

Chihaya Koriyama August 21th, 2019



Why do we need multivariable analysis?

"Treatment (control)" for the confounding effects at analytical level

Stratification by confounder(s)
 Multivariable / multiple analysis

Prediction of individual risk

Regression models for multivariable analysis

| Paired? | Outcome variable | Proper model |
|---------|------------------------------------|---|
| No | Continuous | Linear regression model |
| | Binomial | Logistic regression model |
| Yes | Categorical (≥3) | Multinomial (polytomous) logistic regression model |
| | Binomial (event) with censoring | Cox proportional hazard model |
| | Continuous | Mixed effect model, Generalized estimating equation |
| | Categorical (≥3) | Generalized estimating equation |

LINEAR REGRESSION ANALYSIS



Original data: Doll and Hill Br Med J 1956

Height explaining mathematical ability!!??

| Source SS + | df MS | Number of obs = 32 F(1, 30) = 726.87 |
|---|--|---|
| Model 412.7743 Residual 17.0365 + | 1 412.774322 30 .567882354 | Prob > F = 0.0000 R-squared = 0.9604 Adj R-squared = 0.9590 |
| Total 429.8108 Ability score of maths | 31 13.8648643 | Root MSE = .75358 |
| ama Coef. | Std. Err. t P> t | [95% Conf. Interval] |
| height .4118029 _cons -42.82525 | .0152743 26.96 0.000 2.191352 -19.54 0.000 | .3806086 .4429973 -47.30059 -38.34992 |

Association between height and score of maths





Both height and ability of maths increase with age



Age is a confounding factor in the association between height and ability of maths.



How age itself influences the association between height and the ability of maths?

Let's see the equation Ability of maths (AM) = α + β 1(Height) \rightarrow AM = -42.8 + 0.41(Height)

AM = α + β 1(Height) + β 2(Age) \rightarrow AM = 1.48 - 0.01(Height) + 2.02 (Age)

Significant association between height and the ability of maths was gone after adjusting for the effect of age

| Source | SS | df | MS | Number of obs = 32 F(2, 29) = 851,23 |
|-------------------|-----------------------|---------|--------------------------|---|
| Model Residual | 422.6119 7.19885 | 2 29 | 211.305972 .248236138 | Prob > F = 0.0000 R-squared = 0.9833 Adi R-squared = 0.9821 |
| Total 4 | 429.81079 | 31 | 13.8648643 | Root MSE = $.49823$ |

| ama + | Coef. | Std. Err. | t | P> t | [95% Conf | f. Interval] |
|------------|----------------|-----------------|--------------|--------------|-----------|-----------------|
| height | 0121303 | .0680948 | -0.18 | 0.860 | 1513998 | .1271393 |
| age | 2.02461 | .3216095 | 6.30 | 0.000 | 1.366845 | 2.682375 |
| _cons | 1.483038 | 7.185946 | 0.21 | 0.838 | -13.21387 | 16.17995 |

No association between height and age-adjusted score of maths



Interpretation of coefficients

Let's see the equation Ability of maths (AM) = α + β 1(Height) \rightarrow AM = -42.8 + 0.41(Height)

0.41 points increase by 1cm increase of height

 $AM = \alpha + \beta 1(\text{Height}) + \beta 2(\text{Age})$

 \rightarrow AM = 1.48 - 0.01(Height) + 2.02 (Age)

0.01 points decrease by 1cm increase of height

Confounding effect: magnitude and direction of the association



Interpretation of coefficients in general

To simplify, the explanatory variable is binomial one: 1=exposed or 0=unexposed

- Exposed: Ye = α + β (Exp=1) = α + β Unexposed: Yu = α + β (Exp=0) = α Difference: Ye – Yu = β
- Coefficient estimate: difference in dependent value

Interpretation of coefficients after log-transformation of dependent variable

The explanatory variable is binomial one: 1=exposed or 0=unexposed

Exposed: In (Ye) = α + β (Exp=1) = α + β Unexposed: In (Yu) = α + β (Exp=0) = α Difference: In(Ye) – In (Yu) = β Ratio: Ye / Yu = e β

Coefficient estimate: ratio of dependent value (after exponentiating)

Control of confounding with regression model

- Compared to stratified analysis, several confounding variables can be <u>easily</u> <u>controlled simultaneously</u> using a multivariable regression model.

