UNIVERSITY OF COPENHAGEN

### Faculty of Health Sciences



# Binary outcomes and frequency tables Basic Statistics for health researchers

Alessandra Meddis

Section of Biostatistics, University of Copenhagen

March 10th, 2025

# Outline/Intended Learning Outcomes (ILOs)

#### Preliminaries

ILO: calculate 95% CIs for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% CIs

#### Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

#### Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

#### Confounding

ILO: to exemplify confounding and its potential to be misleading ILO: to name two commonly used remedies

#### Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs

LO: to restate which association measure(s) can be used for each design

#### Screening: jargon

ILO: to recognize some jargon

#### Paired binary data (if time allows)

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-CI and p-values



## Binary outcome

$$Y = \begin{cases} 1 & \text{event } / \text{ positive } / \text{ disease} \\ 0 & \text{no event } / \text{ negative } / \text{ non-disease} \end{cases}$$

### **Binary outcome**

$$Y = \begin{cases} 1 & \text{event / positive / disease} \\ 0 & \text{no event / negative / non-disease} \end{cases}$$

### Parameters

Prevalence: proportion of the population with event at fixed time point.

How many have the disease right now?

Risk: probability that event occurs in given time period: How likely will a subject acquire the disease within 1-year?



## Statistical inference

### Estimating risks and prevalence

$$\hat{p} = \text{Relative frequency} = \frac{\text{Number of events}}{\text{Number of subjects}} = \frac{x}{n}$$

### Confidence limits: normal approximation ("large" $n^1$ )

$$\left[\widehat{p} - 1.96\sqrt{\frac{\widehat{p}(1-\widehat{p})}{n}}; \widehat{p} + 1.96\sqrt{\frac{\widehat{p}(1-\widehat{p})}{n}}\right]$$

Confidence limits: "exact" (any n)

binom.test(x,n)

<sup>1</sup>rule of thumb: when both  $x \ge 5$  and  $n - x \ge 5$ .

## Exact confidence intervals

**Exact:** No approximation. Example:

- x = 7 (number of events)
- ▶ n = 43 (number of subjects)  $\rightarrow \hat{p} = 7/43 = 16.3\%$

We want to be sure at 95% that the true value falls inside the confidence interval  $[p_L;p_U]$ 

## Exact confidence intervals

**Exact:** No approximation. Example:

- x = 7 (number of events)
- ▶ n = 43 (number of subjects) →  $\hat{p} = 7/43 = 16.3\%$

We want to be sure at 95% that the true value falls inside the confidence interval  $[p_L;p_U]$ 

To obtain the exact confidence interval, we look for:

 $\begin{array}{ll} {\sf p}_U & {\sf s.t.P}(X \le 7) = 0.025 \\ {\sf p}_L & {\sf s.t.P}(X \ge 7) = 0.025 \end{array}$ 

**Binomial Distribution:** probability of having x success among n tries, knowing that the probability of success is p

$$\mathbf{P}(X \le x) = \sum_{i=0}^{x} \binom{n}{i} p^{i} (1-p)^{n-i}$$

## Exact confidence intervals

**Exact:** No approximation. Example:

- x = 7 (number of events)
- ▶ n = 43 (number of subjects) →  $\hat{p} = 7/43 = 16.3\%$

We want to be sure at 95% that the true value falls inside the confidence interval  $\left[p_L;p_U\right]$ 

We create a sequence of possible p and we look for  $p^{\ast}% (p)=p^{\ast}(p)$  that solves

$$\mathbf{P}(X \le 7) = 0.025$$

and  $p^{\prime\ast}$  that solves

$$\mathbf{P}(X \ge 7) = 0.025$$





















## Outline/Intended Learning Outcomes (ILOs)

#### Preliminaries

ILO: calculate 95% Cls for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% Cls

#### Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

#### Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

#### Confounding

ILO: to exemplify confounding and its potential to be misleading

#### Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs

ILO: to restate which association measure(s) can be used for each design

#### Screening: jargon

ILO: to recognize some jargon

#### Paired binary data (if time allows)

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-CI and p-values



## Case: clinical trial on Dalteparin <sup>3</sup>

Data: n = 85 diabetic patients with peripheral arterial occlusive disease and chronic foot ulcers, randmomized (double-blind) to:

- ▶ Placebo (n = 42)
- ▶ Dalteparin (n = 43)

#### Outcome:

Category <sup>2</sup>	Label
intact skin	healed
decreased ulcer area $\geq 50\%$	improved
increased ulcer area $\geq$ 50%	impaired
decreased or increased ulcer area $<$ 50%	unchanged
amputation above/below ankle	amputation

**Research question**: Does Dalteparin improve the outcome, when injected once daily until ulcer healing or for a maximum of 6 months?



<sup>8/62</sup> <sup>3</sup>Kalani et al. *Diabetes Care* **26**: 2575-2580, 2003



<sup>&</sup>lt;sup>2</sup>mutually exclusive.

## Frequency table

	Dalteparin	Placebo
Healed	14 (33%)	9 (21%)
Improved	15 (35%)	11 (26%)
Unchanged	7 (16%)	9 (21%)
Impaired	5 (12%)	5 (12%)
Amputation	2 (5%)	8 (19%)
total (100%)	43	42

- Summarizes the outcome data.
- Prepare/Format data for analyzes.

