

Extensions to Local Estimation

Overview

1. Assessing fit using Pseudo- R^2 s
2. GLMM Example
3. GAM Example

1.2. Pseudo- R^2 s

1.2. Pseudo- R^2 s. Omnibus test

- Fisher's C/χ^2 is the global fit statistic for local estimation but has many shortcomings:
 - Sensitive to the number of d-sep tests and the complexity of the model (easier to reject as the complexity increases)
 - Sensitive to the size of the dataset (e.g., high n leads to low P)
 - Fails symmetry when dealing with d-separated non-normal intermediate variables
 - Cannot be computed for saturated models

1.2. Pseudo- R^2 s. Local tests

- How do we infer the confidence in our SEM?
 - Examine standard errors of individual paths, qualitatively assess cumulative precision
 - Explore variance explained (i.e., R^2), qualitatively assess cumulative precision

1.2. Pseudo- R^2 s. General linear regression

- Coefficient of determination (R^2) = proportion of variance in response explained by fixed effects
- For OLS regression, simply $1 -$ the ratio of unexplained (error) variance (e.g., SS_{error}) over the total explained variance (e.g., SS_{total})
- Ranges $(0, 1)$, independent of sample size
- Not good for model comparisons since R^2 monotonically increases with model complexity (go to AIC which is penalized for complexity)

1.2. Pseudo- R^2 s. Generalized linear regression

- Likelihood estimation is not attempting to minimize variance but instead obtain parameters that maximize the likelihood of having observed the data
- In a likelihood framework, equivalent $R^2 = 1 - \frac{\text{log-likelihood of the null (intercept-only) model}}{\text{log-likelihood of the full model}}$
- Leads to identical R^2 as OLS for normal (Gaussian) distributions, not so for GLM – need to use likelihood-based pseudo- R^2 (e.g., McFadden, Nagelkerke)

1.2. Pseudo-R²s. Generalized mixed models

- Becomes even worse for mixed models because variance is partitioned among levels of the random factor, so what is the error variance?
- Need a new formulation of R²:
 - Marginal R² = variance explained by fixed effects only

$$R_{\text{GLMM}(m)}^2 = \frac{\sigma_f^2}{\sigma_f^2 + \sum_{l=1}^u \sigma_l^2 + \sigma_e^2 + \sigma_d^2}$$

The diagram illustrates the components of the Marginal R² formula for Generalized Linear Mixed Models (GLMMs). The formula is $R_{\text{GLMM}(m)}^2 = \frac{\sigma_f^2}{\sigma_f^2 + \sum_{l=1}^u \sigma_l^2 + \sigma_e^2 + \sigma_d^2}$. Four orange arrows point from descriptive labels to specific terms in the formula:

- An arrow points from "Fixed effects variance" to σ_f^2 in the numerator.
- An arrow points from "Fixed effects variance" to σ_f^2 in the denominator.
- An arrow points from "Random effects variance" to the summation term $\sum_{l=1}^u \sigma_l^2$ in the denominator.
- An arrow points from "Residual variance" to σ_e^2 in the denominator.
- An arrow points from "Distribution-specific variance" to σ_d^2 in the denominator.

1.2. Pseudo-R²s. Generalized mixed models

- Conditional R² = variance explained by both the fixed and random effects

The diagram illustrates the formula for the conditional R-squared value in a Generalized Linear Mixed Model (GLMM). The formula is presented as a fraction, with the numerator representing the variance explained by both fixed and random effects, and the denominator representing the total variance. Orange arrows point from descriptive labels to the corresponding mathematical terms in the formula.

$$R_{\text{GLMM}(c)}^2 = \frac{\sigma_f^2 + \sum_{l=1}^u \sigma_l^2}{\sigma_f^2 + \sum_{l=1}^u \sigma_l^2 + \sigma_e^2 + \sigma_d^2}$$

Labels and their corresponding terms in the formula:

- Fixed effects variance** (top left) points to σ_f^2 in the numerator.
- Random effects variance** (top right) points to $\sum_{l=1}^u \sigma_l^2$ in the numerator.
- Fixed effects variance** (bottom left) points to σ_f^2 in the denominator.
- Random effects variance** (bottom middle) points to $\sum_{l=1}^u \sigma_l^2$ in the denominator.
- Residual variance** (bottom center-right) points to σ_e^2 in the denominator.
- Distribution-specific variance** (bottom right) points to σ_d^2 in the denominator.

1.2. Pseudo- R^2 s. Generalized mixed models

- Comparison of marginal and conditional R^2 can lead to roundabout assessment of 'significance' of the random effects (e.g., if conditional R^2 is larger relative to marginal R^2)
- Best to report both and allow readers to determine how their magnitude affects the inferences

1.2. GLMM Example

1.2. SEM Example. Shipley 2009

- Hypothetical dataset: predicting latitude effect on survival of a tree species
- Repeated measures on 5 subjects at 20 sites from 1970-2006
- Survival (0/1) influenced by phenology (degree days until bud break, Julian days until bud break), size (stem diameter growth)



1.2. SEM Example. Shipley 2009

- Two distributions: normal, binary (survival)
- Random effects:
 - Site-only: latitude
 - Site and year: degree days, date
 - Site, year, and subject: diameter, survival



1.2. SEM Example. What is the basis set?



- $\text{Date} \perp \text{Lat} \mid (\text{Degree days})$
- $\text{Growth} \perp \text{Lat} \mid (\text{Date})$
- $\text{Survival} \perp \text{Lat} \mid (\text{Growth})$
- $\text{Growth} \perp \text{Degree days} \mid (\text{Date}, \text{Lat})$
- $\text{Survival} \perp \text{Degree days} \mid (\text{Growth}, \text{Lat})$
- $\text{Survival} \perp \text{Date} \mid (\text{Growth}, \text{Degree days})$

