# TESTS ON CATEGORICAL DATA I

Sorana D. BOLBOACĂ

# OUTLINE

2×2 Tables: Contingency Tables
Risks and Odds in Medical Decisions
2×2 Tables: Tests of Associations

# 2×2 CONTINGENCY TABLES

- Nominal scale: dichotomiale: 2-by-2 contingency Table)
- Ordinal scale: r-by-c contingency table
- Absolute frequency (number of events per category)
- 2  2 contingency table: 4 categories
  - TP = true pozitive
  - FP = false pozitive
  - FN = false negative
  - TN = true negative

# 2×2 Contingency Tables

|  | Caries + | Caries - | Total |
|---|---|---|---|
|  | Caries + | Caries - | Total |
| Fluoridated + | TP = 77 | FP = 29 | = 77+29 = 106 |
| Non-Fluoridated - | FN = 95 | TN = 31 | = 95+31 = 126 |
| Total | =77+95=172 | =29+31=60 | = 77+29+95+31 =232 |

- Degree of freedom (df) = the minimum number of values in cells necessary to compute the values from other cells
  - 2 2 contingency table: if we have the totals on rows and column the values in all 4 cells could be computed
  - df = (r - 1)(c - 1); r = number of rows, c = number of columns

# RISK

- Risk means the same thing as probability to a mathematician. Clinicians and epidemiologists tend to use the word risk in a particular way but we still calculate risks in the same way as any other probability.
- The risk of an outcome is the number of times the outcome of interest occurs divided by the total number of possible outcomes.

# RISK

- In the paper *Caries prevalence in northern Scotland before and 5 years after, water defluoridation* (Stephen et al., 1987, BDJ 163: 324-326) the researchers studied two groups of children in Wick; one group whilst the water was fluoridated and one group after defluoridation. Out of 106 children examined whilst the water was fluoridated 77 had caries.
- We get the risk of caries whilst the water was fluridated by the following calculation:

**Risk = 77 / 106 = 0.73**

# RISK RATIO (RELATIVE RISK)

|  | Caries + | Caries - | Total |
|---|---|---|---|
| Fluoridated + | TP = 77 | FP = 29 | = 77+29 = 106 |
| Non-Fluoridated - | FN = 95 | TN = 31 | = 95+31 = 126 |
| Total | =77+95=172 | =29+31=60 | = 77+29+95+31 =232 |

- If we want to compare the effects of fluoridated and non-fluoridated water we could calculate the risk of having caries for each group:

- Risk of having caries when water is fluoridated = 77 / 106 = 0.73

- Risk of having caries when water is not fluoridated = 95 / 126 = 0.75

# RISKS COMPARISON

- We can compare the risk for each of the groups using the risk ratio. The risk ratio for being caries free when water is fluoridated compared to when it is not fluoridated is:
- **(Risk when fluoridated) / (Risk when not fluoridated) = 0.73 / 0.75 = 0.96**
- The risk of having caries when the water is fluoridated is only 0.96 that of when the water is not fluoridated.
- An risk ratio of 1 means there is no difference between the groups
- The 95% CI includes 1 so we have a (statistically) non-significant result.

# ODDS

|  | Caries + | Caries - | Total |
|---|---|---|---|
| Fluoridated + | TP = 77 | FP = 29 | = 77+29 = 106 |
| Non-Fluoridated - | FN = 95 | TN = 31 | = 95+31 = 126 |
| Total | =77+95=172 | =29+31=60 | = 77+29+95+31 =232 |

- The odds in favor of a particular outcome is the number of times the outcome occurs divided by the number of times it does not occur.

- 77 children who had caries and 29 who didn't:

**Odds = 77 / 29 = 2.66**

# Odds

- Odds of less than 1 mean the outcome occurs less than half the time
- Odds of 1 mean the outcome occurs half the time
- Odds of more than 1 mean the outcome occurs more than half the time

# ODDS RATIO

- If we want to compare the effects of fluoridated and non-fluoridated water we could calculate the odds for each group:
    - Odds for having caries when water is fluoridated = 77 / 29 = 2.66
    - Odds for having caries when water is not fluoridated = 95 / 31 = 3.06

# ODDS RATIO

- We can compare the odds using the odds ratio. The odds ratio for having caries when water is fluoridated compared to when it is not fluoridated is:

**(Odds when fluoridated) ÷ (Odds when not fluoridated)**
**= 2.66 / 3.06 = 0.87**

- So, the odds of having caries when the water is fluoridated are about 90% those of when the water is not fluoridated.
- An odds ratio of 1 means there is no difference between the groups

