# Mixed models
and why sociolinguists should use them
Daniel Ezra Johnson

| VARBRUL / GoldVarb | other |
|---|---|
| dependent variable (DV) | DV, response, y |
| factor group, independent variable (IV) | IV, factor (categorical), predictor, x |
| factor | level |
| factor weight | coefficient, effect, estimate, $\beta$ |
| factor weight range | similar to 'effect size' |
| input probability | intercept |
| applications / total | (response) proportion |

| lmer | other |
|---|---|
| mixed model | mixed-effects, hierarchical, or multilevel model |
| fixed effect | main effect |
| (all) fixed-effects model | flat model |
| conditional modes of random effects | random effect estimates, random effect BLUPs |

# Terminological 'translations'

| PROPERTIES OF DATA | GoldVarb | Rbrul | R | POSSIBLE ANALYSIS |
|---|:---:|:---:|:---:|---:|
| response / DV: 2 categories | ✔ | ✔ | ✔ | logistic regression |
| response: 3+ categories | | | ✔ | ordinal, multinomial logistic |
| response: count | | | ✔ | Poisson regression, etc. |
| response: continuous | | ✔ | ✔ | linear regression |
| predictor(s) / IV(s) : categorical | ✔ | ✔ | ✔ | (any) |
| predictor(s): continuous | | ✔ | ✔ | (any) |
| predictor(s): have interactions | *hard* | | ✔ | (any) |
| random intercept(s) | *?* | ✔ | ✔ | **mixed model** |
| random slope(s) | *??* | | ✔ | mixed model |
| lots of data (need for speed) | | ✔ | ✔ | |
| | | *hard* | ✔ | plots and graphics |
| | | | ✔ | other statistical methods |
| | ✔ | | | "slash" operator |
| | *?* | *?* | | user friendly |

## Comparing Software Tools

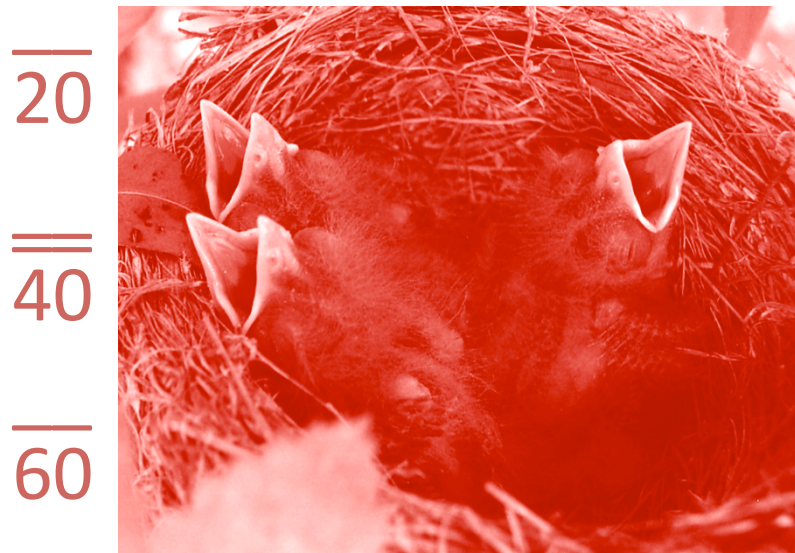GoldVarb          Rbrul                              R

Finding the right tool for the job

- mixed models: both fixed effects and random effects
- fixed effect: ordinary regression predictor (IV)
- random effect: theoretically sampled from a population
    - est. population variance (s.d.) is the real parameter
    - individual estimates (BLUPs) "shrunk" towards mean
    - residual random effects should be normally distributed

- random intercept: individuals "high" or "low" (input prob.)
- random slope: individuals differ w.r.t. predictors (constraints)

- in model fitting, there is a penalty on the random effects
    - as much variance as possible assigned to fixed effects
    - only the left-over variance is assigned to random effects
- this random effect penalty allows nested models to fit
    - sometimes fixed vs. random (or separate runs) is a valid choice
    - but nested predictors must be random effects in a mixed model

