# *Session S:*

# Lecture – SEM Trees

## ATI Staff

Arizona State University

Summer 2016

# Motivating Example

## Development and Validation of Empirically Derived Frequency Criteria for NSSI Disorder Using Exploratory Data Mining

Brooke A. Ammerman
Temple University
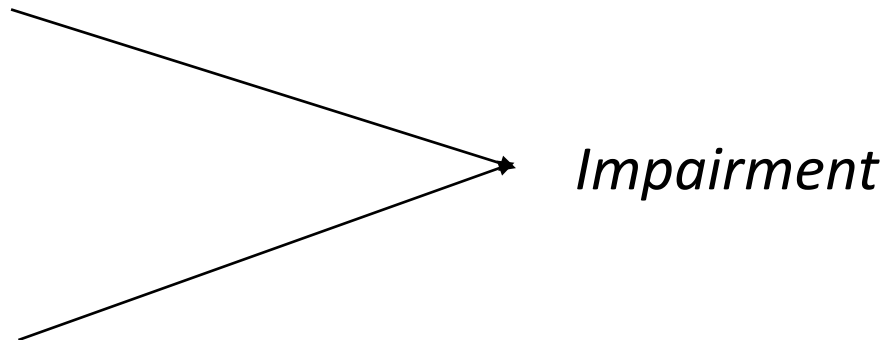
Ross Jacobucci
University of Southern California

Evan M. Kleiman
Harvard University

Jennifer J. Muehlenkamp
University of Wisconsin, Eau Claire

Michael S. McCloskey
Temple University

# Goal of Analysis

- To find clinically meaningful sub groups of non-suicidal self-injury (NSSI)
  - Cutoffs for DSM-V Criteria
  - Current cutoff = 5 NSSI acts in past year
    - Workgroup for disorder stated this is an "arbitrary cutoff"
- Have measure of how many times participants self-injured in last year
  - "Have you ever, intentionally or on purpose, hurt yourself in the following ways, without the intention of killing yourself?"
- Outcome
  - One-factor Factor Analysis consisting of:
    - Suicidality
    - Emotion Dysregulation
    - Emotion reactivity
    - Borderline personality
    - Disordered eating
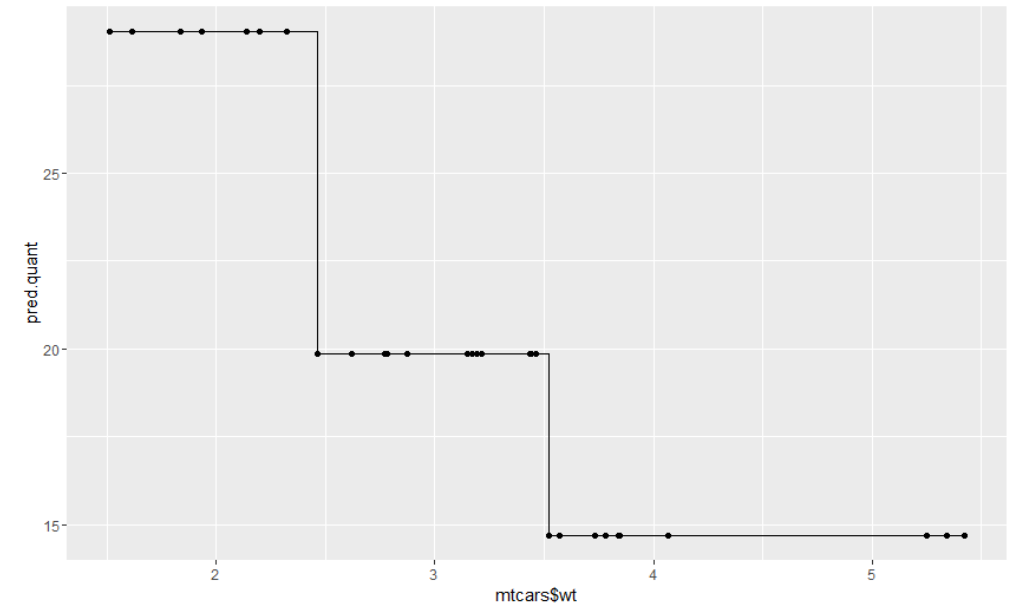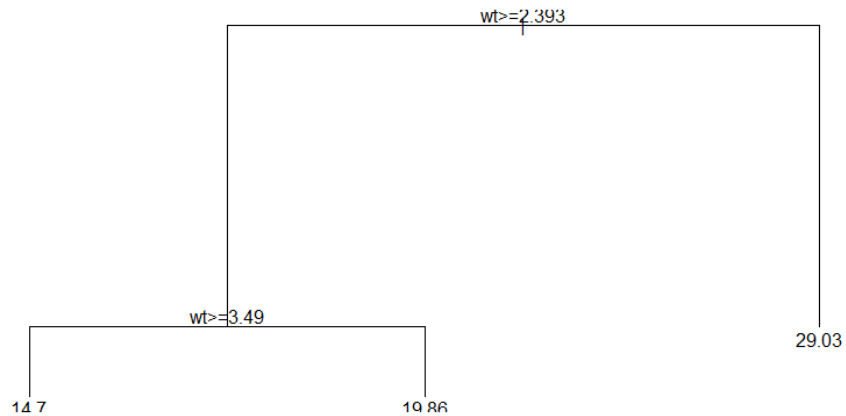    - Anxiety
    - Depression

*Impairment*

# Decision Trees for Groups

# Decision Trees for Groups

- As mentioned in previous presentations, Decision Trees (DTs) are non-linear models for predicting categorical or continuous outcomes

```
library(rpart);library(ggplot2); quant.out <- rpart(mpg ~ wt, data=mtcars); plot(quant.out);text(quant.out)

pred.quant <- predict(quant.out); qplot(mtcars$wt,pred.quant,geom=c("step","point"))
```
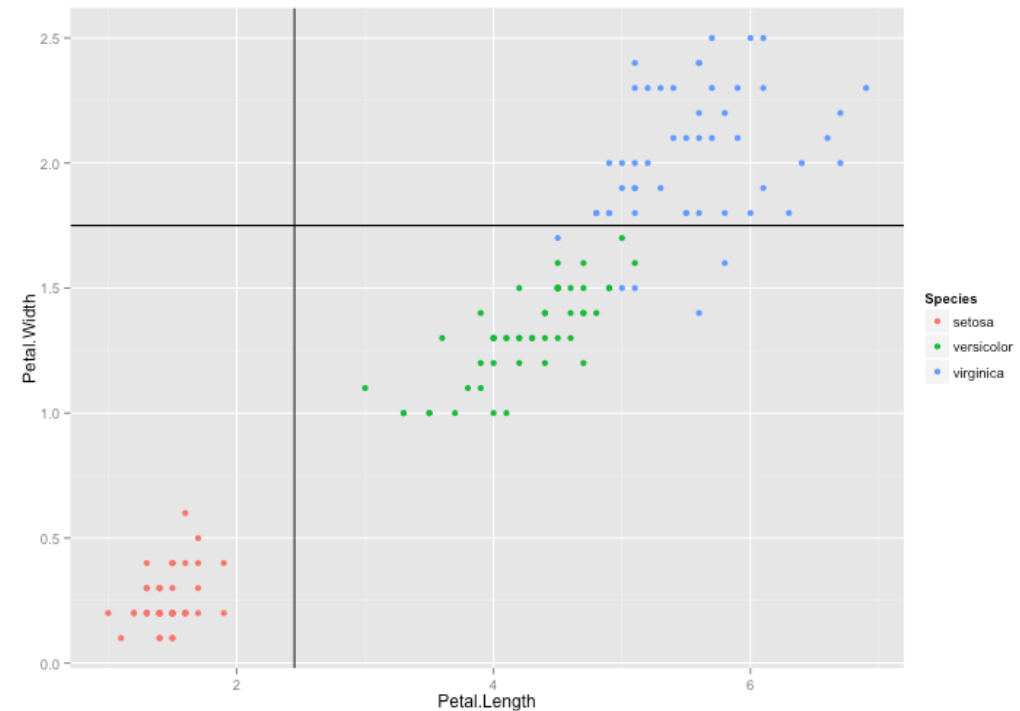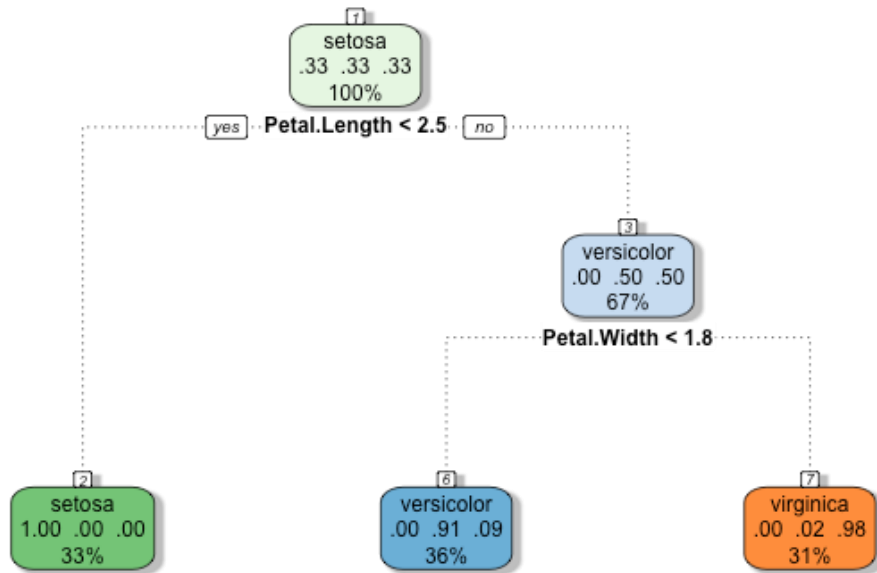
# DTs for Groups Cont'd

- Most applications in Data Mining don't involve people
  - The stepwise prediction for DT applies to groups as well
- Variable(s) you want to create groups with:
  - Predictor (one or more; categorical or continuous)
    - Categorical: Testing whether the levels of the variable are significantly different
    - Continuous: Looking for groups of people based on values of variable(s)
- Variable(s) you care about relationship with groups
  - Outcomes (one or more; categorical or continuous)

# DT for Groups: Traditional Way

- Not people, but flowers!