## Barplot (frequencies)



## Barplot (proportions<sup>4</sup>)

![](_page_20_Figure_3.jpeg)

<sup>4</sup>often better when sample sizes are not equal in both groups.

### Here we pool the outcome categories as follows

Category	Dichotomized outcome
intact skin	better
ulcer area decreased $\geq 50\%$	Detter
decreased or increased ulcer area $< 50\%$	
increased ulcer area $\geq$ 50%	worse
amputation above/below ankle	

**Important:** this dichotomization should be prespecified (i.e. decision made before seeing the data).  $^{5}$ 

<sup>&</sup>lt;sup>5</sup>For an illustration of why prespecification matters, see e.g. Austin & Goldwasser. "Pisces did not have increased heart failure: data-driven comparisons of binary proportions between levels of a categorical variable can result in incorrect statistical significant tevels." Journal of clinical applemiology 613, 2(2008): 295–500.

# Group comparison

Placebo group

Risk of worse outcome 
$$=rac{22}{42}=\widehat{p}_1$$

Dalteparin group

Risk of worse outcome 
$$=rac{14}{43}=\widehat{p}_2$$

<sup>6</sup>whenever possible, we prefer using risk ratios or risk differences to odds ratios.

# Group comparison

### Placebo group

Risk of worse outcome 
$$=rac{22}{42}=\widehat{p}_1$$

Dalteparin group

Risk of worse outcome 
$$=$$
  $\frac{14}{43} = \widehat{p}_2$ 

Association measures<sup>6</sup>

Relative risk: 
$$\frac{\widehat{p}_1}{\widehat{p}_2}$$
 Odds ratio:  $\frac{\widehat{p}_1}{1-\widehat{p}_1}$  Risk difference:  $\widehat{p}_1 - \widehat{p}_2$ 

<sup>6</sup>whenever possible, we prefer using risk ratios or risk differences to odds ratios.

### 2x2 contingency table

![](_page_24_Figure_3.jpeg)

Response

### **Risk estimates**

$$\widehat{p}_1 = \frac{a}{a+b} \qquad \widehat{p}_2 = \frac{c}{c+d}$$

![](_page_24_Picture_7.jpeg)

### Relative risk

![](_page_25_Figure_3.jpeg)

Standard error of  $\log(\widehat{RR})$  and confidence interval of RR  $^7$ 

$$\widehat{\sigma} = \sqrt{\frac{1}{a} - \frac{1}{a+b} + \frac{1}{c} - \frac{1}{c+d}}$$
$$\log(RR) : CI_{95\%} = \left[\log(\widehat{RR}) - 1.96\,\widehat{\sigma}) ; \, \log(\widehat{RR}) + 1.96\,\widehat{\sigma})\right]$$
$$RR : CI_{95\%} = \left[\widehat{RR} \cdot \exp(-1.96\,\widehat{\sigma}) ; \, \widehat{RR} \cdot \exp(1.96\,\widehat{\sigma})\right]$$

![](_page_25_Picture_6.jpeg)

### Relative risk: placebo versus dalteparin

![](_page_26_Figure_3.jpeg)

Standard error of  $\log(\widehat{RR})$  and confidence interval

$$\hat{\sigma} = \sqrt{\frac{1}{22} - \frac{1}{42} + \frac{1}{14} - \frac{1}{43}} = 0.264$$
$$CI_{95\%} = [0.959; 2.7]$$

![](_page_26_Picture_6.jpeg)

### Relative risk: placebo versus dalteparin

![](_page_27_Figure_3.jpeg)

 $CI_{95\%} = [0.959; 2.7]$  (does include 1)

The risk in the placebo group is 1.6 times higher then the risk in the dalteparin group and the risk among patients on placebo could be between 0.9 times lower and 2.7 higher compared with patients on dalteparin.

![](_page_27_Picture_6.jpeg)

### Relative risk: placebo versus dalteparin

![](_page_28_Figure_3.jpeg)

 $1/1.609 = 0.625, CI_{95\%} = [0.37; 1.04]$ 

The risk in the dalteparin group is reduced by a factor 0.622 compared to the placebo group....

![](_page_28_Picture_6.jpeg)

### Risk difference

![](_page_29_Figure_3.jpeg)

Standard error of  $\widehat{\Delta}$  and confidence interval  $^{8}$ 

$$\widehat{\sigma} = \sqrt{ab/(a+b)^3 + cd/(c+d)^3}$$
$$CI_{95\%} = \left[\widehat{\Delta} - 1.96\,\widehat{\sigma}\ ;\ \widehat{\Delta} - 1.96\,\widehat{\sigma}\right]$$

<sup>28</sup> 8 This method is "good enough" with "large enough" sample sizes.