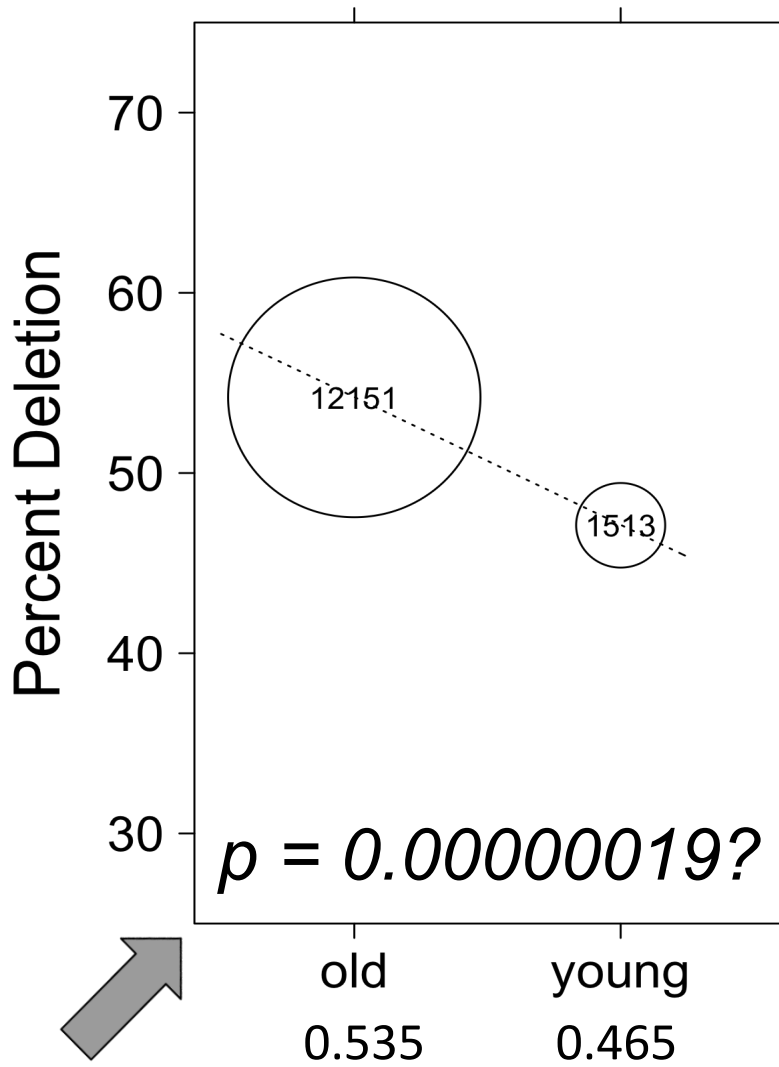# What are mixed models?

