

Introduction to Medical Statistics 2026

Oxford University Clinical Research Unit
March 23-27, 2026

Ronald Geskus and the biostatistics crew
Oxford University Clinical Research Unit
Hospital for Tropical Diseases,
Ho Chi Minh City, Viet Nam



Part III

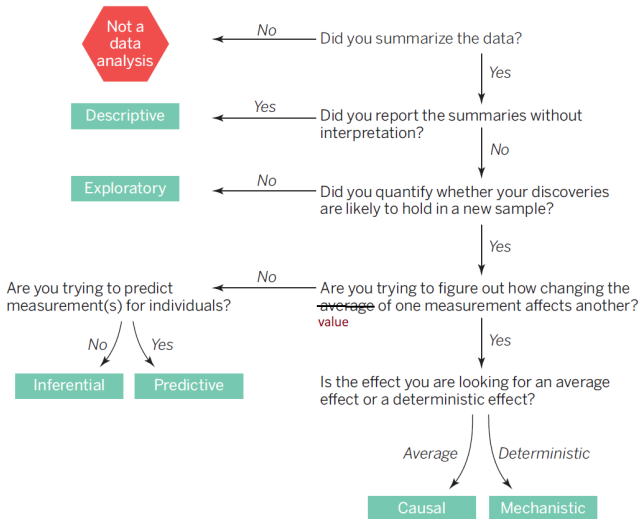
Statistical Analysis: Main Concepts and Principles; Binomial Distribution



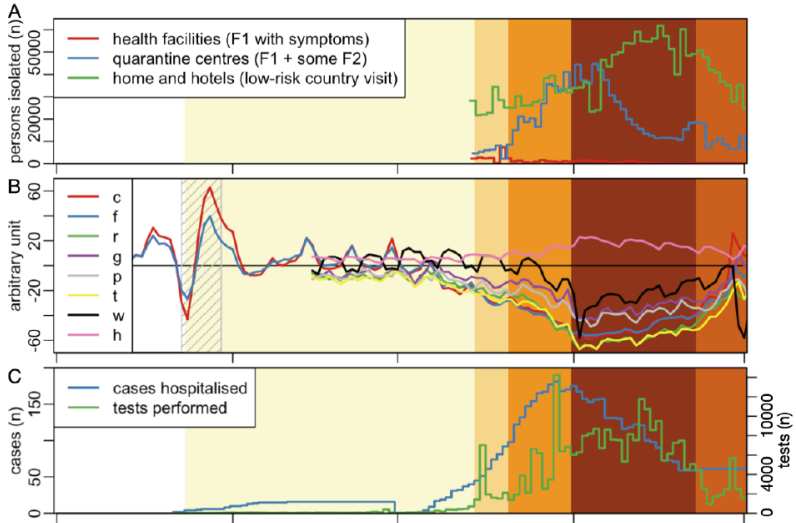
Contents:

1. Types of study questions
2. Sampling variation
3. Binomial distribution
4. Testing hypotheses

Data analysis flowchart

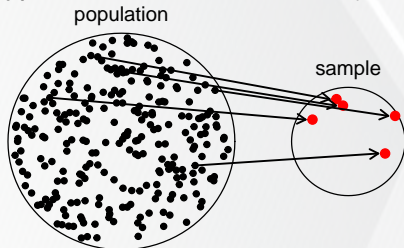


Exploratory: first 100 days of SARS-CoV-2 in Vietnam



“Discovery likely to hold in a new sample”?

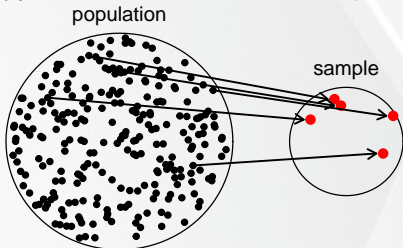
- Data (=sample) from (much) larger target population (often hypothetical, of “infinite size”)



- Assumption: **representative** sample, reflects population. Result from sample **can be generalized to population**, i.e.
 - variable characteristics: sample \approx population
 - relation between variables: in sample \approx in population