- **Petals = Predictor**

```
library(rpart);library(ggplot2);library(rattle);attach(iris)
fit2 <- rpart(Species ~ Petal.Length + Petal.Width, data =
iris,control=rpart.control(maxdepth=2))
fancyRpartPlot(fit2)
p <- qplot(Petal.Length,Petal.Width,colour=Species); q <- p +
geom_vline(xintercept=2.45)
q + geom_abline(intercept=1.75,slope=0)
```
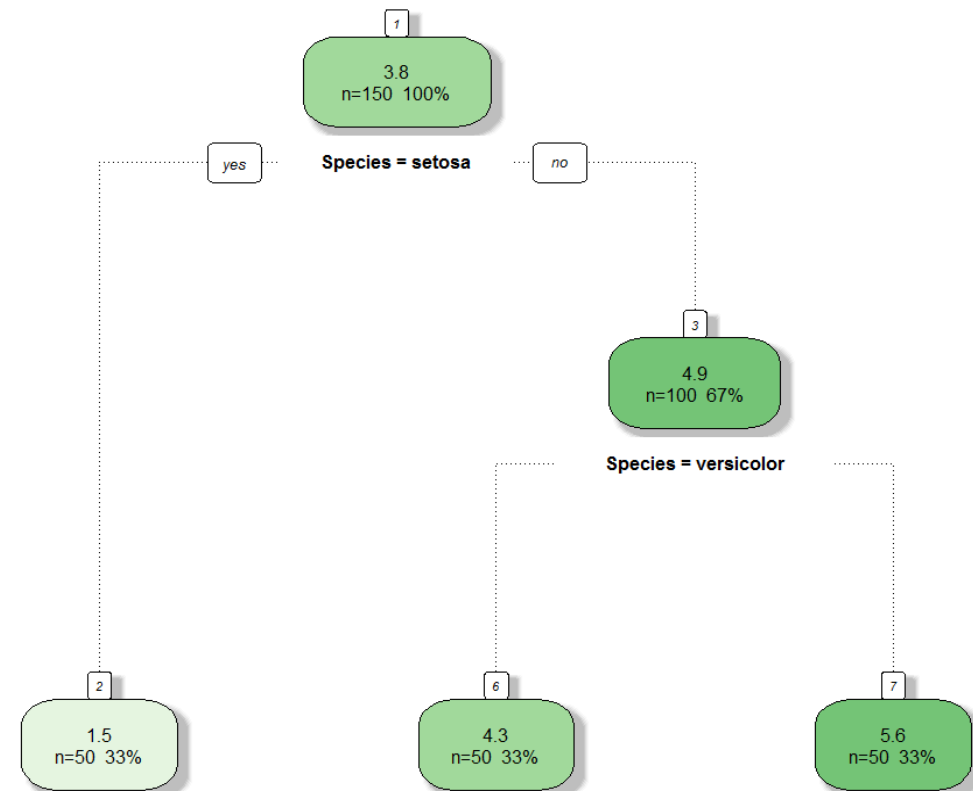


Rattle 2014-Dec-04 11:03:19 RJacobucci

Now, instead of predicting group members, we are searching for groups *within* the predictors.

# DT for Groups: Switch Predictors & Outcomes

- **Flower = Predictor**

- Test whether the groups are

  significantly different

```
library(rpart);library(ggplot2);library(rattle);attach(iris)
fit2 <- rpart(Petal.Length ~ Species, data =
iris,control=rpart.control(maxdepth=2))
```
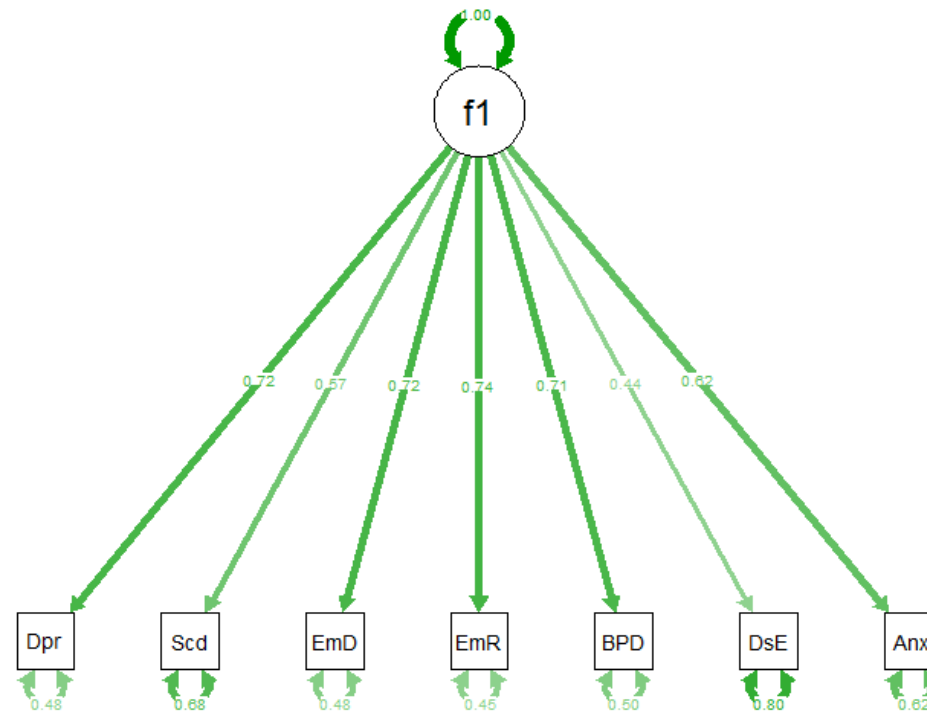
# DT for Groups Cont'd

- Many times we have variables that we might want to create subgroups based on
  - Income → find cutoffs for low, middle, high SES groups
  - Have continuous test scores and find cutoff for admission criteria
    - i.e. only accept those > 160 on GRE
- Which group has biggest differences on outcome?
  - Gender or Ethnicity a more "important" grouping in relation to supporting Clinton or Trump?
- Want to find interactions between groups to create further subgroupings
  - Females in low SES prefer Trump while females in high SES prefer Hillary
- Often times limited by only using 1 variable as an outcome
  - Where multivariate DT methods come in
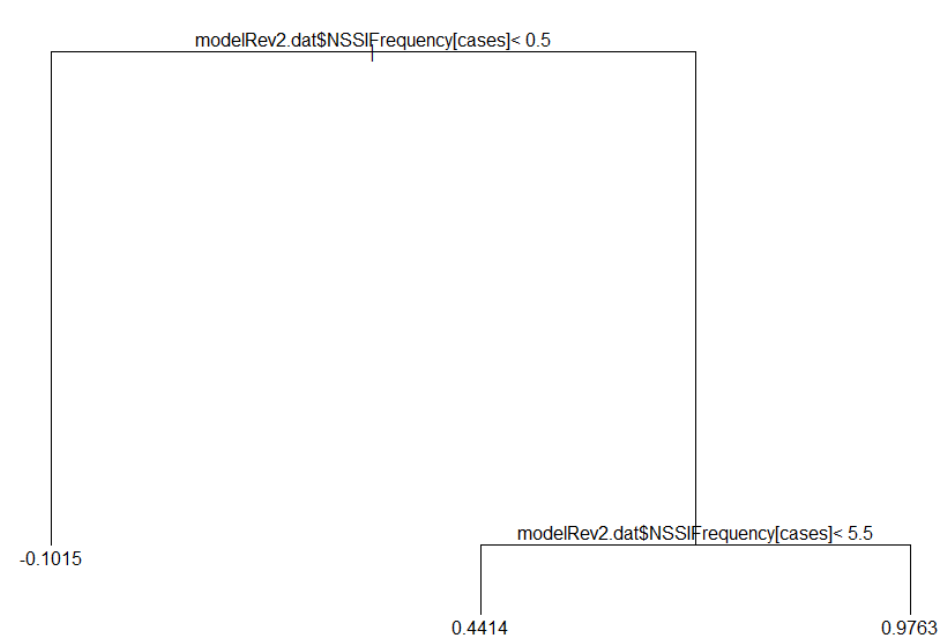    - E.g. multivariate boosting or SEM Trees

