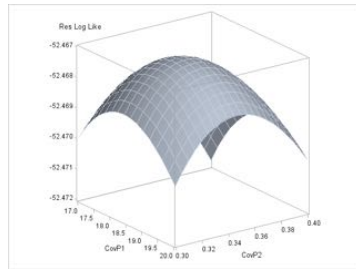# Introduction to Likelihood Methods for SEM

Jarrett E. K. Byrnes
University of Massachusetts Boston



$$\Sigma = \Sigma(\Theta)$$

---

"There are no routine statistical questions, only questionable statistical routines"
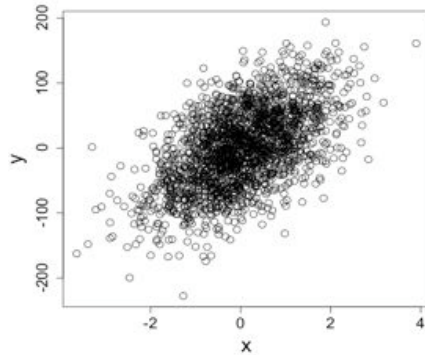
- Sir David Cox

---

## What is Covariance-Based SEM Estimation with Likelihood?

- Estimation of parameters given covariance of the data

- Equivalent to Linear Regressions, but…

- Estimation of each parameter influences the others

- Can accomodate unobserved (latent) variables and feedbacks

---

## A Likely Outline

1. What is SEM using likelihood and covariance matrices?

2. Model Identifiability

3. Sample Size for SEM

4. Standardized Coefficients

5. Introduction to `lavaan`

## Covariance and Correlation



$$COV_{xy} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{n-1}$$

$$r_{xy} = \frac{COV_{xy}}{SD_x SD_y}$$

## Covariance and Correlation

We often use covariances to calculate slopes, but standardized covariances – i.e. correlations – for interpretation.

Raw Covariance Matrix

|       | $x_1$ | $x_2$ | $y_1$ |
|-------|-------|-------|-------|
| $x_1$ | 0.81  |       |       |
| $x_2$ | 0.87  | 1.63  |       |
| $y_1$ | 0.88  | 1.80  | 4.98  |

variance          covariance

Standardized Covariance Matrix

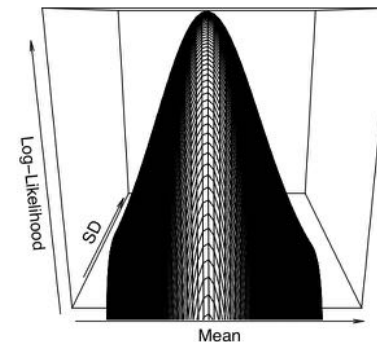|       | $x_1$ | $x_2$ | $y_1$ |
|-------|-------|-------|-------|
| $x_1$ | 1.0   |       |       |
| $x_2$ | 0.76  | 1.0   |       |
| $y_1$ | 0.44  | 0.63  | 1.0   |

correlation

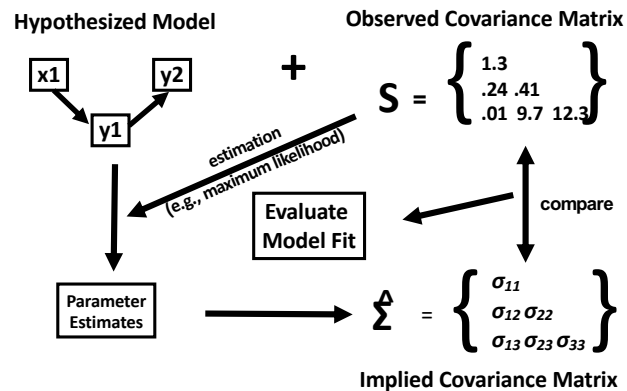## Maximizing Likelihood with One Parameter



Iteration over possible values simple

## Likelihood with Two Parameters



- Algorithms used to search parameter space
- Integrate answer over all data points
  - difficult computationally!

## How does ML Estimation Work?

**Hypothesized Model**          **Observed Covariance Matrix**



## What we're used to with ML

Data Generation: $\mu_i = a + bX_i$

Likelihood Function: $F_r = Y_i \sim dnorm(\mu_i, \sigma)$

**We minimize the likelihood function, $F_r$**

## It's…More Complicated with SEM

Data Generation:

$$\Sigma = \begin{pmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{pmatrix} = \begin{pmatrix} \Lambda_y (I-B)^{-1} (\Gamma\Phi\Gamma' + \Psi)(I-B)^{-1'} \Lambda_y' + \theta_\varepsilon & \Lambda_y (I-B)^{-1} \Gamma\Phi\Lambda_x' \\ \Lambda_x \Phi\Gamma'(I-B)^{-1'} & \Lambda_x\Phi\Lambda_x' + \theta_\delta \end{pmatrix}$$

Likelihood Function:

$$F_{ML} = \log\left|\hat{\Sigma}\right| - \log\left|\mathbf{S}\right| + tr\left(\mathbf{S}\hat{\Sigma}^{-1}\right) - (p+q)$$

## The Maximum Likelihood Fitting Function

$$F_{ML} = \log\left|\hat{\Sigma}\right| - \log\left|\mathbf{S}\right| + tr\left(\mathbf{S}\hat{\Sigma}^{-1}\right) - (p+q)$$

S = Sample covariance matrix
S = Fit covariance matrix
p = endogenous variables
q = exogenous variables

*Linear Algebra Review*

*Det(A) = scalar number*

*A\*A$^{-1}$ = Diagonal matrix of ones*

- If S = $\Sigma$, term 1 - 2 = 0 and terms 3 - 4 = 0.
- $F_{ML}$ = 0 with perfect fit

## Assumptions Behind $F_{ml}$

- Multivariate normality
  - Fairly robust (non-normality of residuals bigger problem)
  - Test with multivariate Shapiro-Wilk's Test (library mvnormtest)
  - In particular, no skew
  - Severe violations bias parameter error and tests of model fit

- No missing data in calculation of S
  - Biases your estimates with pairwise corrections

- No redundant variables
  - S must be positive definite

- Sample size is "large" (more soon)

## A Likely Outline

1. What is SEM using likelihood and covariance matrices?

2. Model Identifiability

3. Sample Size for SEM

4. Standardized Coefficients

5. Introduction to `lavaan`

## Identifiability

1. To fit a model, it must be <u>identified</u>

2. We need as much unique information as parameters

3. What can make a model non-identified?
   - Too many paths relative to # of variables
   - Certain model structures
   - High multicollinearity (r>0.9)
   - Complex model & small sample
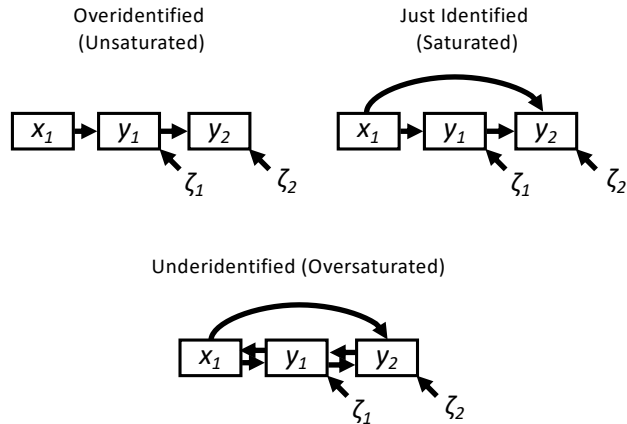
4. How do I know if my model is identified?

