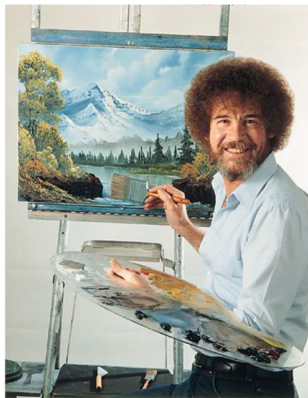


Subgroup detection in GLMMs and GAMs

glmertrees, splinetrees and gamtrees

Marjolein Fokkema

Trees



Code and data:

<https://github.com/marjoleinF/Speech-prosody-workshop-trees>

Or: <https://tinyurl.com/mpphwc2h>

Short history of trees

Trees recursively partition the observations in a dataset based on the values of covariates, in order to find subgroups that have increasingly similar values on the response variable.

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)

Unbiased recursive partitioning:

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)
- ▶ CART: Classification and Regression Trees (Breiman et al., 1984)

Unbiased recursive partitioning:

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)
- ▶ CART: Classification and Regression Trees (Breiman et al., 1984)
- ▶ ID3 (Quinlan, 1986)

Unbiased recursive partitioning:

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)
- ▶ CART: Classification and Regression Trees (Breiman et al., 1984)
- ▶ ID3 (Quinlan, 1986)
- ▶ C4.5 (Quinlan, 1993)

Unbiased recursive partitioning:

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)
- ▶ CART: Classification and Regression Trees (Breiman et al., 1984)
- ▶ ID3 (Quinlan, 1986)
- ▶ C4.5 (Quinlan, 1993)

Unbiased recursive partitioning:

- ▶ GUIDE: Generalized unbiased interaction detection and estimation (Loh, 2002)

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)
- ▶ CART: Classification and Regression Trees (Breiman et al., 1984)
- ▶ ID3 (Quinlan, 1986)
- ▶ C4.5 (Quinlan, 1993)

Unbiased recursive partitioning:

- ▶ GUIDE: Generalized unbiased interaction detection and estimation (Loh, 2002)
- ▶ ctree: Conditional inference trees (Hothorn, Hornik & Zeileis, 2006)

Short history of trees

Early tree methods:

- ▶ AID; Automated Interaction Detection (Morgan & Sonquist, 1963)
- ▶ CART: Classification and Regression Trees (Breiman et al., 1984)
- ▶ ID3 (Quinlan, 1986)
- ▶ C4.5 (Quinlan, 1993)

Unbiased recursive partitioning:

- ▶ GUIDE: Generalized unbiased interaction detection and estimation (Loh, 2002)
- ▶ ctree: Conditional inference trees (Hothorn, Hornik & Zeileis, 2006)
- ▶ MOB: Model-based recursive partitioning (Zeileis, Hothorn & Hornik, 2008)

Short history of trees

Model-based recursive partitioning (MOB)

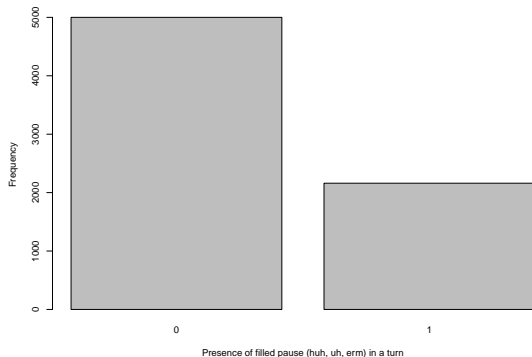
Rationale:

A single global parametric model may not fit all observations well.

- ▶ E.g., (G)LM: $y_i = x_i^\top \beta + \epsilon_i$

Example: Filled pauses

Dataset from Gardner et al. (2021), subset of Switchboard Corpus of American English.



Example: Filled pauses

```
gmod <- glm(NFP.bi ~ 1, data = df, family = "binomial")  
coef(gmod)
```

```
## (Intercept)
```

```
## -0.8388668
```

Model-based recursive partitioning (MOB)

Rationale:

A single global parametric model may not fit all observations well.

- ▶ E.g., (G)LM: $y_i = x_i^\top \beta + \epsilon_i$

When additional covariates are available, it may be possible to partition the dataset into subgroups, and obtain better-fitting models in each of the subgroups.

- ▶ (G)LM tree: $y_i = x_i^\top \beta_j + \epsilon_i$

Example: Filled pauses

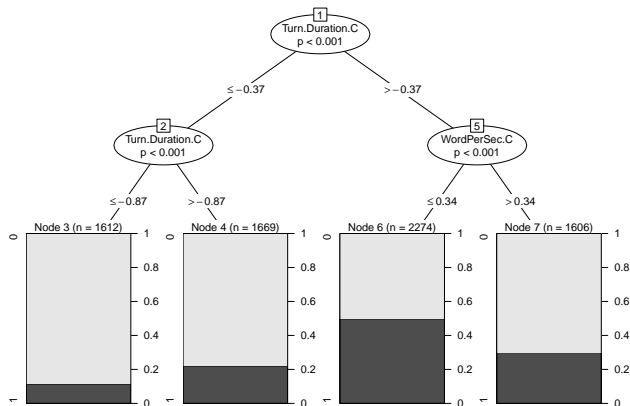
Additional covariates:

