

## *basic generalized linear models*

*Ben Bolker*

```
library(ggplot2)
theme_set(theme_bw())
library(ggExtra)
library(cowplot)

##
## Attaching package: 'cowplot'

## The following object is masked from 'package:ggplot2':
##
##      ggsave

library(dotwhisker)
```

### *Linear models*

- foundation for (G)LM(M)s, other complex models
- flexible, robust, computationally efficient, standard
- includes (multiple) regression, ANOVA, ANCOVA, ...
- natural ways to express dependence, interactions

### *Linear models: assumptions*

- response variables:
  - Gaussian (normally distributed)
  - independent
  - *conditionally* homoscedastic (equal variance)
  - univariate
- predictor variables
  - numeric or categorical (nominal)

### *Linear models: math*

$$z = a + bx + cy + \epsilon$$

or (more predictor variables)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$$

or (more flexible distribution syntax)

$$y \sim \text{Normal}(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots, \sigma^2)$$

or (more complex sets of predictors)

$$\mu = \mathbf{X}\mathbf{f}$$

$$y_i \sim \text{Normal}(\mu_i, \sigma^2)$$

what does “linear” mean?

- $y$  is a linear function of the *parameters*  
( $\partial^2 y / \partial^2 \beta_i = 0$ )
- e.g. polynomials:  $y = a + bx + cx^2 + dx^3$
- or sinusoids:  $y = a \sin(x) + b \cos(x)$
- but **not**: power-law ( $ax^b$ ), exponential ( $a \exp(-bx)$ )

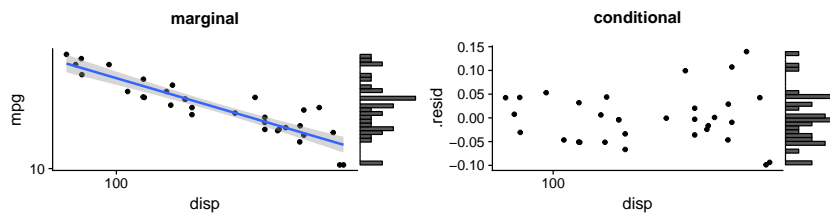
marginal vs. conditional distributions

- common mistake: worry about the overall distribution of the response,  
rather than the *conditional* distribution (i.e., residuals)
- if only categorical predictors, can mean-correct each group, then look at residuals
- otherwise have to fit the model first!

example

MPG vs displacement for cars

```
cars_lm <- lm(log10(mpg) ~ log10(displ), mtcars)
```



(We'll come back to how to judge this later)

categorical predictors

- how do categorical predictors fit into this scheme?
- *dummy variables*: convert to 0/1 values
- R does this automatically with formula syntax
- e.g. for two levels:

```
dd <- data.frame(flavour = rep(c("chocolate",  
  "vanilla"), c(2, 3)))  
print(dd)
```

```
##      flavour  
## 1 chocolate  
## 2 chocolate  
## 3  vanilla  
## 4  vanilla  
## 5  vanilla
```

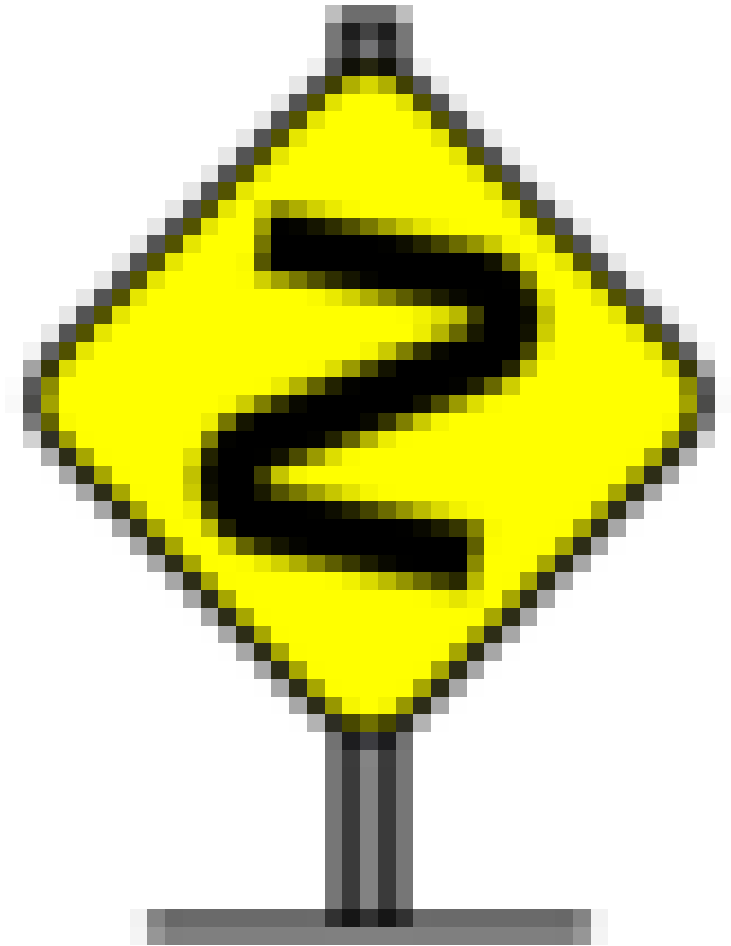
```

model.matrix(~flavour, dd)

## (Intercept) flavourvanilla
## 1          1          0
## 2          1          0
## 3          1          1
## 4          1          1
## 5          1          1
## attr(,"assign")
## [1] 0 1
## attr(,"contrasts")
## attr(,"contrasts")$flavour
## [1] "contr.treatment"