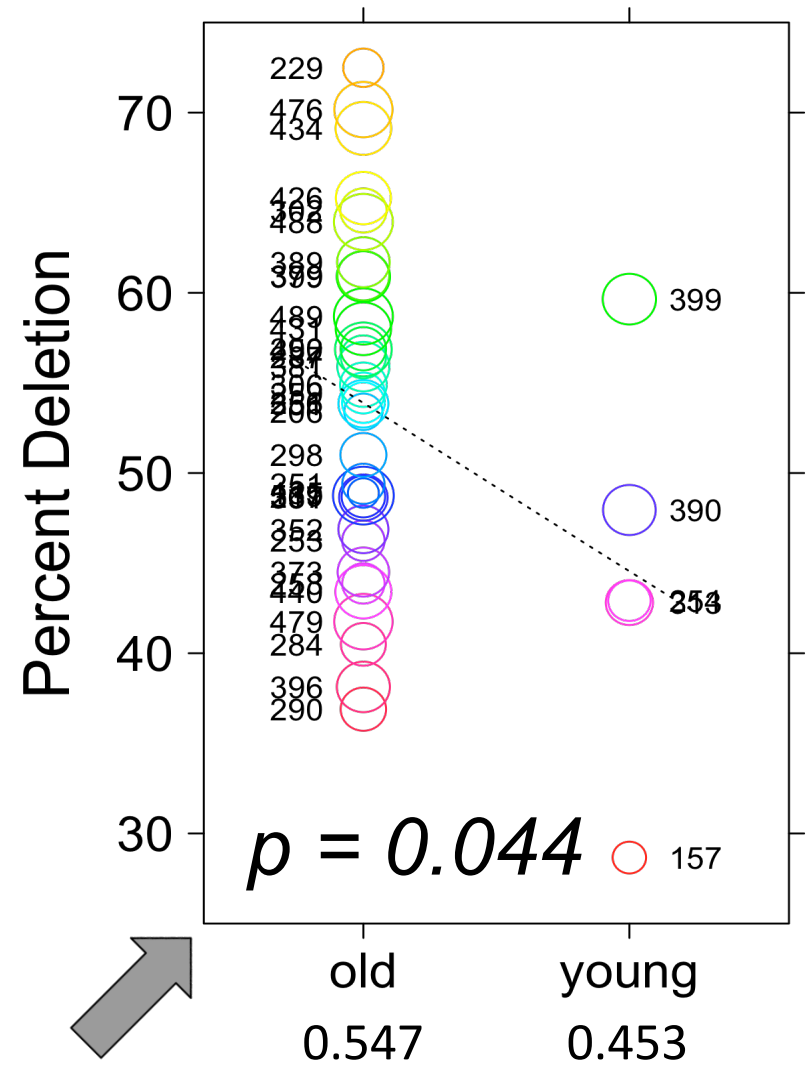Mixed models for nested data

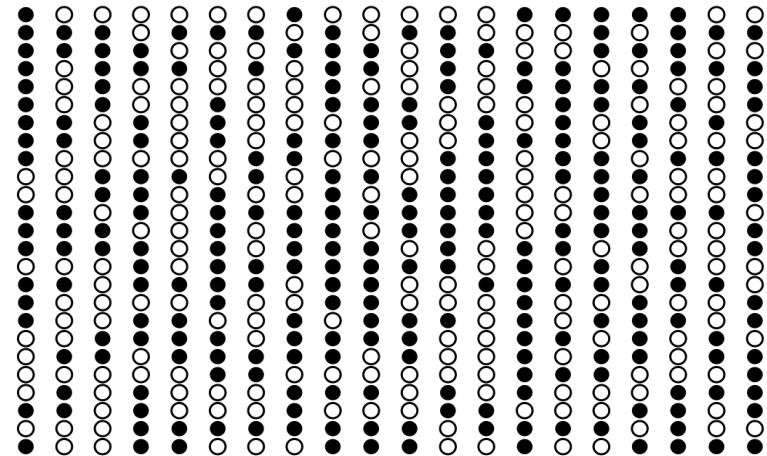