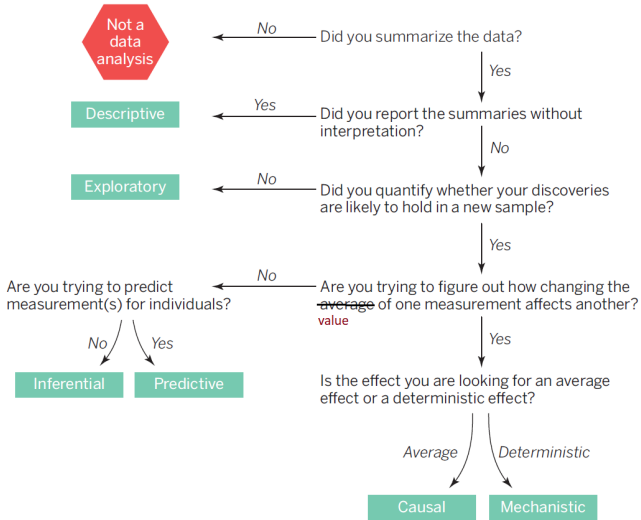
“Discovery likely to hold in a new sample”?

- Data (=sample) from (much) larger target population (often hypothetical, of “infinite size”)



- Assumption: representative sample, reflects population. Result from sample can be generalized to population, i.e.
 - variable characteristics: sample \approx population
 - relation between variables: in sample \approx in population
- Similar result for a new representative sample, but some differences due to sampling variation

Data analysis flowchart



Inferential, predictive or causal?

Relation between variable X and outcome Y

Inferential, predictive or causal?

Relation between variable X and outcome Y

Language: association, risk factor, predictor, effect, cause?

- *Smoking is associated with lung cancer*
- *Smoking is a risk factor for lung cancer*
- *Smoking is a predictor of lung cancer*
- *Smoking has an effect on lung cancer*
- *Smoking can cause lung cancer*

THE FAMILY CIRCUS



8-5

©1988 M. J. J. & Co.
and N. Condit Egan, Inc.

Walt

"I wish they didn't turn on that seatbelt sign so much! Every time they do, it gets bumpy."

Example: tuberculous meningitis (TBM)

Relation between one or more variables and an outcome

- **Causal** (well-defined hypothesis; “effect”)
Does dexamethasone decrease mortality? (intervention)
Role of Leukotriene A4 hydrolase (LTA4H) genotype in disease process (etiology)

Example: tuberculous meningitis (TBM)

Relation between one or more variables and an outcome

- Causal (well-defined hypothesis; “effect”)
Does dexamethasone decrease mortality? (intervention)
Role of Leukotriene A4 hydrolase (LTA4H) genotype in disease process (etiology)
- Inferential (risk factors; “association/correlation”)
What are the risk factors for mortality? No formal causal structure specified
- **Predictive** (personalized medicine; “prognostic/diagnostic value”)
Prognostic: probability of dying within 12 months based on individual characteristics
Diagnostic: probability to have TBM based on individual characteristics

Sample versus population: randomness

- Statistical science: develop methods to obtain **estimate** from sample that approximates true value in population

Sampling variability

- \hat{p} not exactly equal to π
 \hat{x} not exactly equal to μ
- new random sample from target population
 - new set of individuals
 - new value of estimate \hat{p} or \hat{x}

Sampling variability

- \hat{p} not exactly equal to π
 \hat{x} not exactly equal to μ
- new random sample from target population
 → new set of individuals
 → new value of estimate \hat{p} or \hat{x}
- How accurate is our (single) estimate?
sampling distribution: variation in estimate if we repeat the experiment (random sampling and computing estimate) many times

Sampling distribution

- Standard deviation describes variation/spread of any distribution
- Here: distribution of statistic if we repeatedly draw random samples from the population
- Standard deviation of statistic is called **standard error**

Sampling distribution

- Standard deviation describes variation/spread of any distribution
- Here: distribution of statistic if we repeatedly draw random samples from the population
- Standard deviation of statistic is called **standard error**
- If sample size is large then
 - statistic often follows a normal distribution (next class)
 - the amount of variation decreases and the estimate approaches population value

Formalisation of the problem

Probability model for the data

- The number of tumour responses $X \sim B(n = 20, \pi)$
 - Population success probability π
 - Variation based on $n = 20$ individuals in sample follows binomial distribution

Null hypothesis H_0

- $H_0 : \pi = 0.5$ (no effect)
 - $\hat{p} = 15/20$ successes occurred by chance. How likely is this?