 $\widehat{\Delta} =$ 

18/62

Outcome

### Risk difference: placebo versus dalteparin

$$\frac{22}{42} - \frac{14}{43} = 0.198$$
Treatment
$$\frac{1}{14} = 0.198$$
Treatment
$$\frac$$

### Risk difference: placebo versus dalteparin

![](_page_31_Figure_3.jpeg)

Standard error of  $\widehat{\Delta}$  and confidence interval

$$\hat{\sigma} = \sqrt{22 \cdot 20/42^3 + 14 \cdot 29/43^3} = 0.105$$
  
 $CI_{95\%} = [-0.008 \ ; \ 0.404]$ 

![](_page_31_Picture_6.jpeg)

## Risk difference: placebo versus dalteparin

![](_page_32_Figure_3.jpeg)

### $CI_{95\%} = [-0.008 ; 0.404]$ (does include 0)

The risk among patients on placebo is 19.8 % higher compared to risk in the deltaparin group, the risk in the placebo group could be between 0.8% lower and 40.4% higher compared with patients on dalteparin.

# Odds Ratio (OR)

Odds: ratio of the probability of success by probability of failure

$$\mathsf{odds} = p/(1-p) \ ,$$

and the risk can be computed back from the odds, p = odds/(1 + odds). Odds are difficult to interpret, but if risks are small, then risks  $\approx$  odds.

The Odds ratio (OR) is defined as the ratio of the odds

$$OR = \frac{\mathsf{odds}_1}{\mathsf{odds}_2} = \frac{p_1/(1-p_1)}{p_2/(1-p_2)}$$

Concept needed for

- case-control studies (stay tuned!)
- logistic regression (next week)

![](_page_33_Picture_11.jpeg)

$$OR = \frac{\mathsf{odds}_1}{\mathsf{odds}_2} = \frac{p_1/(1-p_1)}{p_2/(1-p_2)} \ .$$

OR are difficult to interpret, but from the equation...

$$RR = \frac{OR}{\left\{1 - p_2\right\} + p_2OR},$$

...and further conclude that

...we can first conclude:

 $\blacktriangleright OR > 1 \Leftrightarrow RR > 1$ 

$$\triangleright$$
  $OR = 1 \Leftrightarrow RR = 1$ 

 $\blacktriangleright OR < 1 \Leftrightarrow RR < 1$ 

▶ the OR is sufficient to deduce whether a risk increases or decreases.
 ▶ if p<sub>2</sub> is small (e.g. rare disease), then OR ≈ RR.

## When is $OR \approx RR$ ?

![](_page_35_Figure_3.jpeg)

![](_page_35_Picture_4.jpeg)
### Odds ratio



Standard error of  $\log(\widehat{OR})$  and confidence interval<sup>9</sup>

$$\widehat{\sigma} = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$
$$CI_{95\%} = \left[\widehat{OR} \cdot \exp(-1.96\,\widehat{\sigma}); \widehat{OR} \cdot \exp(1.96\,\widehat{\sigma})\right]$$

<sup>22/62</sup> 9 This method is "good enough" with "large enough" sample sizes.

# Odds ratio: placebo versus dalteparin



Outcome

# Odds ratio: placebo versus dalteparin

$$\widehat{OR} = \frac{22 \cdot 29}{14 \cdot 20} = 2.279$$
Treatment
Treatment

 worse
 better
 total

 placebo
 22
 20
 42

 dalteparin
 14
 29
 43

 total
 36
 49
 85

Standard error of  $\log(\widehat{OR})$  and confidence interval

$$\widehat{\sigma} = \sqrt{\frac{1}{22} + \frac{1}{20} + \frac{1}{14} + \frac{1}{29}} = 0.449$$
$$CI_{95\%} = [0.946; 5.491]$$



# Odds ratio: placebo versus dalteparin



### $CI_{95\%} = [0.946; 5.491]$ (does include 1)

The placebo group has 2.3 times higher odds of experiencing the worse outcome compared to the dalteparin group.

# Reporting results

The association measures of group 1 (placebo) versus group 2 (Dalteparin) are estimated as

 $RR = 1.609, \quad RD = 0.198$ 

### Equivalent statements:

- The risk in group 1 is 1.609 times higher than in group 2.
- ► The risk in group 1 is 60.9% higher than in group 2.<sup>10</sup>
- ▶ The Risk in group 1 is increased by 19.8% points

# Reporting results

The association measures of group 1 (placebo) versus group 2 (Dalteparin) are estimated as

 $RR = 1.609, \quad RD = 0.198$ 

### Equivalent statements:

- The risk in group 1 is 1.609 times higher than in group 2.
- ▶ The risk in group 1 is 60.9% higher than in group 2.<sup>10</sup>
- ▶ The Risk in group 1 is increased by 19.8% points

Note: RR shows relative change, which depends on baseline risk. RD shows absolute change, which is more informative for public health decisions.

Percentage points increase is different then saying x% higher. The latter refers to the proportional growth relative to the starting percentage.

# When to use?

Metric	When to use it	Interpretation
Risk Differ-	to measure the ab-	increase/decrease in
ence (RD)	solute impact of an	percentage points
	exposure on an out-	
	come	
Relative Risk (RR)	Cohort Studies, to assess how many times more likely one group is to ex- perience the event, rare event	relative increase/de- crease, depend on the baseline risk
Odds Ratio (OR)	Case-control studies, logistic regression	compare odds in- stead of risk



 $\Leftrightarrow$ 

# Testing independence in a randomized clinical trial

Null hypothesis  $H_0$ : the treatment has no effect.