1.2. SEM Example. List of equations



```
library(piecewiseSEM)
library(nlme)
library(lme4)

# Load data
data(shipley); shipley <- na.omit(shipley)

# Create list of structural equations
shipley.sem <- psem(
  lme(DD ~ lat, random = ~1|site/tree, na.action = na.omit,
    data = shipley),
  lme(Date ~ DD, random = ~1|site/tree, na.action = na.omit,
    data = shipley),
  lme(Growth ~ Date, random = ~1|site/tree, na.action = na.omit,
    data = shipley),
  glmer(Live ~ Growth + (1|site) + (1|tree),
    family = binomial(link = "logit"), data = shipley)
)
```



1.2. SEM Example. D-sep tests



```
# Get summary  
summary(shipley.sem)
```

```
Structural Equation Model of shipley.sem
```

```
Call:  
  DD ~ lat  
  Date ~ DD  
  Growth ~ Date  
  Live ~ Growth
```

```
      AIC  
21745.782
```

```
---
```



1.2. SEM Example. D-sep tests



Tests of directed separation:

Independ.Claim	Test.Type	DF	Crit.Value	P.Value
Date ~ lat + ...	coef	18	-0.0798	0.9373
Growth ~ lat + ...	coef	18	-0.8929	0.3837
Live ~ lat + ...	coef	1431	1.0280	0.3039
Growth ~ DD + ...	coef	1329	-0.2967	0.7667
Live ~ DD + ...	coef	1431	1.0046	0.3151
Live ~ Date + ...	coef	1431	-1.5617	0.1184

--

Global goodness-of-fit:

Chi-Squared = NA with P-value = NA and on 6 degrees of freedom

Fisher's C = 11.536 with P-value = 0.484 and on 12 degrees of freedom

Warning message:

check model convergence: log-likelihood estimates lead to negative chi-squared!



1.2. SEM Example. D-sep tests



```
# Look at problematic model & variance components
Live.model <- glmer(Live ~ Growth + Date + DD + lat + (1|site) +
(1|tree), family = binomial(link = "logit"), data = shipley)
boundary (singular) fit: see ?issingular
```

```
VarCorr(Live.model)
```

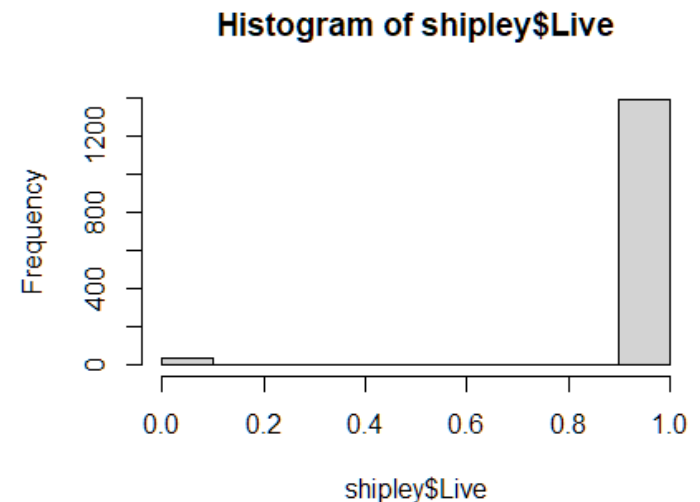
Groups	Name	Std.Dev.
tree	(Intercept)	0
site	(Intercept)	0



1.2. SEM Example. D-sep tests



- Re-specify random structure
- Still no positive χ^2 statistic ☹️
- Consider other distributions (e.g., negative binomial)
- Revert to d-sep test



1.2. SEM Example. D-sep tests



Coefficients:

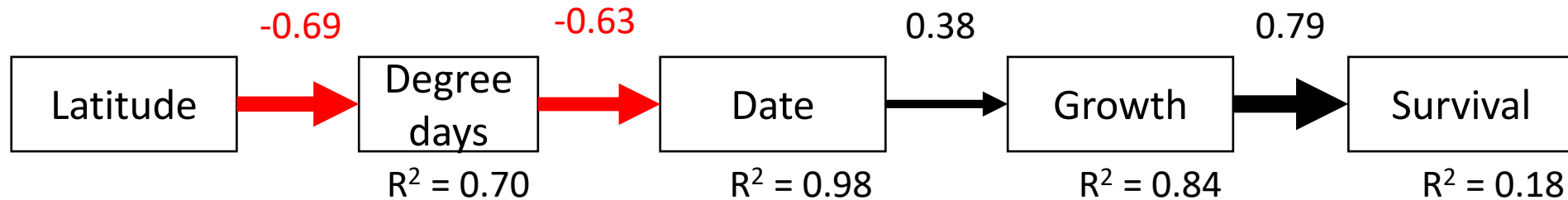
Response	Predictor	Estimate	Std.Error	DF	Crit.Value	P.Value	Std.Estimate
DD	lat	-0.8355	0.1194	18	-6.9960	0	-0.6877 ***
Date	DD	-0.4976	0.0049	1330	-100.8757	0	-0.6281 ***
Growth	Date	0.3007	0.0266	1330	11.2.917	0	0.3824 ***
Live	Growth	0.3479	0.0584	1431	5.9552	0	0.7866 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05

Individual R-squared:

Response	method	Marginal	Conditional
DD	none	0.49	0.70
Date	none	0.41	0.98
Growth	none	0.11	0.84
Live	delta	0.16	0.18

1.2. SEM Example. Populate final model



Coefficients:

