

# Generalised Linear Models

## for data of multiple species

Bert van der Veen

Department of Mathematical Sciences, NTNU

# Outline

---

- ▶ Data collection and common data types
- ▶ Generalised Linear Models background
  - ▶ Assumption checking
- ▶ “Vector” models
- ▶ Building on material from [the GLM workshop](#)

## The ecological process

---

What do we know of the processes that generate these data?

- ▶ Meta-community theory
- ▶ Assembly processes (filtering)
- ▶ Ecological gradient theory

# The ecological process

---

What do we know of the processes that generate these data?

- ▶ Meta-community theory
- ▶ Assembly processes (filtering)
- ▶ Ecological gradient theory

Multispecies models provide a statistical connection to these ecological frameworks. We do not just use a fancy tool, we use a fancy tool because we believe it aligns well with our understanding of the ecological process.

# The ecological process (2)

---

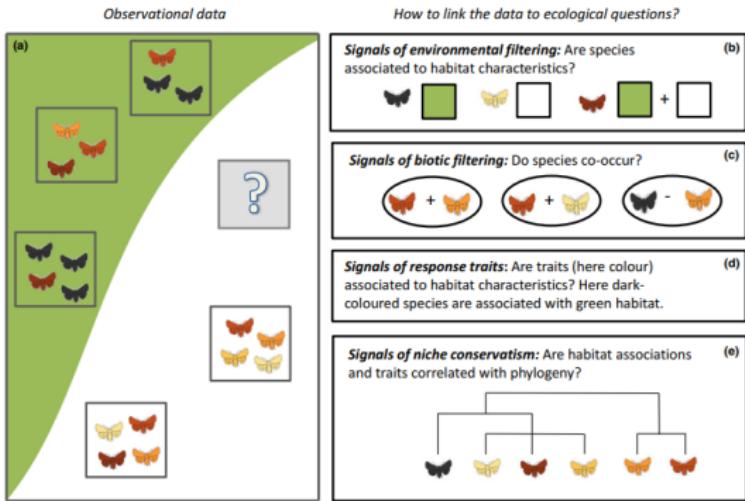


Figure 2 A conceptual illustration of some key questions in community ecology. The green and white colours represent differences in the environmental

Figure 1: Figure 2 from Ovaskainen et al. (2017)

## On ecological communities

---

The concept of an ecological community is of limited use. By definition:

**An ecological community is a group or association of two or more species occupying the same geographical area at the same time**

- ▶ We often think of ecological communities as groups
- ▶ We can also think of a community as a continuum that changes along a gradient (Austun 1985)
- ▶ We can also think of them as the species in our data

## On ecological communities

---

The concept of an ecological community is of limited use. By definition:

**An ecological community is a group or association of two or more species occupying the same geographical area at the same time**

- ▶ We often think of ecological communities as groups
- ▶ We can also think of a community as a continuum that changes along a gradient (Austun 1985)
- ▶ We can also think of them as the species in our data

Connecting model outputs to ecological concepts requires some deep thoughts

## An ecological gradient

### VEGETATION PATTERNS IN THE GREAT SMOKY MOUNTAINS

Change of vegetation along the moisture gradient at lower and higher elevations

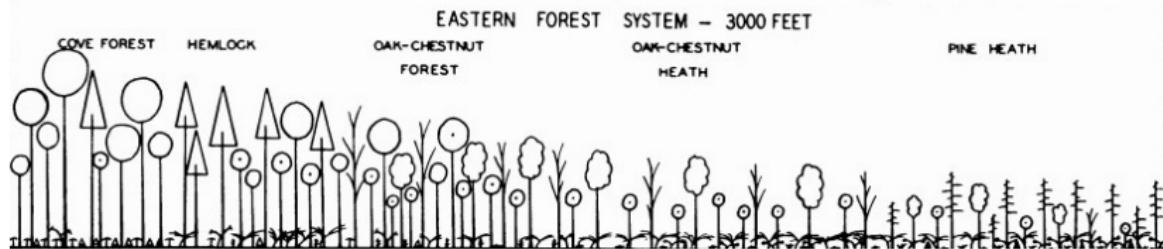


Figure 2: Whittaker (1956)