Prob(worse given dalteparin) = Prob(worse given placebo)

$$p_1 - p_2 = 0 \qquad \text{(Difference =0)}$$

$$\Leftrightarrow \qquad \frac{p_1}{p_2} = 1 \qquad \text{(Relative risk =1)}$$

$$\Leftrightarrow \qquad \frac{p_1/(1-p_1)}{p_2/(1-p_2)} = 1 \qquad \text{(Odds ratio =1)}$$

Popular tests of independence between the treatment group and the outcome groups:

- $\chi^2$  test (normal approximation)<sup>11</sup>
- Fisher's exact test: recommended as the default choice! <sup>12</sup>

<sup>11</sup>This method is "good enough" with "large enough" sample sizes. <sup>2/02</sup>12Recommended because: Why approximate when you can get the exact?



### Observed counts

### Expected counts

Response						Response				
- Exposure -		yes	no	total			yes	no	total	
	yes	а	b	a+b	– Exposure –	yes	E <sub>11</sub>	E <sub>12</sub>	a+b	
	no	с	d	c+d		no	E <sub>21</sub>	E <sub>22</sub>	c+d	
	total	a+c	b+d	N		total	a+c	b+d	N	

- The expected counts are calculated under the null hypothesis of independence between exposure and response
- in a population of size n, for a given risk of event p, we expect to see (on average) np events in this population



### Observed counts

### Expected counts

Response						Response				
Exposure		yes	no	total			yes	no	total	
	yes	а	b	a+b	Exposure .	yes	E <sub>11</sub>	E <sub>12</sub>	a+b	
	no	с	d	c+d		no	E <sub>21</sub>	E <sub>22</sub>	c+d	
	total	a+c	b+d	N		total	a+c	b+d	N	

- The expected counts are calculated under the null hypothesis of independence between exposure and response
- in a population of size n, for a given risk of event p, we expect to see (on average) np events in this population

**Example:** Expected counts for *Exposed*=yes, *Response*=yes ( $\mathbf{E}_{11}$ ):





### Observed counts

### Expected counts

Response						Response				
Exposure		yes	no	total			yes	no	total	
	yes	а	b	a+b	Exposure -	yes	E <sub>11</sub>	E <sub>12</sub>	a+b	
	no	с	d	c+d		no	E <sub>21</sub>	E <sub>22</sub>	c+d	
	total	a+c	b+d	N		total	a+c	b+d	N	

- The expected counts are calculated under the null hypothesis of independence between exposure and response
- in a population of size n, for a given risk of event p, we expect to see (on average) np events in this population

**Example:** Expected counts for *Exposed*=yes, *Response*=yes ( $\mathbf{E}_{11}$ ):

$$p = \mathbf{P}(Exposed = yes, Response = yes) = \mathbf{P}(Exposed = yes) \cdot \mathbf{P}(Response = yes)$$



### Observed counts

### Expected counts

Response						Response			
- Exposure -		yes	no	total			yes	no	total
	yes	а	b	a+b	Exposure -	yes	E <sub>11</sub>	E12	a+b
	no	с	d	c+d		no	E <sub>21</sub>	E22	c+d
	total	a+c	b+d	N		total	a+c	b+d	N

- The expected counts are calculated under the null hypothesis of independence between exposure and response
- in a population of size n, for a given risk of event p, we expect to see (on average) np events in this population

**Example:** Expected counts for *Exposed=yes*, *Response=yes* (**E**<sub>11</sub>):

$$p = \mathbf{P}(Exposed = yes, Response = yes) = \mathbf{P}(Exposed = yes) \cdot \mathbf{P}(Response = yes)$$

$$p = \frac{a+b}{N} \cdot \frac{a+c}{N}$$



### Observed counts

### Expected counts

Response						Response			
Exposure .		yes	no	total			yes	no	total
	yes	а	ь	a+b	Exposure -	yes	E <sub>11</sub>	E <sub>12</sub>	a+b
	no	c	d	c+d		no	E <sub>21</sub>	E <sub>22</sub>	c+d
	total	a+c	b+d	N		total	a+c	b+d	N

- The expected counts are calculated under the null hypothesis of independence between exposure and response
- in a population of size n, for a given risk of event p, we expect to see (on average) np events in this population

**Example:** Expected counts for *Exposed=yes*, *Response=yes* ( $E_{11}$ ):

$$p = \mathbf{P}(Exposed = yes, Response = yes) = \mathbf{P}(Exposed = yes) \cdot \mathbf{P}(Response = yes)$$

$$p = \frac{a+b}{N} \cdot \frac{a+c}{N}$$
$$\rightarrow \mathbf{E}_{11} = N \cdot \frac{a+b}{N} \cdot \frac{a+c}{N} = \frac{(a+b) \cdot (a+c)}{N}$$

. .

