



Composites: Modeling with Formative Measures

Jim Grace

U.S. Department of the Interior
U.S. Geological Survey

1

In this module I consider an important “third” variable type, the composite. Composites are similar to latent variables, but with some fundamentally important differences. This is currently a very brief introduction to the topic, so going to the primary reference below will definitely be helpful.

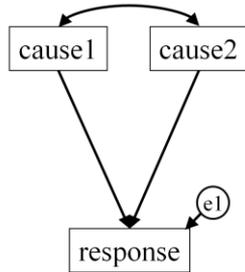
An appropriate citation for this material is

Grace, J.B. and Bollen, KA. 2008. Representing general theoretical concepts in structural equation models: the role of composite variables. *Environmental and Ecological Statistics* 15:191-213.

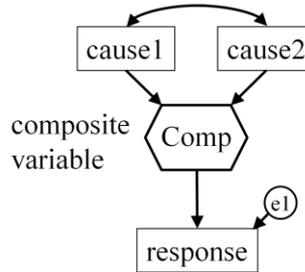
([http://www.odum.unc.edu/content/pdf/Bollen%20Grace%20Bollen%20\(preprint%202008\)%20Environ%20and%20Ecol%20Stats.pdf](http://www.odum.unc.edu/content/pdf/Bollen%20Grace%20Bollen%20(preprint%202008)%20Environ%20and%20Ecol%20Stats.pdf))

Notes: IP-056512; Support provided by the USGS Climate & Land Use R&D and Ecosystems Programs. I would like to acknowledge formal review of this material by Jesse Miller and Phil Hahn, University of Wisconsin. Many helpful informal comments have contributed to the final version of this presentation. The use of trade names is for descriptive purposes only and does not imply endorsement by the U.S. Government. Last revised 20141216. Questions about this material can be sent to sem@usgs.gov.

1. In the simple case, we can use composites to represent collective effects of a set of variables.



First step in compositing process.



Second step in compositing process.



2

Composites are a type of latent variable, but they differ from conventional LVs in some important ways. Certainly they are abstractions, but in our models their values are computed from other variables. In this example we are essentially using a composite variable “Comp” to capture the collective effects of a set of causes on some response. Thus we are translating the model on the left into the model on the right so we can represent the joint effects of cause1 and cause2 with a single path (the path from Comp to response).

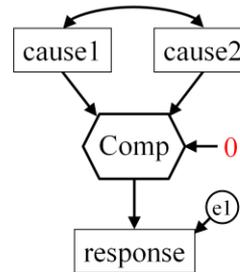
In this situation we refer to the observed indicators cause1 and cause2 as “cause indicators” or sometimes referred to as “formative indicators”.

I like to use the hexagon shape to represent the composite. However, note that often it is represented in presentations by ovals, since it is a form of latent variable.

2. Composites are abstractions, and we must declare them.

Composites have an implied set of scores, similar to regular latent variables.

To specify composites in lavaan, we need to conceptualize a composite as a latent variable with no error (so, error variance is set = 0).



Like with regular latent variables, introducing a composite variable into a model potentially adds two parameters to our model, one for the variance of the composite and one for its scale (intercept in this case). So, we need to fix some parameters to specific values to identify our model.



3

Composites are technically latent variables without variance. The absence of variance is modeled by setting the error variance to zero.

Another way of thinking about this is that the composite variable is one with a predicted set of values, one for each case in the dataset. For this model, the Comp scores are equivalent to

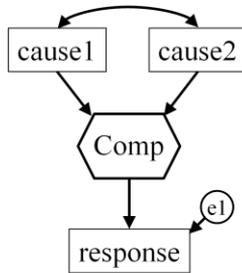
$$\hat{y} = b_0 + b_1 * \text{cause1} + b_2 * \text{cause2}$$

where the b s come from the model

$$\text{responses} \sim \text{lm}(\text{cause1} + \text{cause2}).$$

We will demonstrate this as a method of compute composites by hand later in the presentation.

3. We can create composites with lavaan syntax.



```
# we declare composite as effects of two variables  
Cmod.1 <- `Comp <~ 1*cause1 + cause2  
           response ~ Comp`
```

We use “1*cause1” to set scale of the composite.

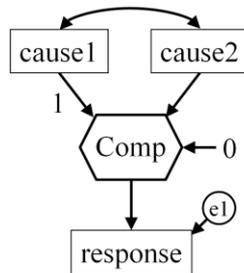


4

Lavaan has a special operator for composites. Just like with latent variables, we have to give the program some information. Here we explicitly indicate that the composite is on the same scale as the first indicator by pre-multiplying cause1 by 1.