# RISKS AND ODDS: OTHER MEASURES OF ASSOCIATION

| Denumire | Formula | Definiție |
|----------|---------|-----------|
| False positive rate | =FP/(FP+TP) | Probability of a false positive test (α) |
| False negative rate | =FN/(FN+TP) | Probability of a false negative test (β) |
| Sensibility | =TP/(TP+FN) | Probability of a true positive test (1- β) |
| Specificity | =TN/(TN+FP) | Probability of a true negative test (1- α) |
| Accuracy | =(TP+TN)/n | General probability of a correct decision |
| Positive predictive value | =TP/(TP+FP) | Probability of a correct positive test |
| Negative predictive value | =TN/(TN+FN) | Probability of a correct negative test |
| Relative risk | =[TP(FP+TN)]/[FN(TP+FP)] | |
| Odd ratio | =(TP·TN)/(FN·FP) | |
| Attributable risk | =TP/(TP+FP)-FN/(FN+TN) | |

# RISKS AND ODDS: MEASURES OF ASSOCIATION

|  | Caries + | Caries - | Total |
|---|---|---|---|
| Fluoridated + | TP = 77 | FP = 29 | = 77+29 = 106 |
| Non-Fluoridated - | FN = 95 | TN = 31 | = 95+31 = 126 |
| Total | =77+95=172 | =29+31=60 | = 77+29+95+31 =232 |

| Name | Formula |
|---|---|
| False positive rate | = 29/(77+29) = 0.2736 |
| False negative rate | = 95/(95+31) = 0.7540 |
| Sensibility | = 77/(77+95) = 0.4477 |
| Specificity | = 31/(31+29) = 0.4833 |
| Accuracy | = (77+31)/232 = 0.4655 |
| Positive predictive value | = 77/(77+29) = 0.7264 |
| Negative predictive value | = 31/(31+95) = 0.2460 |
| Relative risk | = 77(29+31)/95(77+29) = 0.4588 |
| Odd ratio | = (77·31)/(95·29) = 0.8664 |
| Attributable risk | = 77/(77+29)-95/(95+31) = -0.0275 |

# TESTING ASSOCIATION IN CONTINGENCY TABLE

- We can perform a hypothesis test on a contingency table. The test we will use most often is the χ2 test.
- $\chi^2$ Test
  - Is proper to be applied if the sample size is large
  - The test is valid if the expected frequency of each cell is at least equal to 1 and the observed frequency is of 5
  - If the above-described conditions are not meet, the Fisher exact test is the proper test

# χ² TEST

- Indicate if that the two variables are or are not independent BUT DO NOT quantify the power of association between them.
- Steps:

1. Define the hypotheses
2. Define the parameter of the test
3. Define the significance level
4. Define the critical interval
5. Calculate the observed value of the parameter of the test
6. Make a decision

# χ² Test: Problem

- The association between *Streptococcus mutans* (as risk factor) and dental caries was studied. A sample of 620 patients was investigated. The sample contains: 150 patients with caries and *Streptococcus mutans*, 230 patients without caries and without *Streptococcus mutans* and 60 patients with caries but without *Streptococcus mutans.* The presence of *Streptococcus mutans* is asscoiated with dental caries? (df=1; α=0.05; $\chi^2_{critical}$ = 3.84).

# χ² TEST: 1. HYPOTHESES

- $H_0$:
  - There is no association between *Streptococcus mutans* and dental caries.
  - The presence of *Streptococcus mutans* and dental caries are independent.
- $H_1/H_a$:
  - There is an association between *Streptococcus mutans* and dental caries.
  - The presence of *Streptococcus mutans* and dental caries are not independent.

# χ² Test: 2. Parameter of the test

$$\chi^2 = \sum_{i=1}^{r \cdot c} \frac{(f_i^0 - f_i^t)^2}{f_i^t}$$

Follow a distribution law with (r-1)(c-1) degree of freedom

where

- χ² = the parameter of χ² test
- $f_i^o$ = observed frequency
- $f_i^t$ = expected (theoretic) frequency

# χ² Test: 3. Significance level

- Let $\alpha = 0.05$ (5%) be the significance level.

# χ² Test: 4. Critical region

- Critical region: $[\chi_\alpha^2, \infty)$
- For $\alpha = 0.05$:
  - $\chi_\alpha^2 = 3.84$
  - $[3.48, \infty)$

# χ² TEST: 5. PARAMETER OF THE TEST

| **observed** | DC+ | DC- | Total |
|---|---|---|---|
| SP + | TP = **150** | FP = **180** | 330 |
| SP - | FN = **60** | TN = **230** | 290 |
| Total | 210 | 410 | 620 |

| **expected** | DC+ | DC- | Total |
|---|---|---|---|
| SP + | **= 330  210/620** | **= 330  410/620** | 330 |
| SP - | **= 290  210/620** | **= 290  410/620** | 290 |
| Total | 210 | 410 | 620 |

