

①

χ^2 Tests

Pronounced "Khi squared"

χ^2 tests aim at testing values of parameters / independence of categorical variables (as opposed to continuous variables and T-tests you saw in Math 181A for quantitative ideas)

Roughly speaking, a variable is "categorical" (or quantitative) if it has a discrete (or even finite) set of possible outcomes.

→ Think to bins

Examples where χ^2 tests are useful:

- With a sample of 100 San Diegans, you wonder whether the distribution of ages (grouped by categories) in S.D. is the same as in the U.S. overall
- With a sample of 100 San Diegans and 100 New Yorkers, you wonder if the age distribution in S.D. is the same as in N.Y.
- With a sample of 100 San Diegans, you try to figure out if there is a link between sex (δ / φ) and eye color (Brown, Blue, Green, ...).

(2) Before carrying on, we need to model "binned variable" generically, when there are $t \geq 2$ bins.

This is a generalisation of the binomial distribution, that can be seen as $t=2$ bins "0" and "1".

10.2

Def: (Multinomial Distribution)

The i^{th} bin

Let X_i denote the number of times the outcome r_i occurs ($1 \leq i \leq t$) in a series of n independent Bernoulli trials with $P(r_i) = p_i$.

Then the vector (X_1, \dots, X_t) has a Multinomial distribution, and

$$P(X_1=k_1, \dots, X_t=k_t) = \frac{n!}{k_1! \times \dots \times k_t!} p_1^{k_1} \cdots p_t^{k_t},$$

Sometimes denoted
by $\binom{n}{k_1 \dots k_t}$: "Multinomial coefficient"
for all $0 \leq k_1, \dots, k_t \leq n$ with $\sum_{i=1}^t k_i = n$.

Rb: If $\sum_{i=1}^t k_i \neq n$, then

$$P(X_1=k_1, \dots, X_t=k_t) = 0$$

Condition required because we put n elements (exactly) in the t bins.

(3)

Example: The positions on a roulette wheel are divided into 3 colors: red, black and green:

$$\left. \begin{array}{l} 18 \text{ red} \\ 18 \text{ black} \\ 2 \text{ green} \end{array} \right\}$$

If the wheel is fair, determine

- 1) The probability that in 7 independent spins of the wheel, the ball lands in a red slot 4 times and in a black slot 3 times.
- 2) The chance the ball in a red slot at least three times
(in 7 spins)

Sol. 1) $m = 7, k_1 = 4, k_2 = 3, k_3 = m - k_1 - k_2 = 0$

$$P_1 = \frac{18}{38} = P_2, P_3 = \frac{2}{38} \cdot \text{ We find } 0.1873$$

2) This is modelled by a binomial ($t=2$) with parameters
 $m = 7$ and $p = \frac{18}{38}$. We find 0.728

(4)

Prop.: If $X = (X_1, \dots, X_t)$ has a multinomial distribution, then all the marginals of X have binomial distribution and

$$E(X) = \begin{pmatrix} m p_1 \\ \vdots \\ m p_t \end{pmatrix}$$

$$\text{Var}(X) = \begin{pmatrix} m p_1(1-p_1) & -m p_2 p_1 & \cdots & -m p_t p_1 \\ -m p_1 p_2 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ -m p_1 p_t & \cdots & -m p_{t-1} p_t & m p_t(1-p_t) \end{pmatrix}$$

Proof: For all $1 \leq i \leq t$, $(X_i) = X$ can be seen as the sum $X^{(1)} + \dots + X^{(n)}$ of independent multinomial variables with parameters $m=1$, and p_1, \dots, p_t . Hence, $E(X) = m E(X^{(1)}) = \begin{pmatrix} m p_1 \\ \vdots \\ m p_t \end{pmatrix}$ and

$$\text{Var}(X) = m \text{Var}(X^{(1)}). \text{ But for all } 1 \leq i \leq t,$$

$$\text{Cov}(X_i^{(1)}, X_j^{(1)}) = \begin{cases} 0 & \text{if } i \neq j \\ p_i & \text{if } i = j \end{cases}$$

so that $\begin{cases} \text{Var}(X_i) = p_i - p_i^2 = p_i(1-p_i) \\ \text{Cov}(X_i, X_j) = 0 - p_i p_j \end{cases}$, which is the result.

(5)

10.3 Goodness-of-Fit χ^2 Test

To check if a sample made of counts (categorical)
 "fits" some multinomial (binned) you think it should fit.