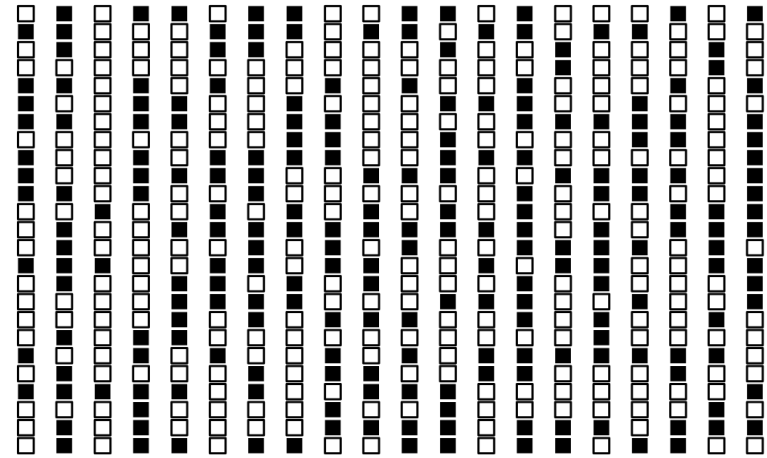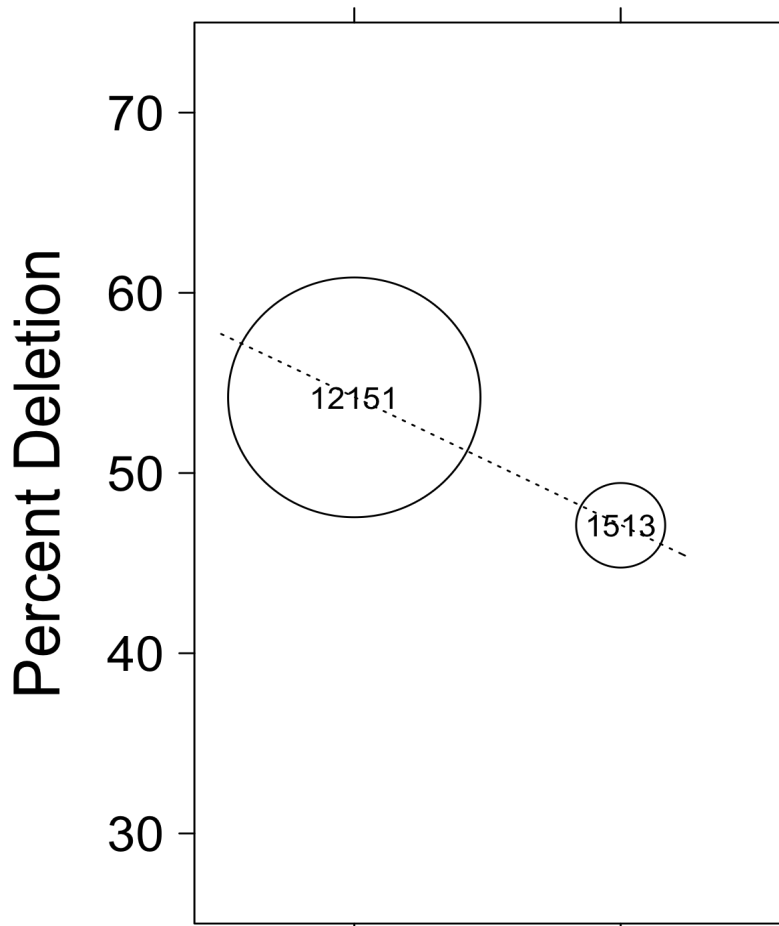When we don't need mixed models