Results from the regression model are readily <u>susceptible to bias</u> if the model is not a good fit to the data.



Epidemiology (Rothman KJ, Oxford University Press)



Epidemiology (Rothman KJ, Oxford University Press)

CORRELATION = REGRESSION ANALYSIS?

Correlation coefficient

- Strength of the correlation between two continuous variables ranging from -1 to 1
- Correlation is a linear association between two variables
- NOT to prove the causal association; x and y variable are interchangeable.

Examples of correlation



What does "r=0" mean?

- No association between x and y?
- No linear association between x and y



Correlation coefficient is not the magnitude of "slope"



Correlation coefficients

Pearson's CC (r): parametric method
 At least, one of the two variables should follow the normal distribution.

Non-parametric methods
 Spearman's CC (ρ)
 Kendall's CC (τ)

r (correlation coefficient) and R-squared



R squared, coefficient of determination, is the proportion of the variance in the dependent variable that is predictable from the independent variable(s).

| Number of obs | = | 759 |
|---------------|---|--------|
| F(1, 757) | = | 190.24 |
| Prob > F | = | 0.0000 |
| R-squared | = | 0.2008 |
| Adj R-squared | = | 0.1998 |
| Root MSE | = | 5.4705 |

$$R^2 = r^2$$

$$\mathsf{R}^2 = 1 - \frac{\sum_i (yi - fi)^2}{\sum_i (yi - \bar{y})^2}$$

r (correlation coefficient) and R-squared



| Number of 005 | | 102 |
|---------------|---|--------|
| F(1, 757) | = | 190.24 |
| Prob > F | = | 0.0000 |
| R-squared | = | 0.2008 |
| Adj R-squared | = | 0.1998 |
| Roo VSE | = | 5.4705 |
| | | |

75Q

Number of the

Adjusted R squared value by the sample size and the number of variable(s) Better to use when you have more variables or small sample size

The R² coefficient of determination, ranging 0-1, is a statistical measure of how well the regression predictions approximate the real data points.

r (correlation coefficient) and regression coefficient



 $\sqrt{0.6084 \times 0.3301} = 0.4481$

ADJUSTMENT OF CORRELATION

Calculation skill and physical development

| | calculation | weight |
|-------------------|-------------|--------|
| | 134 | 24.8 |
| | 136 | 25.6 |
| Is weight related | 117 | 15.9 |
| to calculation | 124 | 16.1 |
| skill? | 137 | 24.2 |
| | 135 | 28.9 |
| | 135 | 31.8 |
| | 137 | 22.3 |
| | 131 | 18.9 |
| | 100 | 15 |

Estimation of age-adjusted correlation coefficient

- Correlation coefficient between weight and calculation skill was 0.79.
- Age is related to both variables, weight and calculation skill : age is a confounder.



partial correlation coefficient

-0.23 • • • after adjusting the effect of age

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)
$$y = \alpha + x_1\beta_1 + x_2\beta_2 + \ldots + x_k\beta_k + \varepsilon$$

where y is a continuous dependent variable, x is a single predictor in the simple regression model, and x_1 , $x_2, ..., x_k$ 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

Logistic regression analysis

Logistic regression is used to model <u>the</u> <u>probability of a binary response</u> as a function of a set of variables thought to possibly affect the response (called covariates).

Y =
$$\begin{cases} 1: \text{ case (with the disease} \\ 0: \text{ control (no disease)} \end{cases}$$

One could imagine trying to fit <u>a linear model</u> (since this is the simplest model !) for the probabilities, but often this leads to problems:



In a linear model, fitted probabilities can fall <u>outside</u> of 0 to 1. Because of this, linear models are seldom used to fit probabilities. In a logistic regression analysis, the **logit** of the probability is modeled, rather than the probability itself.