3 (cont.). Creating composites with lavaan syntax - explained.

When we say “ $\text{Comp} \sim 1*\text{cause1} + \text{cause2}$ ”, we are specifying one of the paths linking the composite with an indicator (in this case, “cause1”).



Lavaan automatically creates the variable “Comp” in this case and sets its error variance to zero.



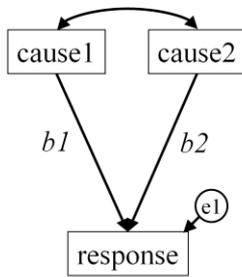
5

Here is a little more behind-the-scenes information. This slide shows more explicitly what the syntax in the previous slide does.

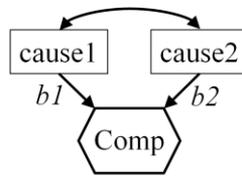
Note that we could have set the scale for cause2 instead of cause1. Also, we could use a different value from 1. For example, we could use the value that came from running the model without the composite.

Further, we could set both values from causes to COMP, but to do that we would need to use the two exact values obtained from running the model without the composite (slide 2 left figure).

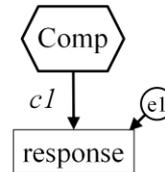
4. We can also compute composite scores by hand.



First, run model without composite and obtain coefficients.



Use coefficients to calculate composite scores.



Use composite scores as a variable.

The coefficient $c1$ represents collective effects and is only meaningful in standardized units (normally, its raw value is 1.0).



It is useful to understand how to compute composite scores by hand.

One reason to do this is because lavaan can have problems solving models containing composites sometimes. Various tricks for helping lavaan include (a) premultiplying cause1 by the exact coefficient found from the model that omitted composites. In rare cases, using the composite scores compute by hand may be the only way to successfully model.

This figure indicates how we can have the values of the composite variable and model with it directly. An example of this follows.

5. An illustration.



```
## lavaan model
library(lavaan)

# model without composite
mod.1 <- 'response ~ cause1 + cause2'

mod.1.fit <- sem(mod.1, data=sim.dat)

summary(mod.1.fit, rsq= T)
```

	Estimate	Std.err	Z-value	P(> z)
Regressions:				
response ~				
cause1	0.780	0.108	7.206	0.000
cause2	0.749	0.230	3.249	0.001
Variances:				
response	2.649	0.530		
R-Square:				
response	0.731			

7

```
### Simulate some data for an example
library(MASS)
mu = c(10, 20)
Sigma <- matrix(c(10, 3, 3, 2), 2, 2)
set.seed(1, kind = NULL, normal.kind = NULL)
xs <- mvrnorm(n = 50, mu, Sigma)
cause1 = xs[,1]
cause2 = xs[,2]
set.seed(3, kind = NULL, normal.kind = NULL)
yerror = rnorm(50, mean=0, sd=2)
response <- 0.812*x1 + 0.673*x2 + yerror
sim.dat <- data.frame(cause1, cause2, response)
```

Now you can use the R code on the slide to generate the results shown.

5. An illustration (cont.).

```
# lavaan model with composite
Cmod.1 <- 'Comp <~ 1*cause1 + cause2
          response ~ Comp'

Cmod.1.fit <- sem(Cmod.1, data=sim.dat)

summary(Cmod.1.fit, rsq=T, standardized=T)
```

	Estimate	Std.err	Z-value	P(> z)	Std.all
Composites:					
Comp <~					
cause1	1.000				0.761
cause2	0.960	0.389	2.471	0.013	0.343
Regressions:					
response ~					
Comp	0.780	0.108	7.206	0.000	0.855
R-Square:					
response	0.731				

interpretable coef

8

Again, you can use the “sim.dat” data created from the R script in the notes of slide 7 and the lavaan code in the slide here to generate the output.



5. An illustration (cont.).

```
# compute composite scores "by hand"
Comp.a <- 0.780*cause1 + 0.749*cause2

# create new data set
dat2 <- data.frame(cause1, cause2, response, Comp.a)

# model direct effect of composite on response
Cmod.2 <- 'Comp.a ~ cause1 + cause2
          response ~ Comp.a'

Cmod.2.fit <- sem(Cmod.2, data=dat2)
summary(Cmod.2.fit, rsq=T, standardized=T)
```

	Estimate	Std.err	Z-value	P(> z)	Std.all
Regressions:					
response ~					
Comp.a	0.999	0.086	11.663	0.000	0.855
R-Square:					
response	0.731				



we get same R-square and std. coef.

