# Chi-square test of independence

**Set up:** Count data on two categorical variables (or factors) $A$ and $B$ obtained from a sample of $n$ subjects. Suppose the categories of $A$ are $i = 1, \ldots, k$, and the categories of $B$ are $j = 1, \ldots, m$. The data are arranged in a $k \times m$ table. Let $O_{ij}$ = observed count in $(i, j)$-th cell.



$$P[X = i, Y = j] = p_{ij}$$

prob. that a randomly selected individual falls in category of $A$ and $j$th category of $B$

**Hypotheses:** $H_0$ : $A$ and $B$ are independent (i.e., are not associated), vs., $H_1$ : $A$ and $B$ are not independent (i.e., are associated). If there is an association, the value one variable depends (at least to some extent) on the value of the other.

**Example:** The table below shows 695 children under 15 years of age are cross-classified according to ethnic group and hemoglobin level. Is hemoglobin level associated (related) to ethnicity?

$P[X=1, Y=1]$ ? $= P[X=1] \cdot P[Y=1]$ ? indep.

| Ethnic Group | Hemoglobin Level (g/100 ml) | | | Total | Proportion |
|---|---|---|---|---|---|
| | ≥10 | 9.0 – 9.9 | <9.0 | | |
| A | 80 | 100 | 20 | 200 | $\frac{200}{695} = \hat{P}(X=1)$ |
| B | 99 | 190 | 96 | 385 | $385/695 = \hat{P}(X=2)$ |
| C | 70 | 30 | 10 | 110 | $110/695 = \hat{P}(X=3)$ |
| Total | 249 | 320 | 126 | 695 | |
| Proportion | $\frac{249}{695}$ | $\frac{320}{695}$ | $\frac{126}{695}$ | | |

$\hat{P}[X=1]$

- If He level is not associated to ethnicity, then the proportion of subjects in population that fall a He group does not depend on ethnicity, i.e., it is the same for each ethnicity group, and vice versa.

$P[\text{He group} = 1 | \text{ethnicity} = A] = P[\text{He group} = 1]$

13/10

$H_0$: The two variables are indep., $H_1$: NOT indep.

To do a chi-square test, we need the expected counts $E_{ij}$ assuming that $H_0$ is true. Let $X$ and $Y$ indicate respective categories of $A$ and $B$ in which a randomly selected subject from the population falls. When $A$ and $B$ are independent,

$$\overset{H_0}{\underset{\downarrow}{}}$$

$$P(X = i, Y = j) = P(X = i)P(Y = j) \text{ for all } i, j.$$

- $P(X = i)$ is estimated as
$$\frac{i\text{-th row total}}{n} = \frac{\sum_\ell O_{i\ell}}{n}$$

- $P(Y = j)$ is estimated as
$$\frac{j\text{-th column total}}{n} = \frac{\sum_\ell O_{\ell j}}{n}$$

- Assuming independence, $P(X = i, Y = j)$ is estimated
$$\longrightarrow \hat{P}(X=i) \cdot \hat{P}(Y=j)$$

- Assuming independence, $E_{ij}$ is estimated as
$$\hat{E}_{ij} = (n) \hat{P}(X=i) \cdot \hat{P}(Y=j)$$
$$= \frac{(i\text{-th row total})(j\text{-th col. total})}{n}$$

**Test statistic:**

$$\chi^2 = \sum_i \sum_j \frac{(O_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}} \sim \chi_\nu^2 \text{ when all } \hat{E}_{ij} \geq 5, \text{ when } H_0 \text{ is true.}$$

**Degrees of freedom:**

Read the notes and this:
$$\nu = (k-1)(m-1)$$
$$\begin{array}{cc} \uparrow & \uparrow \\ \text{\# categories} & \text{\# categories of } B. \\ \text{of } A & \end{array}$$

14/19

**Example (continued):** The expected counts for all cells (in parenthesis below next to the observed counts) are:

| Ethnic Group | Hemoglobin Level (g/100 ml) | | | |
|---|---|---|---|---|
| | $\geq 10$ | 9.0 - 9.9 | < 9.0 | Total |
| A | 80 ( ? ) | 100 (92.09) | 20 (36.26) | (200) |
| B | 99 (137.94) | 190 (177.27) | 96 (69.80) | 385 |
| C | 70 (39.41) | 30 (50.65) | 10 (19.94) | 110 |
| Total | 249 | 320 | 126 | (695) |

$$? = \frac{(200)(249)}{695} = ?$$

$$\chi^2 = 67.8$$

Hmm

# R code:

```
> x <- c(80, 100, 20, 99, 190, 96, 70, 30, 10)
> xmat <- matrix(x, byrow=T, ncol=3)
> xmat
     [,1] [,2] [,3]
[1,]   80  100   20
[2,]   99  190   96
[3,]   70   30   10
> chisq.test(xmat)

        Pearson's Chi-squared test
```

$$(k-1)(m-1)$$

Take $\alpha = 0.05$ ✓
$\alpha = 0.01$ ✓

```
data: xmat
X-squared = 67.8015, df = 4, p-value = 6.606e-14
```

conclusion: Reject $H_0$

=> Conclude that its the two variables are associated.

