

QUESTIONS FOR BIAS

Key concepts





- We want to compare the mean of blood pressure levels between two groups.
- > The blood pressure checker has a problem and <u>always gives 5mmHg-higher</u>than true values.
- > <u>All subjects</u> were examined <u>by the same</u> <u>blood pressure checker</u>.





Proper comparison between groups :

1) Comparison using accurate data

2) Comparison using (<u>in)</u>accurate data

As long as the magnitude of nondem error and bias of What would be the ong problem in this study?

FOR DISCRETE VARIABLES, MEASUREMENTS ERROR IS CALLED CLASSIFICATION ERROR OR MISCLASSIFICATION

Two types of misclassification Non-differential misclassification Misclassification of a study variable that is independent of other study variables Systematic error may not be a critical issue as long as it occurs in all comparison groups.

Differential misclessification

If the error occurrence of the error occurrence of

This is a problem!!

Non-differential Misclassification with Two Exposure Categories

Correct Data Cases Controls	Exposed 240 240	Unexposed 200 600
	OR	= 3.0
Sensitivity = <mark>0.8</mark>		20% of exposed subjects were misclassified
Specificity = 1.0	Exposed	Unexposed
Cases	192	248
Controls	192	648

OR = 2.61







BIAS IN EPIDEMIOLOGIC STUDY

Different types of bias







If parents of cases with leukemia, living in the neighborhood of power lines, suspect the association and tend to agree on participation to the study,

> the association may become than what it should be.

What is this bias? How do you solve it?

In a hospital-based case-control study, the researchers excluded subjects with CVD, whom "Reserpine" was likely to be prescribed, from control group.

They found that "Reserpine was a significant risk factor of breast cancer".





- If you agree, why do you think so?
- If you don't agree, why do you think so? How do you solve it?

A doctor may examine the patient's chest X-ray more carefully if he knew the patient is a heavy smoker but not for non-smoking patients.







What is this bias? Detection bias

How do you avoid detection bias?

Suppose, you conducted a casecontrol study on relationship of prenatal infections and congenital malformations.

You asked mothers regarding prenatal episode of infections by interview / questionnaire.

Cases (mothers of babies with defect)



Controls (mothers of healthy babies)





What is the possible bias? Recall bias

How do you avoid / minimize the bias? Consider using a hospital control

Controlling for misclassification

Blinding

prevents investigators and interviewers from Knowing case/control or exposed/non-exposed status of a given participant

- Form of survey
- mail may impose less "white coat tension" than a phone or face-to-face interview
- Questionnaire
- \square use multiple questions that ask same information
- Accuracy
- Multiple checks in medical records & gathering diagnosis data from multiple sources

Lecture note of Dr. Dorak (http://www.dorak.info/epi)

CONFOUNDING

3 conditions of Confounding

- 1. Confounders are risk factors for the outcome.
- 2. Confounders are related to exposure of your interest.
- 3. Confounders are NOT on the causal pathway between the exposure and the outcome of your interest.



"Prevention" at study design

- Limitation
- Randomization in an intervention study

Matching in a cohort study

Notice: Matching does not always prevent the confounding effect in a case-control study.

How can we solve the problem of confounding?

'Treatment " at statistical analysis

Stratification by a confounder

Multivariable / multiple analysis

Mantel-Haenszel odds ratio

Stratification by confounding factor

- □ After stratification by confounding factor, common OR, OR_{MH}, among all strata should be calculated.
- □ Assumption: there is a common OR among all strata → there is no significant difference in ORs among all strata by homogeneity test.

An example of Mantel-Haenszel estimation 1

Calculate the common OR among all strata

smoking	Case	Control	
+ -	a _i c _i	b _i d _i	M _{1i} M _{0i}
Total	N _{1i}	N _{0i}	T _i

 $OR_c = \Sigma W_i OR_i / \Sigma w_i$

i: "i" th stratum, W_i : weight of "i" th stratum

Practice 1 Mantel-Haenszel odds ratio(1)

- Open the "tsunagi_v]" data by excel
 Please refer Appendix1 for the explanation of each variable.
- Import this data set by your statistical software (STATA, R, and SPSS ···)

Mantel-Haenszel odds ratio(2)

3. Suppose, you want to examine the cancer risk by habitual alcohol drinking.

Please create a contingency table of cancer and alcohol drinking.

STATA command: tab alc cancer, row Please calculate an odds ratio.