# DTs and NSSI

- We want to find subgroups of people based on the number of times they self-injures

- In relation to:
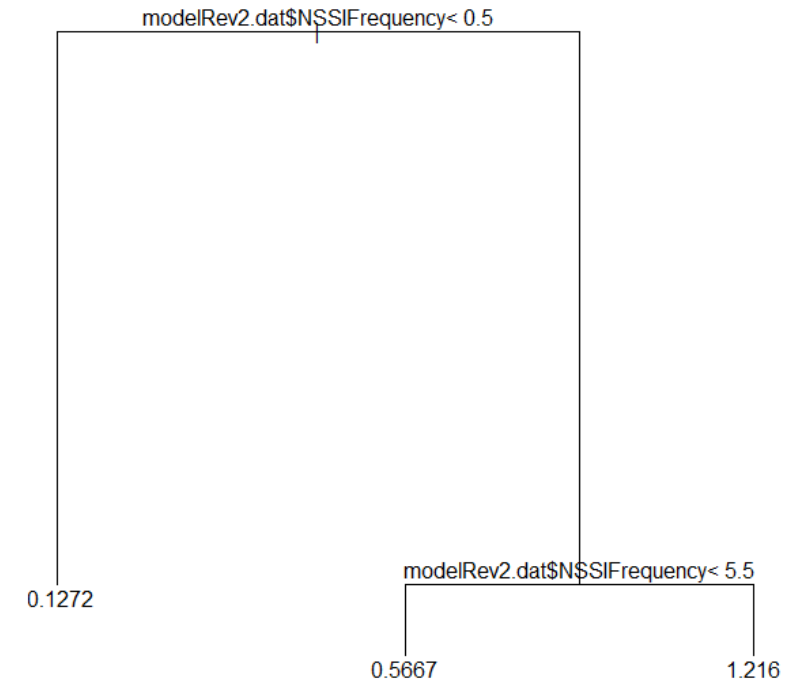  - A summed score
  - A factor score

# DTs and NSSI cont'd

NSSI ➡ Summed Score

NSSI ➡ Factor Score

modelRev2.dat$NSSIFrequency[cases]< 0.5

-0.1015

modelRev2.dat$NSSIFrequency[cases]< 5.5

0.4414

0.9763

modelRev2.dat$NSSIFrequency< 0.5

0.1272

modelRev2.dat$NSSIFrequency< 5.5

0.5667

1.216

# So What Did We Learn?

- Cutoffs for NSSI Frequency Criteria should be:
  - Between 0 & 1
  - Between 5 & 6
- The DT results give us three groups with those cutoffs
- But what does this actually mean?
  - People w/ NSSI value of 0 have a deviation between their actual and predicted factor score that is significantly smaller than between their actual scores and predictions for any other group, including assuming every is in one group
  - More technically, putting people into groups decreases the RMSEA

# SEM Trees & NSSI

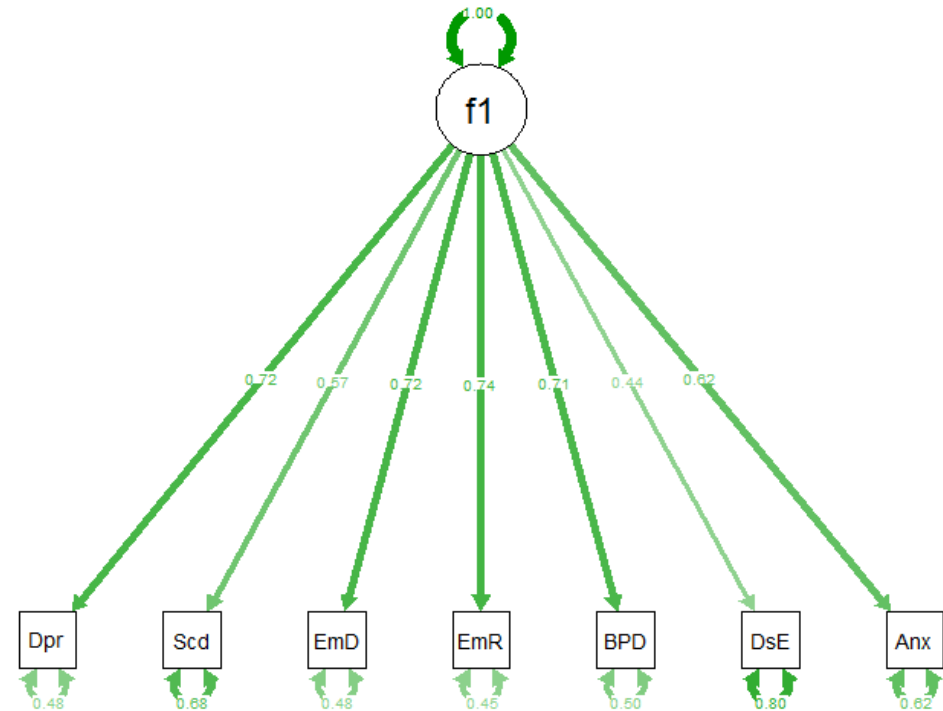# When We Can't Reduce the Outcome to a Single Value

- Need SEM Trees or other multivariate model (Multivariate Boosting or others)
- Instead of reducing our outcome to a single variable we can use:
  - Means of multiple variables
  - A confirmatory factor model
  - Latent growth model
  - Autoregressive model
  - Etc.

# SEM Trees

**Predictor**

**Outcome**

# SEM Tree Cont'd

- Different than just including a covariate in the model

- In SEM Trees, the covariates predict the *model fit*
  - Not just the latent variable

- In predicting the fit of the model, you are indirectly predicting differences in each of the model parameters
  - i.e. factor variance, mean, loadings ….
  - In other words: If you change the model parameters then you change the model fit
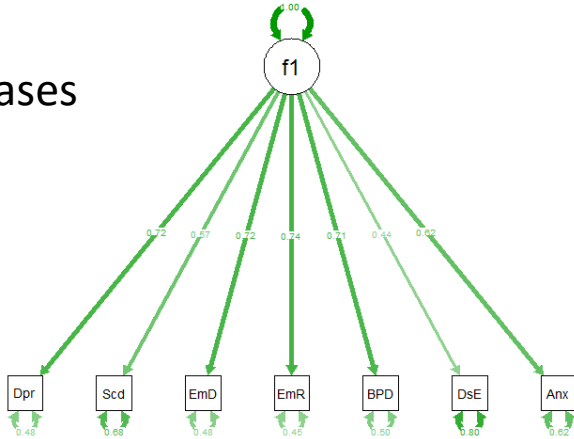
# SEM Trees Algorithm

- In our NSSI example, the NSSI variable has integer values ranging from 0 to 1000

- If a covariate is an integer in SEM Trees, the model tests groups at every value:
  - 1. Fit of multiple group model with groups of people with values of 0 and 1-1000
    - Get model fit
  - 1. Fit of multiple group model with groups of people with values of 0 or 1 and 2-1000
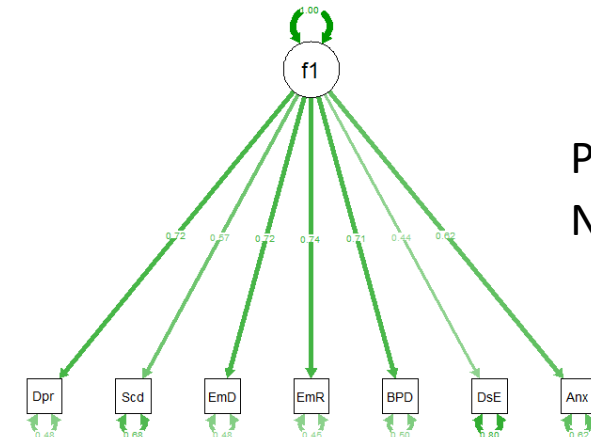    - Get model fit

# SEM Trees Algorithm Cont'd

- At each tested split, the model becomes

People w/ NSSI = 0



All cases



**VERSUS**

**+**

People w/ NSSI > 0
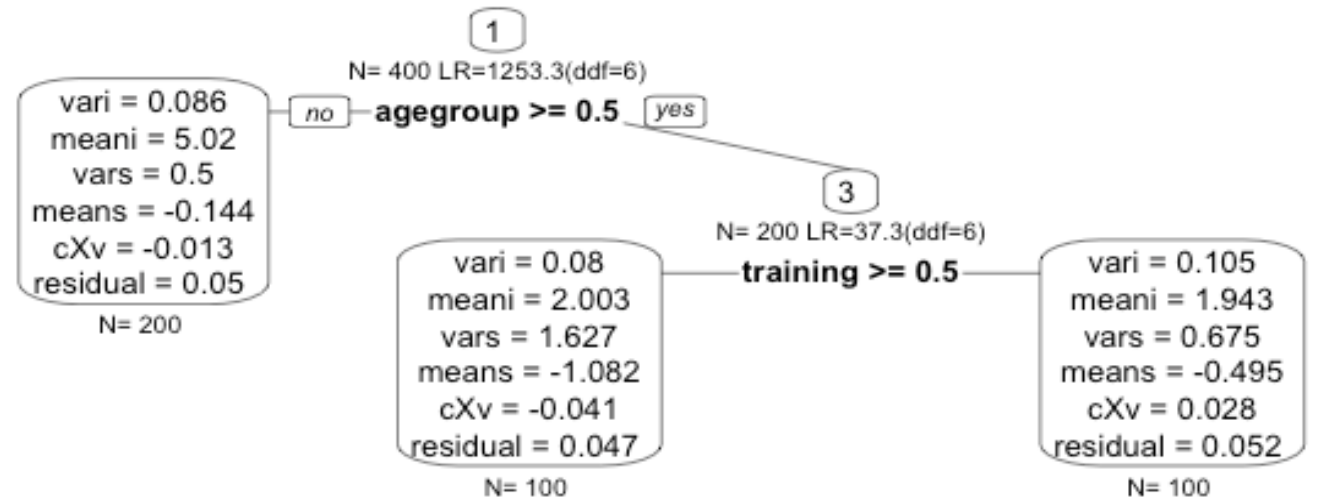
# SEM Trees Algorithm Cont'd