27 / 62

$$\chi^2 = \sum \frac{(\text{observed counts} - \text{expected counts})^2}{\text{expected counts}}$$

### Observed counts

	Response				
		yes	no	total	
Exposure	yes	а	b	a+b	
Exposure	no	с	d	c+d	
	total	a+c	b+d	N	

### Expected counts

28/62

		Response				
		yes	no	total		
Exposure	yes	(a+b)(a+c)/N	(a+b)(b+d)/N	a+b		
Exposure	no	(c+d)(a+c)/N	(c+d)(b+d)/N	c+d		
	total	a+c	b+d	N		

under the null hypothesis the groups are identical, hence data can be merged into a single group

in a population of size n, for a given risk of event p, we expect to see (on average) np events in this population

### under the null hypothesis.

Rule of thumb: a valid analysis requires that all expected counts are  $\geq 5$ .



# Test results

### Null hypothesis:

dalteparin treatment has no effect for chronic foot ulcers.

Test	p-value
Fisher's exact test	0.0808
Pearson's $\chi^2$ test	0.0644
Pearson's $\chi^2$ test with Yates' continuity correction <sup>13</sup>	0.1032

### R code:

```
tab <- rbind(c(22,20),c(14,29))
fisher.test(tab)  # always works (default choice!)
chisq.test(tab,correct=FALSE) # fine with large samples
chisq.test(tab,correct=TRUE) # no longer useful</pre>
```

 $_{_{20/\omega}}^{13}$ Expected to be more precise than the usual Pearson's  $\chi^2$  test when the sample size is versional. NOT RECOMMENDED, with small sample sizes, use Fisher's test instead.

# A note of caution

Because the (simple) formulas for the 95% CI (of the previous slides) are based on large sample size approximations, they are not necessarily consistent with the result of the Fisher's exact test, especially with "very small" sample sizes.

		event	no event
Example:	exposed	5	12
	non-exposed	8	3

• 
$$\hat{p}_1 = 8/11 = 0.73$$
,  $\hat{p}_2 = 5/17 = 0.29$ .

• 
$$\widehat{\Delta} = 0.43 \ (0.09 \ ; \ 0.77)$$

- $\widehat{RR} = 2.47 (1.09; 5.62)$
- $\widehat{OR} = 6.40 (1.18; 34.61)$
- ▶ p-values from Fisher's exact test and Pearson's  $\chi^2$  (with and without Yates correction) are 0.051, 0.063 and 0.025, respectively.

Here the confidence intervals show a significant result, but not Fisher's test.



Advanced methods and software<sup>14</sup> are available to avoid running into this kind of inconsistency between hypothesis test and confidence intervals.

Fortunately, it is rare that we run into this problem.... and even rarer that it matters for the interpretation.

# Larger contigency tables (1/2)

If the table is not 2x2 but, e.g., 3x4 or 2x4, the  $\chi^2$  test and Fisher's exact test are testing an "ANOVA-like" null hypothesis similarly to what the F-test does to compare several means.

### First example:

	underweight	normal	overweight	obese
no SCD	9	51	20	8
SCD	23	61	3	1

### R code:

```
fisher.test(table(d$SCD,d$BMIgroup))
```

```
returns a p-value <0.001, for the null hypothesis
```

 $\mathsf{H}_0$ : "the prevalence of SCD is the same in all groups of BMI"

that is, "no association between BMI group and SCD".

# Outline/Intended Learning Outcomes (ILOs)

#### Preliminaries

ILO: calculate 95% Cls for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% Cls

#### Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

### Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

#### Confounding

ILO: to exemplify confounding and its potential to be misleading ILO: to name two commonly used remedies

#### Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs

LO: to restate which association measure(s) can be used for each design

#### Screening: jargon

ILO: to recognize some jargon

### Paired binary data (if time allows)

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-Cl and p-values

# Sample size and power calculation

Sample size and power calculation is mostly useful for designing clinical trials to determine the appropriate sample size needed to detect the expected effect size with sufficient statistical power.

However, this could be a useful tool in observational studies to understand what is possible to achieve with the available data.

# Sample size and power calculation

Sample size and power calculation is mostly useful for designing clinical trials to determine the appropriate sample size needed to detect the expected effect size with sufficient statistical power.

However, this could be a useful tool in observational studies to understand what is possible to achieve with the available data.

**Textbook formula** ("large *n*" approximation)

$$n = \frac{\left\{z_{\alpha/2}\sqrt{2\bar{p}(1-\bar{p})} + z_{\beta}\sqrt{p_1(1-p_1) + p_2(1-p_2)}\right\}^2}{(p_1 - p_2)^2}$$

 $\blacktriangleright$   $z_{\gamma}$  is the  $\gamma$ -quantile of a standard normal distribution <sup>15</sup>

▶  $\bar{p} = (p_1 + p_2)/2.$ 

▶ *n*: number of observations in **each** group.

 $z_{\alpha/2} = -1.96$  for  $\alpha = 5\%$  and  $z_{\beta} = 0.84$  is 1.28 for  $1 - \beta = 80\%$ 

When calculating the sample size we need to specify:

- $\blacktriangleright$  expected  $p_1, p_2 
  ightarrow$  expected size effect
- the desired power  $(1 \beta)$  and Type I error  $(\alpha)$

Reverse the formula to compute:

- Power for a given sample size: for expected values of p<sub>1</sub> and p<sub>2</sub> and desired n and α.
- ► Least detectable difference (or ratio):  $\delta = p_1 p_2$  (or  $r = p_1/p_2$ ) for given n, expected  $p_1$ , desired  $\alpha$  and minimal power  $(1 \beta)$ .