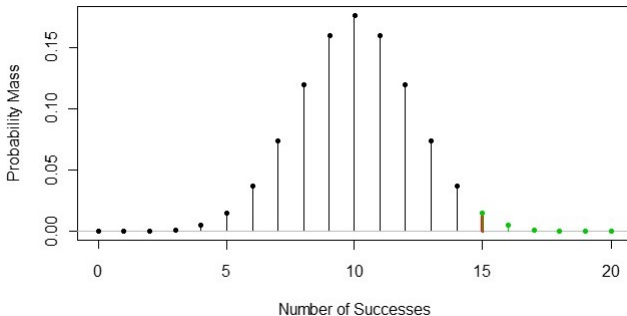
Alternative hypothesis H_A

- $H_A : \text{True } \pi > 0.5$ (one-sided) - “better”
- $H_A : \text{True } \pi \neq 0.5$ (two-sided) - “different”

P-value for one-sided alternative $H_A : \pi > 0.5$

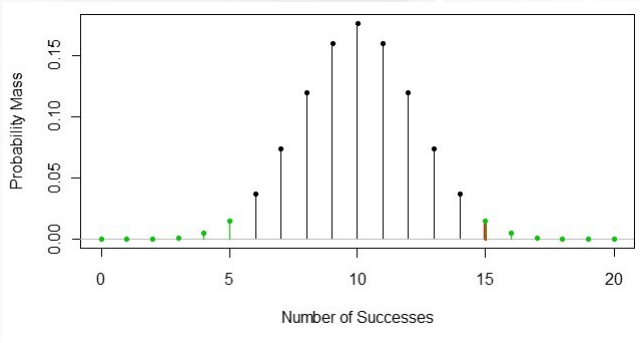
- How likely are 15/20 successes or more extreme (“**better**”) if null hypothesis is true?
- If $\pi = 0.5$:

$$\begin{aligned}
 \text{p-value} &= P(X \geq 15) \\
 &= P(X = 15) + P(X = 16) + \dots + P(X = 20) \\
 &= 0.014 + 0.004 + \dots \\
 &= 0.0207
 \end{aligned}$$



P-value for two-sided alternative $H_A : \pi \neq 0.5$

- If the null hypothesis is true (which on average gives 10 successes), how likely to have deviation of 5 or more extreme
- “More extreme” : $X \leq 5$ or $X \geq 15$
- If $\pi = 0.5$: $p\text{-value} = P(X \leq 5) + P(X \geq 15) = 0.041$



P-value for the example with R

```
> prop.test(x=15,n=20,p=0.5,alternative="two.sided")
```

```
1-sample proportions test with continuity correction
```

```
data: 15 out of 20, null probability 0.5  
X-squared = 4.05, df = 1, p-value = 0.04417  
alternative hypothesis: true p is not equal to 0.5  
95 percent confidence interval:  
 0.5058845 0.9040674  
sample estimates:  
      p  
0.75
```

- For one-sided tests, use alternative="greater" or alternative="less"
- Reject null hypothesis if p-value is below a certain **significance level** α

One- and two-sided tests

- With symmetric distribution:
 - $p\text{-value (two-sided test)} = 2 \times p\text{-value(one-sided test)}$

One- and two-sided tests

- With symmetric distribution:
 - $p\text{-value (two-sided test)} = 2 \times p\text{-value (one-sided test)}$
- Two-sided alternatives are generally preferred
 - We cannot a priori exclude harm of a new treatment
 - One-sided p-value always smaller than 0.5

One- and two-sided tests

- With symmetric distribution:
 - $\text{p-value (two-sided test)} = 2 \times \text{p-value(one-sided test)}$
- Two-sided alternatives are generally preferred
 - We cannot a priori exclude harm of a new treatment
 - One-sided p-value always smaller than 0.5
- Hence
 - Either two-sided tests at significance level α
 - Or one-sided tests at significance level $\alpha/2$
 - Often, $\alpha = 0.05$ is chosen

Hypothesis testing, basic structure

Effect exposure **E** on disease **D**

E → **D** in population

Hypothesis testing, basic structure

Effect exposure **E** on disease **D**

E \rightarrow **D** in population

1. Null hypothesis H_0 : no effect (often " $\beta = 0$ ")
2. Alternative hypothesis H_1 : there is an effect (often " $\beta \neq 0$ ")