Cannot deduce any causal relationship — just that the two variables are related.

16/19
>

# Chi-Square test of Homogeneity

Often we are interested in comparing different populations with respect to a variable of interest, e.g., are the populations of carriers and non-carriers of a certain antigen _homogeneous_ with respect to blood type?

_The dist. of blood group type is same for the two populations_

**Example:** A sample of 150 carriers of a certain antigen and a sample of 500 non-carriers showed the following blood group distributions:

_with 3 populations_
_H0: $p_1 = p_2 = p_3$_

| Blood Group | Carriers | Non-Carriers | Total |
|---|---|---|---|
| O | 72 | 230 | 302 |
| A | 54 | 192 | 246 |
| B | 16 | 63 | 79 |
| AB | 8 | 15 | 23 |
| Total | 150 | 500 | 650 |

$P[BG = O \mid carrier]$
$= P[BG = O \mid non\ carrier]$
$= P[BG = O]$
this equality should also hold for all other BG types

Are carriers and non-carriers similar with respect to blood group distributions?

$\longrightarrow$ Do a chi-square test of independence we know how to do — R&C clark 19

(for P value)

NOTE: Test of homogeneity is a special case of independence...

# Test of Homogeneity vs. Test of Independence

Comparing the layout of this table with the table for the test of independence, we see that the two layouts are Thus, mathematically the tests of homogeneity and independence are exactly the same. So, the same formulas apply. However, there are some key conceptual differences.

**Sampling procedure:**

- *Test of independence:* <u>one</u> overall sample is collected first and then each observation is classified by levels of the two variables. So, neither row nor column totals are fixed in advance.

- *Test of homogeneity:* <u>several</u> samples are collected from several populations with each sample size fixed in advance. After collecting these pre-determined # of observations, each is classified by *various levels of one variable. So, in the above example, ............. column totals are fixed.*

# Number of variables:

- *Test of independence:* **two** variables. → Blood group type.

- *Test of homogeneity:* **one** variable. The column/row representing "population" is fixed due to the sampling process.

# Hypotheses:

- *Test of independence:* H₀: A and B are indep.

- *Test of homogeneity:* H₀: 100 The populations are same w.r.t. the variable $Y$ of interest.

  ↑ Question of independence doesn't arise because only one variable being measured.

Recall: X and Y are indep. if

$$P[X=i, Y=j] = P[X=i] \cdot P[Y=j] \quad \text{for all } i,j$$

This condition is equivalent to:

$$P[X=i \mid Y=j] = P[X=i] \quad \text{for all } i,j$$

$$\cong$$

$$P[Y=j \mid X=i] = P[Y=j] \quad \text{for all } i,j$$

# Nonparametric Tests

**Issue:** We would like test hypothesis on **center** of a distribution (one-sample problem) or compare centers of two distributions (two-sample problem). But the distributions are not normal — e.g., they are skewed or data has outliers.

*[handwritten: mean is not a good meas. of center.]*

**Q:** Why not simply use large-sample $z$ test?

*[handwritten: Fine if mean is a meas. of center. → median of dist.]*

**Alternative measure of "center":**

*[handwritten: → median of dist.]*

**Nonparametric procedures:**

- Typically they don't assume a specific distributional form (e.g., normal); only that the distribution is continuous *[handwritten: X is a cont. r.v.]*. Some procedures assume that the distribution is symmetric.

- More broadly applicable than **parametric procedures** that assume specific distributional form.

- Use these when the distributional assumption behind a parametric procedure is clearly violated.
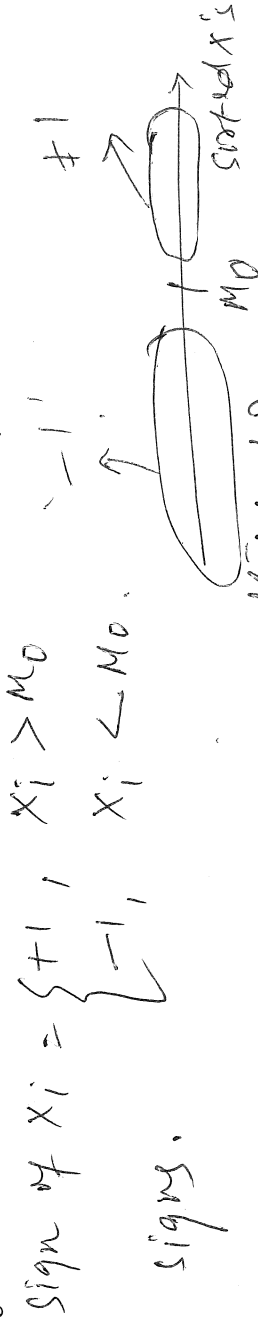
# Sign test