```

- first alphabetical level (chocolate) used as default (use `relevel()` or `factor(..., levels=...)` to change default)



•

*ordered factors* are handled differently

### *R formulas*

- Wilkinson and Rogers (1973)
- `response ~ predictor1 + predictor2 + ...`
- numeric variables used “as is”
- categorical variables (factors) converted to dummy variables
- intercept added automatically (`1+ ...`)
- interaction: `:` *multiplies* relevant columns
- `a*b`: main effect plus interactions
- `model.matrix(formula, data)`

### *Formulas, continued*

- `y~f`: 1-way ANOVA
- `y~f+g`: 2-way ANOVA (additive)
- `y~f*g`: 2-way ANOVA (with interaction)
- `y~x`: univariate regression
- `y~f+x`: ANCOVA (parallel slopes)
- `y~f*x`: ANCOVA (with interaction, non-parallel slopes)
- `y~x1+x2`: multivariate regression (additive)
- `y~x1*x2`: multiv. regression with interaction

If confused, (1) try to write out the equation; (2) `model.matrix()`

### *Contrasts*

- Machinery for translating categorical variables to dummy (0/1) variables
- **treatment** contrasts (default):
  - $\beta_1$  = intercept = expected value of first level (by default, “aardvark”)
  - $\beta_i$  = difference between level  $i + 1$  and baseline
- **sum-to-zero** contrasts:
  - $\beta_1$  = intercept = unweighted mean of all levels
  - $\beta_i$  = difference between level  $i$  and mean; last level not included (!)

too many ways to change contrasts (globally via `options()`; as attribute of factor; contrasts argument in `lm()`)

### *Example 1 (treatment contrasts)*

Data on ant colonies from Gotelli and Ellison (2004):

```

ants <- data.frame(place = rep(c("field", "forest"),
  c(6, 4)), colonies = c(12, 9, 12, 10, 9, 6,
  4, 6, 7, 10))
aggregate(colonies ~ place, data = ants, FUN = mean)

##      place colonies
## 1 field 9.666667
## 2 forest 6.750000

pr <- function(m) printCoefmat(coef(summary(m)),
  digits = 3, signif.stars = FALSE)
pr(lm1 <- lm(colonies ~ place, data = ants))

##              Estimate Std. Error t value
## (Intercept)    9.667      0.958   10.09
## placeforest   -2.917      1.515   -1.92
##              Pr(>|t|)
## (Intercept)    8e-06
## placeforest    0.09

```

*Ants: sum-to-zero contrasts*

```

pr(lm2 <- update(lm1, contrasts = list(place = contr.sum)))

##              Estimate Std. Error t value
## (Intercept)    8.208      0.758   10.83
## place1         1.458      0.758    1.92
##              Pr(>|t|)
## (Intercept)   4.7e-06
## place1        0.09

data(lizards, package = "brglm")

```

*Interactions: example*

- Bear road-crossing
- Predictor variables: sex (categorical: M/F), road type (categorical: major/minor), road length (continuous)
- **Two-way interactions**
  - sex  $\times$  road length: “are females more sensitive to amount of road than males?”
  - sex  $\times$  road type: “do females prefer major over minor roads more than males?”
  - road type  $\times$  road length: “does amount of road affect crossings differently for different road types?”
- **Three-way interaction:** does the difference of the effect of road length between road types differ between sexes?