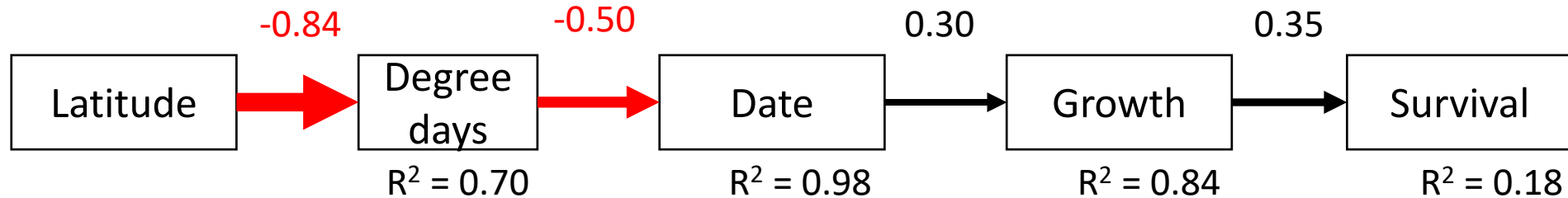
Response	Predictor	Estimate	Std.Error	DF	Crit.Value	P.Value	Std.Estimate
DD	lat	-0.8355	0.1194	18	-6.9960	0	-0.6877 ***
Date	DD	-0.4976	0.0049	1330	-100.8757	0	-0.6281 ***
Growth	Date	0.3007	0.0266	1330	11.2.917	0	0.3824 ***
Live	Growth	0.3479	0.0584	1431	5.9552	0	0.7866 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05

Individual R-squared:

Response	method	Marginal	Conditional
DD	none	0.49	0.70
Date	none	0.41	0.98
Growth	none	0.11	0.84
Live	delta	0.16	0.18

1.2. SEM Example. Refit using *lavaan*



...

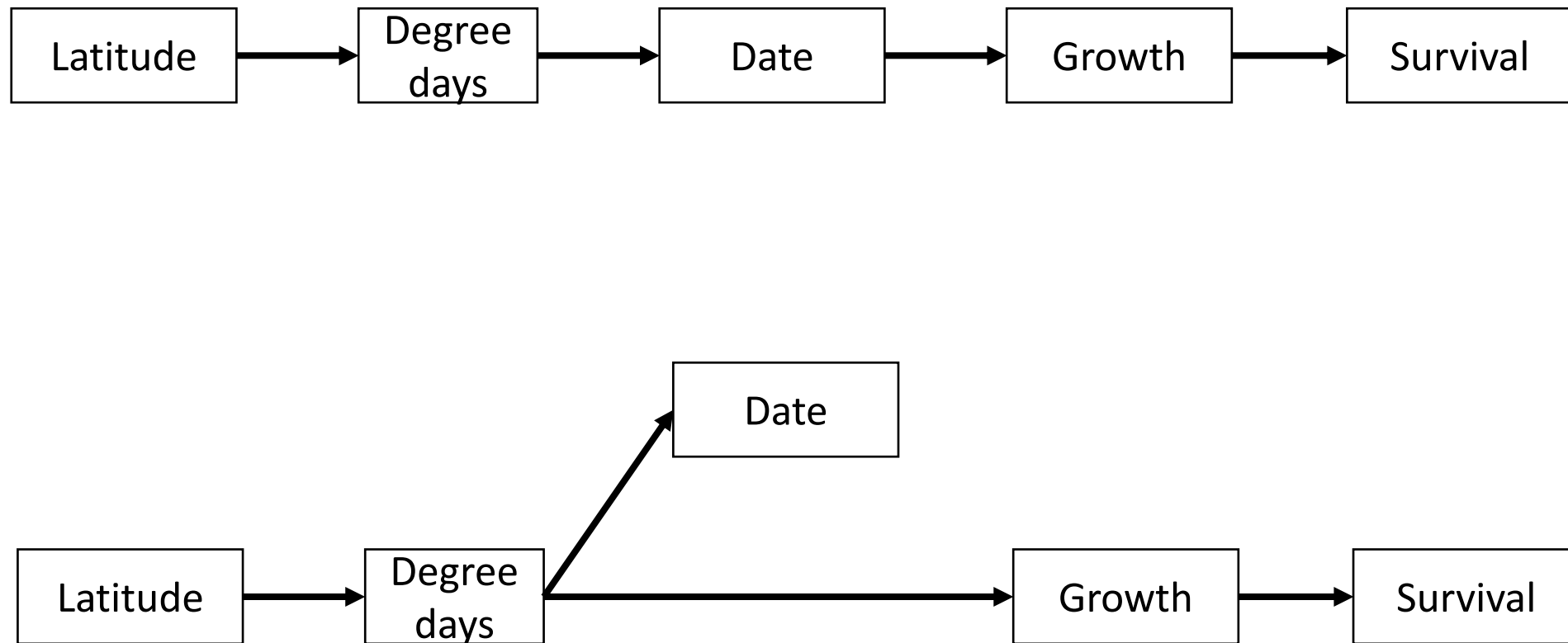
Estimator	ML
Model Fit Test Statistic	38.433
Degrees of freedom	6
P-value (Chi-square)	0.000

...