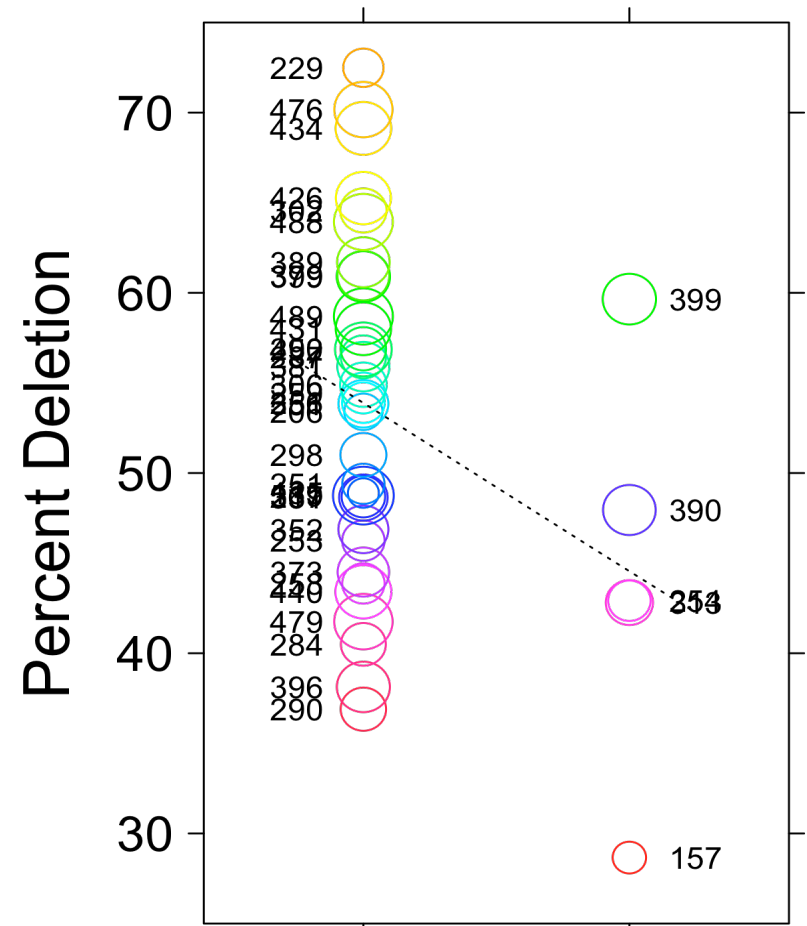And when we might need them

Random effects and significance

large effect size: 0.167 vs. 0.833
small significance: p = 0.08

small effect size: 0.45 vs. 0.55
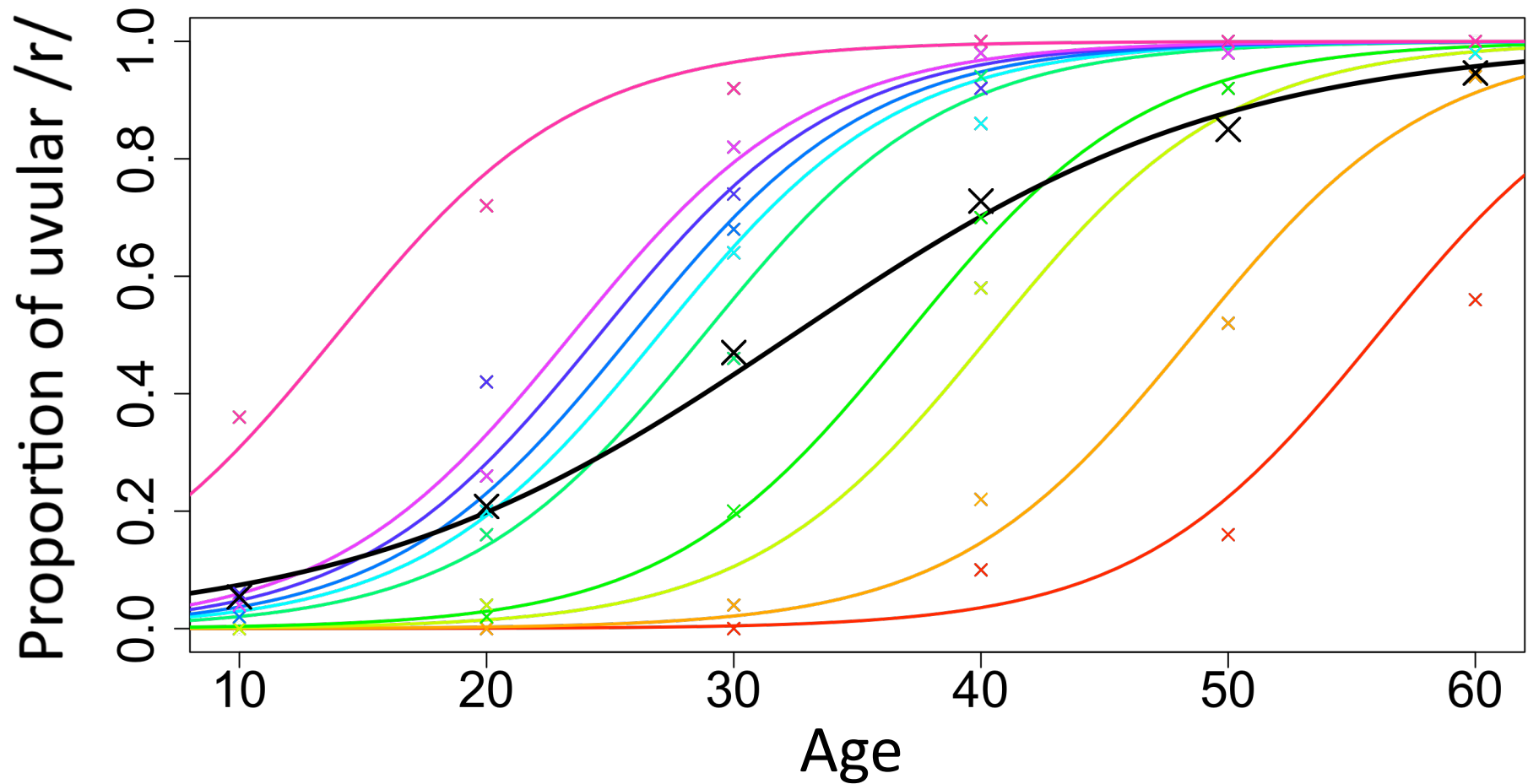larger significance: p = 0.002