Setting: Independent samples, sample size large enough.

Hypotheses: $H_0: p_1 = p_{1_0}, p_2 = p_{2_0}, \dots, p_t = p_{t_0}$

p_i ↑
 The true probability
 of bin i . p_{i_0} ↑
 The expected probability
 of bin i .

$H_1: p_i \neq p_{i_0}$ for at least one i

Idea: • If p_i truly is equal to p_{i_0} (i.e. under H_0), you expect to count approximately $m \times p_{i_0}$ elements, among m samples, ending up in bin i .

• If not, this count should deviate from mp_{i_0} .

Take away formula to encode this (i.e. test statistic)

$$D_n = \sum_{\text{bins}} \left(\frac{\text{Expected counts} - \text{Observed count}}{\text{Expected count}} \right)^2$$

To take bins into account altogether → bins Expected counts → Renormalization factor

For two bins
 not to compensate

- ⑥ Under H_0 , D should be small
Under H_1 , D should be large (and go to infinity)

We formalize this idea through the following result, that states how multinomials behave.

Thm: If $X = (X_1, \dots, X_t)$ has multinomial distribution with parameters n , and p_1, \dots, p_t , and that $m p_i \geq 5$ ($1 \leq i \leq t$), then

$$D = \sum_{i=1}^t \frac{(X_i - m p_i)^2}{m p_i} \stackrel{\text{has approximate distribution}}{\sim} \chi^2_{t-1}$$

χ^2 statistic Observed counts Expected counts

Proof: From the central limit theorem,

$$\sqrt{n} \left\{ \begin{pmatrix} X_1 \\ \vdots \\ X_t \\ \hline n \end{pmatrix} - \begin{pmatrix} p_1 \\ \vdots \\ p_t \end{pmatrix} \right\} \xrightarrow[n \rightarrow \infty]{D} N_t(0, \Sigma)$$

where $\Sigma = \begin{pmatrix} p_1(1-p_1) & & & & \\ & \ddots & & & \\ & & p_j p_{j+1} & & \\ & & & \ddots & \\ & & & & p_t(1-p_t) \end{pmatrix}$

⑦ Δ Notice that Σ is not an invertible matrix (sum of columns is equal to 0), so that we cannot renormalize the left-hand side properly.

Let us consider $Y = \begin{pmatrix} X_1 \\ \vdots \\ X_{t-1} \end{pmatrix}$ (subvector of X). Then we have

$$\sqrt{m} \left(\frac{Y}{m} - \begin{pmatrix} p_1 \\ \vdots \\ p_{t-1} \end{pmatrix} \right) \xrightarrow[m \rightarrow \infty]{\mathcal{D}} N_{t-1}(0, \Sigma^*)$$

where Σ^* is the upper-left $(t-1) \times (t-1)$ submatrix of Σ .

One may check that Σ^* is invertible, and that

$$(\Sigma^*)^{-1} = \begin{pmatrix} \frac{1}{p_1} + \frac{1}{p_T} & \frac{1}{p_T} & \cdots & \frac{1}{p_T} \\ \frac{1}{p_T} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \frac{1}{p_T} \\ \frac{1}{p_T} & \cdots & \frac{1}{p_T} & \frac{1}{p_{T-1}} + \frac{1}{p_T} \end{pmatrix}$$

$$\text{Furthermore, } D = \left\| \sum_{n=1}^{t-1} \sqrt{n} \left(\frac{Y}{m} - \begin{pmatrix} p_1 \\ \vdots \\ p_{t-1} \end{pmatrix} \right) \right\|^2 = m \left(\frac{Y}{m} - \begin{pmatrix} p_1 \\ \vdots \\ p_{t-1} \end{pmatrix} \right)^T (\Sigma^*)^{-1} \left(\frac{Y}{m} - \begin{pmatrix} p_1 \\ \vdots \\ p_{t-1} \end{pmatrix} \right)$$

And since $Z_n \xrightarrow[m \rightarrow \infty]{\mathcal{D}} N_{t-1}(0, I_{t-1})$, we get the result.

(8)

Rk: The proof only prove an asymptotic convergence. Condition " $mp_i \geq 5$ for all $1 \leq i \leq t$ " is here to account for this. Always have in mind that this is an approximate statement.

From this, we derive a goodness-of-fit test:

Coro: (χ^2 goodness-of-fit test)

Let k_1, \dots, k_t be the observed counts for the outcomes π_1, \dots, π_t .