STATA command: cc cancer alc

or

cs cancer alc, or

Same OR but 95%CI is slightly different



Case-control study

. cc cancer alc				Proportion	
	Exposed	Unexposed	Total	Exposed	
Cases Controls	79 316	77 738	156 1054	0.5064 0.2998	
Total	395	815	1210	0.3264	
	Point	estimate	[95% Conf.	Interval]	
Odds ratio	2.3	96104	1.678901	3.416628	(exact)
Attr. frac. ex. Attr. frac. pop	. 58	26558 150629	.4043721	.7073138	(exact)
		chi2(1) =	26.38 Pr>ch	2 = 0.0000	

Cohort study

cs cancer alc, or							
	alc Exposed	Unexposed	Total				
Cases Noncases	79 316	77 738	156 1054				
Total	395	815	1210				
Risk	. 2	.0944785	.1289256				
	Point	estimate	[95% Conf.	Interval]			
Risk difference Risk ratio Attr. frac. ex. Attr. frac. pop	.1055215 2.116883 .5276074 .2671858		.0612577 1.584051 .3687071	.1497852 2.828946 .6465115			
Odds ratio	2.3	96104	1.706524	3.364381	(Cornfield)		
		chi2(1) =	26.38 Pr>chi	2 = 0.0000			

Mantel-Haenszel odds ratio(3)

4. Since we know that cancer risk increases with age, you may want to confirm the association between alcohol drinking and cancer risk by age group (<60, 60-69, ≥ 70).</p>

Please create contingency tables of cancer
STATA : by age_gp, sort: tab alc cancer, row

Please calculate odds ratios for each age group.

An example of Mantel-Haenszel estimation 1

	age	alcohol	Case	Control	OR
1	<60	+	13	129	1.54
		-	14	214	1 (ref)
2	60-69	+	32	105	3.95
		-	19	246	1 (ref)
3	≥70	+	34	82	2.62
		-	44	278	1 (ref)
Total		+	79	316	2.40
		-	77	738	1



Mantel-Haenszel odds ratio(5)

You can also calculate OR_{MH} by yourself.

$$OR_{MH} = \sum (a_i^*d_i/T_i) / \sum (b_i^*c_i/T_i)$$

$$OR_{MH} = \frac{(13^*214/370) + (32^*246/402) + (34^*278/438)}{(129^*14/370) + (105^*19/402) + (82^*44/438)}$$

= 2.69

Practice 2

 Using the "tsunagi_v1" data set, please examine the association between habitual alcohol drinking and cancer risk by sex stratification.



	mare	OR	[95% Conf.	Interval]	M-H Weight	
	0	.9455128	.3977241	2.01076	7.399209	(exact
	1	1.992308	1.179265	3.441301	11.52993	(exact
	Crude	2.396104	1.678901	3.416628		(exact
M-H C	ombined	1.583126	1.061639	2.360773		
st of h	omogeneity	(M-H)	chi2(1) =	2.68 Pr>ch	ni2 = 0.1014	

Q1. Is this OR_{MH} statistically significant? Q2. Is it OK to report OR_{MH} when the homogeneity test is statistically significant?

How can we solve the problem of confounding?

"Treatment " at statistical analysis

Stratification by a confounder
 Multivariable / multiple analysis

Multivariate ≠ Multivariable (Multiple)

Am J Public Health. 2013 January; 103(1): 39–40. Published online 2013 January. doi: 10.2105/AJPH.2012.300897 PMCID: PMC3518362 NIHMSID: NIHMS514677

Multivariate or Multivariable Regression?

Bertha Hidalgo, PhD, MPH^{III} and Melody Goodman, PhD, MS

Author information
Article notes
Copyright and License information

See letter "Hidalgo and Goodman Respond" in volume 10 on page e1.

This article has been cited by other articles in PMC

Abstract

Go to: 🕑

The terms multivariate and multivariable are often used interchangeably in the public health literature. However, these terms actually represent 2 very distinct types of analyses. We define the 2 types of analysis and assess the prevalence of use of the statistical term multivariate in a 1-year span of articles published in the American Journal of Public Health. Our goal is to make a clear distinction and to identify the nuances that make these types of analyses so distinct from one another.

Multivariable (Multiple) analysis

A multivariable model can be thought of as a model in which multiple variables are found on the right side of the model equation. This type of statistical model can be used to attempt to assess the relationship between a number of variables; one can assess independent relationships while adjusting for potential confounders.

This is the model to control the effects of confounders!

By contrast, a multivariable or multiple linear regression model would take the form

(2) $y = \alpha + x_1\beta_1 + x_2\beta_2 + \ldots + x_k\beta_k + \epsilon$

where y is a continuous dependent variable, x is a single predictor in the simple regression model, and x_1 , x_2, \ldots, x_k are the predictors in the multivariable model.

As is the case with linear models, logistic and proportional hazards regression models can be simple or multivariable. Each of these model structures has a single outcome variable and 1 or more independent or predictor variables.

Multivariate analysis

Multivariate, by contrast, refers to the modeling of data that are often derived from longitudinal studies, wherein an outcome is measured for the same individual at multiple time points (repeated measures), or the modeling of nested/clustered data, wherein there are multiple individuals in each cluster. A multivariate linear regression model would have the form

(3) $Y_{n \times p} = X_{n \times (k+1)} \beta_{(k+1) \times p} + \varepsilon$

where the relationships between multiple dependent variables (i.e., *Ys*)—measures of multiple outcomes—and a single set of predictor variables (i.e., *Xs*) are assessed.