## Identification: An Algebra-Eye View

| | | |
|---|---|---|
| 3 = a + b | | |
| 4 = 2a + b | *a* and *b* have unique solutions | **Just identified** |
| | | |
| 3 = a + b + c | *a*, *b*, and *c* have no unique | |
| 4 = 2a + b + 3c | solution | **Underidentified** |

3 = a + b
4 = 2a + b
7 = 3b + a
*a* and *b* have unique solutions, more knowns than unknowns
**Overidentified**

## Different Model States

Overidentified
(Unsaturated)

$x_1$ → $y_1$ → $y_2$

$\zeta_1$ $\zeta_2$

Just Identified
(Saturated)

$x_1$ → $y_1$ → $y_2$

$\zeta_1$ $\zeta_2$

Underidentified (Oversaturated)

$x_1$ ⇄ $y_1$ ⇄ $y_2$

$\zeta_1$ $\zeta_2$

---

## The T-Rule
### # of Parameters v. Covariance Matrix

x1

$\gamma_{12}$ $\delta_2$

y1

$\beta_{12}$ $\zeta_1$

y2

$\zeta_2$

|  | x1 | y1 | y2 |
|---|---|---|---|
| x1 | 0.5 |  |  |
| Cov(x,y1,y2)= y1 | 0.7 | 0.5 |  |
| y2 | 0.2 | 0.8 | 0.3 |

• # Parameters ≤ # Unique Entries in a Covariance Matrix

**T-rule: t ≤ (p+q)(p+q+1)/2**

• t=# params, p = # endogenous variables, q = # exogenous variables

---

## How Do I Count the Number of Parameters?

x1  Yes, there is a variance here

$\gamma_{12}$

y1

$\beta_{12}$ $\zeta_1$

y2

$\zeta_2$

If variance and covariances among exogenous variables is not shown
either draw them or use modified formula:
**T-rule: t* ≤ (p+q)(p+q+1)/2 - q(q+1)/2**

---

## You will see path diagrams drawn many ways…

$\delta_1$ $\delta_2$

x1 x2

y1

$\zeta_1$

y2

$\zeta_2$
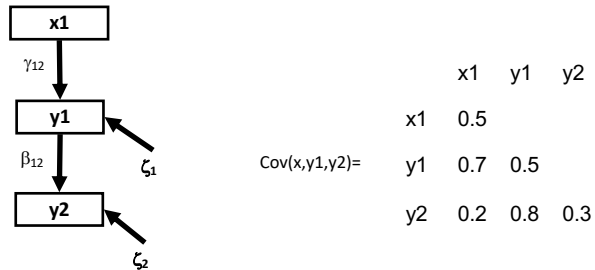
x1 x2

y1

$\zeta_1$

y2

$\zeta_2$

x1 x2

y1

$\zeta_1$

y2

$\zeta_2$

Check what researcher is doing with exogenous variables!
DF of all of these models = 4*5/2 – 8 = 2

Model Degrees of Freedom
DF = $t_{max}$ - t

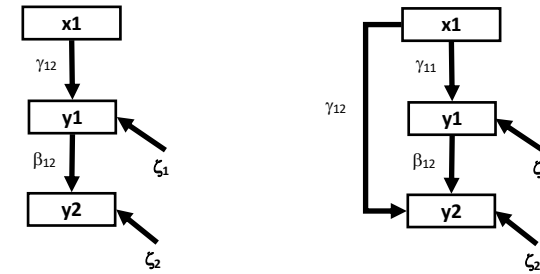Estimating 5 parameters from 6 variance/covariance relationships

**DF=1**
**Model Is _Overidentified_**

|     | x1  | y1  | y2  |
|-----|-----|-----|-----|
| x1  | 0.5 |     |     |
| y1  | 0.7 | 0.5 |     |
| y2  | 0.2 | 0.8 | 0.3 |

Cov(x,y1,y2)=

Identification in SEM
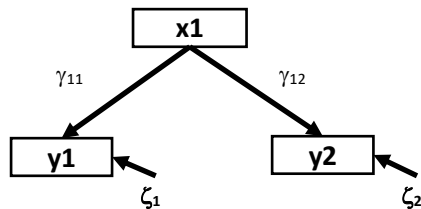_# of Parameters v. Covariance Matrix_

**Overidentified**   **Just Identified**

_Just Identified models have no DF to evaluate fit_
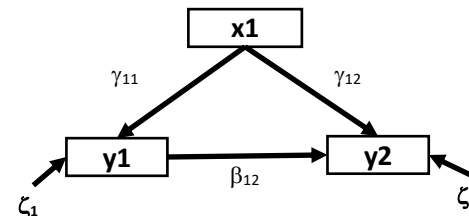
Identification in SEM
_Many Regressions_

**Yes**: There are no relationships between endogenous variables
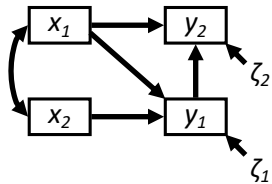**SUFFICIENT CONDITION**

Identification in SEM
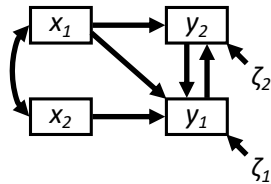_No Feedbacks_

**Yes**: Model is Recursive
**SUFFICIENT CONDITION**

6

Feedbacks and SEM

Recursive — Non-recursive

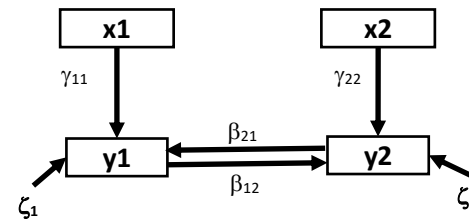**Recursive** = each item in a series is directly determined by the preceding item