- ▶ NVarbs: Variable context (present or not)
- ▶ Turn.Duration.C: Standardized turn duration
- ▶ CharPerWord.C: Mean word length
- ▶ WordPerSec.C: Speech rate

```
library("partykit") ## Implements, amongst others, MOB  
gt <- glmtree(NFP.bi ~ 1 | NVarbs + Turn.Duration.C +  
              CharPerWord.C + WordPerSec.C,  
              data = df, family = "binomial",  
              minsize = 1500)
```


Example: Filled pauses

```
plot(gt)
```



Example: Filled pauses

```
coef(gt)
```

```
##           3           4           6           7  
## -2.05524402 -1.26629497 -0.01055419 -0.85264070
```

MOB algorithm (Zeileis et al., 2008)

GLM-based recursive partitioning:

- a) Fit a GLM to all observations in the current subgroup.
- b) Test for instability of the GLM parameters with respect to each of the partitioning variables.
- c) If there is some overall parameter instability, split the subgroup with respect to the partitioning variable associated with the highest instability.
- d) Repeat Steps (a) through (c) in each of the resulting subgroups.

Step b): Parameter stability tests

##

\$'1'

##

	NVarbs	Turn.Duration.C	CharPerWord.C	WordPerSe
## statistic	24.926	519.185	59.89	115
## p.value	0.000	0.000	0.00	0

##

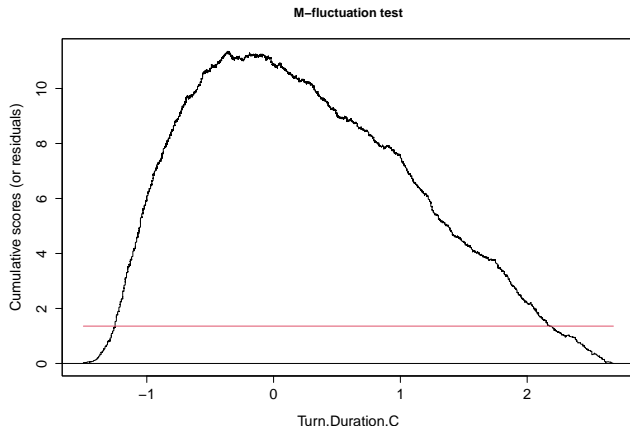
\$'2'

##

	NVarbs	Turn.Duration.C	CharPerWord.C	WordPerSe
## statistic	5.700	68.284	42.452	26
## p.value	0.164	0.000	0.000	0

Step b): Parameter stability tests

Computation of test statistic for `Turn.Duration.C` in the first node:



Generalized linear mixed-effects model tree (Fokkema et al., 2018; Fokkema & Zeileis, in press)

- ▶ LMM trees extend LM trees with random effects, much like the LMM extends the LM:

$$y_i = X_i\beta + Z_ib_i + \epsilon_i$$

- ▶ LMM tree:

$$y_i = X_i\beta_j + Z_ib_i + \epsilon_i$$

- ▶ Allows to account for and quantify dependence between observations within the same subjects or clusters i .

GLMM trees: Estimation

0. Initialize: Set step $r = 0$ and all random-effect estimates $\hat{b}_{i,(r)} = 0$.
1. Estimate subgroups: Set $r = r + 1$. Fit an LM tree $X_i \hat{\beta}_{j,(r)}$ using $Z_i \hat{b}_{i,(r-1)}$ as an offset. Extract the partition or subgroup memberships $j_{(r)}$.
2. Estimate full mixed-effects model: Fit the mixed-effects model $\mu_i = X_i \beta_{j,(r)} + Z_i b_{i,(r)}$ with the subgroups $j_{(r)}$ from Step 1. Extract the random-effect estimates $\hat{b}_{i,(r)}$ from the fitted model.
3. Repeat Steps 1 and 2 until convergence.

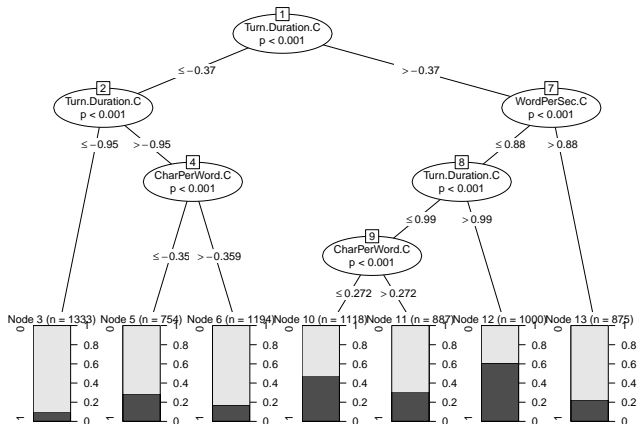
This procedure can easily be generalized to non-Gaussian responses within the GLM(M).