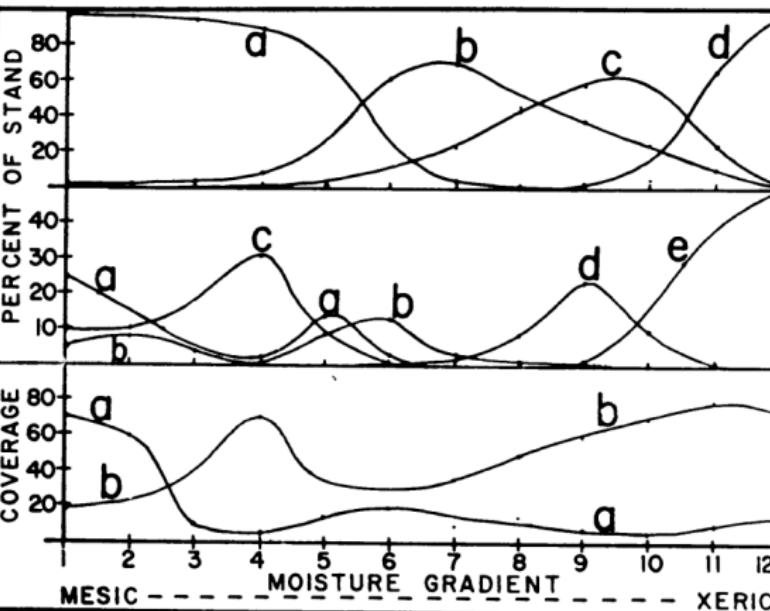


FIG. 4. Transect of the moisture gradient, 3500-4500 ft. Top—curves for tree classes; a, mesic; b, submesic; c, subxeric; d, xeric. Note expansion of mesic stands, compared with Figs. 2 and 3. Middle—curves for tree

## Multispecies models

---

There are multiple statistical frameworks for studying the processes:

- ▶ Generalised Linear Models
- ▶ Generalised Linear Mixed-effects Models
- ▶ Generalised Additive Models (not covered here)
- ▶ Generalised Linear Latent Variable Models

## Multispecies models

---

There are multiple statistical frameworks for studying the processes:

- ▶ Generalised Linear Models
- ▶ Generalised Linear Mixed-effects Models
- ▶ Generalised Additive Models (not covered here)
- ▶ Generalised Linear Latent Variable Models

Or ecologically:

- ▶ Species distribution models
- ▶ Joint Species Distribution Models
- ▶ Ordination

## Multispecies models

---

There are multiple statistical frameworks for studying the processes:

- ▶ Generalised Linear Models
- ▶ Generalised Linear Mixed-effects Models
- ▶ Generalised Additive Models (not covered here)
- ▶ Generalised Linear Latent Variable Models

Or ecologically:

- ▶ Species distribution models
- ▶ Joint Species Distribution Models
- ▶ Ordination

and more. Each method has its limitations (assumptions). It is up to us to assess which are appropriate.

## Generalised linear models (GLMs)

---

GLMs as a framework were introduced by Nelder and Wedderburn (1972) uniting many different models. With a special focus on teaching statistics.

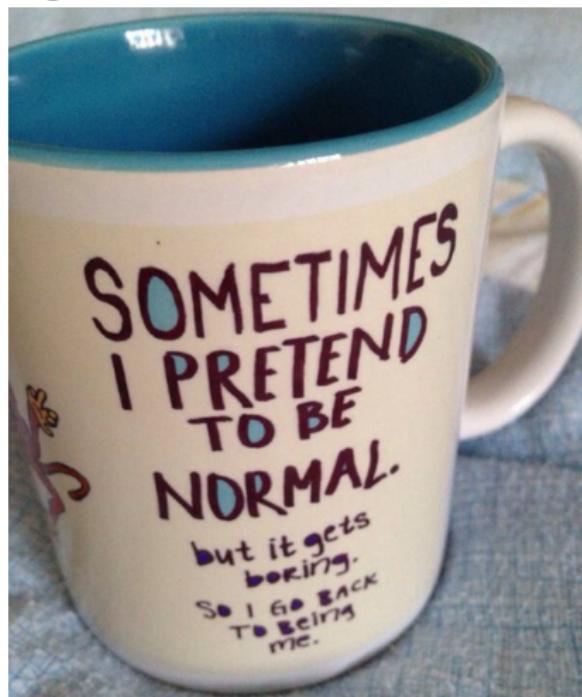
- ▶ Linear regression
- ▶ Logistic regression
- ▶ Probit regression
- ▶ Complementary log-log regression
- ▶ Log-linear regression
- ▶ Gamma regression

# Generalised Linear Models

---

For when the assumptions of linear regression fail.