**Non-recursive** = there is bidirectionality (feedbacks) implicit in the model
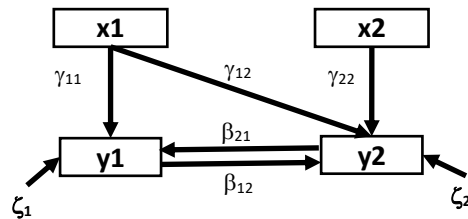


Identification in SEM
*Feedbacks with Different Causes*

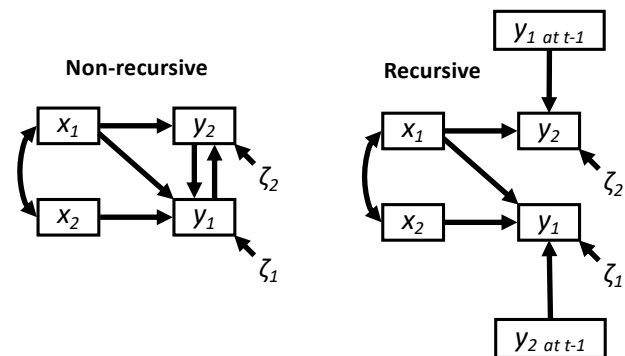**YES:** Model is Non-recursive, but y's have unique information
**NECESSARY CONDITION**
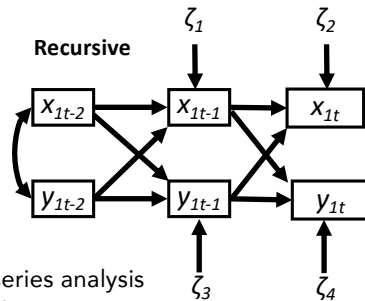


Identification in SEM
*Is this model identified?*

**NO!** Model is Non-recursive
AND not enough information for unique solution



Break Feedbacks with Time

Non-recursive — Recursive

## Cross-Lagged Panel Models to Solve Feedbacks!



- Time series analysis
- BACI designs
- Etc…

## Identification: Can I Fit My Model?

- **NECESSARY**: Fewer parameters than entries in covariance diagonal matrix (T-Rule)

- **SUFFICIENT**: And my model is recursive

- If you have feedbacks, then…
  - Break your model into time-lags (this is easy)
  - Or, ensure you have unique information for all variables (this can be hard!)

## A Likely Outline

1. What is SEM using likelihood and covariance matrices?

2. Model Identifiability

3. Sample Size for SEM

4. Standardized Coefficients
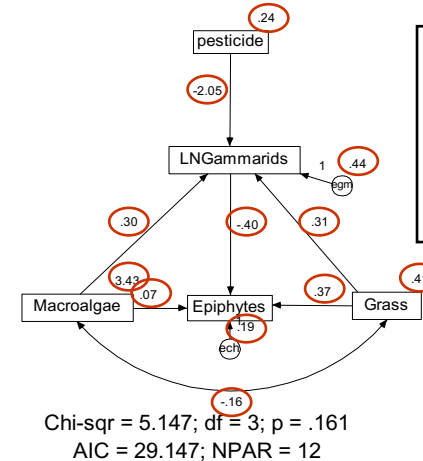
5. Introduction to `lavaan`

## Sample Size

1. The further you are in a model from an exogenous data-generating, the weaker it's influence.

2. Our ability to detect these tapering effect sizes is proportional to our information (especially sample size) and the number of parameters being estimated.

3. Sample size sets an upper limit for the complexity of the model we can obtain.

4. Sample Size influences our ability to detect lack of model fit
   - This might not be a benefit…

## So…What's my Sample Size?

1. Rules of thumb for sample size - at least 5 samples per estimated parameter
   - prefer 20 samples per parameter
   - Really, $p^{3/2}/n$ should approach 0 (Portnoy 1988)

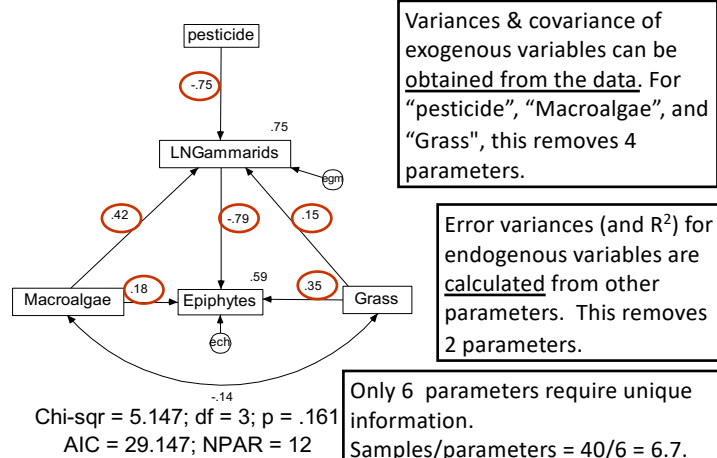2. Path coefficients add to our parameter list, not the variances

---

## Number of Estimated Parameters



There are a total of 12 parameters shown.

However, only 6 of these require unique information…

Chi-sqr = 5.147; df = 3; p = .161
AIC = 29.147; NPAR = 12

---

## Parameters Needing Unique Information



Variances & covariance of exogenous variables can be obtained from the data. For "pesticide", "Macroalgae", and "Grass", this removes 4 parameters.

Error variances (and $R^2$) for endogenous variables are calculated from other parameters. This removes 2 parameters.

Chi-sqr = 5.147; df = 3; p = .161
AIC = 29.147; NPAR = 12

Only 6 parameters require unique information.
Samples/parameters = 40/6 = 6.7.

---

## A Likely Outline

1. What is SEM using likelihood and covariance matrices?

2. Model Identifiability

3. Sample Size for SEM
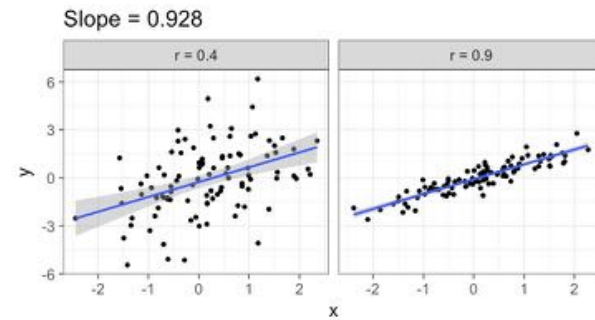
4. Standardized Coefficients

5. Introduction to `lavaan`

## Standardization

- *Unstandardized coefficient* = absolute strength of the pathway
  - " An 1 unit change in *X* results in some unit change in *Y* "

$$\beta_{xy\ std} = b_{xy} * sd_x/sd_y$$

- *Standardized coefficient* = relative strength of the pathway
  - " A 1 standard deviation change in *X* results in some standard deviation change in *Y* "
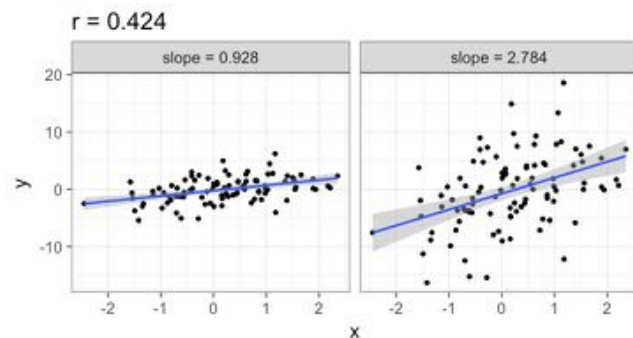  - Path Coefficient

