

university of copenhagen

Combined effects of two or more variables, statistical considerations

Niels Keiding & Anne Helby Petersen
Section of Biostatistics

Scientific meeting on the use of composite variables in medical research

Danish Epidemiological Society & Danish Society for Occupational and Environmental Health

September 20, 2018 - Bispebjerg Hospital, Copenhagen

Slide 1/16

university of copenhagen section of biostatistics **Example: PRISME-data**

Data courtesy of Sigurd Mikkelsen, see e.g. Vammen et al. (2016) *JOEM* **58**, 994-1001.

Public employees from Aarhus, Denmark were recruited in 2007 and followed up in 2009. At baseline they filled out a questionnaire and the

follow-up focused on diagnosis of clinical depression (individuals with clinical depression at baseline were excluded). Scored according to demands and control (dichotomized at median values) the following results were recorded:

Depression: yes Depression: no Total
n = 3035 n = 59 n = 3094

Low demands, high control 892 (29.4) 13 (22.0) 905 (29.3) Low demands, low control 668 (22.0) 14 (23.7) 682 (22.0) High demands, high control 795 (26.2) 7 (11.9) 802 (25.9) High demands, low control 680 (22.4) 25 (42.4) 705 (22.8)

We model

$$P(\text{clinical depression}) = \frac{\exp(-\mu_1 z_1 + \mu_2 z_2 + \dots + \mu_k z_k)}{1 + \exp(-\mu_1 z_1 + \mu_2 z_2 + \dots + \mu_k z_k)}$$

which is called $\text{logit}(-\mu_1 z_1 + \mu_2 z_2 + \dots + \mu_k z_k)$.

In the beginning, we use the principal exposure variables

;

$$z_1 = z_2 = \begin{cases} 1, & \text{if demands} \\ & \text{demands} < \text{median} \\ & > \text{median} \\ 0, & \text{if control} \\ & < \text{median} \end{cases}$$

$z_3 = z_1 \cdot z_2 =$
 0, if control > median ; median ('strain') 0, otherwise
 1, if demands > median and control <

Later additional covariates z_4, z_5, \dots may be added to handle confounding (examples: sex, age, socio-economic status).

Statistical models for occurrence of depression *Standard full statistical model*

demands
 low high
 control high — — + ”
 low — + “ —+“+”+—

Estimates 95% Conf. interv. P

—: intercept -4.229 (-4.829, -3.725)

”: main effect of demands -0.504 (-1.487, 0.394) 0.285 “: main effect of control 0.363

(-0.404, 1.137) 0.350 —: interaction effect of demands and control ('strain') 1.066

(-0.046, 2.248) 0.066

Tendency to interaction effect: the combination of high demands and low

control increases probability of clinical depression.

Slide 4/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence of depression *Standard full statistical model*

Net contrast to baseline: High control, low demands

demands

low high
 control high *Baseline* $\neq 0.504$
 low 0.363 0.925^u

Estimates 95% Conf. interv. P

u : main effect of demands -0.504 ($-1.487, 0.394$) 0.285 c : main effect of control 0.363
 ($-0.404, 1.137$) 0.350 —: interaction effect of demands and control ('*strain*') 1.066
 ($-0.046, 2.248$) 0.066

(*) $0.925 = 0.363 \neq 0.504 + 1.066$

Statistical models for occurrence of depression *Model IPD ('strain')*

demands

low high

control high – –

low – – + —

Main effects of control and demands are assumed to be zero.

Estimates 95% Conf. interv. P
-: intercept -4.229 (-4.829, -3.725)

—: interaction effect of demands and control ('strain') 0.935 (0.402, 1.454) 0.000

Strong 'strain' effect, partly generated by main effects that are not modelled.

Slide 5/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics **The**

'hierarchical principle' of

regression analysis

If a model contains an interaction between categorical variables, then we must keep lower order interactions and main effects associated with this interaction in the model as well.

Violated by the '*strain*' model.

Slide 6/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence of depression *Model IPD ('strain')*

*Net contrast to baseline: High control
and/or low demands*

demands

low high

control high *Baseline Baseline*

low *Baseline* 0.935

Main effects of control and demands are
assumed to be zero.

Estimates 95% Conf. interv. P

—: interaction effect of demands and control ('strain') 0.935 (0.402, 1.454) 0.000

The contrasts to baseline are heavily driven

*by the model
assumptions of vanishing main effects.*

Slide 7/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence
of depression *Model with only main
effects*

demands
low high
control high — — + ”

low – + “ – + “ + ”

Interaction effect of demands and control
(*'strain'*) assumed to be zero.

Estimates 95% Conf. interv. P

–: intercept -4.530 (-5.081, -4.045)

Ⓜ: main effect of demands 0.193 (-0.324, 0.718) 0.464

Ⓛ: main effect of control 0.885 (0.353, 1.448) 0.001

***Strong estimated main effect of control
(low control related to increased risk of
depression).***

Statistical models for occurrence
of depression *Model with only main
effects*

*Net contrast to baseline: High control,
low demands*

demands

low high

control high *Baseline* 0.193

low 0.885 1.078^ú

Interaction effect of demands and control ('strain') assumed to be zero.