- ▶ Linearity (straight line)
- ▶ Independence of errors
- ▶ Homoscedasticity (same variance for all errors)
- ▶ Normality (distribution of errors)



## Generalised linear models (2)

GLMs extend the linear model framework to address:

- ▶ Variance changes with the mean
- ▶ Range of  $y$  is bounded



**The basis of many statistical models in Biology**

## Components of a GLM

---

- ▶ Systematic component:  $\eta$
- ▶ Random component: data/distribution
- ▶ The link function: connects these components
  - ▶ This is not a data transformation
- ▶ The variance function

**But no explicit error term**

## GLM Likelihood

---

- ▶ We use MLE for estimation
- ▶ With a distribution in the “exponential family” (for fixed  $\phi$ )

All GLMs have the likelihood:

$$\mathcal{L}(y_i; \Theta) = \exp\left\{\frac{y_i \eta_i - b(\eta_i)}{a(\phi)} + c(y_i, \phi)\right\} \quad (1)$$

## Generalised linear model

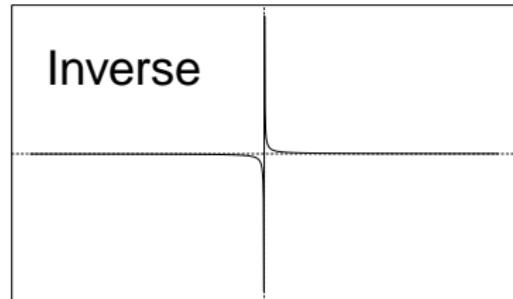
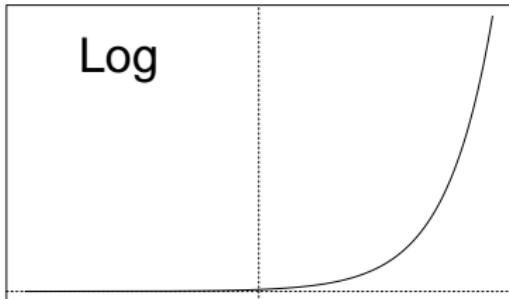
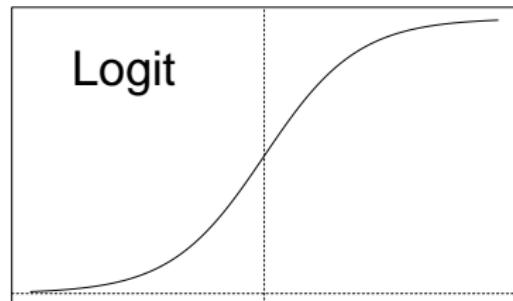
---

$$\begin{aligned} g\{\mathbb{E}(y_i|x_i)\} &= \eta_i = \alpha + x_i\beta \\ \mathbb{E}(y_i|x_i) &= g^{-1}(\eta_i) = g^{-1}(\alpha + x_i\beta) \end{aligned} \tag{2}$$

$g(\cdot)$  is the **link function**

## The link function

- ▶ Is a smooth/monotone function
- ▶ Has an inverse  $g^{-1}(\cdot)$
- ▶ Restricts the scale
- ▶  $g(\cdot)$  can be e.g.



## Variance function

---

Perhaps most critically, variance changes with the mean:

$$\text{var}(y_i; \mu_i, \phi) = \frac{\partial^2 b(\eta_i)}{\partial \eta_i^2} a(\phi)$$

- ▶  $\phi$ : the dispersion parameter, constant over observations
  - ▶ Fixed for some response distributions
- ▶  $a(\phi)$  is a function of the form  $\phi/w_i$  (McCullagh and Nelder 1989)