## Same Slope, Different Correlation



**Regression**: y = a + bx

**Standardized Coefficient**s: r = b * sd(x)/sd(y)

## Different Slope, Same Correlation



**Regression**: y = a + bx

**Standardized Coefficient**s: r = b * sd(x)/sd(y)

## Which to Use?

| Unstandardized | Standardized |
|---|---|
| Good for prediction: coefficients are in raw units | Good for ranking: coefficients are in equivalent units |
| Has direct real world meaning | Less clear real world meaning |
| Can be compared across pathways or models that have identical units | Can be compared across all pathways in all models |

## A Likely Outline

1. What is SEM using likelihood and covariance matrices?

2. Model Identifiability

3. Sample Size for SEM

4. Standardized Coefficients

5. Introduction to `lavaan`

## What is `lavaan`?

- Stands for LAtent VAriable Analaysis

- Written by Yves Roseel in 2010

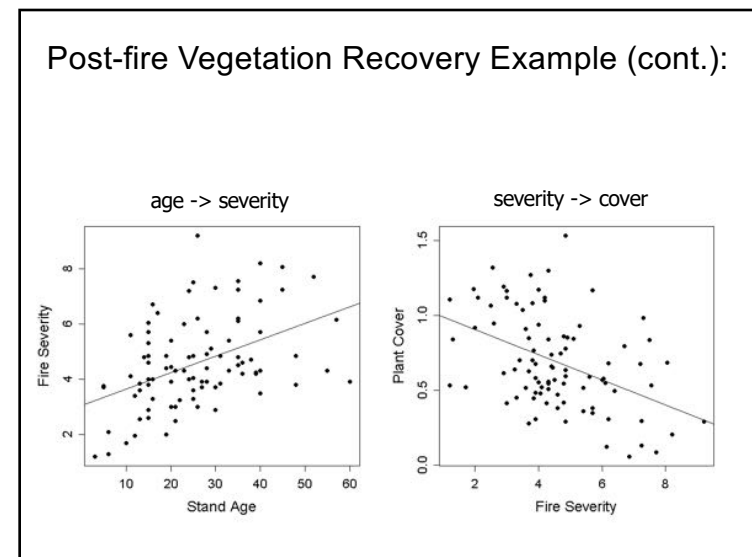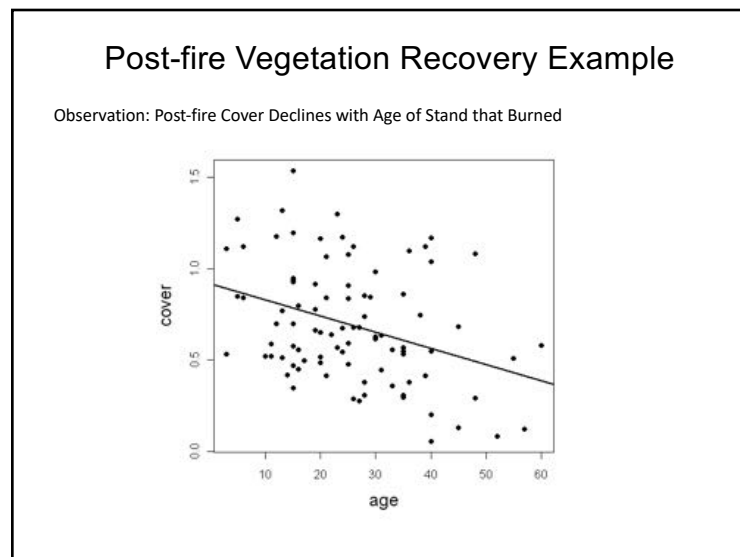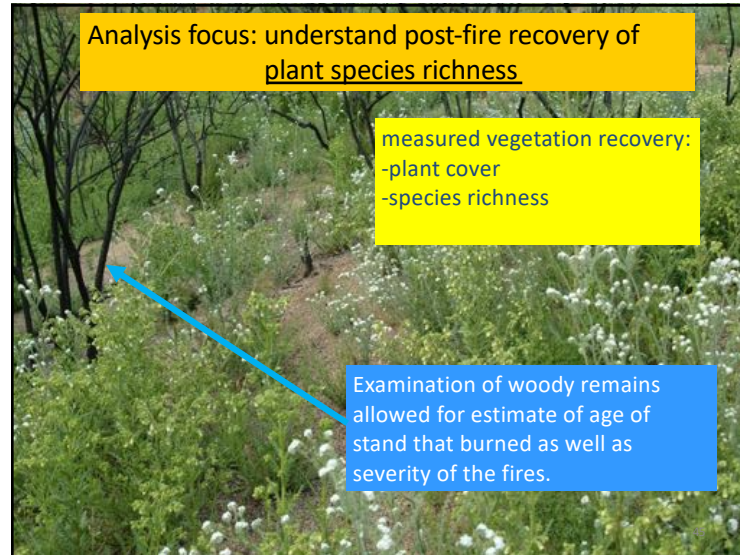- Currently in version 5, but 6 coming soon
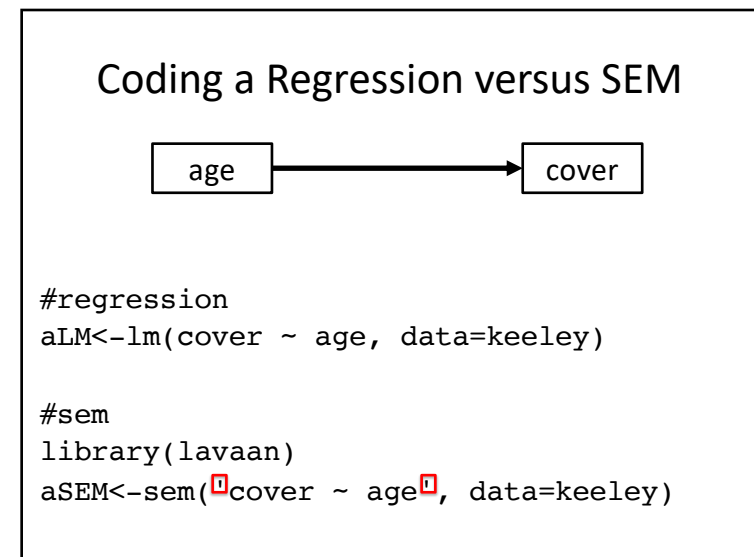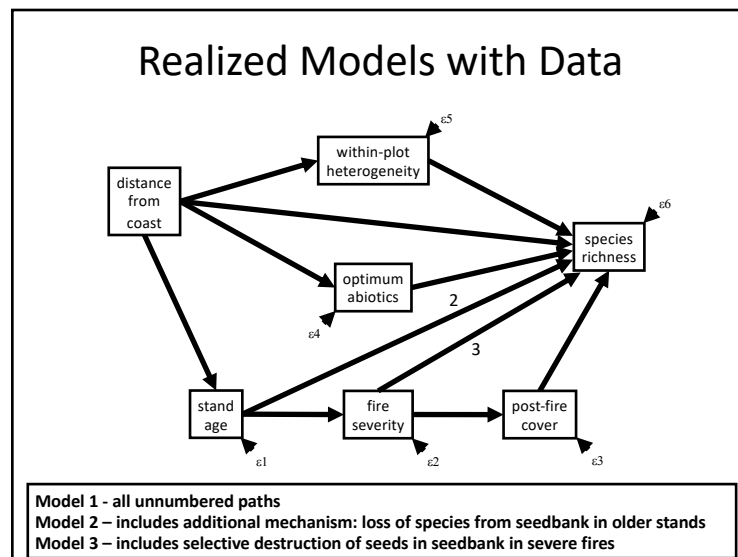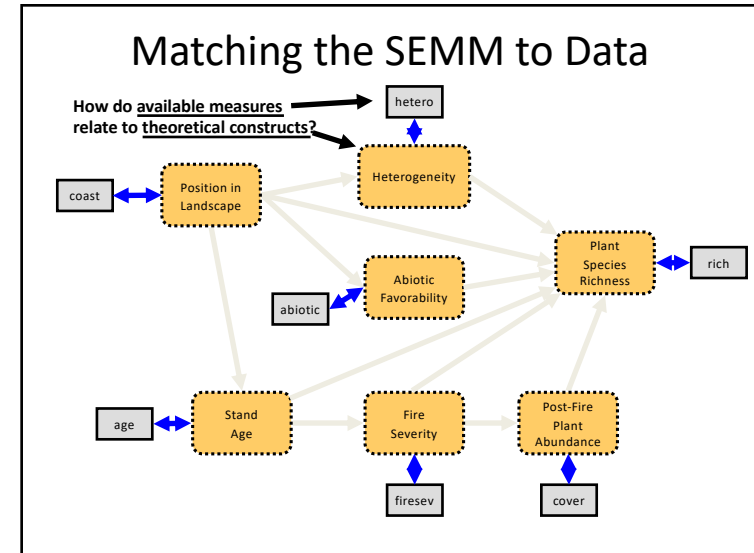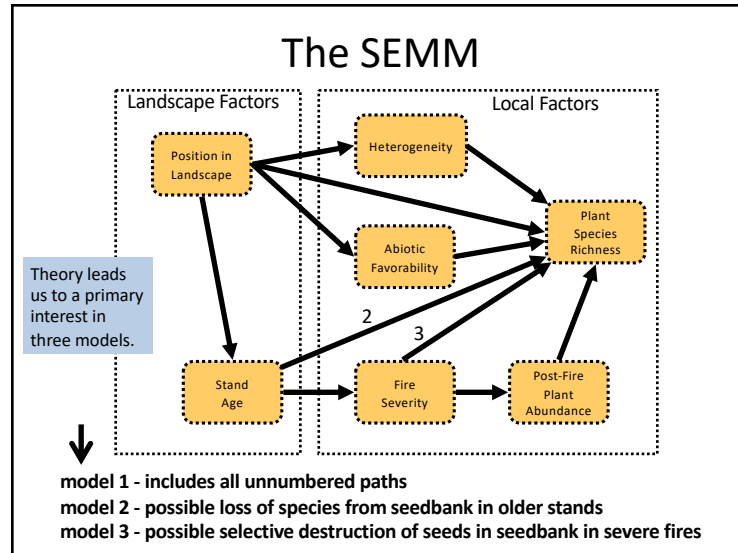
- Uses R `lm` syntax