**Data:** $X_1, \ldots, X_n$ — i.i.d. sample from $X$.

**Hypotheses:** $H_0 : M = M_0$ vs. $H_1 :$ one of three possibilities,
$M > M_0$ or $M < M_0$ or $M \neq M_0$

*(pop. median / null value)*

*If $X$ is cont.*
$$P[X = M_0] = 0.$$

**Signs:** Remove the $X$'s that are equal to $M_0$ and reduce the sample size accordingly.

$n^* = \#$ obs. that are not equal to $M_0$.

$$\text{sign of } x_i = \begin{cases} +1, & x_i > M_0 \\ -1, & x_i < M_0 \end{cases}$$

**Test statistic:**

$S = \#$ positive signs.

**Null distribution:**

If $H_0$ is true:
sign of $x_i \sim \text{Bernoulli}(p = \tfrac{1}{2})$.

$\Rightarrow S \sim \text{Bin}(n^*, p = \tfrac{1}{2})$.

**When to reject $H_0$?**

- $H_1 : M > M_0$:   $S$ is too large (compared with $\tfrac{n^*}{2}$)
  
  $P[\text{Bin}(n^*, \tfrac{1}{2}) \geq S_{obs}]$   $\leftarrow$ p-value

- $H_1 : M < M_0$:   $S$ is too small

  $P[\text{Bin}(n^*, \tfrac{1}{2}) \leq S_{obs}]$

- $H_1 : M \neq M_0$:   $S$ is either too large or too small.

  $2 \min \{ \text{the two} \}$   2-sided p-value

# R code:

```
# Time between keystrokes data from Example 10.9

x <- c(0.24, 0.22, 0.26, 0.34, 0.35, 0.32, 0.33, 0.29,
       0.19, 0.36,
       0.30, 0.15, 0.17, 0.28, 0.38, 0.40, 0.37, 0.27)

# Histogram and boxplot

par(mfrow=c(1,2)) # 2 plots in 1 row

hist(x)
qqnorm(x)
qqline(x)

library(nortest)
```

```
> shapiro.test(x)

        Shapiro-Wilk normality test

data:  x
W = 0.9611, p-value = 0.6233
>
```
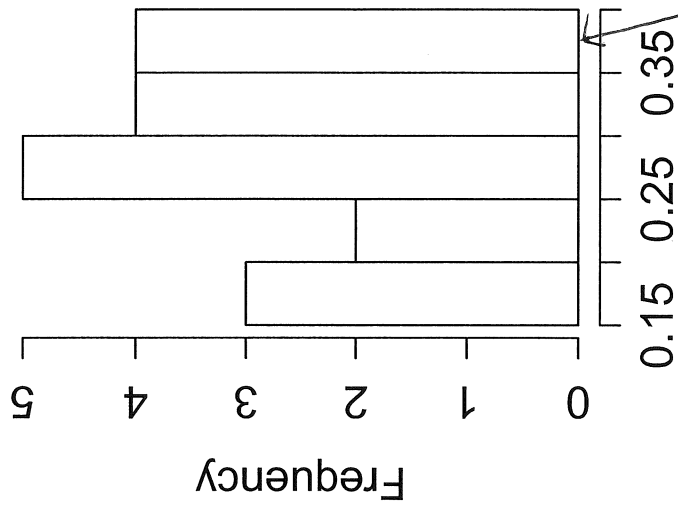
Normality appears reasonable from this as well Q-Q plot.

Histogram of x

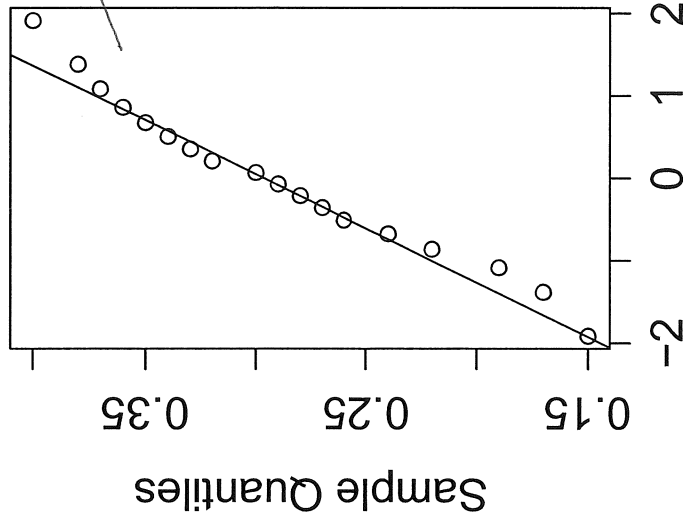Normal Q-Q Plot

normality appears reasonable.

suggests distribution with ... quantiles

POB  Both models may be reasonable for these data. With a larger sample size, we may be able to distinguish b/w its two distributions, with a better distribution with a ...