Regressions:

	Estimate	Std.Err	z-value	P(> z)
DD ~				
lat	-0.860	0.023	-37.923	0.000
Date ~				
DD	-0.517	0.016	-32.525	0.000
Growth ~				
Date	0.173	0.020	8.508	0.000
Live ~				
Growth	0.006	0.001	9.854	0.000

1.2. SEM Example. Compare these models

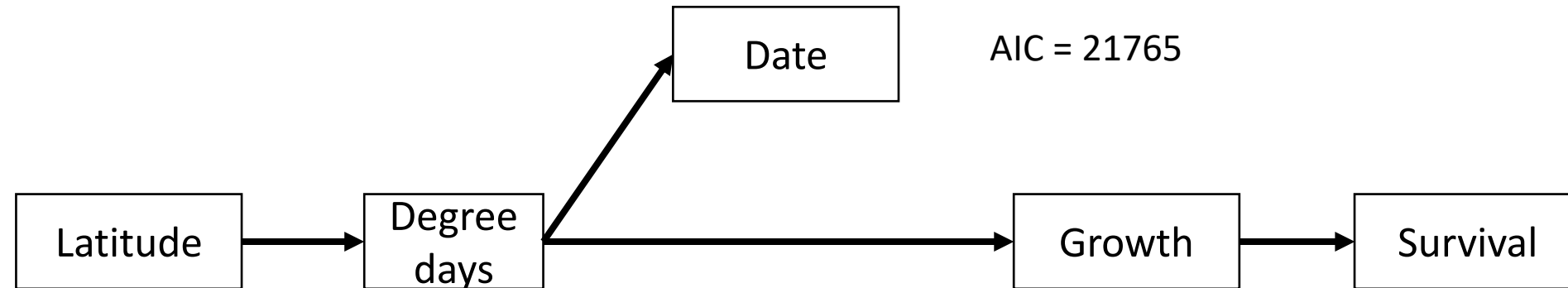


1.2. SEM Example. Compare these models

AIC = 21745



AIC = 21765