## *A Reminder*

*1. SOFTWARE IS A TOOL*

*2. IT IS NOT PERFECT*

*3. ALWAYS MAKE SURE IT IS DOING WHAT YOU THINK IT IS DOING!*

**Mediation in Analysis of Post-Fire Recovery of Plant Communities in California Shrublands***



*Five year study of wildfires in Southern California in 1993. 90 plots (20 x 50m), (data from Jon Keeley et al.)

Analysis focus: understand post-fire recovery of plant species richness

measured vegetation recovery:
-plant cover
-species richness

Examination of woody remains allowed for estimate of age of stand that burned as well as severity of the fires.



Other factors measured included:
- local abiotic conditions (aspect, soils)
- spatial heterogeneity
- landscape-level conditions (location, elevation)

## Post-fire Vegetation Recovery Example

Observation: Post-fire Cover Declines with Age of Stand that Burned



## Post-fire Vegetation Recovery Example (cont.):

age -> severity          severity -> cover

The SEMM



Matching the SEMM to Data



Realized Models with Data



Coding a Regression versus SEM

## summary(aSEM)

**The model converged!**

```
lavaan (0.5-23.1097) converged normally after  10 iterations

  Number of observations                          90
                              Model is saturated
  Estimator                   so, χ2 test has no df  ML
  Minimum Function Test Statistic              0.000
  Degrees of freedom                               0

Parameter estimates:

  Information                              Expected
  Standard Errors                          Standard

                    Estimate  Std.err  Z-value  P(>|z|)
Regressions:
  cover ~
    age               -0.009    0.002   -3.549    0.000

Variances:
    .cover             0.087    0.013
```

## Compare to Regression

```
                Estimate  Std.err  Z-value  P(>|z|)
Regressions:
  cover ~
    age           -0.009    0.002   -3.549    0.000

Variances:
    .cover         0.087    0.013
```
**Compare to Residual SE sqrt(0.087)=0.295**

```
> summary(aLM)

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.917395   0.071726   12.79  < 2e-16 ***
age         -0.008846   0.002520   -3.51  0.00071 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2988 on 88 degrees of freedom
```