- Fitting the SEM in each group results in a separate fit
    - LL (log-likelihood) for each model
    - Model parameters are now group specific
- This can be compared using the likelihood ratio test
    - Is $LL_{all} <^* LL_{NSSI = 0} + LL_{NSSI = 1-1000}$
        - * Fit will always be better in multiple group model, but about <u>significance</u>
        - Larger LL is better (less negative)
        - 2 times the difference between models is a chi square distribution
            - Df = number of parameters of multiple group model minus parameters of no group model

# SEM Trees Algorithm Cont'd

- If the covariate is an integer (& ordinal), tests every possible group in sequence
  - 0 vs. 1-1000, 0-1 vs 2-1000, 0-2 vs 3-1000 etc…
- Same thing if numeric (continuous; i.e. 0.324)
  - Will test every value in sequence, so may be better to round first
- If categorical (factor), one vs. the rest scheme
  - If 5 categories, 15 possible splits
    - Computationally intensive
- Once a best split is determined:
  - Move down one level to start over searching for an additional split
  - Continues until:
    - No longer improves model fit
    - Reaches other stopping criterion
      - Too small N in a node
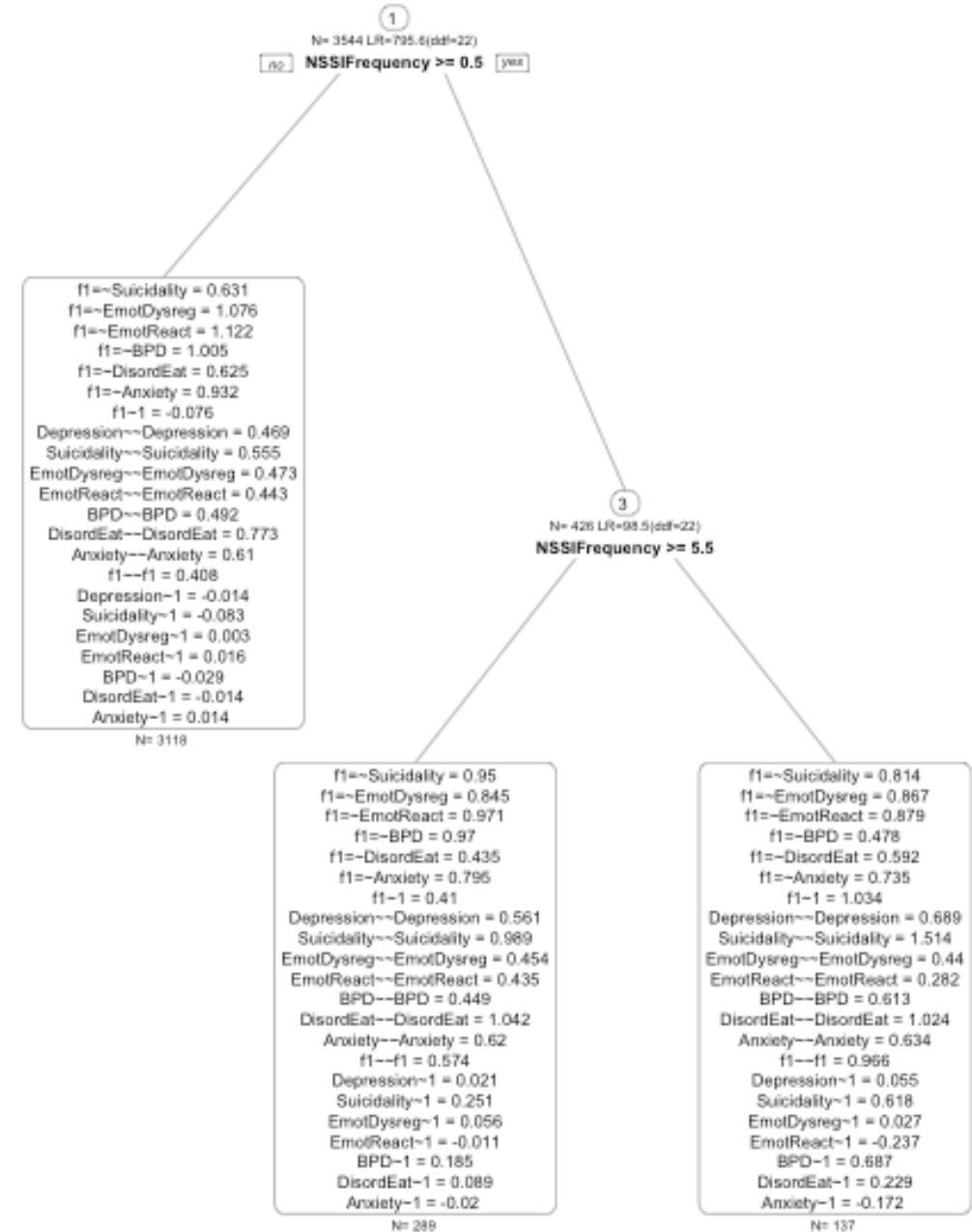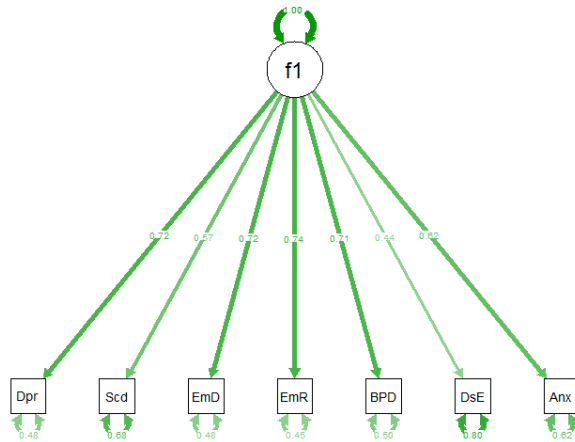      - A priori set maximum depth (# of groups)

# SEM Trees Algorithm Cont'd

- Just like DT, SEM Trees can create a large tree

- Say the first split is between 0 & 1

- Next level, test splits on other variables

- Example:
  - Second split occurs at agegroup = 1 & between training 0-1
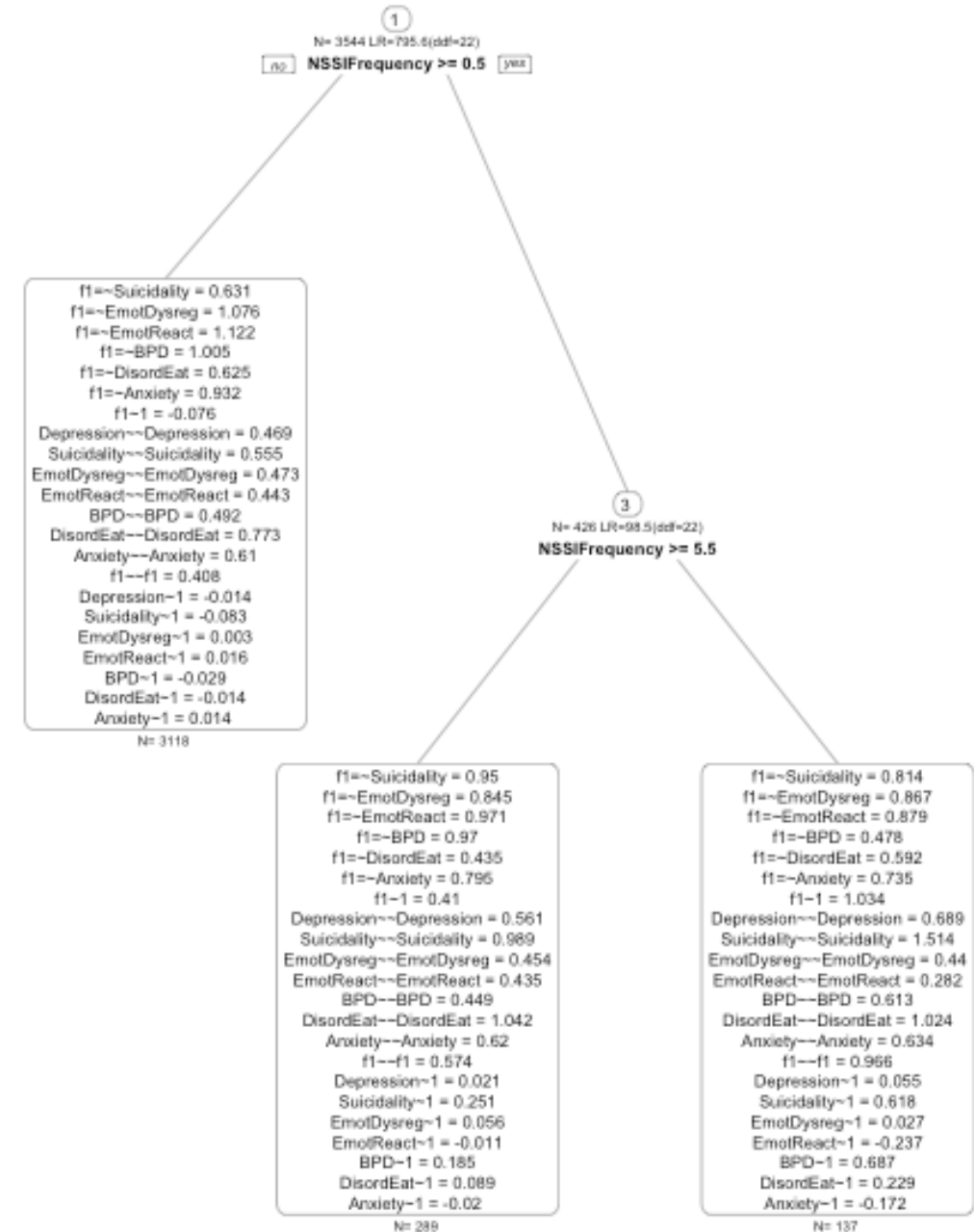  - No further split for agegroup = 0