# Fitting GLMs

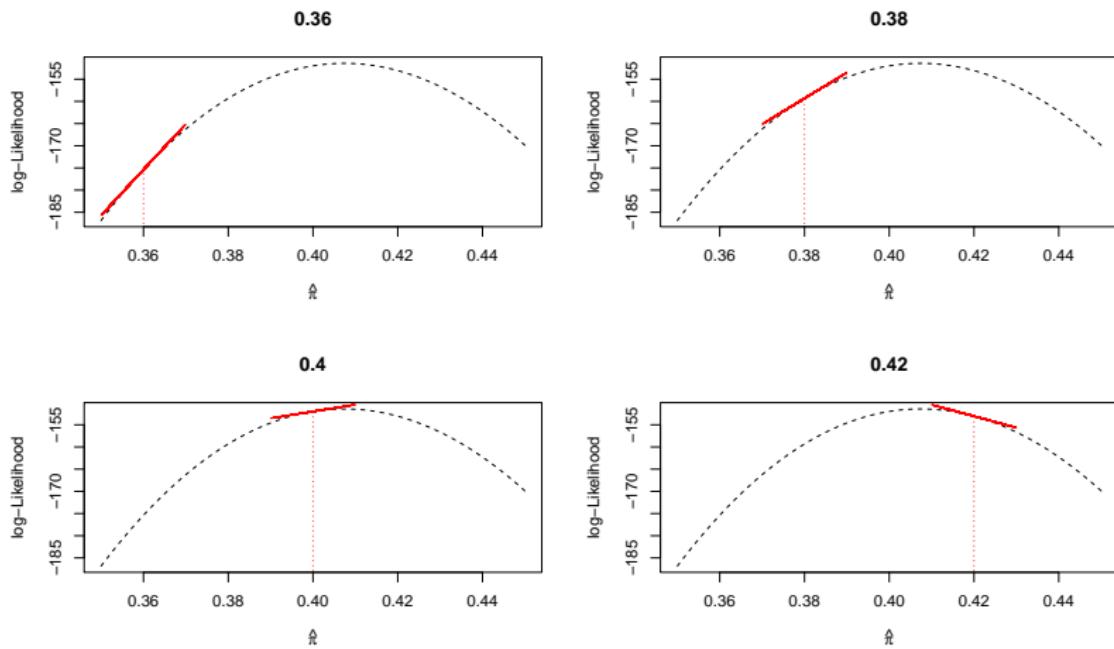
---

Parameters in GLMs need to be estimated **iteratively**.

- ▶ More difficult to fit
- ▶ Requires numerical optimisation
- ▶ Susceptible to local convergence

Holds for GLLVMs too

# Estimating GLMs



We need a good algorithm to find the maximum!

## Why is this important?

---

- 1) A basic (mathematical) understanding helps apply methods correctly.
- 2) GLMMs/GLVMs may not always converge to the best solution immediately.
- 3) This can help to diagnose your model.

## Often used distributions in ecology

---

- ▶ Binomial: occurrence/counts. Presence of species, number of germinated seeds out of a total
- ▶ Poisson: counts. Abundance
- ▶ Negative binomial (fixed dispersion): counts. Number of species or abundance
- ▶ Gamma: (positive) continuous. Body size or biomass
- ▶ Ordinal (cumulative link). Cover classes
- ▶ Beta (logit link). Cover (note: not a GLM)

## Example: Swiss bird occurrence

---

**Observation** process: see if a bird is present (or we might hear it)

**Alternatively:** The proportion of a species in a place

**Alternatively:** Count of birds in the forest

There are often many ways to observe the same ecological process.

We need to **disentangle** this from the ecological process.

# Example: Swiss birds

- ▶ Data by Schmid et al. (1998): the Swiss breeding bird atlas
- ▶ Occurrence of 56 species at 2524 locations recorded over a 4-year period