**But what about the intercept?**

## Intercepts Estimated with Mean Structure

```
> aMeanSEM<-sem('cover ~ age',
 data=keeley, meanstructure=T)
```

```
> summary(aMeanSEM)
...
                Estimate  Std.err  Z-value  P(>|z|)
Regressions:
  cover ~
    age           -0.009    0.002   -3.549    0.000

Intercepts:
    .cover         0.917    0.071   12.935    0.000

Variances:
    .cover         0.087    0.013
```

## Intercepts Estimated with Mean Structure

```
> aMeanSEM<-sem('cover ~ age',
 data=keeley, meanstructure=T)
```
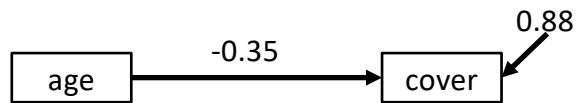
## Standardized Coefficients

```
>standardizedSolution(aSEM)
    lhs op   rhs est.std    se       z pvalue
1 cover  ~   age  -0.350 0.090 -3.912      0
2 cover ~~ cover   0.877 0.063 13.973      0
3   age ~~   age   1.000 0.000     NA     NA
```
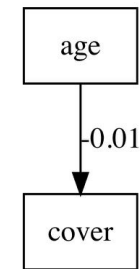


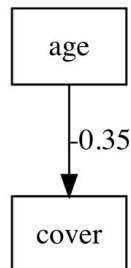**Also:** `summary(aSEM, standardized=T, rsq=T)`

## Can I See It?

```
library(lavaanPlot)
lavaanPlot(model = aSEM, coefs = TRUE)
```
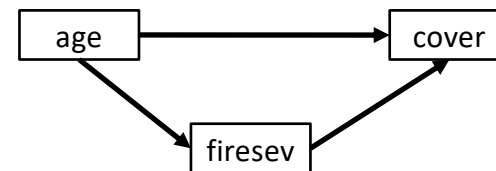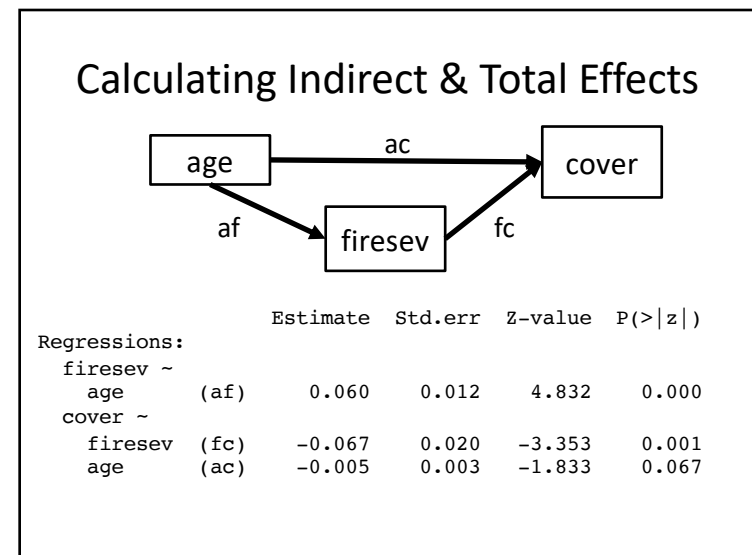


## Can I See It?

```
lavaanPlot(model = aSEM, coefs = TRUE,
           stand=TRUE)
```
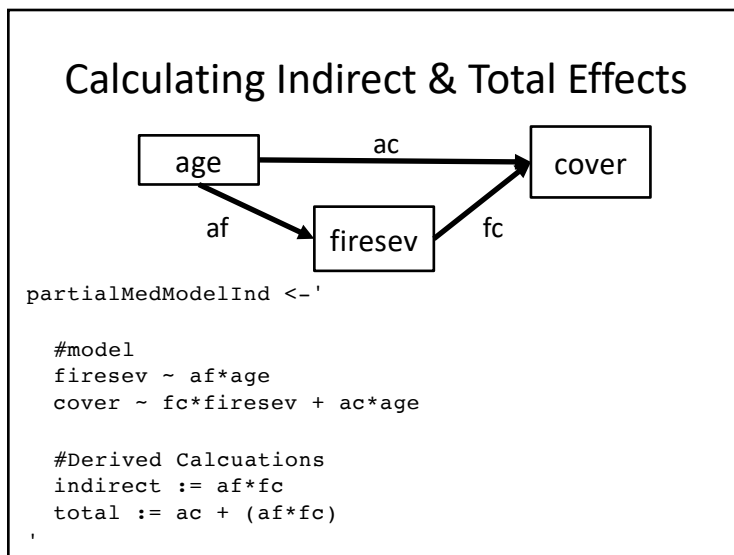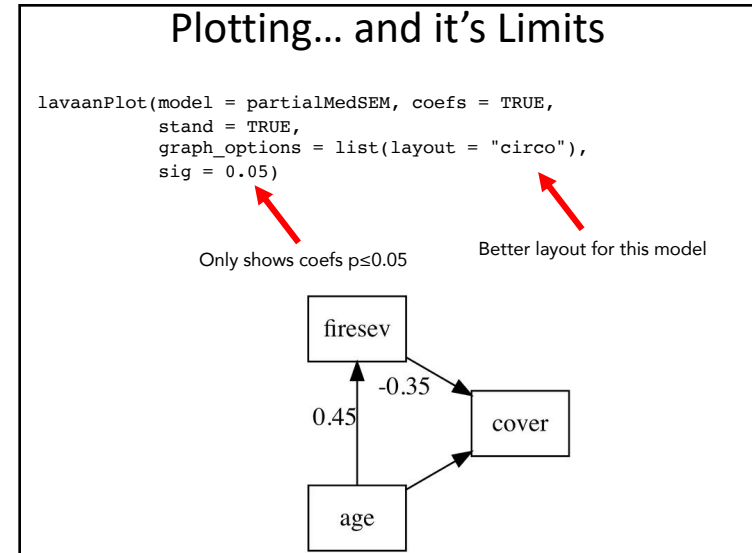