At the α level of confidence for the test

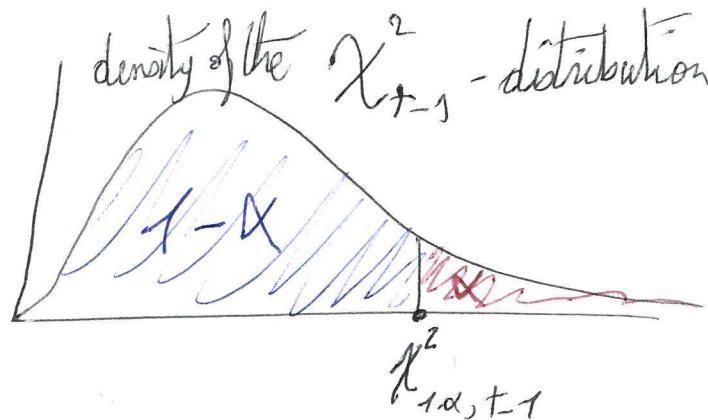
$$\left\{ \begin{array}{l} H_0: p_i = p_{i_0} \quad \text{True proba. of bin } i \\ \text{vs} \\ H_1: p_i \neq p_{i_0} \quad \text{for at least one } i. \end{array} \right. \quad \begin{array}{l} \xrightarrow{\text{tested proba. of bin } i} \\ \xleftarrow{\text{vs}} \end{array}$$

with $mp_i \geq 5$ for all $i \leq t$, reject H_0 when

$$D = \sum_{i=1}^t \frac{(k_i - mp_{i_0})^2}{mp_{i_0}} \geq \chi^2_{1-\alpha, t-1}$$

Quantile of order $1-\alpha$ of the χ^2_{t-1} distribution

Rk: The quantile of order $1-\alpha$ of a distribution is the value that has $\begin{cases} 1-\alpha & \text{area to its left} \\ \alpha & \text{--- right} \end{cases}$



Example: Testing Mendel's theory on alleles for peas.

Mendel's theory for peas:

- An allele R/r codes roundness
- An allele Y/y codes yellowness
- R and Y are dominant
- r and y are recessive

In a cross, this theory predicts that observed phenotypes come with frequencies

- * $\frac{9}{16}$: Yellow Round \rightarrow Phenotype #1
- * $\frac{3}{16}$: non-yellow Round \rightarrow #2
- * $\frac{3}{16}$: Yellow non-round \rightarrow #3
- * $\frac{1}{16}$: non-yellow non-round \rightarrow #4

(10)

After 100 peas bred, you get:

Phenotype	Observed counts
# 1	62
# 2	24
# 3	9
# 4	5

Do these data provide evidence that Mendel's theory is wrong?

→ This amounts to test if $H_0: (P_1, P_2, P_3, P_4) = \left(\frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16}\right)$ is true or false. (here, $n=100$)

We can display our data in the following way:

Phenotype	Observed count	Expected count under H_0
# 1	62	$100 \times \frac{9}{16} = 56.25$
# 2	24	$100 \times \frac{3}{16} = 18.75$
# 3	9	$100 \times \frac{3}{16} = 18.75$
# 4	5	$100 \times \frac{1}{16} = 6.25$

⚠ DO NOT round the expected counts.

$$\text{We get } \chi^2 = D = \frac{(62 - 56.25)^2}{56.25} + \frac{(24 - 18.75)^2}{18.75} + \frac{(9 - 18.75)^2}{18.75} + \frac{(5 - 6.25)^2}{6.25}$$

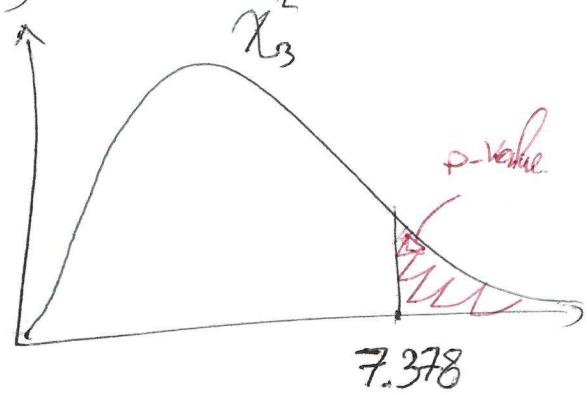
$$\approx 7.378.$$

Since all the expected counts ($n P_{i_0}$) are ≥ 5 we can apply the χ^2 -test, with $df = 4 - 1 = 3$