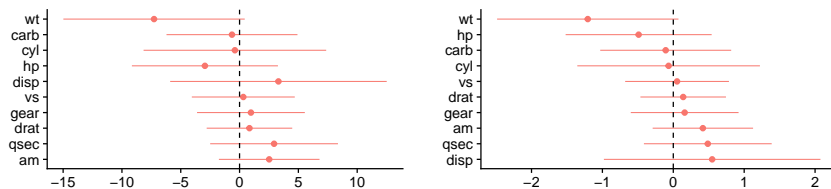
*Centering (Schielzeth 2010)*

- in interaction models, interpretation of main effects **depends on the center-point of the predictors**
- *centering* makes main effects much more interpretable
  - numeric predictors (subtracting the mean by default; other choices could be sensible)
  - categorical predictors: sum-to-zero (weighted or unweighted)
- e.g. if Gregorian year is a predictor, the intercept is at year 0 (!)
- also improves model stability, decorrelates coefficients

*Scaling (Schielzeth 2010)*

- scaling parameters improves interpretability
- standard deviation scaling:  
parameter magnitudes = importance

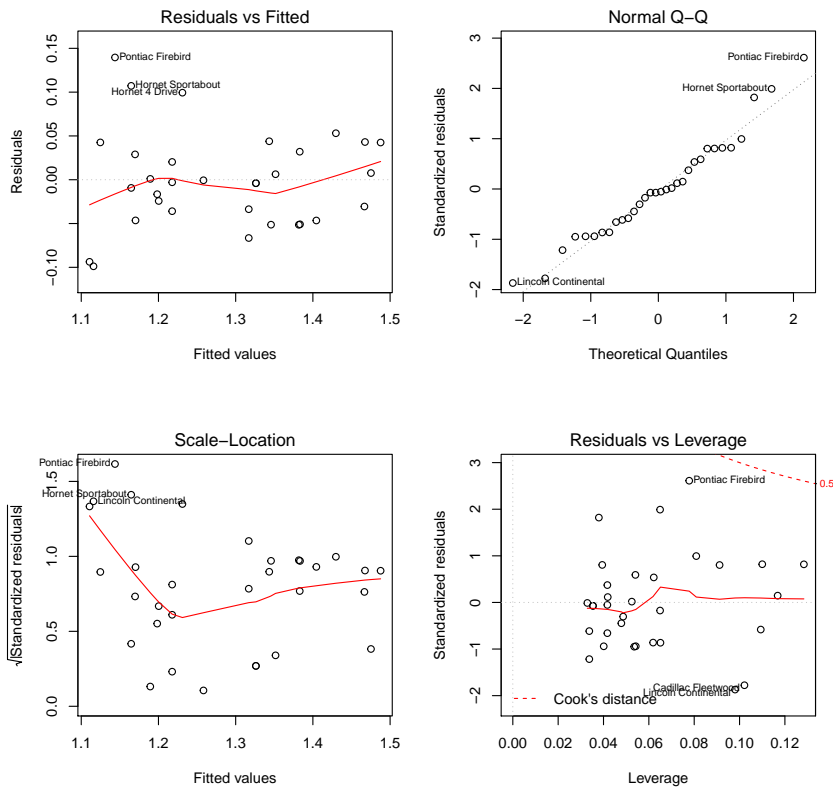
```
mtcars_big <- lm(mpg ~ ., data = mtcars)
mtcars_big_sc <- lm(mpg ~ ., data = as.data.frame(scale(mtcars)))
dwfun <- function(.) {
  dwplot(., order_vars = names(sort(coef(.)))) +
    geom_vline(xintercept = 0, linetype = 2)
}
plot_grid(dwfun(mtcars_big), dwfun(mtcars_big_sc))
```

*LM diagnostics*

- fitted vs. residual: pattern in mean? (linearity)
- scale-location: pattern in variance? (homoscedasticity)
- Q-Q plot: Normality of **residuals**
- leverage/Cook's distance: influential points?
- independence is often hard to test
- Normality is the **least important** of these assumptions

*LM diagnostics*

```
par(mfrow = c(2, 2))
plot(cars_lm)
```



- smooth lines help interpretation
- highlighted points are 3 most extreme (`id.n` argument)

### Diagnostics

- statisticians: “don’t use p-values to evaluate LM assumptions”
- everyone else: “so what should I do?”
- statisticians: “look at pictures”
- everyone else: “how do I decide whether to worry?”
- statisticians: “...”