# SEM Trees & NSSI

- Use a one factor CFA as the outcome
- Came up with the same splits as with DT

# SEM Trees & NSSI Cont'd

- Can investigate individual parameters to get a better clue to how groups differ
  - Factor mean increased w/ NSSI frequency
  - Most variable means increased
  - Factor variance increased

# SEM Trees & NSSI Conclusion

- Agreement in cutoffs across methods
  - Important, as SEM Trees hasn't been studied that much
- SEM Trees allows for a more comprehensive evaluation
  - While being able to make specific comparisons across groups
- Was able to incorporate more information into the model than if we used mixture models
  - This is not always the case in comparing mixtures and SEM Trees
  - We had a **very** informative covariate
- Not that difficult to program and didn't take long to run
  - Demonstrated in Lab Session T

# SEM Trees Options

# SEM Trees Options

- Because there is an SEM for each group
  - Should set minimum number of cases per node of tree
    - `min.N` in `semtree.control`
- Large number of covariates?
  - Need to prevent type-1 errors (overfitting)
    - Also, large, uninterpretable tree
  - Can use Bonferroni correction or use cross-validation
  - `bonferroni = TRUE` or `method="cv"` in `semtree.control`
- Large number of response options for covariate(s)
  - These variables are more likely to be split on
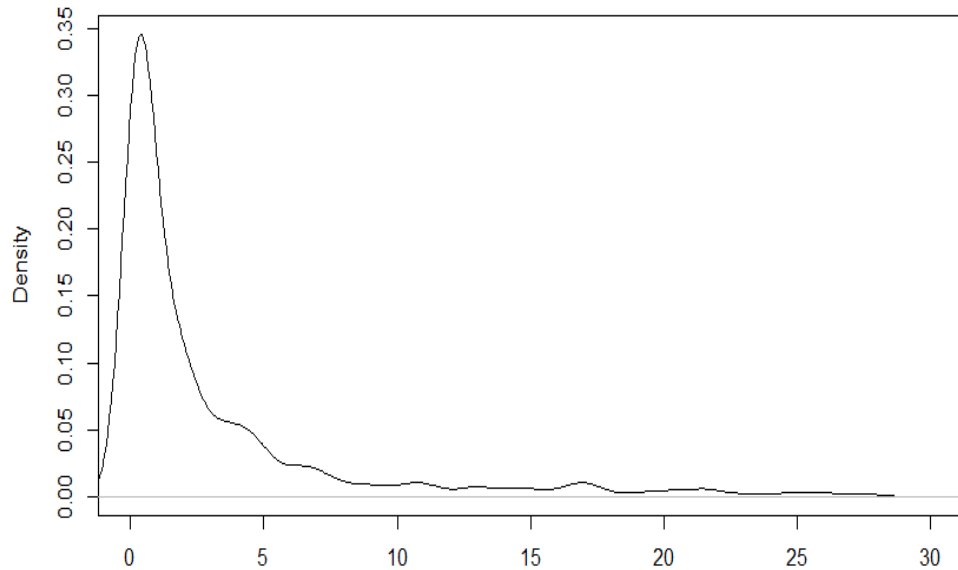  - Correct with changing `method="fair"`

# SEM Trees & Invariance

- In many contexts, it is necessary to ensure that groups are measured on the same construct

- This can either be call differential item functioning (IRT) or measurement invariance (SEM)

- This entails constraining certain parameters to be equal across groups

- Most common: constrain factor loadings

- Easy to do in SEM Trees
  - Set `invariance = ` parameter names in `semtree.control`
    - Demonstrated in Lab Session T

# SEM Trees vs. Mixture Models

# Mixture Models

- Traditional way of doing things
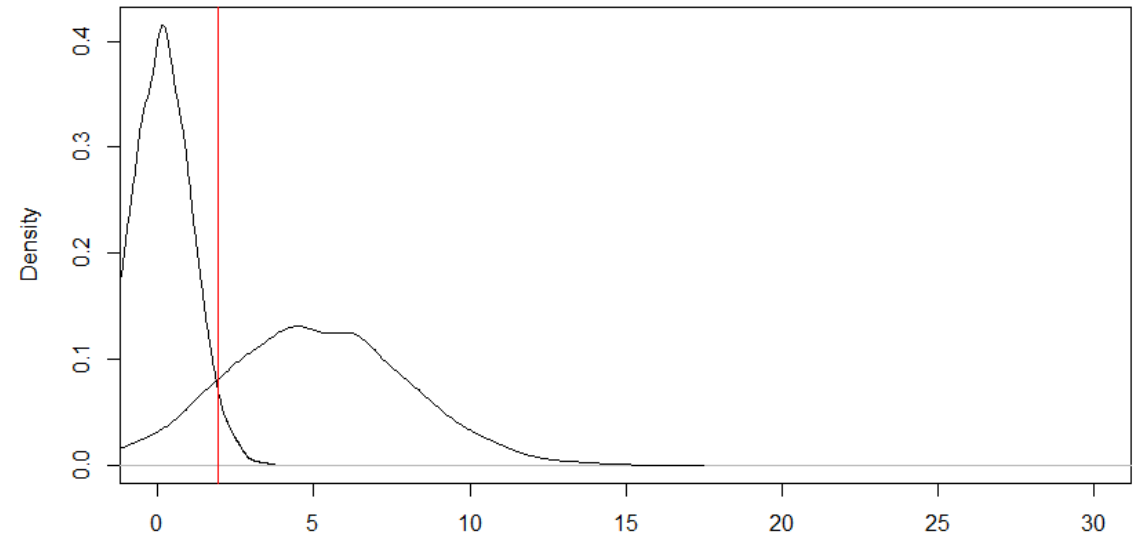  - Look for multiple underlying normal distributions underlying NSSI distribution