### Sample size calculation

Subjects needed to detect significant risk difference with a power of 80%, if the risks in the two groups are 25% and 50%.



- ▶ n = 58 subjects needed in each group (i.e. 116 in total) to detect significant risk difference with a power of 80% and  $\alpha = 0.05$ .
- at fixed p<sub>2</sub> = 0.5, for larger p<sub>1</sub>, namely decreasing the risk difference, we observe a fast increase in the needed sample size.

36 / 62

### Power calculation

**Example:** an initial calculation suggests n = 58 subjects per group (i.e. 116 in total), for detecting a difference of 25% survival between the two groups, assuming 50% survival in the placebo group (with 80% power). But what does the power become if we were too optimistic with the expected treatment effect? E.g. what if the difference in survival probability is only 15%?



### Least detectable difference

**Example:** My grant can finance a total sample size of n = 150 (i.e. 75 per group). What is the smallest survival difference that I can hope to show with a decent power (e.g. 80%), if I expect 80% survival in the "standard of care" (i.e. control) group? And if I expect 85% in the "standard of care" group?



Note: you need to supply a value for p1, not p2, otherwise the software is looking for a lower risk and it returns 0.72.

### One power/sample size calculation is often not enough.

It is good to understand how the needed sample size and power are affected by varying  $p_1 \mbox{ and } p_2$ 

### Discussions on

- Budget and resources allocations
- Ethical implications
- Is it worth continuing with the study knowing that we have small power?



# Outline/Intended Learning Outcomes (ILOs)

#### Preliminaries

ILO: calculate 95% Cls for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% Cls

#### Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

#### Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

### Confounding

# ${\sf ILO}:$ to exemplify confounding and its potential to be misleading ${\sf ILO}:$ to name two commonly used remedies

#### Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs

LO: to restate which association measure(s) can be used for each design

#### Screening: jargon

ILO: to recognize some jargon

### Paired binary data (if time allows)

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-Cl and p-values

# Confounding

"A simple definition of confounding is the confusion of effects. This definition implies that the effect of the exposure is mixed with the effect of another variable, leading to a bias."<sup>16</sup>

Failing to take a confounding variable into account can lead to a **false conclusion** that the outcome are in a **causal relationship** with the predictor variable.

Confounding variables are typically encountered in observational studies, but not in "ideal" randomized experiments.

<sup>&</sup>lt;sup>41/62</sup>16 Rothman (2012), Epidemiology: an introduction.

# Confounding example (birth order and risk of Down syndrome <sup>17</sup>)





# Confounding example (birth order and risk of Down syndrome <sup>17</sup>)



# Confounding example (birth order and risk of Down syndrome <sup>17</sup>)





# When can association mean causation? (1/2)

### We usually say that (statistical) association does not imply causation

- Association: when changes in one variable are observed alongside changes in another variable
- Causation: changes in one variable directly cause changes in another variable.

# When can association mean causation? (1/2)

### We usually say that (statistical) association does not imply causation

- Association: when changes in one variable are observed alongside changes in another variable
- Causation: changes in one variable directly cause changes in another variable.
- Example:
  - Clear association between Being Danish and enjoying licorice-flavored treats. However, being Danish not cause an individual to like licorice, nor does liking licorice cause someone to be Danish.

# When can association mean causation? (1/2)

We usually say that (statistical) association does not imply causation

- Association: when changes in one variable are observed alongside changes in another variable
- Causation: changes in one variable directly cause changes in another variable.

In presence of confounding we might not be able to identify the true causal effect.

We need (among others) that the groups we are comparing are similar with respect to everything except the treatment under study (exchangeability assumption).

When we succeed to correctly control for confounding, conditional exchangeability holds and association can be interpreted as causation.

# When can association mean causation? (2/2)

An example where association implies causation is "ideal" randomized experiments.

The randomization ensures that the two groups that we compare are similar with respect to everything except the intervention / treatment under study. Hence, if a difference in outcome is observed between the two groups, then we can be confident that this is the consequence of this unique difference in exposure / treatment.

In non-randomized (or non "ideally" randomized) experiments the two compared groups will usually differ with respect to more than one characteristic. This generates multiple plausible explanations for the observation of the difference in outcome – some causal and some non causal.
## Adjusted analysis

Suppose that in addition to the outcome and the exposure group a categorical confounder variable (e.g. gender) is measured for each individual.

## Subgroup analysis

Analyze  $2\times 2$  contingency tables separately in each strata defined by the confounder variable.

## Logistic regression (next week)

To compute a "weighted" average of the subgroup analyses, assuming that the exposure-outcome association is the same in all subgroups.<sup>18</sup>.

<sup>&</sup>lt;sup>46/62</sup>18 Applicable also with continuous confounders.