P = probability of getting disease $(0 \sim 1)$



As always, we use the natural log. The logit is therefore **the log odds**, since odds = p / (1-p)

Logistic regression model

Now, we have the same function with linear regression model in the right side.

logit (px) = log
$$\begin{bmatrix} px \\ ---- \\ 1 - px \end{bmatrix} = \alpha + \beta x$$

where px = probability of event for a given value x, and α and β are unknown parameters to be estimated from the data.

 \rightarrow Multivariable analysis is applicable to adjust the effect of confounding factor.

Interpretation of coefficients of logistic regression model

The explanatory variable is binomial one: 1=exposed or 0=unexposed

Exposed: $\log (O_e) = \alpha + \beta (Exp=1) = \alpha + \beta$ Unexposed: $\log (O_u) = \alpha + \beta (Exp=0) = \alpha$ Difference: $\log(O_e) - \log (O_u) = \beta$ Odds ratio: $O_e / O_u = e^{\beta}$

Coefficient estimate: Odds ratio (after exponentiating)

SURVIVAL ANALYSIS

Survival analysis

- Survival time: from the entry point(for example, when the treatment starts) until end point(for example, disease recurrence or the death from the disease)
- Censoring: the follow-up is stopped because of other reason (for example, study period is over or the death from other reason)





Log-rank test:

Statistic test for the difference of survival probability



The statistic follows the chi-square distribution (df=1)

P=0.016

Limitations of Kaplan-Meier method

- Mainly descriptive
- Doesn't control for covariates
- Requires categorical predictors
- Can't accommodate time-dependent variables

Cox proportional hazard model

This model is expressed by the following formula; $\lambda(t | x_1, ..., x_k) = \lambda_0(t) \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)$ where $\lambda(t)$ is the hazard of variable x_k and t is the time until the case is alive, and $\lambda_0(t)$ is the baseline hazard. We assume that the log of hazard ratio is proportional to the variable X.

Log negative-log plot is useful to check

Hazard ratio: $\lambda_1(t)/\lambda_0(t) = \exp(\beta)$

Exposed group

Un-exposed group





STRATEGY FOR CONSTRUCTING REGRESSION MODELS

Basic principles

- 1. <u>Stratified analysis should be first.</u>
- 2. Determine which **confounders to include** in the model.



3. Estimate the shape of the exposuredisease relation.

Dose-response relation

4. Evaluate interaction

How to determine confounders: data-dependent manner

- 1. Start with a set of predictors of outcome based on the strength of their relation to the outcome.
- 2. Build a model by introducing predictor variables one at a time: check the amount of change in the coefficient of the exposure term
 - > 10% change: include it as a confounder

Example of a confounder (age)

Ability of maths (AM) = α + β 1(Height) \rightarrow AM = -42.8 + 0.41(Height) > 10% change AM = α + β 1(Height) + β 2(Age) \rightarrow AM = 1.48 - 0.01(Height) + 2.02 (Age) How to determine confounders: data-independent manner

Some researchers argue that "<u>Without data analysis</u>, decide confounders, important risk factors of the outcome, based on the previous studies."

How can we pick-up "important risk factors"? If there are few studies, how can we know confounders?



How many explanatory variables can we use in a model?

| Model | Number of explanatory variables | Example |
|-------------------------------|---|--|
| Linear regression model | Sample size / 15 | <u>Up to around 6-7</u> 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 a categorical variable in your model, you have to count that 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.

PROPENSITY SCORE

If you cannot recruit enough sample size

Calculate "propensity score" which can be used for adjustment of confounding effects.

Example

Aspirin Use and All-Cause Mortality Among Patients Being Evaluated for Known or Suspected Coronary Artery Disease

A Propensity Analysis

Patricia A. Gum, MD Maran Thamilarasan, MD Junko Watanabe, MD Eugene H. Blackstone, MD Michael S. Lauer, MD

Context Although aspirin has been shown to reduce cardiovascular morbidity and short-term mortality following acute myocardial infarction, the association between its use and long-term all-cause mortality has not been well defined.