# SEM Trees vs. Mixture Models

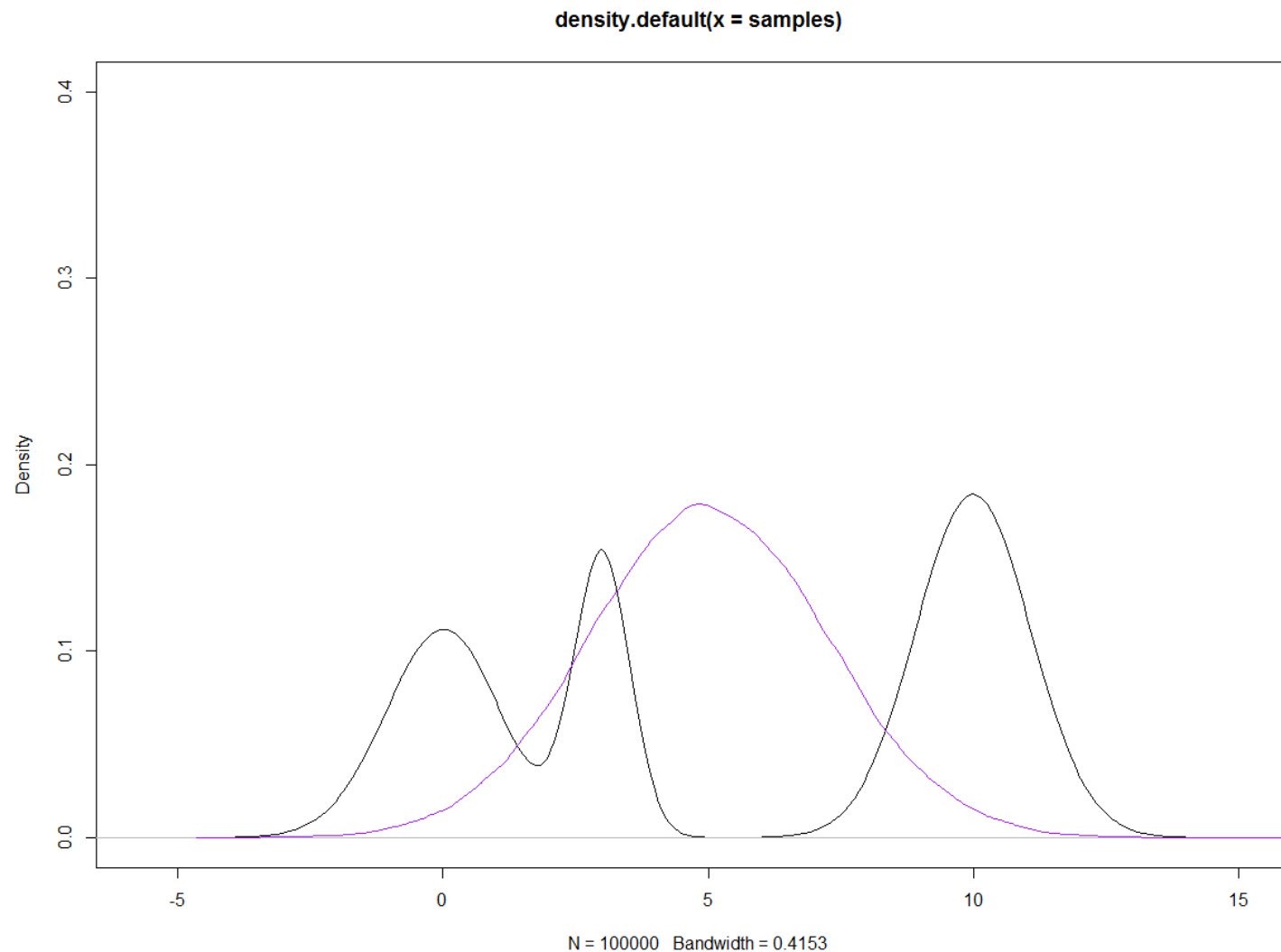- Do very similar things, but find groups (classes) a little differently
- Let's say we have a variable we are trying to find groups in
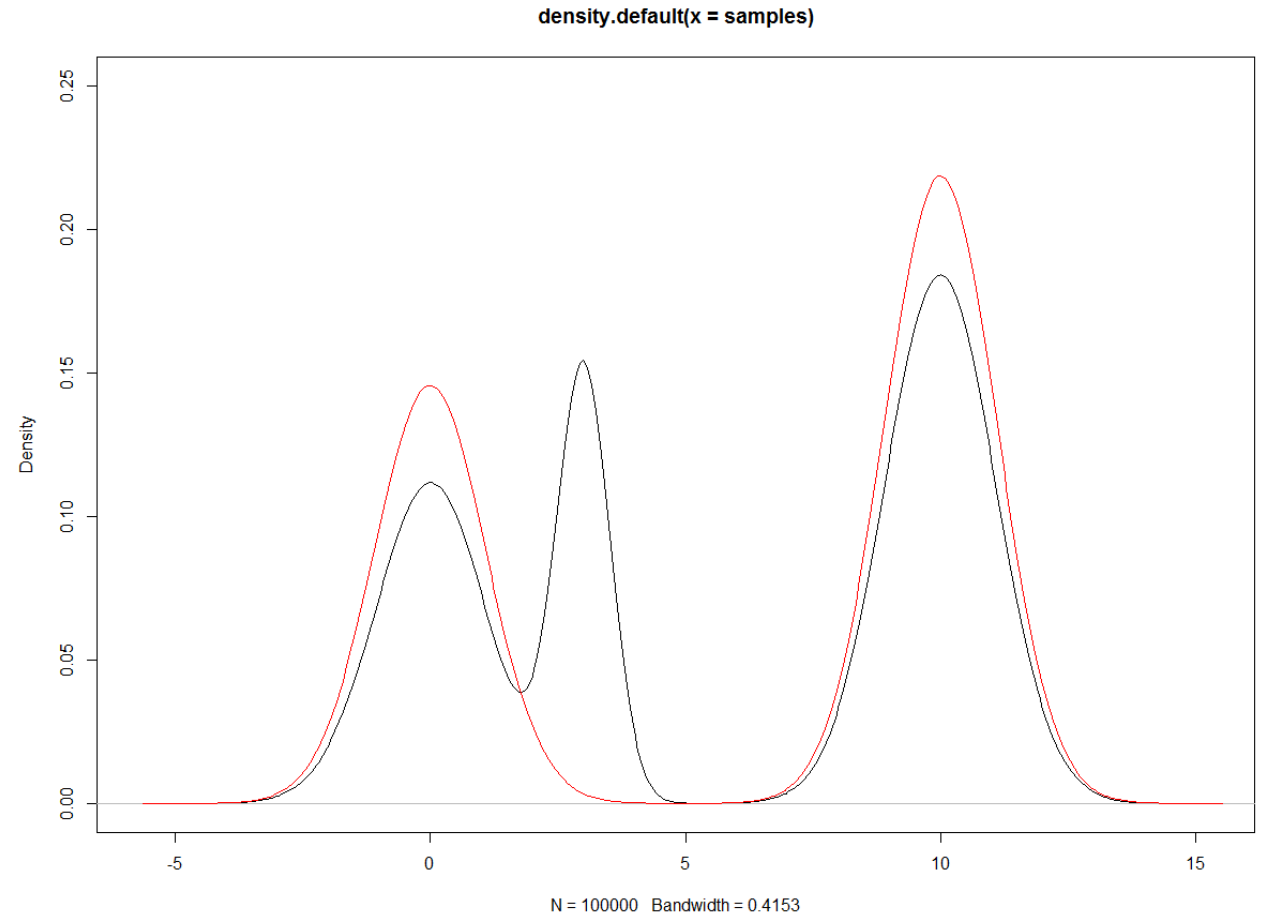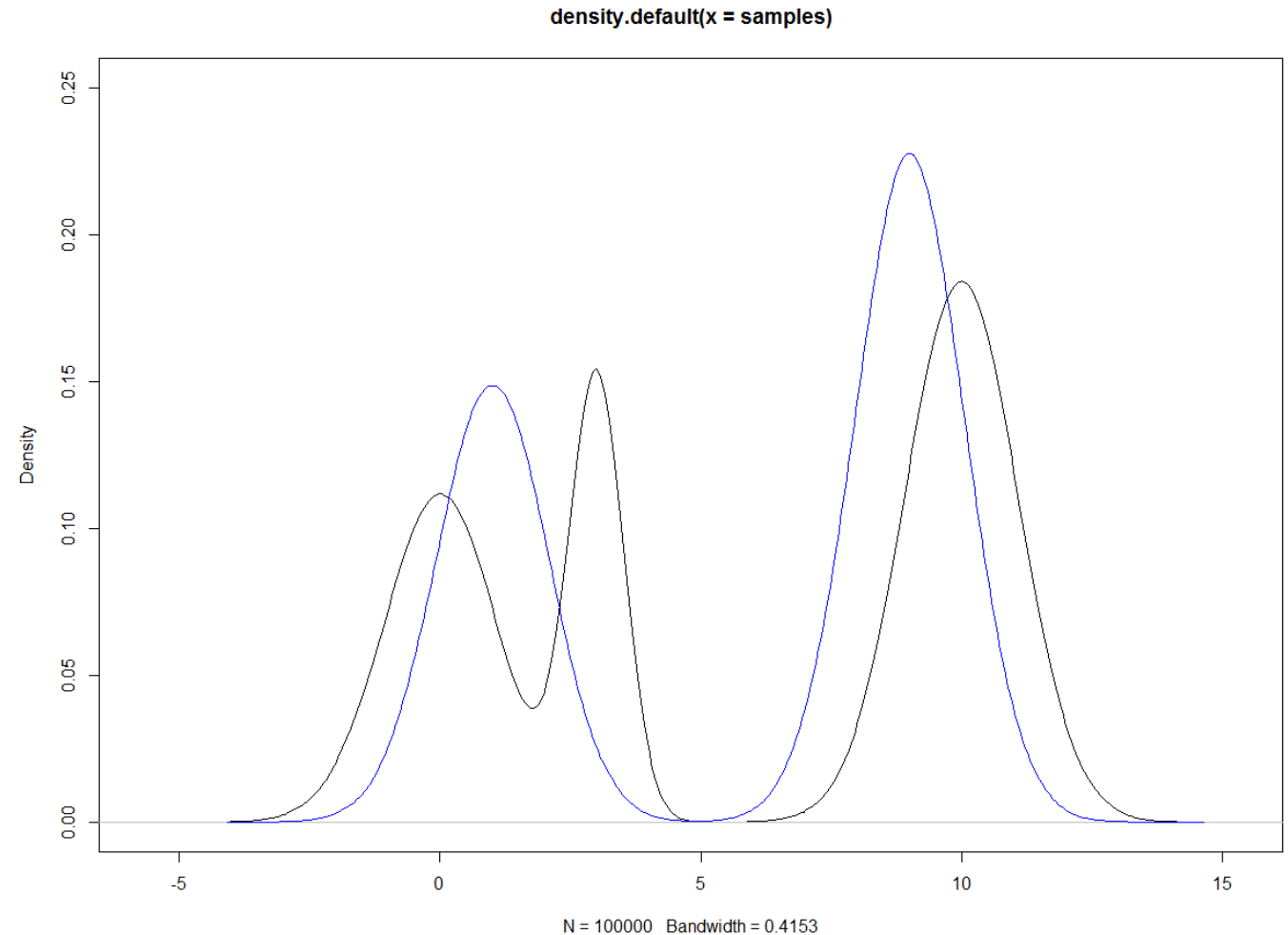- Hypothetical distribution for a cognitive score variable:



density.default(x = samples)

