

Introduction to structural equation modelling

Basic modelling

Frank Pennekamp

04.11.2021

Department of Evolutionary Biology and Environmental Sciences

University of Zurich

Introduction of the Swiss SEM team

- Dr. Frank Pennekamp (main instructor)
- Dr. James Grace (advanced topics and model clinic)
- Dr. Rachel Korn (course development)
- Dr. Noémie Pichon, Dr. Fletcher Halliday, Dr. Eliane Meier, Dr. Hugo Saiz, Dr. Debra Zuppinger-Dingley, Rebecca Oester, Annabelle Constance, Fabienne Wiederkehr (course development)



Schedule

- Day 1:
 - General introduction to SEM to model ecological systems
 - Fitting SEMs to data
 - Model pruning, visualization and reporting
 - Discussion with James Grace
- Day 2:
 - Latent and composite variables
 - Interactions
 - Complex sampling designs
 - Spatial autocorrelation
 - Discussion with James Grace
- Day 3:
 - Self-study with possibility to meet with instructor(s)

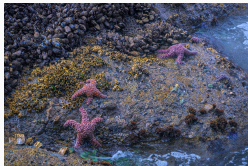
- What the course is about:
 - Global estimation with R package lavaan
 - Hands on exercises and live coding
 - We will work (mostly) with a single, ecological dataset
- What will not be covered
 - Local estimation of SEMs (with piecewiseSEM)
 - Advanced topics like incorporating random effects, feedbacks, temporal autocorrelation

Learning objectives

- Participants understand the benefits and limits of SEMs
- Participants are able to fit, interpret and visualize a SEM
- Participants are able to apply SEM to their own dataset

Getting started with Structural Equation Modeling

Research questions



- Ecology is about the interactions between organisms and their environment.
- We often have ideas how things could be connected in ecological systems.
- To test hypotheses, we need a way to dissect when they occur for a reason versus randomness.
- We use statistics to understand what is signal and what is noise.
- Research questions are often about cause and effect.

Why do we need to use SEM in Ecology

- System thinking
 - “Understanding the whole rather than the parts in isolation”
- SEM unites multiple variables in a single causal network: simultaneous tests of multiple hypotheses.
- Causality is central: SEM implicitly assumes that the relationships among variables represent causal links.

Correlation Vs. Causation



- “Correlation does not imply causation”
- Causation indicates a relation between two variables in which one variable is affected by another.

- We often have multiple observed variables.
- We want to test and evaluate multivariate causal relationship.
- Test direct and indirect effects on assumed causal relationships.
- Incorporate observed and latent variables.
- Include interaction terms can test main effects and interaction effects.

Two goals of SEM:

- 1) Understand the patterns of correlation/covariance among a set of variables.
- 2) Explain as much of their variance as possible with the model specified.

Thinking about the model

A SEM is usually specified based on theory to determine and validate a proposed causal process and/or model.

Which variables to include?

- Supported by theory.
- Ecologically meaningful.
- Garbage in - garbage out (both data quality and ecologically meaningful).

Good practice: Make a table / graph of putative causal relationships before analysis.

- Continuum between theory (hypothesis-driven) to exploratory (data-driven) modelling.
- Important to be explicit about the approach taken.

Introduction to the dataset



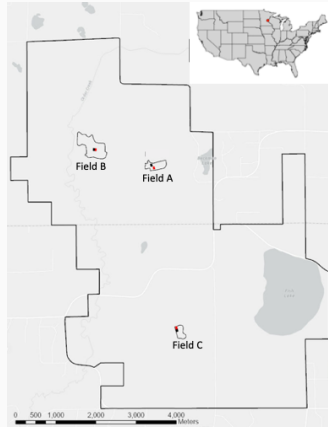
We will use an experimental dataset collected at the Cedar Creek Ecosystem Science Reserve to examine long-term consequences of human-driven environmental changes ecosystem responses to:

- Disturbance
- Nitrogen deposition
- Changes in precipitation

Understanding the recovery of a grassland for two decades following an intensive agricultural disturbance under ambient and elevated nutrient conditions.

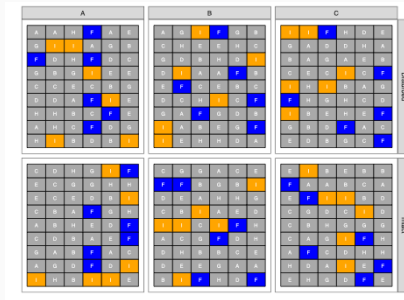
- 1) How has aboveground biomass changed as a function of disturbance (disking) and nutrient addition?
- 2) How are these effects mediated by diversity?