# χ² TEST: 5. PARAMETER OF THE TEST

| observed | DC+ | DC- |
|----------|-----|-----|
| SP + | **150** | **180** |
| SP - | **60** | **230** |

| expected | DC+ | DC- |
|----------|-----|-----|
| SP + | **= 112** | **= 218** |
| SP - | **= 98** | **= 192** |

$$\chi^2 = \frac{(150-112)^2}{112} + \frac{(180-218)^2}{218} + \frac{(60-98)^2}{98} + \frac{(230-192)^2}{192}$$

$$\chi^2 = \frac{38^2}{112} + \frac{(-38)^2}{218} + \frac{(-38)^2}{98} + \frac{(38)^2}{192}$$

$$\chi^2 = \frac{1444}{112} + \frac{1444}{218} + \frac{1444}{98} + \frac{1444}{192} = 12.89 + 6.63 + 14.73 + 7.52 = \boxed{41.77}$$

- If $\chi^2 \in [3.84, \infty)$ $H_o$ is rejected with a risk of error of type I ($\alpha$).

- If $\chi^2 \notin [3.84, \infty)$ $H_o$ is accepted with a risk of error of type II ($\beta$).

- Since $41.77 \in [3.84, \infty)$ $H_o$ is rejected with a risk of error of 5%.

- **There is an association between *Streptococcus mutans* and dental caries.**

# Continuity correction (Yates's correction)

- For small sample sizes the $\chi^2$ test is too likely to reject the null hypothesis (it tends to spot differences where none really exist).

  - A continuity correction can be made to allow for this.

  - Two conditions have to be met:

    - All expected frequencies must be greater than 1
    - 80% of observed frequencies must be greater than 5

# TESTUL X²: CORECȚIA YATES

$$\chi^2 = \sum_{i=1}^{r \cdot c} \frac{|\, f_i^0 - f_i^t \,|^2 - 0.5}{f_i^t}$$

- 0.5 = Yates correction

# FISHER'S EXACT TEST

- Chi-square procedures can be legitimately applied only if all values of **E** are equal to or greater than 5.
- If a 2×2 contingency table fails to meet the conditions required for the $\chi^2$ test then Fisher's exact test can be used.
- It is based on different mathematics to the $\chi^2$ test which are more robust when sample sizes are small.

# FISHER'S EXACT TEST

- $H_o$: there is no association between smoking and dental caries
- If the null hypothesis is true - if any ostensible association between smoking and dental caries were the result of nothing more than mere chance coincidence -how likely is it that we might end up with a result this large or larger?

| observed | DC+ | DC- | Total |
|----------|------|------|-------|
| smoking + | TP = 2 | FP = 7 | 9 |
| smoking - | FN = 8 | TN = 2 | 10 |
| Total | 10 | 9 | 19 |

# FISHER'S EXACT TEST

- Suppose that the initial assessment was performed and the number of subjects who do and do not show characteristics (smoking and dental caries) were counted, but have not yet sorted the subjects according to the correspondences of smoking and dental caries. In this case, all they would have would be the marginal totals shown in the following table/
- Given these marginal totals, there are 10 possible ways in which the specific correspondences between smoking and dental caries.

|  | DC+ | DC- | Total |
|---|---|---|---|
| smoking + |  |  | 9 |
| smoking - |  |  | 10 |
| Total | 10 | 9 | 19 |

# Fisher's exact test

- The p-value is calculated directly from the formula:

$$p = \frac{(a+c)!(b+d)!(c+d)!(a+b)!}{n!a!b!c!d!}$$

- The p-value for the observed contingency table must be added to the p-value of the more extreme contingency table.

# FISHER'S EXACT TEST

|  | DC+ | DC- | Total |
|---|---|---|---|
| smoking + | **6** | **2** | 8 |
| smoking - | **1** | **6** | 7 |
| Total | 7 | 8 | 15 |

|  | DC+ | DC- | Total |
|---|---|---|---|
| smoking + | **7** | **1** | 8 |
| smoking - | **0** | **7** | 7 |
| Total | 7 | 8 | 15 |

# FISHER'S EXACT TEST

- The p-value must be calculated for the two contingency tables:

$$P_1 = \frac{7!8!7!8!}{15!6!2!6!} = 0.0305$$

$$P_2 = \frac{7!8!7!8!}{15!7!0!7!} = 0.0012$$

- Therefore $p = p_1 + p_2 = 0.0305 + 0.0012 = 0.0317$

# FISHER'S EXACT TEST

- The p-value = 0.0317 < $\alpha$ = 0.05 $\Rightarrow$ that smoking is associated with dental caries.

# Summary

**Conditions for the $\chi 2$ test**

All expected values must be greater than 1

80% of expected values must be greater than 5

Online calculator

**r×n: Contingency Table**

Online calculator