# Significance vs. 'effect size'

Unbalanced data and effect size

age coefficient w/ no random effect: 0.113 log-odds/year
age coeff. w/ speaker random effect: 0.205 log-odds/year

Crossed factors and effect size

## speaker-nesting predictors

*constant within (data from) each speaker*
age?  gender  race  class  c.o.p. ...

- significance more accurate:
  p = larger, "no longer significant"?
- effect sizes more accurate with
  unbalanced data: larger/smaller

## speaker-crossed predictors

*vary within (data from) each speaker*
age?  style  phon./gram. context...

- effect sizes more accurate:
  larger (logistic regression only)

# Summary: speaker effect's effects

## speaker-nesting predictors

*constant within (data from) each speaker*

age?  gender  race  class  c.o.p. ...

## word-nesting predictors

*constant within (data from) each word*

frequency  gram. cat.  int. phon. ..

- significance more accurate:
  p = larger, "no longer significant"?
- effect sizes more accurate with
  unbalanced data, larger/smaller

## speaker-crossed predictors

*vary within (data from) each speaker*

age?  style  phon./gram. context...

## word-crossed predictors

*vary within (data from) each word*

stress  style  ext. phon. ...

- effect sizes more accurate:
  larger (logistic regression only)

# Word effect just like speaker effect

## speaker-nesting predictors

*constant within (data from) each speaker*

age?  gender  race  class  c.o.p. ...

## word-nesting predictors

*constant within (data from) each word*

frequency  gram. cat.  int. phon. ..

- significance more accurate:
  p = larger, "no longer significant"?
- effect sizes more accurate with
  unbalanced data, larger/smaller

## speaker-crossed predictors

*vary within (data from) each speaker*

age?  style  phon./gram. context...

## word-crossed predictors

*vary within (data from) each word*

stress  style  ext. phon. ...

word

speaker

- effect sizes more accurate:
  larger (logistic regression only)

# Crossed random effects for speaker & word

- use random effect estimates to identify 'new' fixed effects
  - modeled subject/word variation may include true individual variation, as well as unmodeled fixed effects
  -

- use random effect estimates to (empirically) build groups

- use random effect estimates as predictors in new models

- use random effect population variances to predict behavior of new subjects and words not in the original sample

- can perform an easy transformation into the 'language' of GoldVarb (with some caveats) – this is not a real problem

# Other benefits of mixed models

- cutting-edge statistics, like VARBRUL was in the 1970's
  - follow evolution on R-sig-ME
- double debate over p-values:
  - best way to calculate them
  - should they be used at all?
- convergence problems
  - requires more data (1000's > 100's)
- mixed model tool can be used well or badly, just like any model
  - still need to address multicollinearity
- should not be the only tool
  - mixed models are a better hammer, but everything is still not a nail
- "All models are wrong … but some are useful." – Box

# Drawbacks to mixed models

70     2.  Theory and Computational Methods for LME Models

Substituting (2.16) into (2.15) into (2.6) provides the likelihood as

$$L(\beta, \theta, \sigma^2 | y) = \prod_{i=1}^{M} \frac{\exp\left[-\|c_{0(i)} - R_{00(i)}\beta\|^2/2\sigma^2\right]}{(2\pi\sigma^2)^{n_i/2}} \, \text{abs}\left(\frac{|\Delta|}{|R_{11(i)}|}\right)$$

$$= \frac{\exp\left(-\sum_{i=1}^{M} \|c_{0(i)} - R_{00(i)}\beta\|^2/2\sigma^2\right)}{(2\pi\sigma^2)^{N/2}} \prod_{i=1}^{M} \text{abs}\left(\frac{|\Delta|}{|R_{11(i)}|}\right).$$