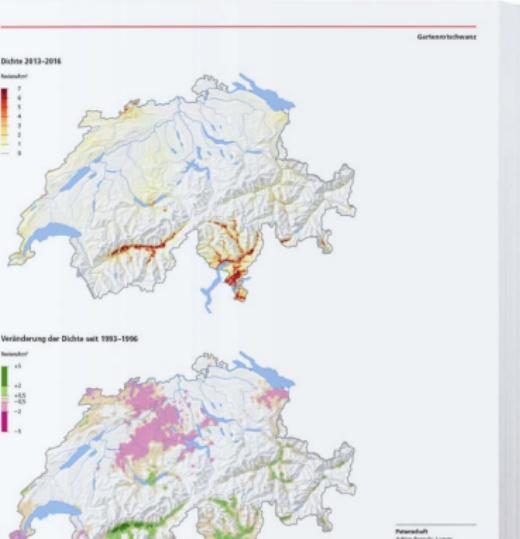


Der Gartenrotschwanz kommt in der ganzen Schweiz vor, oft aber nur in geringer Dichte. Am häufigsten ist er im Süden und in den grossen Tälern des Tessins, im Mäuse GR und im Bergell GR, im Mittel- und Oberrheintal sowie in der Region Basel. Rund 90 % des Bestands leben unter 1500 m. Durch die Zerstörung von natürlichen Lebensräumen, nötiger Böden und Höhlen verhindert wird. Das können eher ländliche Wohngegenden, Parks, Familiengärten, Rebberge mit Blumen, Hochstamm-Obirgäste oder Klettermauern sein. Der Bestand ist in den Jahren 1993 bis 2013 um 45 % gestiegen, von 1993 bis 2016 um 51 %. In den Jahren 2013–2016, bis auf 2250m bei S-Card GR (M. Ernst beobachtet), Der höchste Bruthabungsstandort von 2220m bei Zürich VSM<sup>16</sup>. In La Chaux-de-Fonds ist mit 1700m der höchste Standort. Im Jura zwischen 850 und 950m liegt die maximale Dichte bei 6,1 Bevölkerungen/km<sup>2</sup>.

Für den Schweizer Bestand waren die Dichten der Schätzgruppen in der Satzzeit eingeschleudet<sup>17</sup>. Der Rückgang in der Satzzeit ist mit einem Abfall der Brutrate zu erklären. Seit 2002 zeigt sich eine leichte Erholung. Regrettabel unterscheiden sich die Trends jedoch stark. In den ländlichen Gebieten nördlich der Alpen dominieren die bereits gemeldeten Bestände wieder aus. Die Verluste betreffen im Kanton

Zürich von 1993 bis 2006 851 km<sup>2</sup>, am Bodensee von 1993 bis 2010 und 2013–2016 und im Kanton Basel-Landschaft von 1993 bis 2013 3,4 % und 2013–2016. Den starken Einbußen in Höhen von 300 bis 900 m, die eine Folge des Verschwindens von Hochstamm-Obirgästen, intensiver Grünlanddienstes und der Zerstörung von Naturhaushalten für die Ansiedlung sind<sup>18</sup>, steht ein geringer Bestandeszuwachs in Lagen über 1000m gegenüber. Für die bisher von Erbächen eher verschonten Bestände in trockener bebauter Umgebung ist die Zerstörung durch den Menschen eine zusätzliche Gefahr darstellen<sup>19</sup>. Im Wallis und im Tessin nehmen die Bestände selbst in tiefen Lagen zu<sup>20</sup> und tragen damit wesentlich zum gesamten positiven Landestrend bei. Die Zunahme auf der 3,1 km<sup>2</sup> grossen Walliserfläche bei Leuk VS<sup>21</sup> sorgen die enorme Bestandsgeschwindigkeit dieser Art.

In Italien und Frankreich sind die Trends stark positiiv<sup>22,23</sup>, während in Spanien und Portugal die Bestände regressive Verluste<sup>24</sup> annehmen stabil<sup>25,26</sup>. Die europäische Tendenz ist von 1980 bis 2014 leicht positiv<sup>27,28</sup>.



# The data

---

Falco_subbuteo	Anthus_trivialis	Phylloscopus_bonelli	Tetrao_tetrix	Parus_caeruleus	Dendrocopos_major
0	1	0	1	0	1
1	0	0	0	1	1
0	0	0	0	1	1
0	1	0	0	1	1
0	0	1	1	0	0
1	0	0	0	1	1

# The environmental variables

avg	cov	cv	dns	fhd	p10	p25	p95	rt_p2595	std	uhd
13.549999	34.1	0.6287823	34.1	1.8422778	2.68	5.950000	28.14	0.2114428	8.52	2.381401
13.790000	32.0	0.5547498	32.0	1.8031981	3.96	7.090000	28.20	0.2514184	7.65	2.276368
19.340000	14.3	0.5351603	14.3	2.0021141	4.79	9.139999	34.03	0.2685865	10.35	2.364149
15.460000	34.5	0.4618370	34.5	1.7484871	5.19	9.790000	26.36	0.3713961	7.14	2.369797
2.290000	2.0	0.7292576	2.0	0.3038808	1.06	1.200000	5.74	0.2090592	1.67	1.270930
8.929999	61.3	0.5890258	61.3	1.4213525	2.21	4.570000	18.05	0.2531856	5.26	2.392563