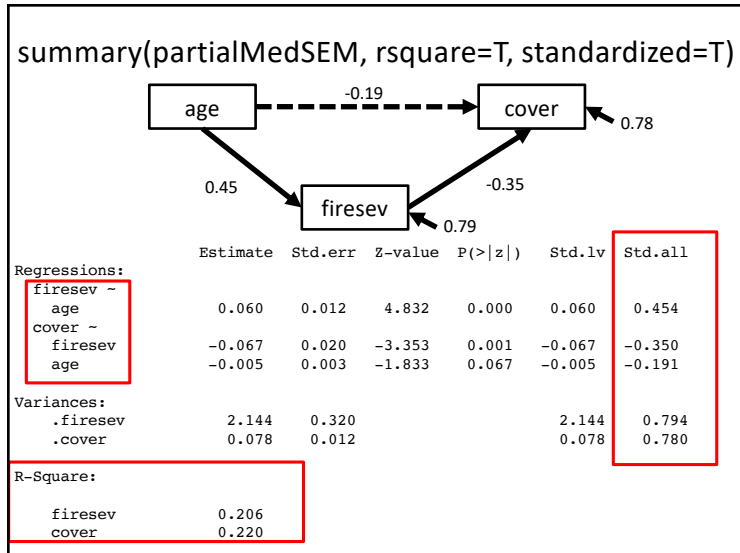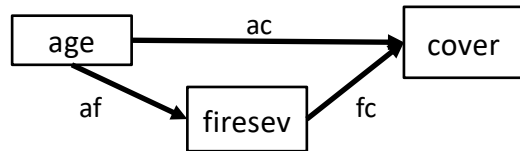## Indirect Effects and Fire



```
partialMedModel<-' firesev ~ age
              cover ~ firesev + age'

partialMedSEM<-sem(partialMedModel,
           data=keeley)
```

## summary(partialMedSEM, rsquare=T, standardized=T)



```
                    Estimate  Std.err  Z-value  P(>|z|)  Std.lv  Std.all
Regressions:
  firesev ~
    age              0.060    0.012    4.832    0.000    0.060    0.454
  cover ~
    firesev         -0.067    0.020   -3.353    0.001   -0.067   -0.350
    age             -0.005    0.003   -1.833    0.067   -0.005   -0.191

Variances:
  .firesev           2.144    0.320                      2.144    0.794
  .cover             0.078    0.012                      0.078    0.780

R-Square:

    firesev          0.206
    cover            0.220
```

## Plotting... and it's Limits

```
lavaanPlot(model = partialMedSEM, coefs = TRUE,
           stand = TRUE,
           graph_options = list(layout = "circo"),
           sig = 0.05)
```

Only shows coefs p≤0.05          Better layout for this model



## Calculating Indirect & Total Effects



```
partialMedModelInd <-'

  #model
  firesev ~ af*age
  cover ~ fc*firesev + ac*age

  #Derived Calcuations
  indirect := af*fc
  total := ac + (af*fc)
'
```

## Calculating Indirect & Total Effects



```
                  Estimate  Std.err  Z-value  P(>|z|)
Regressions:
  firesev ~
    age    (af)     0.060    0.012    4.832    0.000
  cover ~
    firesev (fc)   -0.067    0.020   -3.353    0.001
    age    (ac)    -0.005    0.003   -1.833    0.067
```

## Calculating Indirect & Total Effects



```
                Estimate   Std.err  Z-value  P(>|z|)

...

Defined parameters:
    indirect        -0.004    0.001   -2.755    0.006
    total           -0.009    0.002   -3.549    0.000
```
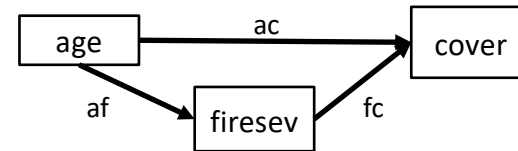
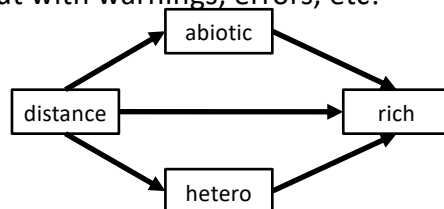## Calculating Indirect & Total Effects



```
> standardizedSolution(partialMedSEMInd)
        lhs op        rhs est.std     se       z pvalue
...
10 indirect :=      af*fc  -0.159 0.054 -2.947   0.003
11    total := ac+(af*fc)  -0.350 0.090 -3.912   0.000
```

## Take Lavaan for a Spin!

1. Fit this model!
2. Fill in Standardized Coefficients and $R^2$ for this model
3. Calculate summed direct and indirect effects of distance on richness
4. Call out with warnings, errors, etc!



## The dreaded variance warning!

Warning message:
In lav_data_full(data = data, group = group, cluster = cluster,  :
  lavaan WARNING: some observed variances are (at least) a factor 1000 times larger than others; use varTable(fit) to investigate
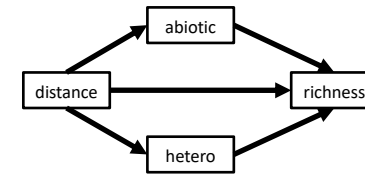
## Diagnosing Error Issues

```
> inspect(distFit, "obs")
$cov
          rich    hetero  abiotc  distnc
rich    225.646
hetero    0.784   0.013
abiotic  58.312   0.241  58.314
distance 77.089   0.347  30.824  77.094
```
***Is this OK?***

1. Does it indicate an outlier or data problem?

2. This is a likelihood algorithm problem – can be fine!

3. If you are worried, rescale by 10s, see if answers change
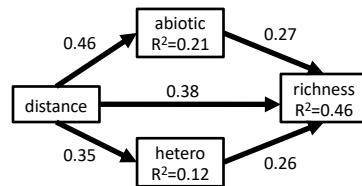
---

## Solution 1: The Model



```
#The Richness Partial Mediation Model
distModel <- 'rich ~ distance + abiotic + hetero
          hetero ~ distance
          abiotic ~ distance'

distFit <- sem(distModel, data=keeley)

standardizedSolution(distFit)
```

---

## Solution 2: Coefficients



```
      lhs op       rhs est.std     se       z pvalue
1    rich  ~ distance   0.377  0.092   4.117  0.000
2    rich  ~  abiotic   0.268  0.087   3.079  0.002
3    rich  ~   hetero   0.256  0.082   3.104  0.002
4  hetero  ~ distance   0.346  0.099   3.498  0.000
5  abiotic ~ distance   0.460  0.094   4.911  0.000
6    rich ~~     rich   0.539  0.080   6.708  0.000
7  hetero ~~   hetero   0.880  0.131   6.708  0.000
8 abiotic ~~  abiotic   0.789  0.118   6.708  0.000
9 distance ~~ distance  1.000    NA      NA     NA
```
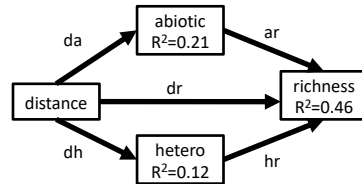
---

## Solution 3: Direct and Indirect



```
distModelEff <- '
rich ~ dr*distance + ar*abiotic + hr*hetero
hetero ~ dh*distance
abiotic ~ da*distance

#The effects
direct := dr
indirect := dh*hr + da*ar
total := direct + indirect
'
```