The term in the exponent has the form of a residual sum-of-squares for $\beta$ pooled over all the groups. Forming another orthogonal-triangular decomposition

$$\begin{bmatrix} R_{00(1)} & c_{0(1)} \\ \vdots & \vdots \\ R_{00(M)} & c_{0(M)} \end{bmatrix} = Q_0 \begin{bmatrix} R_{00} & c_0 \\ 0 & c_{-1} \end{bmatrix} \quad (2.17)$$

produces the reduced form

$$L(\beta, \theta, \sigma^2 | y)$$
$$= (2\pi\sigma^2)^{-N/2} \exp\left(\frac{\|c_{-1}\|^2 + \|c_0 - R_{00}\beta\|^2}{-2\sigma^2}\right) \prod_{i=1}^{M} \text{abs}\left(\frac{|\Delta|}{|R_{11(i)}|}\right). \quad (2.18)$$

For a given $\theta$, the values of $\beta$ and $\sigma^2$ that maximize (2.18) are

$$\hat{\beta}(\theta) = R_{00}^{-1} c_0 \quad \text{and} \quad \hat{\sigma}^2(\theta) = \frac{\|c_{-1}\|^2}{N}, \quad (2.19)$$
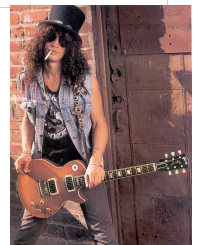
which give the profiled likelihood

$$L(\theta | y) = L\left(\hat{\beta}(\theta), \theta, \hat{\sigma}^2(\theta) | y\right)$$
$$= \left(\frac{N}{2\pi \|c_{-1}\|^2}\right)^{N/2} \exp\left(-\frac{N}{2}\right) \prod_{i=1}^{M} \text{abs}\left(\frac{|\Delta|}{|R_{11(i)}|}\right), \quad (2.20)$$

or the profiled log-likelihood

$$\ell(\theta | y) = \log L(\theta | y)$$
$$= \frac{N}{2}\left[\log N - \log(2\pi) - 1\right] - N \log \|c_{-1}\| + \sum_{i=1}^{M} \log \text{abs}\left(\frac{|\Delta|}{|R_{11(i)}|}\right). \quad (2.21)$$

The profiled log-likelihood (2.21) is maximized with respect to $\theta$, producing the maximum likelihood estimate $\hat{\theta}$. The maximum likelihood estimates $\hat{\beta}$ and $\hat{\sigma}^2$ are then obtained by setting $\theta = \hat{\theta}$ in (2.19).

- it is fixed-effect models that make an assumption:
  - that residual subject and word variances are zero
  - i.e. that word-specific phonology is wrong
- mixed models are agnostic
  - random effects can be zero
  - they do not *assume* a word-specific (or speaker-specific) phonology, they *allow* for it *if it is supported by the data*
- must model speaker/word
  - with random effects, if nested
  - often crossed r. effects for both
- or other results will be wrong
  - maybe not very far wrong?
- as quantitative linguists, we strive for right numbers

Sali Tagliamonte
fellow panelists
Josef Fruehwald
Maryam Bakht
Meghan Armstrong
Kyle Gorman
Kirk Hazen
David Sankoff
Florian Jaeger
Rbrul testers
R developers

Doug Bates `lmer`
Qdoba on Bleecker

Pinheiro, José C. and Douglas M. Bates. 2000. *Mixed-Effects Models in S & S-PLUS.* New York: Springer.

Baayen, R. Harald, Douglas J. Davidson and Douglas M. Bates. 2008. Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language* 59, 390-412.
[I recommend this whole special issue on Emerging Data Analysis.]

Johnson, Daniel Ezra. 2009. Getting off the GoldVarb Standard: introducing Rbrul for mixed-effect variable rule analysis. *Language and Linguistics Compass* 3/1: 359-383.

Rbrul (a work in progress) is at: www.danielezrajohnson.com/Rbrul.R

# Conclusions, thanks, references