# Outline/Intended Learning Outcomes (ILOs)

#### Preliminaries

ILO: calculate 95% Cls for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% Cls

### Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

#### Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

#### Confounding

ILO: to exemplify confounding and its potential to be misleading ILO: to name two commonly used remedies

### Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs

ILO: to restate which association measure(s) can be used for each design

Screening: jargon

ieo. to recognize some jargon

### Paired binary data (if time allows)

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-Cl and p-values



## Observational study design

In a prospective **cohort study**, an outcome or disease-free study population is first identified by an exposure (e.g., onset of diabetes) or other inclusion criteria and followed in time until the disease or outcome of interest occurs.

**Case-control** studies identify subjects by outcome status at the outset of the investigation. First, subjects with outcome are identified and classified as cases. For each case a given number of controls (e.g., 4) are selected. A candidate control is a subject without the outcome but from the same source population.



## **Observational Study Designs: Case Control vs Cohort**



the study. Good for rare diseases. In rare diseases, OR approximates RR. In non-rare diseases, the direction of OR and RR are the same, but the actual number obtained of OR and RR are different. You CANNOT obtain a RR for this. It makes no sense to.



# Cohort study: example from Egerup et al. (2020)<sup>19</sup>

**Research question**: How larger is the 1-year risk of infection (leading to an hospitalization) among newborns of kidney-transplanted women?

		Infection within first year of life			
		yes	no	total	
Kidney- transplanted mother	yes	26	98	124	
	no	133	1098	1231	
	total	159	1196	1355	



The estimated risk ratio is  $\widehat{RR} = 1.94$  (Cl<sub>95%</sub> = [1.33; 2.83]).



 $_{50/62}$  19 Egerup et al. "Increased risk of neonatal complications and infections in children of kidney-transplanted women: A nationwide controlled cohort study." American Journal of Transplantation (2020).

51/6220

# Case-control study: example of Frachon et al.<sup>20</sup>

**Research question**: Is the use of benfluorex associated with unexplained mitral regurgitation?





- Case study described in the movie "150 Milligrams" (2016) (The original title in French is "La fille de Brest")
- France's biggest modern health scandal

Frachon et al. "Benfluorex and unexplained valvular heart disease: a case-control study." PloS one 5.4 (2010).

# Case-control study: example of Frachon et al.<sup>21</sup>





 $\widehat{OR} = 40.4 \ (CI_{95\%} : [9.7; 168])$ 

Risk is defined as the probability (or proportion) of an outcome (such as a disease) occurring over a certain period in a population at risk. here we do not have the population at risk, we know the number of cases and the number of controls (here 2 per case) is defined by the study design.

- The statistic RR depends also on the ratio between controls and cases and should not be used for measuring association in case-control studies
- The statistic  $\widehat{OR}$  works, because it compares the ODDs instead.

<sup>&</sup>lt;sup>52/62</sup>21 Frachon et al. "Benfluorex and unexplained valvular heart disease: a case-control study." PloS one 5.4 (2010).

# Why does $\widehat{OR}$ work? (1/2)



 $OR = \frac{\pi_1/(1-\pi_1)}{\pi_0/(1-\pi_0)}$ 



- 97% of the cases are included in the case-control study and 1% of the "non cases" are selected as controls; all included "blinded" from exposure (i.e. before looking for the information on the exposure).
- Connection to notations of previous slides π<sub>1</sub> = p<sub>1</sub> and π<sub>0</sub> = p<sub>2</sub>.
- E="exposure", F="Fail", S="Survive", D="Disease", H="Healthy".
- source: "Statistical models in Epidemiology", by Clayton and Hills, page 155.

# Why does $\widehat{OR}$ work? (2/2)





source: "Statistical models in Epidemiology", by Clayton and Hills, page 156.

$$\widehat{OR} \approx \frac{\frac{0.1 \times \pi_1 \times 0.97}{0.1 \times (1-\pi_1) \times 0.01}}{\frac{0.9 \times \pi_0 \times 0.97}{0.9 \times (1-\pi_0) \times 0.01}} = \frac{\pi_1 / (1-\pi_1)}{\pi_0 / (1-\pi_0)}$$

# Why does $\widehat{OR}$ work? (2/2)





source: "Statistical models in Epidemiology", by Clayton and Hills, page 156.

$$\widehat{OR} \approx \frac{\frac{0.1 \times \pi_1 \times 0.97}{0.1 \times (1-\pi_1) \times 0.01}}{\frac{0.9 \times \pi_0 \times 0.97}{0.9 \times (1-\pi_0) \times 0.01}} = \frac{\pi_1 / (1-\pi_1)}{\pi_0 / (1-\pi_0)}$$

## but

$$\widehat{RR} \approx \frac{\frac{0.1 \times \pi_1 \times 0.97}{0.1 \times \pi_1 \times 0.97 + 0.1 \times (1 - \pi_1) \times 0.01}}{\frac{0.9 \times \pi_0 \times 0.97 + 0.1 \times (1 - \pi_1) \times 0.01}{0.9 \times \pi_0 \times 0.97 + 0.9 \times (1 - \pi_0) \times 0.01}} \\ = \frac{\pi_1 / (\pi_1 \times 0.97 + (1 - \pi_1) \times 0.01)}{\pi_0 / (\pi_0 \times 0.97 + (1 - \pi_0) \times 0.01)} \\ \neq \frac{\pi_1}{\pi_0}$$