(11)

On the table of the χ^2 distribution, we see that the area to the right of 7.378 (= p-value) satisfies

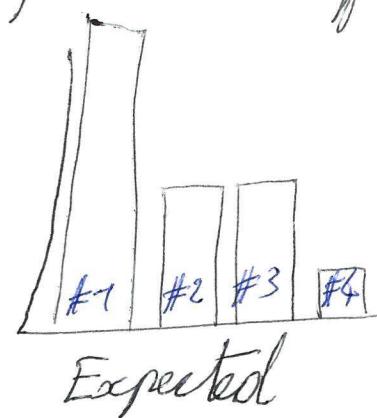
$$0.05 \leq p\text{-value} \leq 0.1.$$



Hence, at the level $\alpha = 5\% = 0.05$, we do not reject H_0 .

Rk: What is this test really doing?

In essence, a χ^2 goodness-of-fit test tells you if the observed data differ from the data expected on some theory by an amount that exceeds the type of variation we expect from random effect.



Said differently, $\chi^2 = D$ tells you if the two histograms (barplots) are "close enough" to each other (null) or not (alternative).

(12)

10.4 Goodness-of-fit χ^2 test: Parameters unknown

In some situations one has reason to believe that some observed counts should be distributed as a multinomial (and wants to test this) but has no a priori parameter values (p_1, \dots, p_t) . Hence, the raw goodness-of-fit χ^2 test cannot be applied.

A strategy to overcome this issue is to replace the expected counts $n p_i$ (not available anymore) by estimated expected counts $n \hat{p}_i$, where \hat{p}_i is the estimated proportion of outcome r_i .

Doing so, we'll have to take into account the extra randomness of \hat{p}_i into the number of degrees of freedom of the χ^2 distribution.

Thm: If $X = (X_1, \dots, X_t)$ has multinomial distribution with $s \leq t$ unknown parameters (= parameters you need to estimate) let \hat{p}_i denote the estimated probability of outcome r_i based on X .

Then, if $n \hat{p}_i \geq 5$ for all $1 \leq i \leq t$

$$\chi^2_D = \sum_{i=1}^t \frac{(X_i - n \hat{p}_i)^2}{n \hat{p}_i} \underset{n \rightarrow \infty}{\sim} \chi^2_{t-1-s}$$

Proof: Not covered here, it has the same flavor as for \hat{p}_i replaced by p_i .

(13)

Rk: As a corollary, we get a test procedure to test if a sample comes from some (partially unspecified) multinomial distribution.

- Next section gives an important application of this result

10.5 Contingency Tables and the χ^2 Independence Test

Generalizing the above, where we studied a single variable X with t possible outcomes, we now move to the study of two categorical variables (X, Y) . The main idea here is to determine (test) whether or not X and Y are independent. (X, Y) describes two traits of a single individual.

Example:

- X : Gender, Y : Eye color
- X : Sexual orientation, Y : State of birth in the US.

In general, say that: X has possible outcomes A_1, \dots, A_n

$$Y \quad = \quad B_1, \dots, B_c$$

We setup notation

n and c have no reason to be equal.

- $P(X = A_i \text{ and } Y = B_j) = p_{ij}$, with $\sum_{i=1}^n \sum_{j=1}^c p_{ij} = 1$
- $P(X = A_i) = p_i$
- $P(Y = B_j) = q_j$

$$\begin{cases} \text{Clearly, } p_i = \sum_{j=1}^c p_{ij} \\ q_j = \sum_{i=1}^n p_{ij} \end{cases}$$

14

Say we want to test if X and Y are independent (H_0). If H_0 is true, then $P(X=A_i, Y=B_j) = P(X=A_i) \cdot P(Y=B_j)$, by definition of independence. In other words, $p_{ij} = p_i q_j$ for all (i, j) .

Say we estimate, from sample, p_1, \dots, p_n and q_1, \dots, q_c with $\hat{p}_1, \dots, \hat{p}_n$ and $\hat{q}_1, \dots, \hat{q}_c$, we can test if $p_{ij} = p_i q_j \approx \hat{p}_i \hat{q}_j$.

Thm: (χ^2 Independence test)

Suppose that n observations are partitioned by the events A_1, \dots, A_n and also by B_1, \dots, B_c . Let k_{ij} denote the number of observations in the sample that belong to $A_i \cap B_j$.