Hypothesis testing, basic structure

Effect exposure **E** on disease **D**

E \rightarrow **D** in population

1. Null hypothesis H_0 : no effect (often " $\beta = 0$ ")
2. Alternative hypothesis H_1 : there is an effect (often " $\beta \neq 0$ ")
3. Calculate value of test statistic *TEST* based on sample data
4. Calculate p-value: probability that *TEST* exceeds some value if H_0 were true
5. If p-value small (often: $< 5\%$): reject H_0
unlikely that observed difference is due to chance

Hypothesis testing, basic structure

Effect exposure **E** on disease **D**

E \rightarrow **D** in population

1. Null hypothesis H_0 : no effect (often " $\beta = 0$ ")
2. Alternative hypothesis H_1 : there is an effect (often " $\beta \neq 0$ ")
3. Calculate value of test statistic *TEST* based on sample data
4. Calculate p-value: probability that *TEST* exceeds some value if H_0 were true
5. If p-value small (often: $< 5\%$): reject H_0
unlikely that observed difference is due to chance
6. Otherwise: do not reject H_0

Does not imply that H_0 is true (power):

No proof of effect \neq proof of no effect

Outline

Types of study questions

Sampling variation

Binary Variable; Estimating Proportions

Testing a hypothesis

Single proportion

Compare proportions between two subgroups

Categorical variables - more than 2 groups

Alternative tests for specific settings

2x2 table; example group and outcome

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

2x2 table; example group and outcome

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

- Percentage diseased in the male group:

$$\frac{32}{49} = 65\%$$

- Percentage diseased in the female group:

$$\frac{118}{245} = 48\%$$

2x2 table; example group and outcome

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

- Percentage diseased in the male group:

$$\frac{32}{49} = 65\%$$

- Percentage diseased in the female group:

$$\frac{118}{245} = 48\%$$

2x2 table; example group and outcome

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

- Percentage diseased in the male group:

$$\frac{32}{49} = 65\%$$

- Percentage diseased in the female group:

$$\frac{118}{245} = 48\%$$

- Does disease probability differ by gender?

Hypothesis test

- H_0 : probability of disease the **same** in both groups
- H_A : probability of disease **differs** per group
- Can the difference be due to chance, i.e. $H_0 : \pi_F = \pi_M$ holds?
- How do we quantify difference from equal probability?
- “Equal disease probability” is same as saying that group and outcome are unrelated/independent: knowing the group does not help in learning the outcome

Chi-squared test for independence

- For each cell: compare **observed number** (O) with **expected number** (E) under null hypothesis of independence
- Large discrepancy between O and E is an indication that probabilities in both groups differ

Expected count under null hypothesis

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

Expected count under null hypothesis

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

- Proportion male: $49/294 = 0.17$
- Proportion diseased: $150/294 = 0.51$

Expected count under null hypothesis

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

- Proportion male: $49/294 = 0.17$
- Proportion diseased: $150/294 = 0.51$

Expected count under null hypothesis

	Disease present	Disease absent	total
male	32	17	49
female	118	127	245
total	150	144	294

- Proportion male: $49/294 = 0.17$
- Proportion diseased: $150/294 = 0.51$
- Expected number diseased and male under H_0 :

$$n \times \hat{P}(\text{male}) \times \hat{P}(\text{Diseased}) = 294 \times (49/294) \times (150/294) = 25$$

Chi-squared statistic calculation (I)

	O	E
Disease present + male	32	25
Disease present + female	118	125
Disease absent + male	17	24
Disease absent + female	127	120
Total	294	294

Chi-squared statistic calculation (II)

	O	E	O-E
Dis. present + male	32	25	7
Dis. present + female	118	125	-7
Dis. absent + male	17	24	-7
Dis. absent + female	127	120	7
Total	294	294	0

Always equals 0

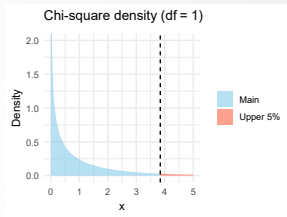
Chi-squared statistic calculation (III)

	O	E	O-E	$(O-E)^2/E$
Dis. present + male	32	25	7	1.960
Dis. present + female	118	125	-7	0.392
Dis. absent + male	17	24	-7	2.042
Dis. absent + female	127	120	7	0.408
Total	294	294	0	$X^2=4.802$