Objectives To determine whether aspirin is associated with a mortality benefit in stable patients with known or suspected coronary disease and to identify patient characteristics that predict the maximum absolute mortality benefit from aspirin.

 Table 1. Baseline and Exercise Characteristics According to Aspirin Use*

| Variable | Aspirin (n = 2310) | No Aspirin (n = 3864) | Р Value |
|--|-----------------------|-----------------------------|------------|
| Demographics | | | |
| Age, mean (SD), y | 62 (11) | 56 (12) | <.001 |
| Men, No. (%) | | 2167 (56) | <.001 |
| Clinical history Diabetes, No. (%) Almost all pro | gnostic | 432 (11) | <.001 |
| Hypertension, No. (%) factors (n=2) | 8) are | 1569 (41) | <.001 |
| Tobacco use, No. (%) | o) aro | 500 (13) | .001 |
| Prior coronary artery c related to asni | rin use! | 778 (20) | <.001 |
| Prior coronary artery by | | 2 (9) | <.001 |
| Prior percutaneous coronary intervention, No. (%) | 667 (29) | 148 (4) | <.001 |
| Prior Q-wave MI, No. (%) | 369 (16) | 285 (7) | <.001 |
| Atrial fibrillation, No. (%) | 27 (1) | 55 (1) | .04 |
| Congestive heart failure, No. (%) | 127 (6) | 178 (5) | .12 |
| Medication use | | | |
| Digoxin use, No. (%) | 171 (7) | 216 (6) | .004 |
| β-Blocker use, No. (%) | 811 (35) | 550 (14) | <.001 |
| Diltiazem/veraparnil use, No. (%) | 452 (20) | 405 (10) | <.001 |
| Nifedipine use, No. (%) | 261 (11) | 283 (7) | <.001 |
| Lipid-lowering therapy, No. (%) | 775 (34) | 380 (10) | <.001 |
| ACE inhibitor use, No. (%) | 349 (15) | 441 (11) | <.001 |
| Cardiovascular assessment and exercise capacity Body mass index, mean (SD), kg/m ² | 29 (5) | 30 (7) | <.001 |
| Ejection fraction, mean (SD), % | 50 (9) | 53 (7) | <.001 |
| Resting heart rate, mean (SD), beats/min | 74 (13) | 79 (14) | <.001 |
| | | | |

Denting Elected economic and a (CDV and 11a

After matching by propensity score, the distribution of prognostic factors are similar between aspirin users and non-users.

Table 3. Selected Baseline and Exercise Characteristics According to Aspn.

| Use in Propensity | | | |
|---|-----------------------|-----------------------------|------------|
| It is just like a RCT (pseud RCT) | Aspirin (n = 1351) | No Aspirin (n = 1351) | P Value |
| Demographics | | | |
| Age, mean (SD), y | 60 (11) | 61 (11) | .16 |
| Men, No. (%) | 951 (70) | 974 (72) | .33 |
| Clinical history | | | |
| Diabetes, No. (%) | 203 (15) | 207 (15) | .83 |
| Hypertension, No. (%) | 679 (50) | 698 (52) | .46 |
| Tobacco use, No. (%) | 161 (12) | 162 (12) | .95 |
| Cardiac variables | | | |
| Prior coronary artery disease, No. (%) | 652 (48) | 659 (49) | .79 |
| Prior coronary artery bypass graft, No. (%) | 251 (19) | 235 (17) | .42 |
| Prior percutaneous coronary intervention, No. (%) | 166 (12) | 147 (11) | .25 |
| Prior Q-wave MI, No. (%) | 194 (14) | 206 (15) | .52 |
| Atrial fibrillation, No. (%) | 21 (2) | 24 (2) | .65 |
| Congestive heart failure, No. (%) | 79 (6) | 89 (7) | .43 |
| | | | |

