test whether the outcome is dependent on the treatment. Use a 5% significance level.
- (iii) If one of the treatments is abandoned, explain briefly which it should be.

(i) This is a goodness-of-fit test. $H_0: p_A = p_B = p_C = \frac{1}{3}$; H_1 : not all equal. Expected frequencies $\frac{90}{3} = 30$ each:

$$X^2 = \frac{1^2 + 5^2 + 4^2}{30} = \frac{42}{30} = 1.4$$

$\nu = 3 - 1 = 2$, critical value $\chi_2^2(5\%) = 5.991$. Since $1.4 < 5.991$, do not reject H_0 : insufficient evidence that the treatments were not equally likely.

(ii) H_0 : outcome is independent of treatment; H_1 : it is not. The contingency table (with the “did not improve” row obtained by subtraction):

| | A | B | C | Total |
|--------------|----|----|----|-------|
| Improved | 14 | 18 | 25 | 57 |
| Not improved | 17 | 7 | 9 | 33 |
| Total | 31 | 25 | 34 | 90 |

Expected frequencies $E_{ij} = \frac{R_i C_j}{90}$: top row 19.63, 15.83, 21.53; bottom row 11.37, 9.17, 12.47 (all ≥ 5). Within each column $|O - E|$ is the same top and bottom (5.63, 2.17, 3.47), so

$$X^2 = 5.63^2 \left(\frac{1}{19.63} + \frac{1}{11.37} \right) + 2.17^2 \left(\frac{1}{15.83} + \frac{1}{9.17} \right) + 3.47^2 \left(\frac{1}{21.53} + \frac{1}{12.47} \right) = 6.74$$

$\nu = (2 - 1)(3 - 1) = 2$, critical value 5.991. Since $6.74 > 5.991$, reject H_0 : there is evidence at the 5% level that the outcome depends on the treatment.

(iii) The proportions improving are $\frac{14}{31} \approx 0.45$, $\frac{18}{25} = 0.72$, $\frac{25}{34} \approx 0.74$: treatment A shows far fewer improvements than expected, while B and C show more. Abandon treatment A.

Textbook Exercises: [CUPS] Ch 6 §1; [S3&4] S3 Ch 5

Yates' Continuity Correction

The chi-squared distribution is *continuous*, but the values O_{ij} are integers, and the E_{ij} (built from integer totals) are fixed rationals — so X^2 can only take certain isolated values. As always when a discrete quantity is approximated by a continuous distribution, we should consider a continuity correction.

For large tables there are so many cells, and so many achievable values of X^2 , that the approximation is fine without one. The worst case is the 2×2 table: there $\nu = 1$, so fixing the margins leaves only *one* free cell, and (check this!) all four cells have the *same* value of $|O - E|$. The achievable values of X^2 are then few and widely spaced.

Definition (Yates' correction). For a 2×2 contingency table, **Yates' continuity correction** replaces the test statistic by

$$X^2 = \sum \frac{\left(|O_i - E_i| - \frac{1}{2}\right)^2}{E_i},$$

i.e. each difference is shrunk towards zero by half before squaring. This correction is used in the special case of a 2×2 table.

Remark. Only apply the correction when $|O - E| > \frac{1}{2}$. If $|O - E| \leq \frac{1}{2}$ the observed values are already as close to expected as integer data can be: the fit is essentially perfect and no test is needed (blindly applying the formula would *increase* the discrepancy, which is nonsense). Some books write the corrected statistic with $\max(|O - E| - \frac{1}{2}, 0)$ for this reason.

Example (2×2 with Yates' correction)

In a trial, 80 patients were randomly assigned a drug or a placebo:

| | Recovered | Did not recover | Total |
|---------|-----------|-----------------|-------|
| Drug | 26 | 14 | 40 |
| Placebo | 17 | 23 | 40 |
| Total | 43 | 37 | 80 |

Test, at the 5% significance level, whether recovery is independent of treatment.

H_0 : recovery is independent of treatment. H_1 : recovery is not independent of treatment.