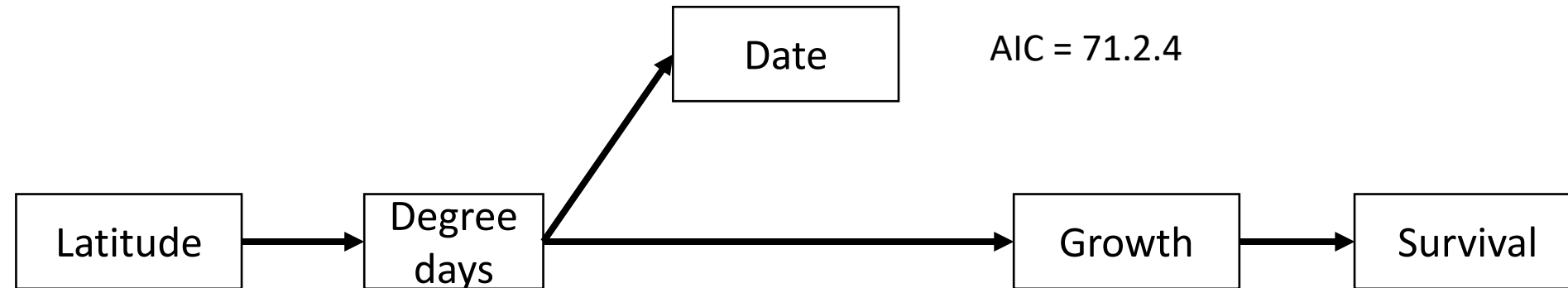


1.2. SEM Example. Compare these models (d-sep)

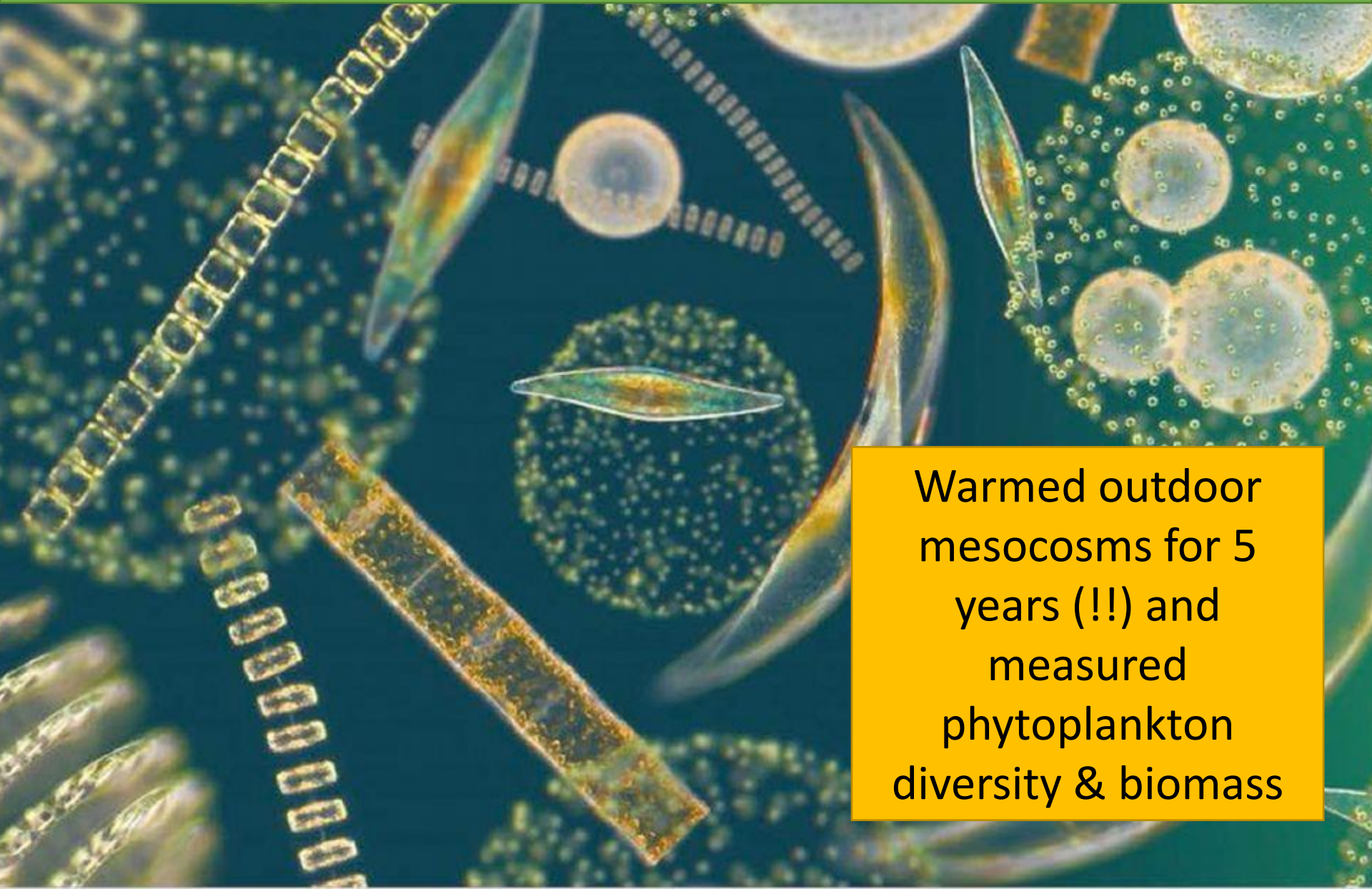
AIC = 49.54



AIC = 71.2.4

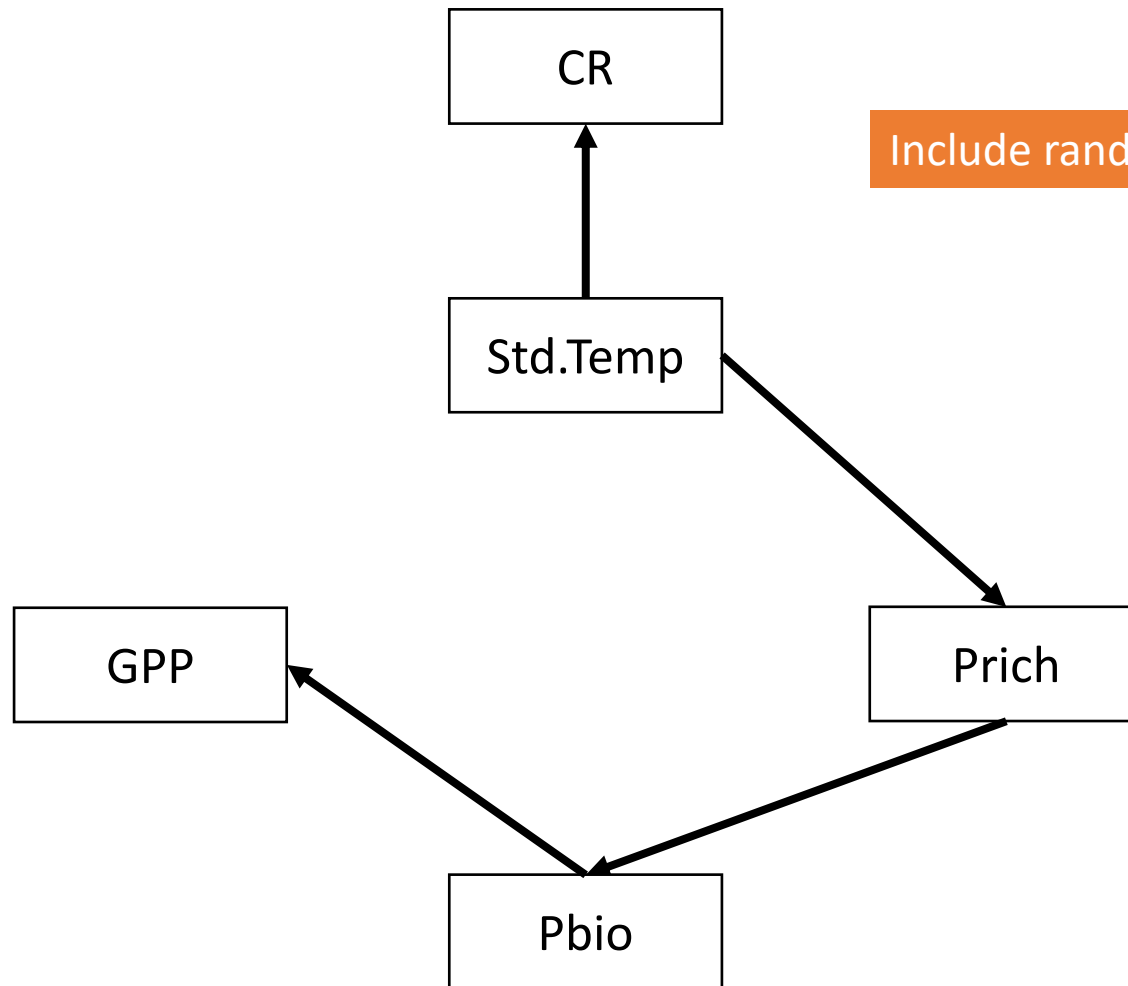


Yvon-Durocher et al (2015): Experimental warming on phytoplankton diversity and biomass



Warmed outdoor mesocosms for 5 years (!!) and measured phytoplankton diversity & biomass