To test H_0 : The A_i 's are independent from the B_j 's

vs

H_1 : They're not independent

The null is rejected at the α level of significance if

$$\chi^2 = D = \sum_{i=1}^n \sum_{j=1}^c \frac{(k_{ij} - n \hat{p}_i \hat{p}_j)^2}{n \hat{p}_i \hat{p}_j} \geq \chi^2_{1-\alpha, (n-1)(c-1)}$$

Rmk: This test is valid only if $n \hat{p}_i \hat{p}_j \geq 5$ for all (i, j) .

(15)

Prob: We are considering the test of multinomial distributions for $n \times c$ possible outcomes. In the process, we estimated :

$\hat{p}_1, \hat{p}_2, \dots, \hat{p}_n \rightarrow$ Which makes only $\boxed{n - 1}$ estimates,
since we have $\sum_{i=1}^n \hat{p}_i = 1 \Rightarrow \hat{p}_n = 1 - \sum_{i=1}^{n-1} \hat{p}_i$

$\hat{q}_1, \dots, \hat{q}_{c-1} \rightarrow \boxed{c-1}$ estimates b.

In total, we estimated $s = (n-1) + (c-1)$ parameters. And at the end of the day, we get the degrees of freedom equal to

$$f - s - 1 = n \times c - (n + c - 2) - 1 \\ = (n-1)(c-1)$$

Example: Testing association between smoking and cancer.

Pyrobenzene is a major component of cigarette smoke. Researchers injected rats with different levels of pyrobenzene, and looked for tumor development. In 230 rats injected, they got these data:

	No tumor	One tumor	≥ 2 tumors	
Control	74	5	7	
Low dose	63	12	5	
High dose	45	15	10	
\uparrow Observed counts				$\downarrow c = 3$

Rk: Usually, it helps to display data in a two-way table, also called contingency table.

(16)

Do these data show evidence of an association between pyrobenzene exposure and tumor development for rats?

We compute the expected counts $m \hat{P}_i \hat{P}_j$.

Rk: Useful formula : $\boxed{\hat{m} \hat{P}_i \hat{P}_j = \frac{(\text{Row total}) \times (\text{Column total})}{m}}$

Here, we get the expected contingency table as follows

	No tumor	One tumor	≥ 2 tumors
Control	63.3	11.13	5.57
Low dose	63.3	11.13	5.57
High dose	55.39	9.74	4.87

Expected counts

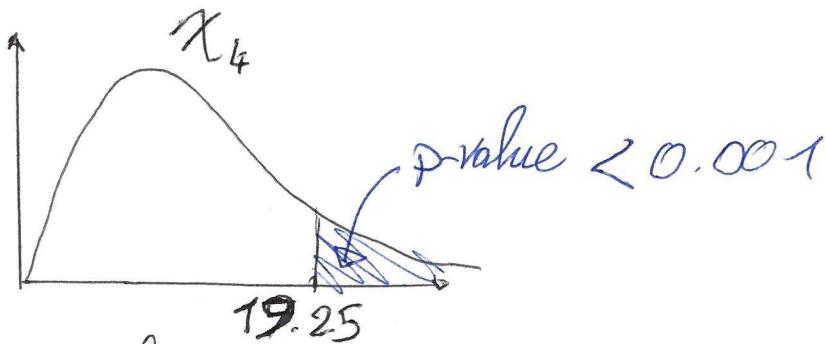
⚠ Do not round
expected counts!

We get $\chi^2 = D = \sum_{\text{cells}} \frac{(O \text{bserved count} - E \text{xpected count})^2}{E \text{xpected count}}$

$$= 19.25$$

looking at a statistical table of the χ^2 distribution with $df = (n-1) \times (c-1) = (3-1) \times (3-1) = 4$, we get p-value < 0.001.

(17)



Hence, as $p\text{-value} < 0.001 < \alpha = 0.05$, we reject the null:

At the level of significance $\alpha = 5\%$, there is evidence of an association between pyrobenzene exposure and tumor development in rats.

⚠ A χ^2 test of independence tells you if there is an association or not, but it doesn't inform you on the sense of this association!! In the example above, pyrobenzene, the χ^2 test might actually have detected that pyrobenzene cures cancer...!

Q: What is this test really doing?

It basically compares how distributions look relative to each other