## Solution 3: Direct and Indirect
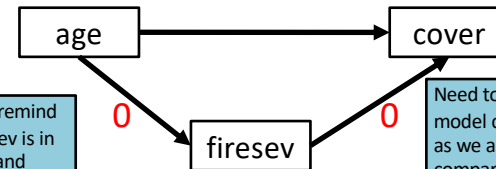


```
> standardizedSolution(distFitEff)
       lhs op          rhs est.std   se      z pvalue
...
10   direct :=          dr  0.377 0.086  4.390  0.000
11 indirect :=   dh*hr+da*ar  0.212 0.055  3.835  0.000
12   total := direct+indirect  0.589 0.062  9.433  0.000
```

**What would you say about direct and indirect effects in this system?**

## What if we know better?



Fill in 0's to remind us that firesev is in the model, and fixed to 0

Need to do this for model comparison, as we are comparing covariance matrices

```
zeroMedModel<-' firesev ~ 0*age
                cover ~ 0*firesev + age'

zeroMedFit<-sem(zeroMedModel,
                data=keeley)
```

## What lavaan sees…

```
> inspect(aSEM, "obs")
$cov
       cover   age
cover  0.100
age   -1.381 156.157
...

> inspect(zeroMedFit, "obs")
$cov
       firesv  cover   age
firesev  2.700
cover   -0.227   0.100
age      9.319  -1.381 156.157
...
```
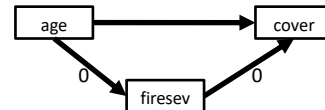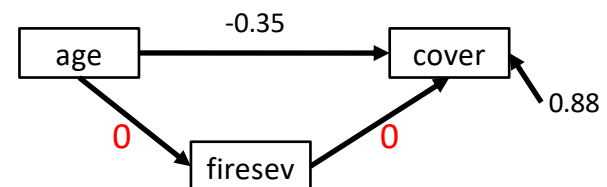


## standardizedSolution(zeroMedFit)



```
      lhs op     rhs est.std   se      z pvalue
1 firesev  ~     age  0.000   NA     NA     NA
2   cover  ~ firesev  0.000   NA     NA     NA
3   cover  ~     age -0.350 0.099 -3.549      0
4 firesev ~~ firesev  1.000 0.149  6.708      0
5   cover ~~   cover  0.877 0.131  6.708      0
6     age ~~     age  1.000   NA     NA     NA
```
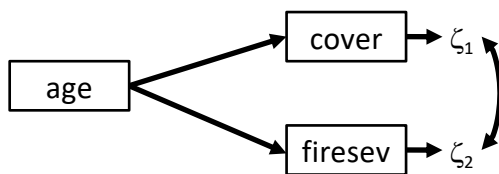
19

## Or… Just use intercepts!

age → cover

firesev

```
zeroMedModel2<-'
  firesev ~ 1
  cover ~ age
'
```

## Or… Just use intercepts!

age → cover

firesev

|   | lhs op | rhs | est.std | se | z | pvalue |
|---|--------|-----|---------|-----|-------|--------|
| 1 | firesev ~1 | | 2.778 | 0.232 | 11.956 | 0 |
| 2 | cover ~ | age | -0.350 | 0.090 | -3.912 | 0 |
| 3 | cover ~~ | cover | 0.877 | 0.063 | 13.973 | 0 |

## What about Correlated Error?

age → cover → $\zeta_1$
age → firesev → $\zeta_2$

```
#what about correlations
corModel <-'firesev ~ age
          cover ~ age
          cover ~~ firesev'

corFit <- sem(corModel, data=keeley)
```

## What about Correlated Error?

age → cover: -0.35, 0.87
age → firesev: 0.45, 0.79
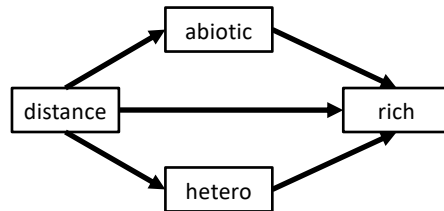cover ~~ firesev: -0.33
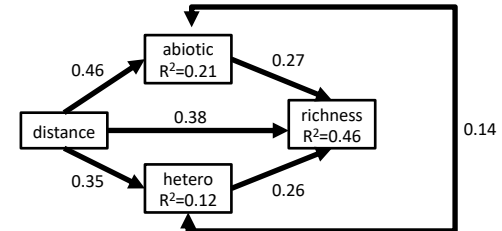
```
> standardizedSolution(corFit)
      lhs op    rhs est.std    se       z pvalue
1 firesev  ~    age   0.454 0.094   4.832      0
2   cover  ~    age  -0.350 0.099  -3.549      0
3 firesev ~~  cover  -0.333 0.094  -3.556      0
4 firesev ~~ firesev  0.794 0.118   6.708      0
5   cover ~~  cover   0.877 0.131   6.708      0
6     age ~~    age   1.000    NA      NA     NA
```

## Final Exercise

1. How does this model differ if the abiotic and hetero error correlate?

2. Fit assuming that there is a 1:1 (think 1 instead of 0) relationship between distance and richness
   – No error correlation please

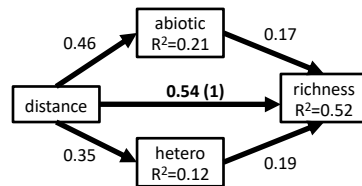

---

## Solution 1: Error Correlation



```
corErrorModel <- '
  rich ~ distance + abiotic + hetero
  hetero ~ distance
  abiotic ~ distance

  abiotic ~~ hetero        Coefficients unaffected
'
```

---

## Solution 2: The New Model



```
oneDistModel <- 'rich ~ 1*distance + abiotic + hetero
          hetero ~ distance
          abiotic ~ distance'

oneFit<-sem(oneDistModel, data=keeley)
summary(oneFit, stdandardized=T, rsquare=T)
```
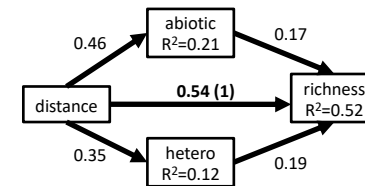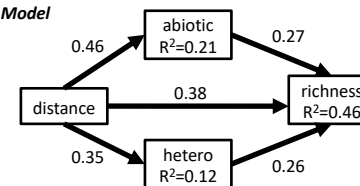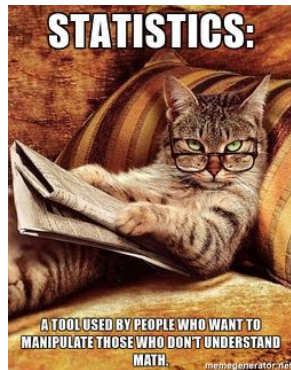
---

## Solution 2: The New Model



*Unconstrained Model*