Introduction to the dataset



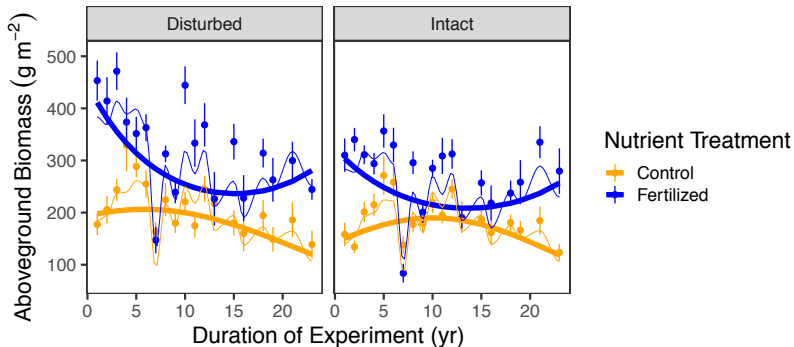
Map showing location of the study site (Cedar Creek Ecosystem Science Reserve), the location of the three study fields within the reserve, and location of the 35 x 55 m intact (black) and disturbed (red) plots within each field.

Introduction to the dataset



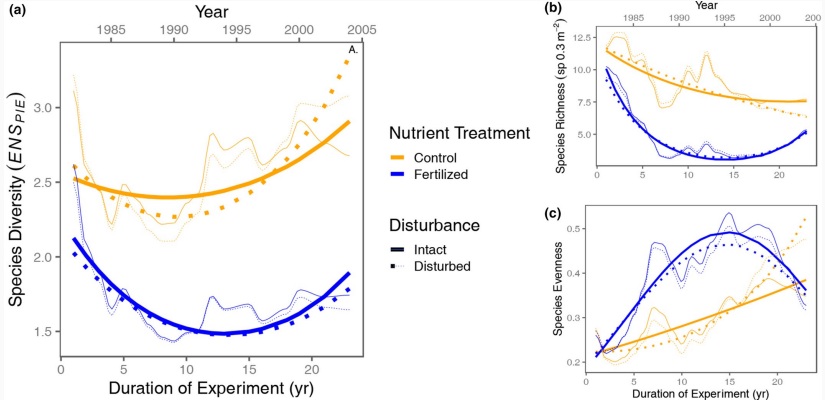
Location of the 4 x 4 m nutrient treatment plots within each 35 x 55 m Intact or Disturbed plot within each of three fields (A, B, and C). Letters indicate the nutrient treatments, and the colored plots are treatments that are the focus of the analyses presented here: Control (orange) and 9.5 g N m⁻² yr⁻¹ (blue).

Introduction to the dataset



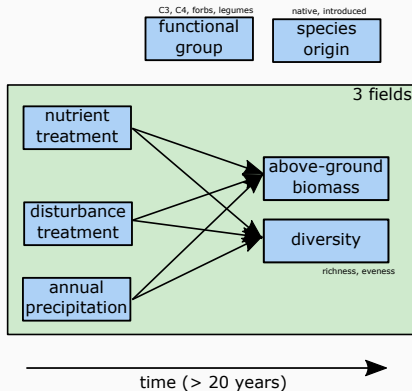
Effect of soil disturbance (disking) and nutrient enrichment on live, aboveground plant biomass. Colors indicate nutrient addition treatment: Control and NPK+ (all nutrients plus $9.5 \text{ g N m}^{-2} \text{ yr}^{-1}$).

Introduction to the dataset

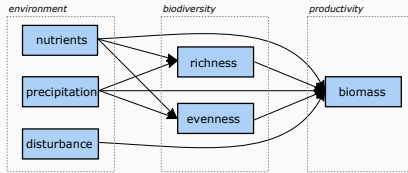


Effect of soil disturbance (disking) and nutrient enrichment on (a) diversity (ENS_{PIE}), (b) richness (S , species $0.3\ m^{-2}$), and (c) evenness ($ENS_{PIE}\ S^{-1}$).

Question of interest: what is the effect of richness on biomass?



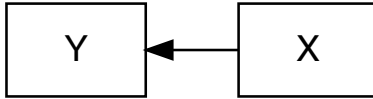
Meta-model



A metamodel summarizes the concept behind a model and links it to theory.

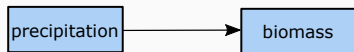
Productivity (biomass) is directly influence on the one hand by the environment (nutrients, disturbance and precipitation) and on the other hand by biodiversity (richness and evenness). Also some elements of the environment influence biodiversity and thus, have an additional indirect effect on productivity via biodiversity.

Graphical models:



- Directed acyclic graphs (no loops).
- Variables are nodes (boxes).
- Edges (one-headed arrows) are causal relationships such as X affects Y.
- Edges (two-headed arrows) are non-causal relationships such as X and Y are correlated.