# Assume Homogeneity

- Doesn't do so well



density.default(x = samples)

N = 100000  Bandwidth = 0.4153

# Comparison Example

- Mixture models will attempt to fit this distribution directly

- An example with 2 classes

density.default(x = samples)

Density

N = 100000   Bandwidth = 0.4153

# Comparison Example Cont'd

- SEM Trees will try and fit this distribution
  - But only through splitting on the covariates
  - Meaning SEM Trees is more constrained
  - Doesn't have access to the whole search space
  - For example, the blue line could represent the distribution for a gender variable



density.default(x = samples)

N = 100000   Bandwidth = 0.4153
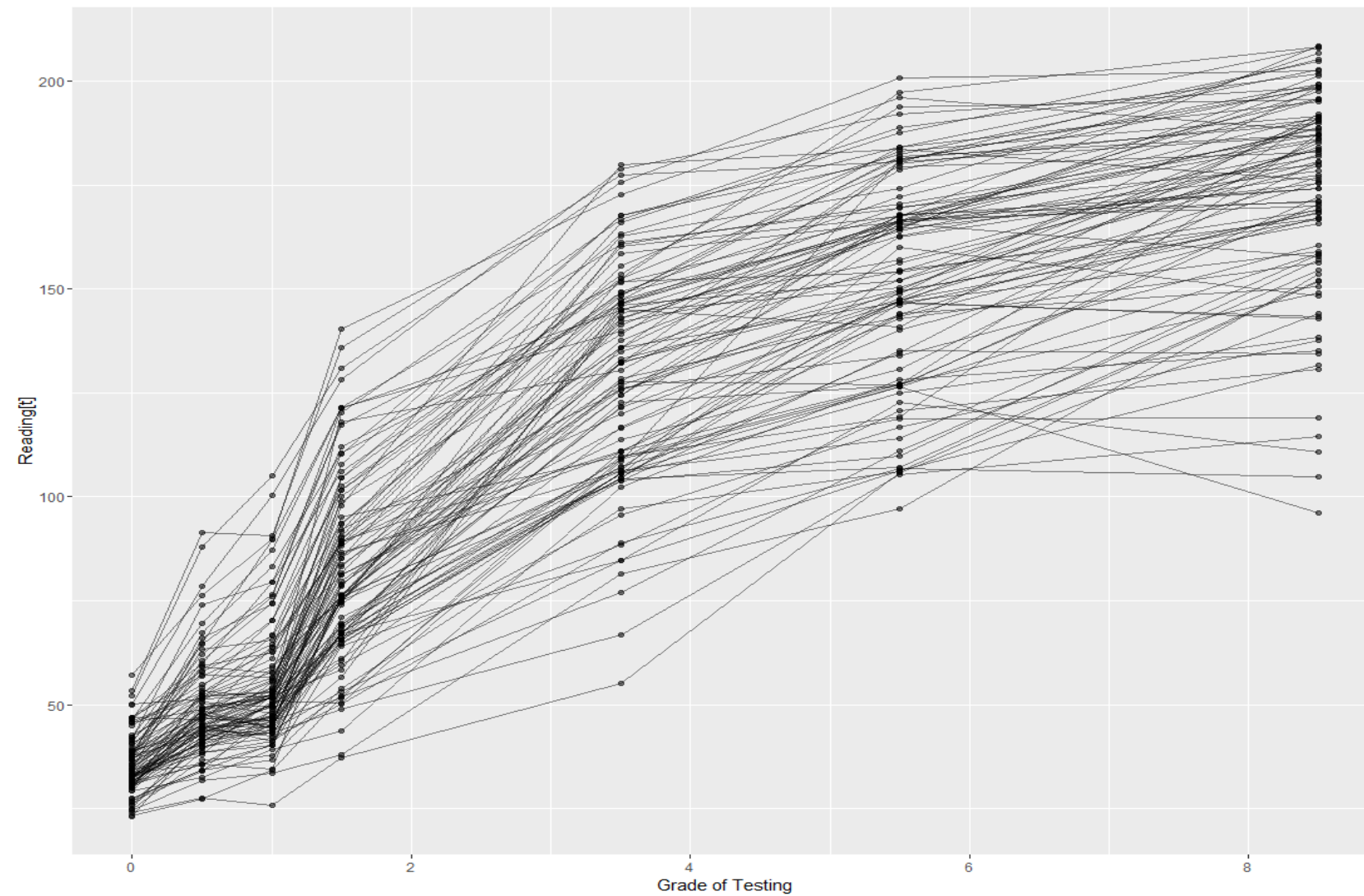
# Comparison Example Summary

- SEM Trees requires informative covariates to split on
  - Mixtures do not
- SEM Trees may be less likely to overfit
- SEM Trees allows interactions and non-linear prediction
  - In mixtures, the classes are all at the same level, not nested or based on interactions
  - In mixtures, the relationship between class membership and auxiliary variables is constrained to be linear
- Can be very useful to test both methods and compare the resultant groups
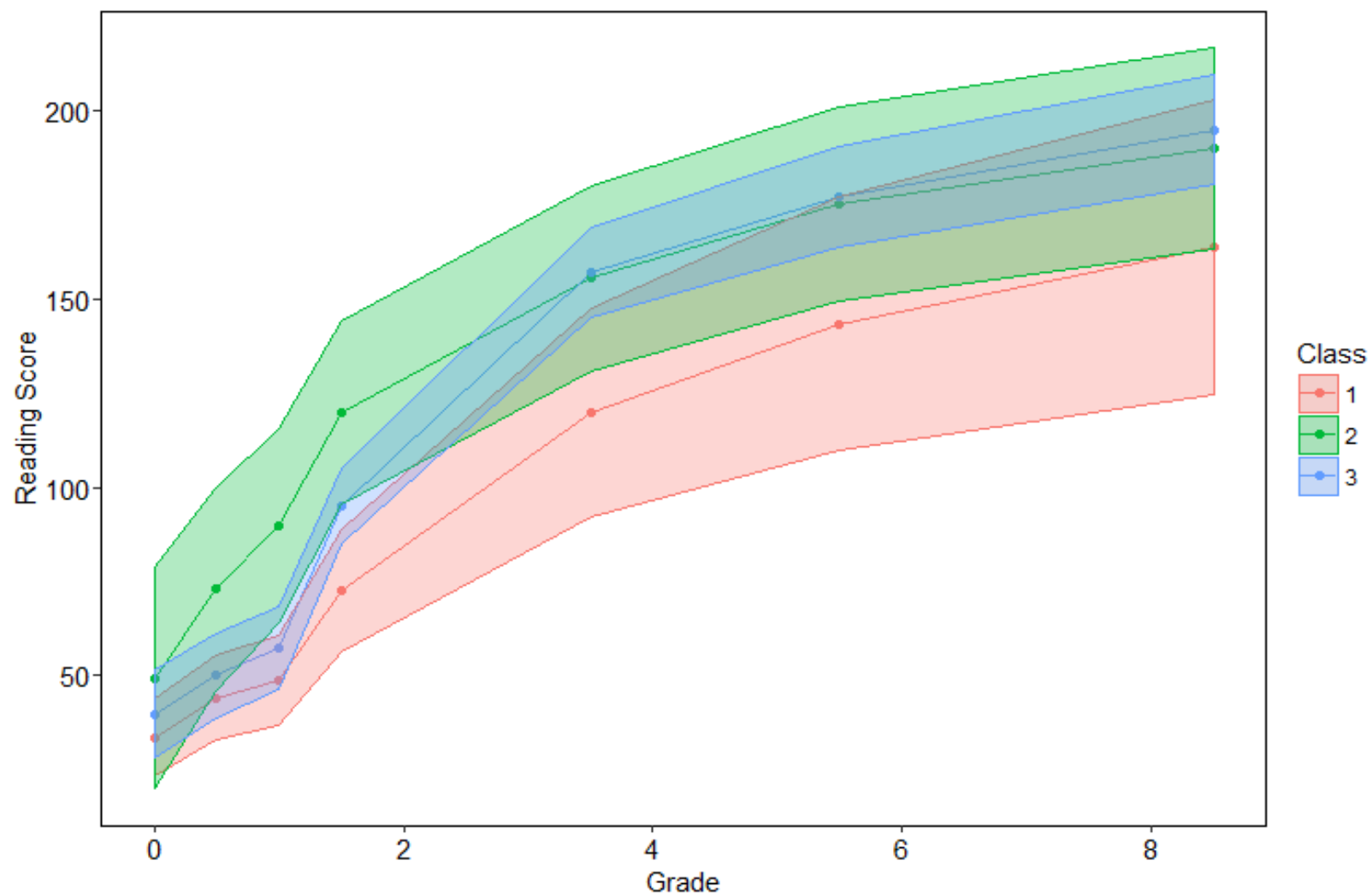
# SEM Tree vs. Growth Mixture Model