This model is to analyze the relationship between "multiple outcomes" and a single set of predictors.

LOGISTIC REGRESSION ANALYSIS

Practice **3** Multivariable analysis

 Let's see the association between habitual alcohol drinking and cancer risk by logistic regression model.

STATA : logistic cancer alc or logit cancer alc, or

 Please examine this association adjusting for the effects of age and sex.



STATA : logistic cancer alc male age

logistic cancer alc male age						
Logistic regression Log likelihood = -431.32695				Number of obs = LR chi2(3) = Prob > chi2 = Pseudo R2 =		1210 67.45 0.0000 0.0725
cancer	Odds Ratio	Std. Err.	z	P> z	[95% Conf.	Interval]
alc male age	1.877452 2.099375 1.041058	.3864541 .4306218 .0093757	3.06 3.62 4.47	0.002 0.000 0.000	1.254174 1.404405 1.022844	2.810476 3.138251 1.059597

STATA : logistic cancer alc male age_gp

logistic cancer alc male age_gp						
Logistic regression Log likelihood = -431.86621					Number of obs = LR chi2(3) = Prob > chi2 = 0 Pseudo R2 = 0	
cancer	Odds Ratio	Std. Err.	z	P> z	[95% Con	f. Interval]
alc male age_gp	1.825007 2.198312 1.669985	.3736455 .4481714 .1951521	2.94 3.86 4.39	0.003 0.000 0.000	1.221779 1.474193 1.328136	2.726067 3.278117 2.099824

STATA : xi: logistic cancer alc male(i.age_gp)

Categorical variable (>2 categories)

. xi: logistic i.age_gp	cancer alc r _Iage_gp_	nale i.age_g _1-3	p (naturall	y coded;	_Iage_gp_1 o	mitted)
Logistic regre Log likelihood	ession d = -431.81689	9		Numbe LR ch Prob Pseud	r of obs = i2(4) = > chi2 = o R2 =	1210 66.47 0.0000 0.0715
cancer	Odds Ratio	Std. Err.	z	P> z	[95% Conf.	Interval]
alc male _Iage_gp_2 _Iage_gp_3	1.825179 2.192979 1.792761 2.84385	.3735059 .4473008 .4573496 .6937492	2.94 3.85 2.29 4.28	0.003 0.000 0.022 0.000	1.222123 1.470332 1.087359 1.763025	2.725812 3.270798 2.955776 4.587276

If there is no linear trend of the cancer risk by age, it would be better to use categorical variable for age.



Practice 4 Regression analysis (1)

Suppose, you want to know predictors of systolic blood pressure in the subjects of "tsunagi_v]" data.

What do you have to check first?



Distribution of systolic blood pressure



<section-header><text><text>

Practice 4 Regression analysis (2)

Age is one of the predictors of systolic blood pressure.

Please conduct regression analysis using "age" as a explanatory variable.

STATA : reg lsbp age



	STATA commands				
. reg isop age					
Source	SS	df	MS		Number of obs = 1210
Model	2 84842128	1	2 94942129		F(1, 1208) = 91.90
Residual	37.4400935	1208	.030993455		R-squared = 0.0707
					Adj R-squared = 0.0699
Total	40.2885149	1209	.033323834		Root MSE = $.17605$
lsbp	Coef.	Std.	Err. t	P> t	[95% Conf. Interval]
age	.0044274	. 0004	618 9.59	0.000	.0035213 .0053334
_cons	4.531922	.0301	.251 150.44	0.000	4.472819 4.591026

SBP= 4.531922 + 0.0044274*age

This indicates that SBP will increase 0.004 per age (year).



SBP= 4.557933 + **0.0431853***age(10)

Std. Err.

.0045294

.0276

P>|t|

0.000

0.000

t

9.53

165.14

[95% Conf. Interval]

.0342989

4.503784

.0520717

4.612083

coef.

.0431853

4.557933

1sbp

age10

_cons

cf. SBP= 4.531922 + **0.0044274***age

Practice 4 Regression analysis (4)

- Suppose, hemoglobin level may be one of the predictors of systolic blood pressure.
- Please pick-up other potential predictors (other than hemoglobin) for systolic blood pressure in this data set based on your Knowledge.
- And, conduct regression analysis.

How many explanatory variables can we use in a model?

Model	Number of explanatory variables	Example
Linear regression model	Sample size / 15	Up to around 6-7 variables in 100 subjects
Logistic regression model	Smaller sample size of outcome / 10	<u>Up to 10 variables if</u> the numbers of cases and controls are 100 and 300, respectively.
Cox proportional hazard model	The number of event / 10	<u>Up to 9 variables if</u> you have 90 events out of 150 subjects

ATTENTION!

- When you include categorical variable in your model, you have to count that variable as (the number of categories - 1).
 - □ For example, the variable of age group used in the previous practice, we have to count it as "two" (=3 categories -1) variables.