Table 4. Cox Proportional Hazards Analyses of Aspirin Use and Mortality Among Propensity-Matched Patients (n = 2702)*

| Model | Hazard Ratio (95% Cl) | <i>P</i> Value |
|---|-----------------------------|-------------------|
| Unadjusted | 0.53 (0.38-0.74) | .002 |
| Adjusted for propensity | 0.53 (0.38-0.74) | <.001 |
| Adjusted for propensity and selected variables† | 0.59 (0.42-0.83) | .002 |
| Adjusted for propensity and all covariates‡ | 0.56 (0.40-0.78) | <.001 |
| *Cl indicates confidence inte | erval. | r |

You need to include only propensity score in the model.

†Selected variables included prior coronary artery disease,

prior coronary artery bypass grafting, prior percutane-

ous intervention, and ejection fraction \leq 40%.

‡For a list of covariates, see Table 2 footnote (†).

WE SHOULD NOT RELY ON P VALUE TOO MUCH

Statistic significance vs. Clinical significance

Statistic significance *≠* Clinical significance

- P value(s) do NOT tell us the significance in clinical practice / biological importance.
- If your sample size is quite large, you may obtain a result with statistic significance. So what?

RCT of donepezil for Alzheimer's disease

Lancet. 2004 Jun 26;363(9427):2105-15.

Long-term donepezil treatment in 565 patients with Alzheimer's disease (AD2000): randomised double-blind trial.

Courtney, C1 Earroll D. Crave P. Hills P. Lupeb L. Sollwood E. Edwards S. Harduman W. Pafton, L. Crame P. London C. Shaw H, Bentham P; AD2000 Collaborative Group. Cognition averaged 0.8 MMSE points better 🕀 Auth (95%CI 0.5-1.2; p<0.0001) and functionality Abstrac BACKG 1.0 BADLS points (0.5-1.6; p<0.0001) with eimer's disease. We aimed to determin ogical symptoms, carers' **a**? psycholo donepezil over the first 2 years. METHO eriod in which they were randomly no completed this period were rerandomised to either donepezil (5 or 10 mg/day) or placebo, with allocated donepezil (5 mg/ double-blind treatment co Donepezil is not cost effective, with benefits defined by loss of either all patier assessments were soud below minimally relevant thresholds. More FINDINGS: Cognition averaged 0.8 MI points better (0.5-1.6; p<0.0001) with d effective treatments than cholinesterase institutionalisation (42% vs 44% at 3 ve bnal care in the donepezil group compared inhibitors are needed for AD. institutional care was 0.96 (95% CI 0.7 and psychological symptoms, carer psy mq donepezil.

INTERPRETATION: Donepezil is not cost effective, with penefits below minimally relevant thresholds. More effective treatments than cholinesterase inhibitors are needed for Alzheimer's disease.

Which description is appropriate?

- In a RCT study, the mortality rates of new drug A and old drug B were 30% and 20%, respectively. And, the <u>p value was 0.6</u> for the difference of them.
- 1. The mortality rate of drug A is equivalent to that of drug B.
- 2. We cannot say "there is a difference in the mortality between drug A and B".
- 3. We cannot say "the mortality of drug A is higher than that of drug B".

We cannot tell the equivalence by p value

- In the previous example, the sample size of each arm was 10.
- The result tells us "you failed to reject the null hypothesis because of small sample size".

□ 293 subjects for each arm are required.

The difference of mortality rate is 10% and its 95% CI is -30%, 50%.

American Statistical Association Releases Statement on Statistical Significance and P-values

Provides Principles to Improve the Conduct and Interpretation of Quantitative Science https://www.amstat.org/newsroom/pressreleases/P-ValueStatement.pdf

- 1. P-values can indicate how incompatible the data are with a specified statistical model.
- 2. P-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.
- Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold.
- 4. Proper inference requires full reporting and transparency.
- 5. A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.
- 6. By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.