9

Here we use some R code to compute composite scores (“Comp.a”) by hand.

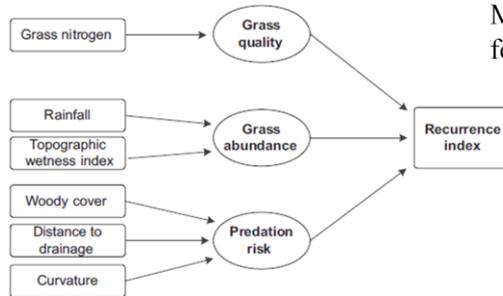
We then create a new dataset with the composite scores = “dat2”

We can model using the composite scores as observed variables.

6. We can use composites to represent general concepts.

Body size and the division of niche space: food and predation differentially shape the distribution of Serengeti grazers

J. Grant C. Hopcraft^{1,2*}, T. Michael Anderson³, Saleta Pérez-Vila⁴, Emilian Mayemba⁵ and Han Olff¹



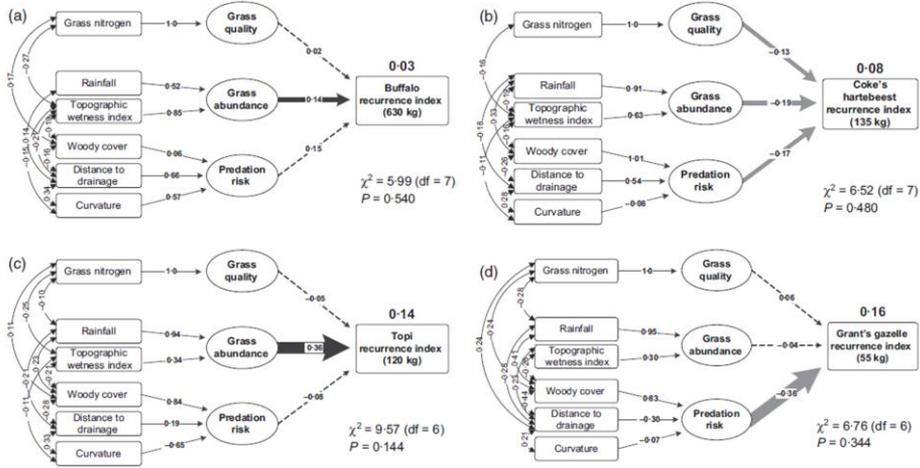
Model of predation risk versus food rewards for grazers.



Fig. 2. The *a priori* structural equation model used to assess the

Composites are very handy for use in addressing ecological questions.

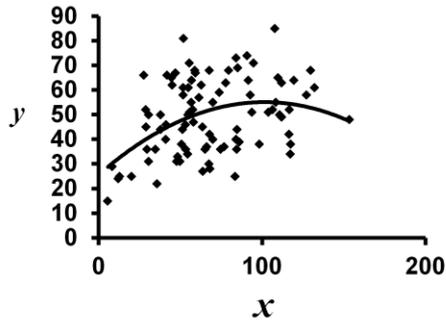
7. And we can use them to compare apples and oranges (with care in interpretation).



There capacity to let us make general comparisons is one of their most appealing features.

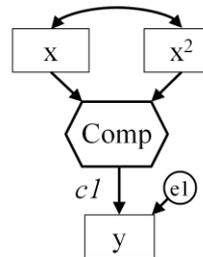
8. We can also use composites as devices.

Consider this nonlinear relationship.



We can model this using polynomial regression:

$$y = \gamma x + \gamma x^2 + \varepsilon$$



Centering a variable before squaring it reduces the autocorrelation between polynomial terms.

```
# center x before squaring  
x <- x - mean(x)  
x2 <- x^2
```



Note: endogenous nonlinearities covered in separate module.

12

But, composite are also handy devices, for example in polynomial modeling of nonlinear relations. What is different here is that one of the cause indicators, x-square, is really a device rather than a separate variable. Otherwise, the compositing process is similar to the case of compositing independent causes. This does not hold, however, when the nonlinearity is endogenous. In that case, the x-square term needs special consideration. This situation is covered in a separate module, “Composites for endogenous nonlinearities”.