### testing hypotheses and interpreting results

- parameter-by-parameter: `summary()` (*t* test)
- multi-parameter comparisons: `anova()`, `car::Anova()` (*F* test)
- order matters
- interactions/main effects matter

### From LM to GLM

#### Why GLMs?

- assumptions of LMs do break down sometimes

- count data: discrete, non-negative
- proportion data: discrete counts,  $0 \leq x \leq N$

<https://twitter.com/thedavidpowell/status/984432764215754753>

### *GLMs in action*

- vast majority of GLMs
  - *logistic regression* (binary/Bernoulli data)
  - *Poisson regression* (count data)
- lots of GLM theory carries over from LMs
  - formulas
  - parameter interpretation (partly)
  - diagnostics (partly)

### *link functions*

- transform *prediction*, not response
- e.g. rather than  $\log(\mu) = \beta_0 + \beta_1 x$ , use  $\mu = \exp(\beta_0 + \beta_1 x)$
- in this case log is the **link function**, exp is the **inverse link function**
- extreme observations don't cause problems (usually)

### *families*

- link function plus variance function
- typical defaults
  - Poisson: log (exponential)
  - binomial: logit/log-odds (logistic)

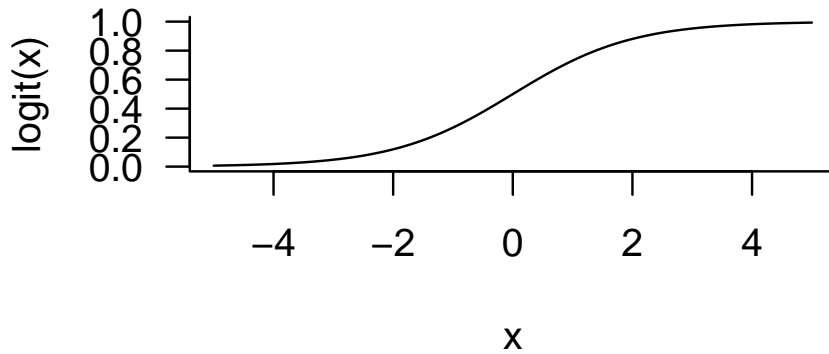
### *log link*

- proportional scaling of effects
- small values of coefficients ( $< 0.1$ )  $\approx$  proportionality
- otherwise change per unit is  $\exp(\beta)$
- large parameter values ( $> 10$ ) mean some kind of trouble

### *logit link*

```
par(las = 1, bty = "l")
curve(plogis(x), from = -5, to = 5, ylab = "logit(x)")
```





- `qlogis()` function (`plogis()` is logistic/inverse-link)
- *log-odds* ( $\log(p/(1-p))$ )
- most natural scale for probability calculations
- interpretation depends on *base probability*
  - small probability: like  $\log$  (proportional)
  - large probability: like  $\log(1-p)$
  - intermediate ( $0.3 < p < 0.7$ ): effect  $\approx \beta/4$

#### *back-transformation*

- confidence intervals are symmetric on link scale
- can back-transform estimates and CIs for  $\log$
- logit is hard (must pick a reference level)
- don't back-transform standard errors!

#### *estimation*

- iteratively re-weighted least-squares
- usually Just Works

#### *inference*

like LMs, but:

- one-parameter tests are usually  $Z$  rather than  $t$
- CIs based on standard errors are approximate (Wald)
- `confint.glm()` computes *likelihood profile* CIs

#### *References*

Gotelli, Nicholas J., and Aaron M. Ellison. 2004. *A Primer of Ecological Statistics*. Sunderland, MA: Sinauer.

Schielzeth, Holger. 2010. "Simple Means to Improve the Interpretability of Regression Coefficients." *Methods in Ecology and Evolution* 1: 103–13. doi:10.1111/j.2041-210X.2010.00012.x.

Wilkinson, G. N., and C. E. Rogers. 1973. "Symbolic Description

of Factorial Models for Analysis of Variance." *Applied Statistics* 22 (3): 392–99. doi:10.2307/2346786.