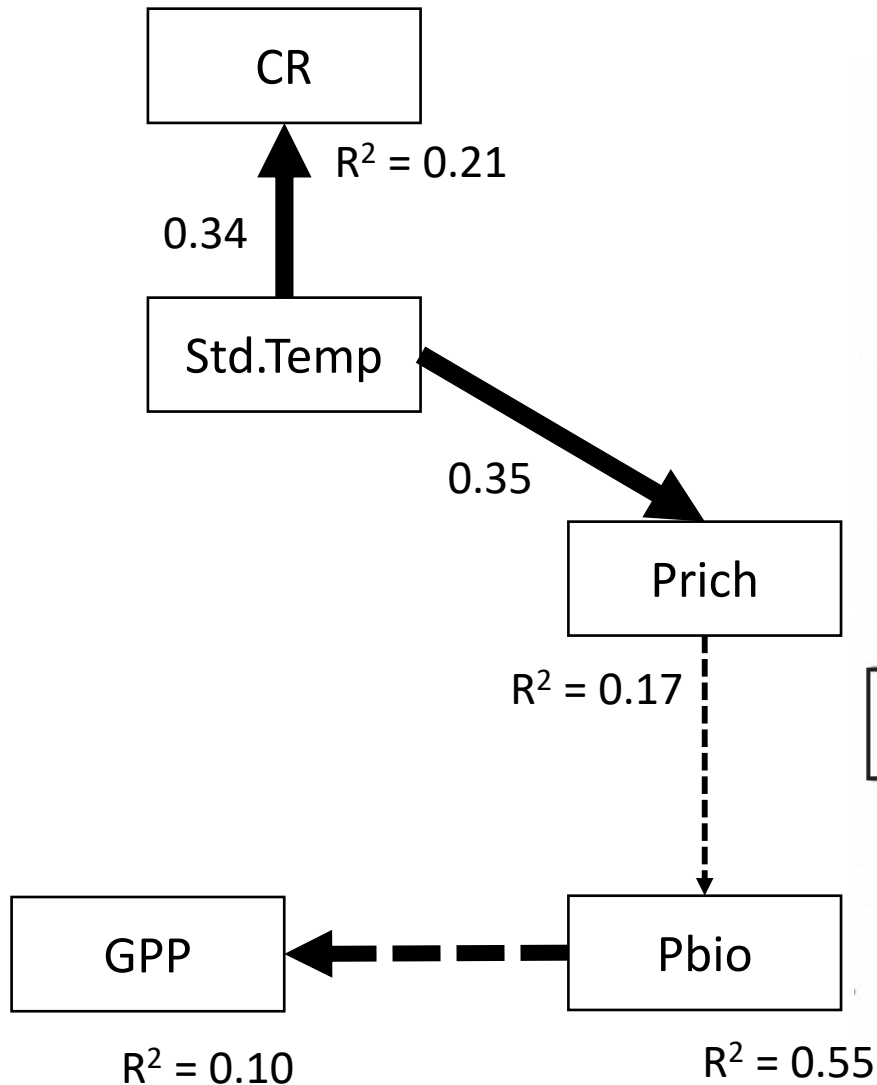
ACTIVITY. Fit Durocher dataset



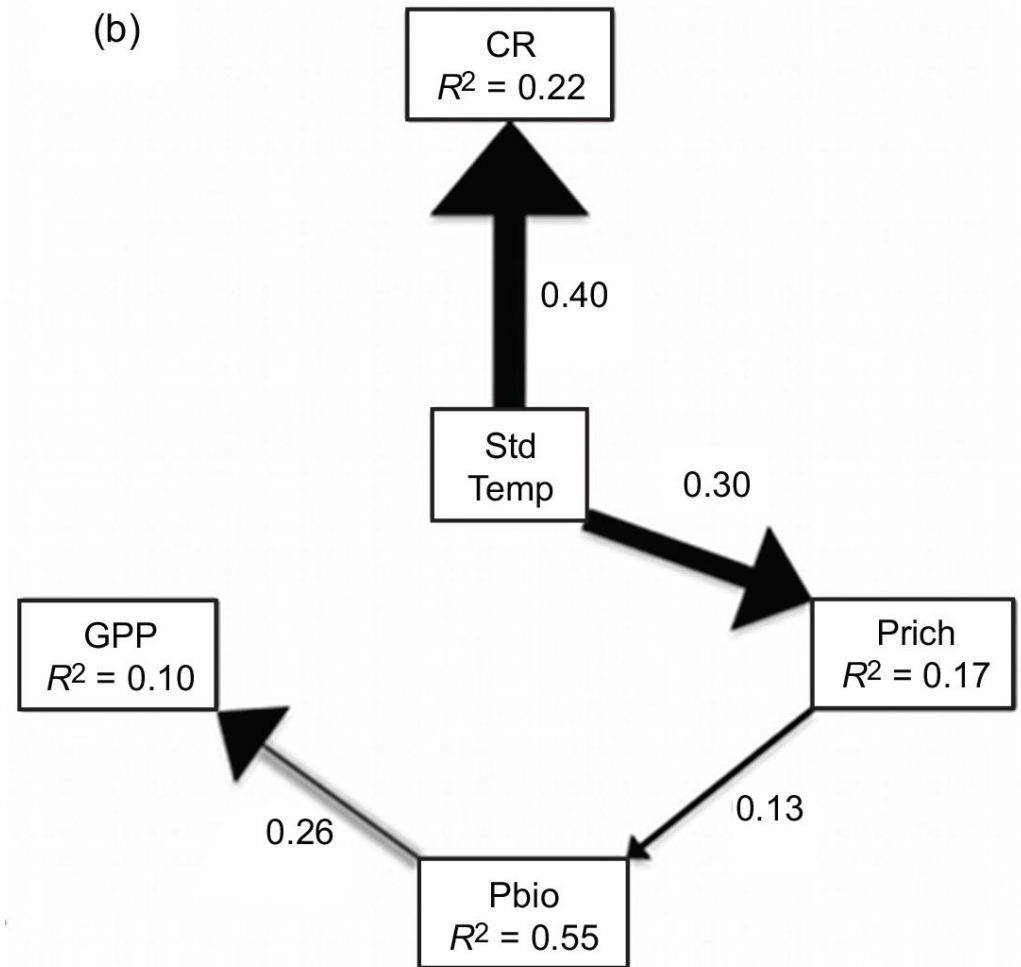
Include random effect of Pond.ID!