Example: Filled pauses

```
library("glmertree")  
glmmt <- glmertree(NFP.bi ~ 1 | (1 | Speaker_Number) |  
                  NVarbs.bi + Turn.Duration.C +  
                  CharPerWord.C + WordPerSec.C,  
                  minsize = 750, data = df,  
                  family = "binomial")
```


Example: Filled pauses

```
plot(glmmt, which = "tree")
```



Example: Filled pauses

```
fixef(glmmt)
```

```
##      (Intercept)
## 3      -2.4054977
## 5      -1.1114517
## 6      -1.7355811
## 10     -0.1797613
## 11     -0.8547658
## 12      0.3968831
## 13     -1.1937829
```

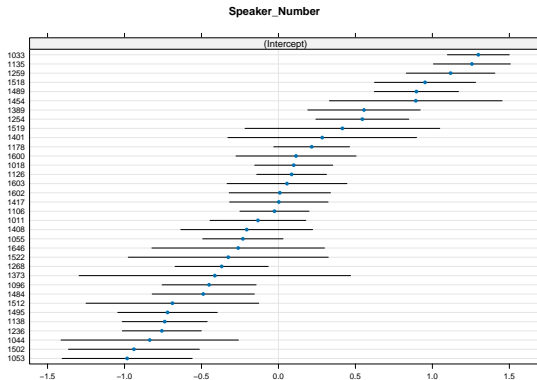
```
VarCorr(glmmt)
```

```
## Groups              Name              Std.Dev.
## Speaker_Number (Intercept) 0.6716
```

Example: Filled pauses

```
plot(glmmt, which = "ranef")
```

```
## $Speaker_Number
```



Example: Articulatory trajectories

Integrating splines

- ▶ The (G)LMM (or GLMM tree model) can incorporate parametric splines:

$$y_i = X_i\beta_j + Z_ib_i + \epsilon_i$$

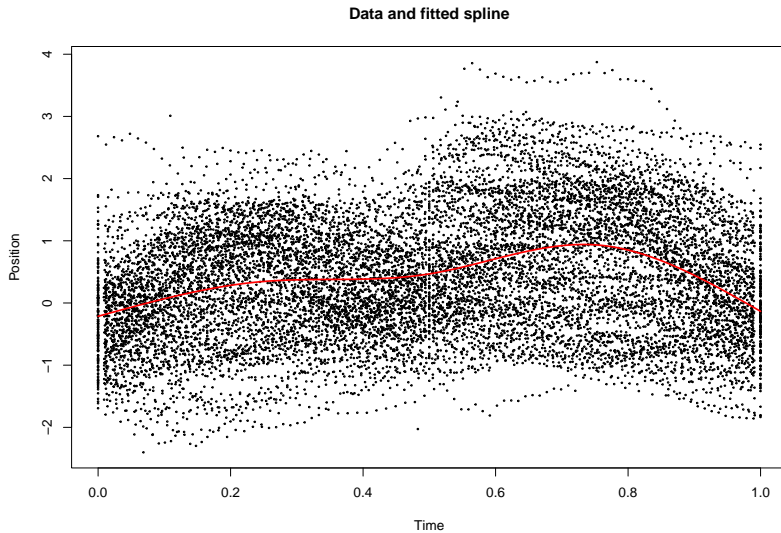
- ▶ “Only” need to add non-linear basis functions to the design matrix X_i .

Example: Articulatory trajectories

Example: Articulatory trajectories

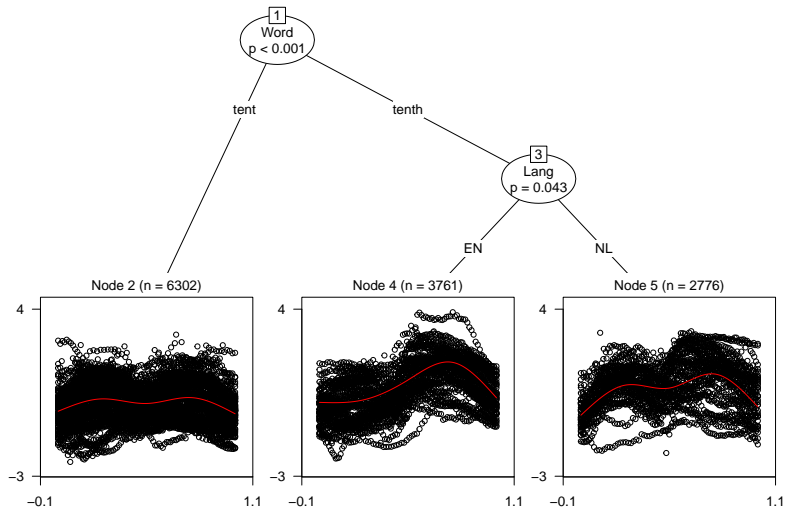
```
## (Intercept)    spl_basis1    spl_basis2    spl_basis3    spl_bas
##   -0.2138743     0.4389011     1.3899242     1.1945478    -0.3325
```

Example: Articulatory trajectories



Example: Articulatory trajectories

Example: Articulatory trajectories



Example: Articulatory trajectories

##	Groups	Name	Std.Dev.
##	Speaker	(Intercept)	0.44710
##	Residual		0.74111