Estimates 95% Conf. interv. P

β: main effect of demands 0.193 (-0.324, 0.718) 0.464

β: main effect of control 0.885 (0.353, 1.448) 0.001

$$(*) 1.078 = 0.193 + 0.885$$

Slide 8/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence of depression *Model with freely*

varying parameters in the four possible combinations of demands and control

demands

low high

control high -- + μ_2

low - + μ_3 - + μ_4

This is just a reparameterization of the standard model:

$$\mu_2 = \text{''}, \mu_3 = \text{''}$$

$$\mu_4 = \text{“ + ”} + \text{—}, \text{ i.e. —} = \mu_4 \neq \mu_2 \neq \mu_3$$

Slide 9/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence
of depression *Model with equal main
effects of demands and control as well as
interaction between them*

demands
 low high
 control high — — + “
 low — + “ — + 2“ + —

Estimates 95% Conf. interv. P

—: intercept -4.229 (-4.829, -3.725)

“: main effect of control = main effect of demand -0.015 (-0.700, 0.707) 0.966 2“ + —:

net contrast to baseline for high demands, low control 0.925 (0.265, 1.633) 0.007

Note: — is still the interaction effect of demands and control
 ('strain')

Strong strain effect, but partly driven by

arbitrary assumption of equal effects of demands and control.

Slide 10/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence of depression *Model with equal main effects of demands and control as well as interaction between them*

Net contrast to baseline: High control, low demands

demands

low high

control high *Baseline* ≠0.015

low ≠0.015 0.925

Estimates 95% Conf. interv. P

“: main effect of control = main effect of demand -0.015 (-0.700, 0.707) 0.966 2^{***} + —:
net contrast to baseline for high demands, low control 0.925 (0.265, 1.633) 0.007

Note: — is still the interaction effect of demands and control
(‘strain’)

Statistical models for occurrence of depression Relations between models

Standard full statistical model:

—, “, ”, —

— = 0 “ = ”

Model with only main effects: —, “, ” Model IPD (strain):

Model with equal main effects and interaction:

—, —

—, “, —

“ = ” = 0

**Statistical models for occurrence
of depression** *Standard full statistical
model*

Minimal confounder control (age, sex, marital status)

Modest confounder control (age, sex, marital status, education)

demands

low high

control high — — + ”

low — + “ —+ “+”+—

Estimates Estimates Estimates

—: intercept -4.229 -5.298 -5.378

”: main effect of demands -0.504 -0.513 -0.508

“: main effect of control 0.363 0.378 0.395

—: interaction effect of demands and control ('strain') 1.066 1.054 1.046

*Small effects of correction for possible confounding, no change in general conclusion:
Tendency to interaction effect: the combination of high demands and low control increases probability of clinical depression.*

Slide 12/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

Statistical models for occurrence of depression *Model IPD ('strain')*

Minimal confounder control (age, sex, marital status)

Modest confounder control (age, sex, marital status, education)

demands

low high

control high — —

low — — + —

Main effects of control and demands are assumed to be zero.

Estimates Estimates Estimates

—: intercept -4.229 -5.259 -5.230

—: interaction effect of demands and control ('strain') 0.935 0.926 0.931

Small effects of correction for possible

*confounding, no change in general conclusion:
Strong 'strain' effect, partly generated by main
effects that are not modelled.*

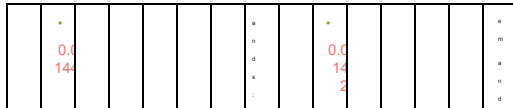
Slide 13/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research

university of copenhagen section of biostatistics

**An overview of predicted
depression probabilities** In models
with no confounder adjustment, we estimate

$$P(\text{clinical depression}) = \frac{\exp(- + \dots)}{1 + \exp(- + \dots)}$$

Control: High Control: Low



a
 n
 d
 s
 :
 n
 i
 g
 n
 a
 n
 t
 s
 :
 144

effects and interaction

0.01 0.02 0.03 0.04 0.01 0.02 0.03 0.04
 Probability of depression

selection

”Stepwise variable selection has been a very popular technique for many years, but if this procedure had just been proposed as a statistical method, it would most likely be rejected because it violates every principle of statistical estimation and hypothesis testing.”

F.E. Harrell (2015). *Regression Modeling Strategies*. Second Edition. Springer, New York, p. 67.

There are two issues in model selection: we hope to *minimize bias* and to *minimize variance*. Much variable selection goes too far in focusing on the second target, getting down to very few variables with bias as a result.

Stay with a model that does cover the important structure in the problem and the data and does not violate the hierarchical principle.

Conclusion

As occasional visitors to this interesting area we feel that the standard mainstream statistical approach would be adequate and relevant for these issues.

Slide 16/16 — Niels Keiding & Anne Helby Petersen — Combined effects of two or more variables, statistical considerations
— Scientific meeting on the use of composite variables in medical research