P = 0.028

Under null hypothesis, statistic X^2 has a chi-squared distribution with one degree of freedom

Upper percentile of χ^2 -distribution



- Starts at zero
- 95-percentile at 3.84

Chi-squared test of independence in R

```
> chisq.test(matrix(c(32,118,17,127),nrow=2),correct=FALSE) # Variant 1
```

Pearson's Chi-squared test

X-squared = 4.802, df = 1, **p-value = 0.02843**

```
> prop.test(x=c(32,118),n=c(49,245),correct=FALSE) # Variant 2, with CI for difference
```

2-sample test for equality of proportions without continuity correction

X-squared = 4.802, df = 1, **p-value = 0.02843**

alternative hypothesis: two.sided

95 percent confidence interval for the difference:

0.024 to 0.32

sample estimates:

prop 1	prop 2
0.65	0.48

- Note: **correct=TRUE** gives “Yates’ continuity correction” (default in R)

Interpreting the chi-squared test of independence

- After identifying the association, you may want to know the strength of this association
- The strength of the association is measured by the difference in proportion, relative risk (RR), or odds ratio (OR) (see Thursday class on logistic regression)

Outline

Types of study questions

Sampling variation

Binary Variable; Estimating Proportions

Testing a hypothesis

Single proportion

Compare proportions between two subgroups

Categorical variables - more than 2 groups

Alternative tests for specific settings

Chi-squared test for independence

- The chi-squared statistic can be used for testing association in any two-way contingency table
- Variables do not need to correspond to group and outcome
- For a table with c column and r rows, the distribution of the statistics (under null hypothesis) is a chi-squared distribution with $(r - 1) \times (c - 1)$ degrees of freedom.

2x3 table

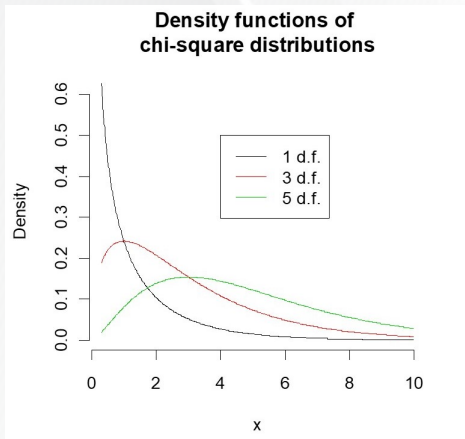
	Body Mass Index (BMI)			
	Low	Normal	Obese	total
male	a	b	c	a+b+c
female	d	e	f	d+e+f
total	a+d	b+e	c+f	

Example with more than two subgroups

TLR2 T597C SNP bacterial genotype frequencies in patients with TBM caused by Beijing genotype isolates compared to chord blood controls (Caws *et al.*; PLoS Pathog 2008)

Lineage/group	Genotype			Row Total
	TT	TC	CC	
Chord blood controls	205 0.544	154 0.408	18 0.048	377
TBM Beijing	24 0.400	25 0.417	11 0.183	60
Column Total	229	179	29	437

Density function of chi-squared distributions



Example with more than two groups in R

```
> chisq.test(matrix(c(205,24,154,25,18,11),ncol=3),  
corr=FALSE)
```

Pearson's Chi-squared test

```
data: matrix(c(205,24,154,25,18,11),ncol=3)  
X-squared = 16.3897, df = 2, p-value = 0.0002761
```

“There is clear evidence for an association between TLR2 T597C genotype and TBM caused by Beijing genotype isolates ($p=0.0003$).”

Outline

Types of study questions

Sampling variation

Binary Variable; Estimating Proportions

Testing a hypothesis

- Single proportion
- Compare proportions between two subgroups
- Categorical variables - more than 2 groups
- Alternative tests for specific settings

Fisher's exact test

- Used in small samples: expected number under null hypothesis is ≤ 1 in at least one cell
 - note: many programs (including R) give a warning if chi-squared test is used and expected frequency under null hypothesis is ≤ 5 in at least one cell

McNemar test for paired dichotomous outcomes

- Two correlated observations per individual
 - example: disease status per individual before and after treatment

	mild after	severe after	total
mild before	100	50	150
severe before	200	100	300
total	300	150	450