Expected frequencies: $E_{11} = \frac{40 \times 43}{80} = 21.5$, and similarly

| E_{ij} | Recovered | Did not recover |
|----------|-----------|-----------------|
| Drug | 21.5 | 18.5 |
| Placebo | 21.5 | 18.5 |

All $E \geq 5$. Every cell has $|O - E| = 4.5$ (as always in a 2×2 table the four differences agree, and $4.5 > 0.5$ so the correction applies). This is a 2×2 table, so use Yates' correction with $|O - E| - \frac{1}{2} = 4$:

$$X^2 = \frac{4^2}{21.5} + \frac{4^2}{18.5} + \frac{4^2}{21.5} + \frac{4^2}{18.5} = 2 \left(\frac{16}{21.5} + \frac{16}{18.5} \right) = 3.218$$

Degrees of freedom: $\nu = (2 - 1)(2 - 1) = 1$. Critical value: $\chi_1^2(5\%) = 3.841$.

Since $3.218 < 3.841$, we do not reject H_0 . There is insufficient evidence to suggest that recovery is associated with the treatment received.

(Without the correction we would have had $X^2 = 2 \left(\frac{20.25}{21.5} + \frac{20.25}{18.5} \right) = 4.07 > 3.841$ and the opposite conclusion — the correction genuinely matters near the critical value.)

Example (OCR S3, June 2007)

The students in a large university department take a trial examination some time before the proper examination. A random sample of 60 students took both examinations during a particular course. 42 students passed the trial examination, 36 passed the proper examination and 13 failed both examinations.

(i) Copy and complete the following contingency table.

| | | Proper | | |
|-------|-------|--------|------|-------|
| | | Pass | Fail | Total |
| Trial | Pass | | | 42 |
| | Fail | | 13 | |
| | Total | 36 | | 60 |

(ii) Carry out a test of independence at the $\frac{1}{2}\%$ level of significance.

(i) $Fail-Pass = 36 - ?$: work from the given entries. Trial Fail total = $60 - 42 = 18$, so $Fail-Pass = 18 - 13 = 5$; $Pass-Pass = 36 - 5 = 31$; $Pass-Fail = 42 - 31 = 11$; Proper Fail total = 24.

| | | Proper | | |
|-------|-------|--------|------|-------|
| | | Pass | Fail | Total |
| Trial | Pass | 31 | 11 | 42 |
| | Fail | 5 | 13 | 18 |
| | Total | 36 | 24 | 60 |

(ii) H_0 : trial and proper examination results are independent; H_1 : they are not. Expected frequencies $E_{ij} = \frac{R_i C_j}{60}$: 25.2, 16.8, 10.8, 7.2 (all ≥ 5). Every cell has $|O - E| = 5.8$. This is a 2×2 table, so apply Yates' correction with $|O - E| - \frac{1}{2} = 5.3$:

$$X^2 = 5.3^2 \left(\frac{1}{25.2} + \frac{1}{16.8} + \frac{1}{10.8} + \frac{1}{7.2} \right) = 9.29$$

$\nu = (2 - 1)(2 - 1) = 1$; critical value $\chi_1^2(0.5\%) = 7.879$. Since $9.29 > 7.879$, reject H_0 : there is evidence, even at the $\frac{1}{2}\%$ level, that trial and proper examination results are not independent. (Reassuring — a trial examination that told you nothing about the real one would be a waste of everyone's time.)

Remark (Fisher's exact test). The chi-squared test for a 2×2 table quietly assumes the row and column totals are fixed. If they really are — pick the cell in the top-left as the variable and the hypergeometric distribution gives the exact probability of each possible table; summing over tables at least as extreme as the one observed is **Fisher's exact test**, no approximation needed. Yates' correction is precisely an attempt to nudge the chi-squared approximation towards these exact values. In real life margins are rarely all fixed: a medical trial fixes the group sizes but not the recovery counts (a *comparative trial*); a survey classifying people by two attributes fixes only the grand total (a *double dichotomy*); fixing everything is an *independence trial*. The chi-squared machinery is used for all three at this level.

Textbook Exercises: [CUPS] Ch 6 §1; [S3&4] S3 Ch 5; [Toller] Ch 9