# Outline/Intended Learning Outcomes (ILOs)

### Preliminaries

ILO: calculate 95% Cls for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% Cls

## Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

## Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

## Confounding

ILO: to exemplify confounding and its potential to be misleading

## Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs ILO: to restate which association measure(s) can be used for each des

## Screening: jargon

ILO: to recognize some jargon

## Paired binary data (if time allows

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-CI and p-values



## Medical test / screening: jargon

Y: Outcome (disease status) E.g. prostate cancer

X: Test result (biomarker). E.g.  $X = \begin{cases} 1 & \text{positive if PSA} > 4.0 \text{ ng/mL} \\ 0 & \text{negative if PSA} \le 4.0 \text{ ng/mL} \end{cases}$ 

	Y = 1	Y = 0
X = 1	True <b>positive</b>	False <b>positive</b>
X = 0	False negative	True negative

- True positive rate (sensitivity): P(X = 1 | Y = 1)
- True negative rate (specificity): P(X = 0 | Y = 0)
- False positive rate (1-specificity): P(X = 1 | Y = 0)

For a good diagnostic we want high TPR and low FPR.

- The positive predictive value: P(Y = 1 | X = 1)
- ▶ The negative predictive value: P(Y = 0 | X = 0)They depend on the disease prevalence .



# Outline/Intended Learning Outcomes (ILOs)

#### Preliminaries

ILO: calculate 95% Cls for population proportions

ILO: distinguish between exact and approximate (asymptotic) 95% Cls

### Group comparison

ILO: to define a suitable association measure and compute its 95% CI ILO: to (correctly) use the  $\chi^2$  test and Fisher's test

#### Sample size and power calculation

ILO: to identify why and how to make power and sample size calculations ILO: to analyse their strengths and limitations

## Confounding

ILO: to exemplify confounding and its potential to be misleading

ILO: to name two commonly used remedies

### Cohort vs case-control study

ILO: to differentiate the cohort and case-control designs

LO: to restate which association measure(s) can be used for each design

#### Screening: jargon

ILO: to recognize some jargon

## Paired binary data (if time allows)

ILO: to exemplify paired binary data

ILO: to calculate appropriate 95%-CI and p-values



## When do we typically meet paired binary data?

## Comparison of diagnostic tests

Example: compare sensitivity (i.e. True Positive Rate) of two diagnostic tests based on either Method 1 (e.g. Blood culture) or Method 2 (e.g. PCR: Polymerase Chain Reaction) using the the same blood samples (i.e. same patients).

## Crossover clinical trials

Example: compare two sedatives, w.r.t. proportions of side effects (e.g. not waking when fire alarm rings), each drug is given to each patient one evening (two evenings separated by one week). The same patients receive the two drugs.

## Why does pairing matter?

## Comparison of diagnostic tests

Example (cont'): blood samples of "heavily" infected patients are easier to test positive than those of "mildly" infected patients. Hence, if one test is positive, the chance that the second test is positive is higher than expected in average.

## Crossover clinical trials

Example (cont'): some people sleep better than others. Some will never wake no matter what. Others are bad sleepers and will always wake. Hence, if a subject wakes the first night, the chance that he/she wakes up the second night is higher than expected in average.

Take home message: we expect less variability between two observations from the same patient than between two observations from two different patients. Appropriate statistical analysis will recognize this smaller variability. Less variability implies less random variation, which further implies more certainty, that is, narrower 95% CI and smaller p-values (than if the pairing was "wrongly" ignored).

## How are paired data often presented?

## Comparison of diagnostic tests<sup>22</sup>

Example (cont'):

		PCR-test		
		Negative	Positive	
RC test	Negative	1	19	
DC-lest	Positive	2	2	

## **Remarks:**

- 1. This 2 by 2 table shows the pairing (and the raw data).
- 2. If the sensitivity of the two diagnostic tests are equally good, we expect (approx.) the same counts in the "upper right" and "lower left" cells (based on the correct definition of Positives)



## Which statistical method with paired binary data?

- ▶ For p-value computation, we often use a McNemar's test
- Modern software can compute an "exact" version of the McNemar's test.
- An exact confidence interval can be computed for each of the two compared specificities (as seen in the first slides of the lecture)<sup>23</sup>

 $_{61/62}^{22}$ Large sample (i.e. "approximate") confidence intervals can be computed for the difference in proportions ( not shown in this course), but no "exact" method exists.

## Which R code and conclusions?

```
library(exact2x2)  # load a useful package
tab <- rbind(c(1,19),c(2,2))  # 2 by 2 table
mcnemar.exact(tab)  # exact McNemar test
binom.test(x=sum(tab[,2]),n=sum(tab)) # sensitivity for PCR-test (95%-CI)
binom.test(x=sum(tab[2,]),n=sum(tab)) # sensitivity for BC-test (95%-CI)
```

## **Conclusions:**

The sensitivity of the PCR test (88%, 95%-CI=[68,97]) was found significantly higher than that of the blood culture test (17%, 95%-CI=[5,37]) among patients with deep-seated candidiasis (p-value<0.001).