# Comparison Example

- Compare differences in trajectories
  - Using a latent growth curve model
  - Change in reading from Kindergarten to 8th grade
  - Early Childhood Longitudinal Study - Kindergarten Cohort (ECLS-K) data
- Questions we are trying to answer
  - Do all children learn at the same rate or are the subgroups?
    - i.e. Late bloomers (start low, but catch up)
    - Those who start high and improve fast
  - Can other measured variables give us insight into these groups
    - Covariates include fine motor skills (*fine*), gross motor skills (*gross*), approaches to learning (*learn*), self-control (*control*), interpersonal skills (*interp*), internalizing behaviors (*int*), and general knowledge (*gk1*)
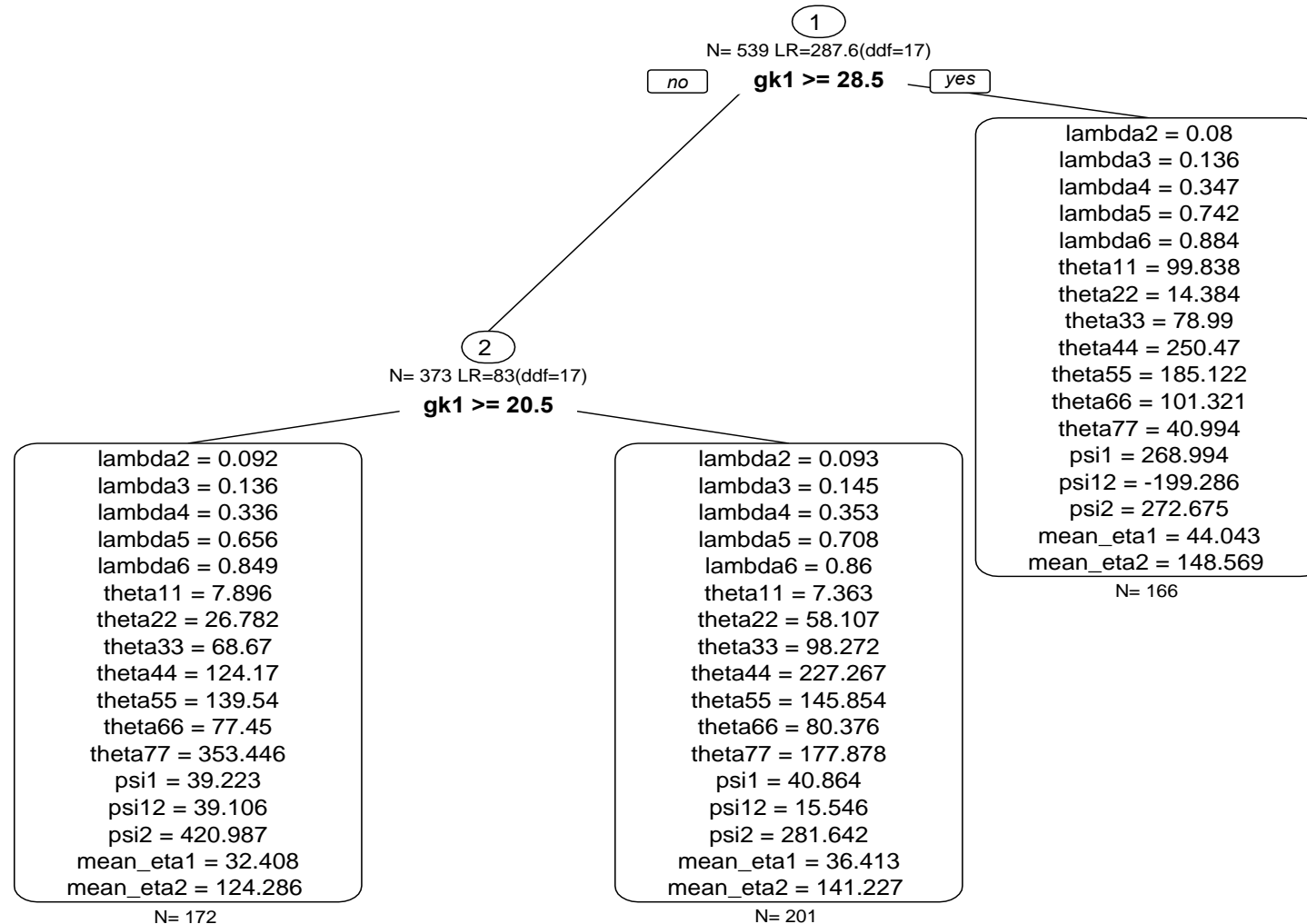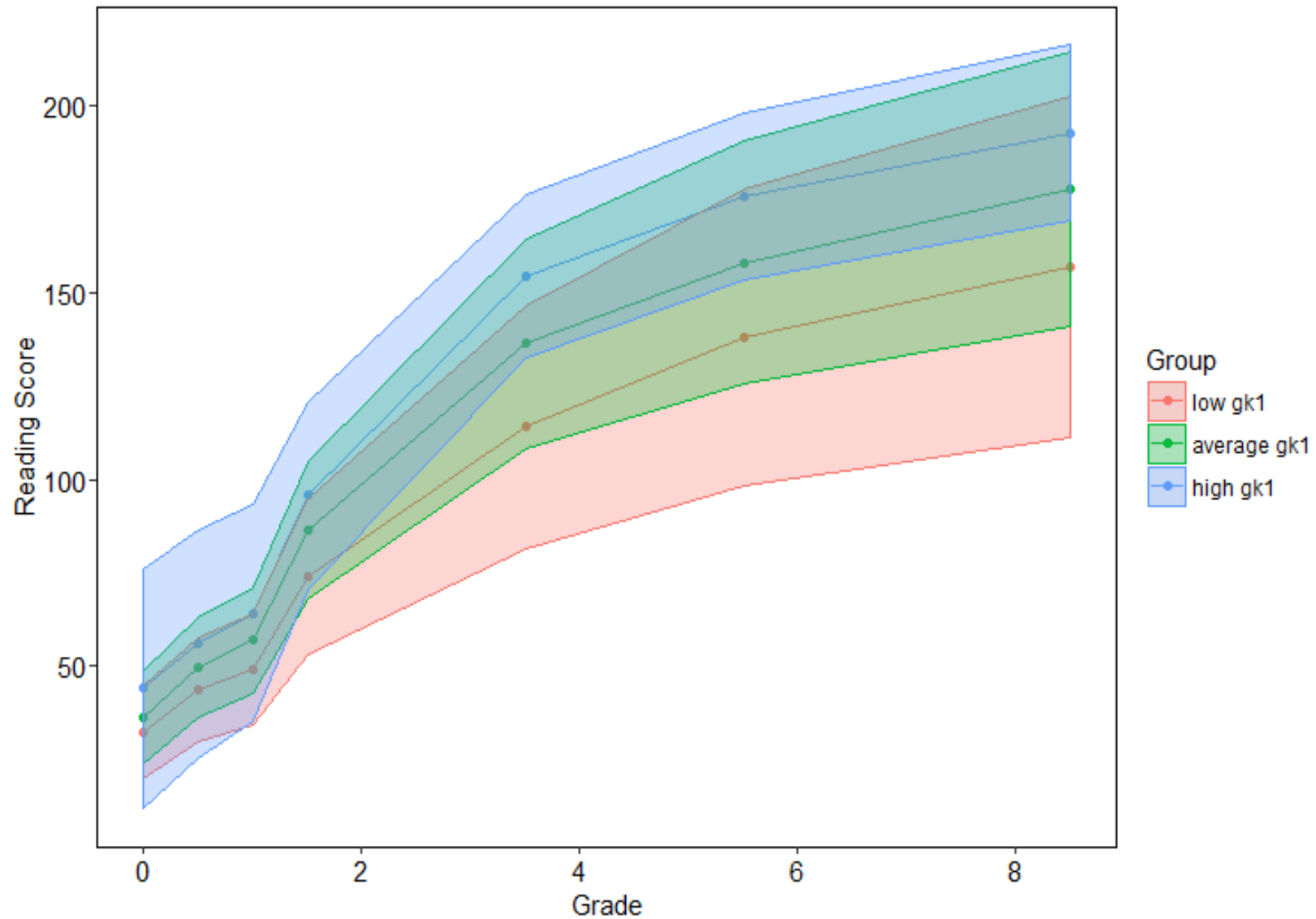
# Trajectory Plot

# Growth Mixture Model

# Growth SEM Trees

# Growth SEM Trees

# Comparison Example Summary

- SEM Trees split twice on the general knowledge variable
  - Was a significant auxiliary variable in Mixture Model
- SEM Trees didn't fit as well
  - More variability within group and overlap across groups
- The Growth Mixture Model had more interesting classes
  - Class 1 caught up to Class 3 in grades 4 and 5
    - Different loadings across groups
  - SEM Trees groups really only differed in Kindergarten

# SEM Trees Conclusion

- Has been extended to allow multiple trees
  - SEM Forests – analogous extension from DT → Random Forests
  - Can take a really long time to run (hours/days)
- Very new method
  - Hasn't been tested to a large extent to set guidelines for best practices
- Can take a long time to run
  - Small model and small # of covariates can be 5 minutes
  - Large model and large # of covariates can be hours
- REMEMBER: Just because you get out a tree, does not mean that it represents real groups
  - Needs to make theoretical sense
  - Could just be the product of overfitting