1.2. SEM Example. Your turn...

Model-wide $P = 0.063$ or $P < 0.001$



(b)



1.2. SEM Example. Your turn...

- Try removing incomplete cases first: `complete.cases`
 - What is their mistake here?
- Methods state: “with multiple measurements of variables made seasonally, nested within replicate mesocosms,” but then, “a path model as a set of hierarchical linear mixed effects models, each of which included hypothesized relationships between a response variable and a set of predictors as fixed effects and mesocosm ID as a random effect on the intercept.”
 - Play with the random structure?
- What about by treatment (Ambient vs. Heated)?
- Can anyone reproduce this result? Is it time to write a response?

1.3. GAM Example

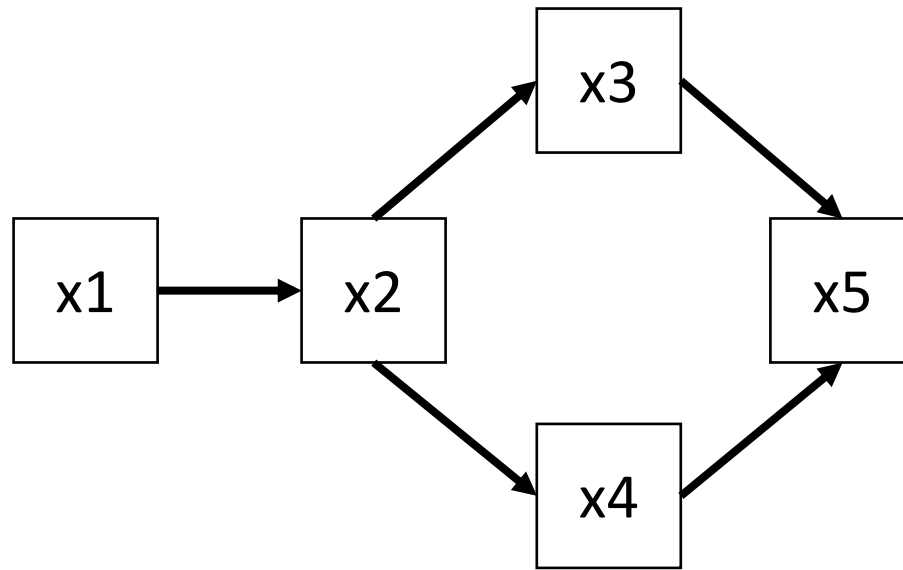
1.3. Generate example data

- Example data from appendix of Shipley and Douma using a mix of non-normal and non-linear variables

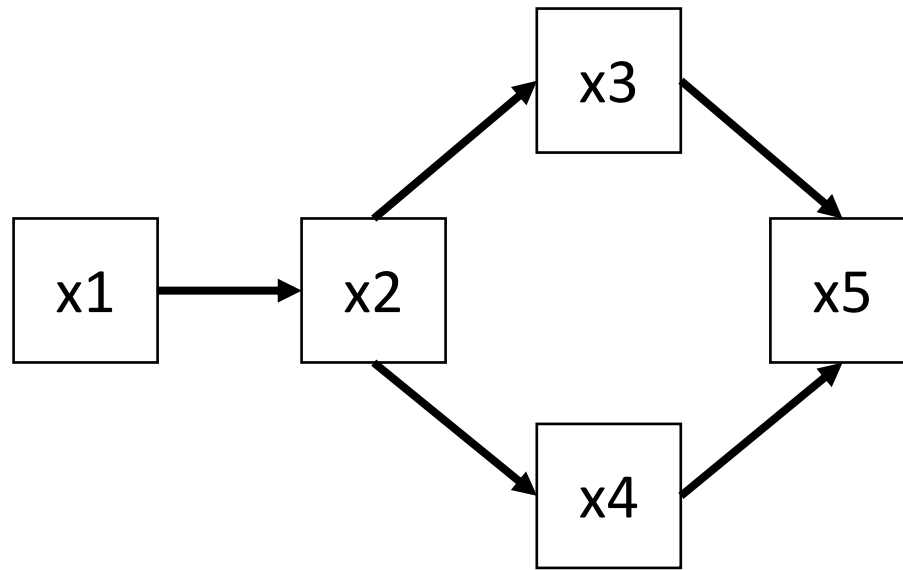
```
# Generate data from paper
set.seed(100)
n <- 100
x1 <- rchisq(n, 7)
mu2 <- 10*x1/(5 + x1)
x2 <- rnorm(n, mu2, 1)
x2[x2 <= 0] <- 0.1
x3 <- rpois(n, lambda = (0.5*x2))
x4 <- rpois(n, lambda = (0.5*x2))
p.x5 <- exp(-0.5*x3 + 0.5*x4)/(1 + exp(-0.5*x3 + 0.5*x4))
x5 <- rbinom(n, size = 1, prob = p.x5)
dat2 <- data.frame(x1 = x1, x2 = x2, x3 = x3, x4 = x4, x5 = x5)
```