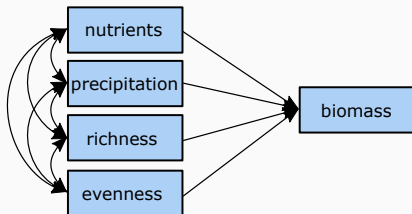
A simple bivariate model.



`lm(biomass ~ precipitation)`

- Linear regression
- Regression coefficient quantifies the strength of relationship
- Change in Y for one unit change in X

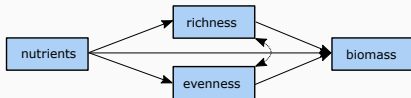
Multiple independent variables



`lm(biomass ~ precipitation + nutrients + ...)`

- Often more than one independent variable important → multiple regression
- Estimates partial regression coefficients (i.e. effect of precipitation on biomass when nutrient addition is fixed)
- Only direct effects.

Multiple independent variables



```
lm(richness ~ nutrients)
```

```
lm(evenness ~ nutrients)
```

```
lm(biomass ~ richness + evenness)
```

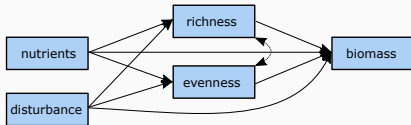
- Indirect effect is the effect of an independent variable on a dependent variable through one or more intervening or mediating variables.
- Indirect effects can be quantified by the product of the compound path.

Mediation



- Tests whether a particular variable has a mediating effect on a path.
- Often used to test underlying mechanisms.
- In our example, we could ask whether the effect of nutrients on biomass is mediated through biodiversity.
- Possibilities: complete mediation, partial mediation, no mediation.

System level approach



- **Exogenous** variables only have paths emanating from them (i.e., do not have arrows going into them).
- **Endogenous** variables have paths directed into them.
- An endogenous variable can also have arrows directing out of it, but the sole condition is that they must be predicted.

Grace's 8 rules of path coefficients

- 1) Unspecified relationships among exogenous variables are the bivariate correlations.
- 2) When two variables are connected by a single path, the coefficient of that path is the regression coefficient.
- 3) The strength of a compound path (one that includes multiple links) is the product of the individual coefficients.
- 4) When variables are connected by more than one pathway, each pathway is the 'partial' regression coefficient.

Grace's 8 rules of path coefficients

- 5) Errors on endogenous variables relate the unexplained correlations or variances arising from unmeasured variables.
- 6) Unanalyzed (residual) correlations among two endogenous variables are their partial correlations.
- 7) The total effect one variable has on another is the sum of its direct and indirect effects.
- 8) The total effect (including undirected paths) is equivalent to the total correlation.

- χ^2 statistic: good fit when failing to reject the null hypothesis that the χ^2 statistic is different from 0 ($P > 0.05$).
- Comparative fit index (CFI): this statistic considers the deviation from a 'null' model. In most cases, the null estimates all variances but sets the covariances to 0. A value > 0.9 is considered good.
- Root-mean squared error of approximation (RMSEA): statistic penalizes models based on sample size. A value < 0.10 is acceptable, and anything < 0.08 is good.
- Standardized root-mean squared residual (SRMR): the standardized difference between the observed and predicted correlations. A value < 0.08 is considered good.

- (Multivariate) Normality
- Global estimation¹
- Directed acyclic relationships²
- Linear relationships³

¹Local estimation possible.

²Causal loops possible.

³Nonlinear relationships possible.

- Underidentified: not enough pieces of information to identify parameters uniquely ($df < 0$)
- Saturated: Just enough information to uniquely identify parameters, but no df to check model fit ($df = 0$)
- Over-identified: parameters can be uniquely identified and positive dfs to test model goodness-of-fit ($df > 0$)

“t-rule” to quickly gauge whether a model is under-, just, or overidentified:

$$t \leq \frac{n(n+1)}{2}$$

t = number of unknowns (parameters to be estimated)

n = number of knowns (observed variables).

The left hand side is how many pieces of information we want to know.
The right hand side reflects the information we have to work with and is equal to the number of unique cells in the observed variance-covariance matrix (diagonal = variance, and lower triangle = covariances).

- Replication should be at least 5x the number of estimated coefficients (not error variances or other correlations).
- To estimate two relationships, we require at least $n = 10$ to fit that model.
- Ideally, replication is 5-20x the number of estimated parameters.
- The larger the sample size, the more precise (unbiased) the estimates.

- 1) Review the relevant theory and research literature to support model specification
- 2) Specify a model (e.g., diagram, equations)
- 3) Determine model identification
- 4) Select measures for the variables represented in the model
- 5) Collect data
- 6) Conduct preliminary descriptive statistical analysis (e.g., scaling, missing data, collinearity issues, outlier detection)
- 7) Estimate parameters in the model
- 8) Assess model fit
- 9) Re-specify the model if meaningful
- 10) Interpret and present results visually

Questions?