1.3. Fit this SEM using `lm` and get GoF

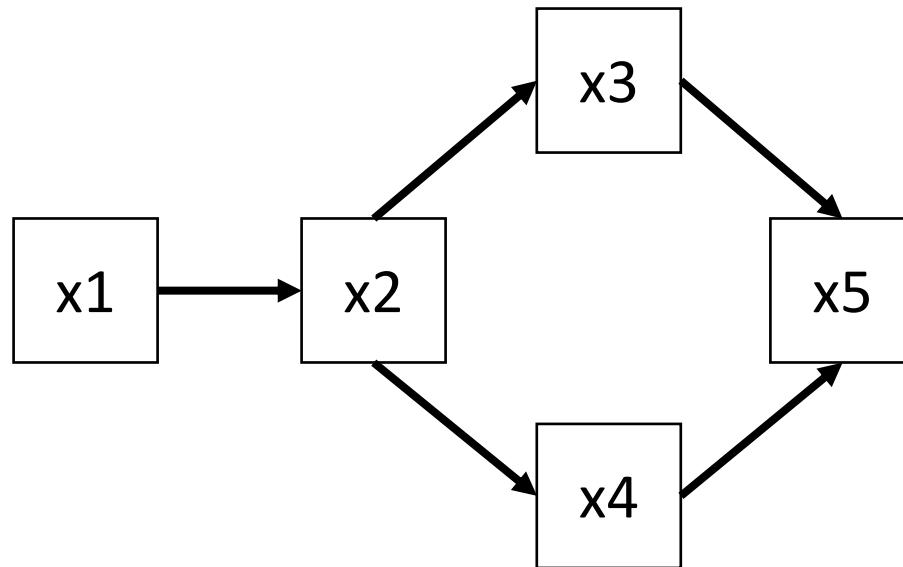


1.3. Fit this SEM using `lm` and get GoF



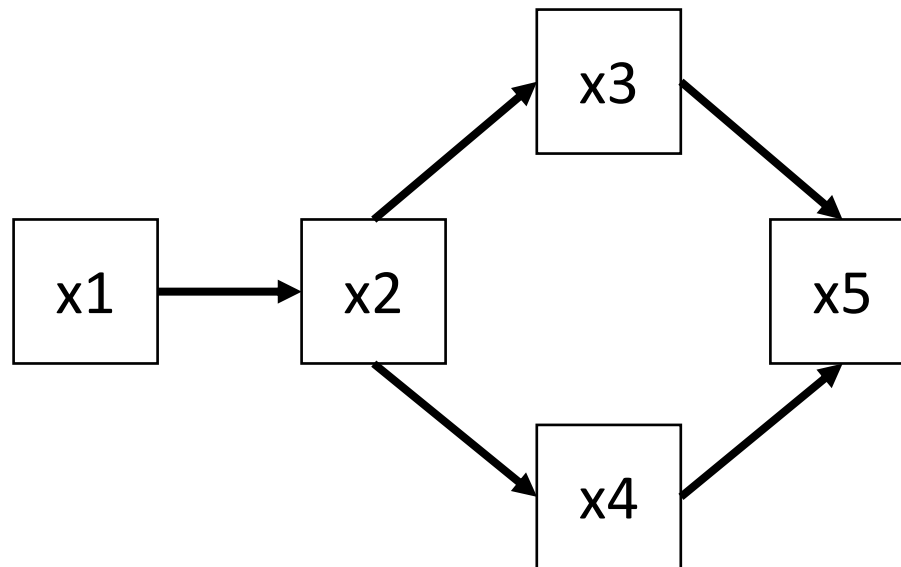
```
LLchisq(shipley_psem2)  
  Chisq df P.value  
1 4.143  5  0.529
```

1.3. Fit using GAM and GLM



```
shipley_psem3 <- psem(  
  gam(x2 ~ s(x1), data = dat2, family = gaussian),  
  glm(x3 ~ x2, data = dat2, family = poisson),  
  gam(x4 ~ x2, data = dat2, family = poisson),  
  glm(x5 ~ x3 + x4, data = dat2, family = binomial)  
)
```

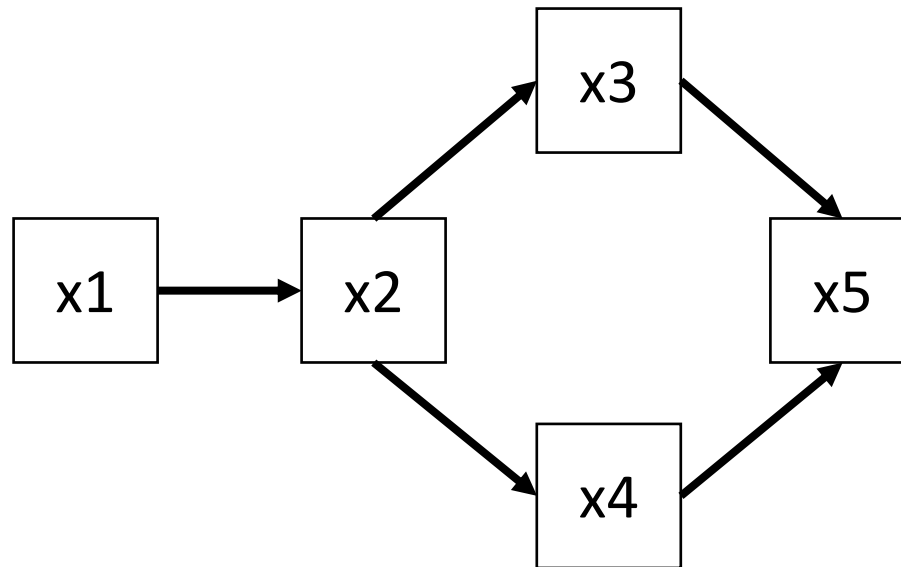
1.3. Fit using GAM and GLM



```
# Get goodness-of-fit  
LLchisq(shipley_psem2)
```

	Chisq	df	P.value
1	4.143	5	0.529

1.3. Fit using GAM and GLM



```
# Compare linear and non-linear models  
AIC(shipley_psem2, shipley_psem3)
```

	AIC	K	n
1	1240.20	13.000	100
2	1190.75	11.563	100

1.3. Truly Non-Linear Implementations

- Possible to compare models with the same typology but different ML fitting functions and forms (or nested models)
- Do not get coefficients returned by `coefs` because smoothed terms are non-linear functions
- How to present this path diagram???

1.3. Truly Non-Linear Implementations

- Piecewise SEM can be extended to many different model types: as long as you can get a P -value or compute a log-likelihood, you can estimate fit
 - Matrix regression (Barnes et al. 2016)
 - Spatially-explicit models