- ▶ Bioclimatic variables (bioclim)
- ▶ Topography (slope, aspect, TPI, TWI) from a DEM
- ▶ Potential evapotranspiration (PET) from solar radiation
- ▶ Moisture index, degree days above zero
- ▶ Vegetation structure from LiDAR

## The binomial GLM

---

$$p(y_{ij} = 1) = p_{ij} = g^{-1}(\eta_{ij}) \quad (3)$$

# The binomial GLM

---

Link functions:

- ▶ Logit:  $\log(\frac{\pi_i}{1-\pi_i})$  and inverse  $\frac{\exp(\eta_i)}{1+\exp(\eta_i)}$  - *the canonical link*
- ▶ Probit:  $\Phi^{-1}(\pi_i)$  and inverse  $\Phi(\eta_i)$
- ▶ Complementary log-log:  $\log(-\log(1-\pi_i))$  and inverse  $1 - \exp(-\exp(\eta_i))$
- ▶ Log-log
- ▶ Logit is canonical and easier to interpret
- ▶ Probit is sometimes easier mathematically than Logit
- ▶ Complementary log-log for counts

## Data format

---

There are two ways to format these data:

**Wide format:** Species as columns (as presented)

**Long format:** Species is one column, and “Site” is another column

The format of the data does not affect the model. Some functions accept long format, other wide format, but the formulation of the model is up to us.

## Swiss birds: to long format

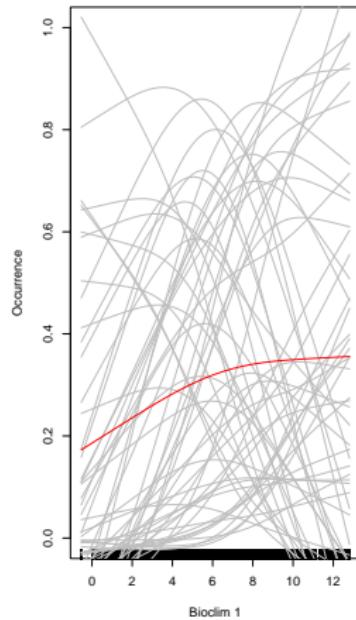
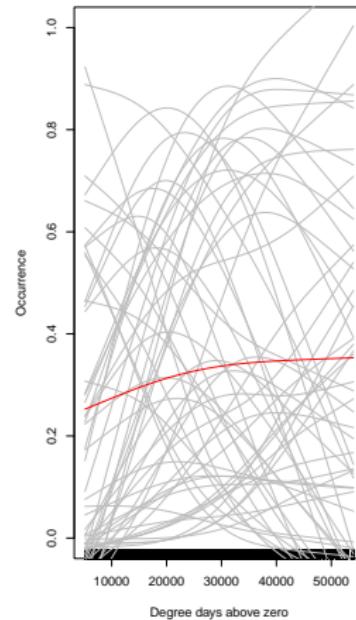
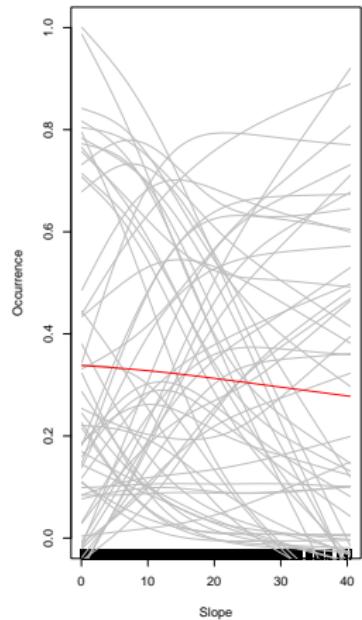
---

```
data <- data.frame(y, X)
datalong <- reshape(data,
                     varying = colnames(y),
                     v.names = "occ",
                     idvar = "Site",
                     timevar = "Species",
                     direction = "long")

datalong$Species <- factor(datalong$Species,
                            labels = colnames(y))
```

# Swiss birds: visually inspect the data

---



Is there a community level trend?

## Swiss birds: fit a model

---

```
model1 <- glm(occ~slp,  
               data = datalong, family="binomial")  
coef(model1)
```

```
## (Intercept) slp  
## -0.6589784 -0.0065449
```

Here we assume that the intercept and slp effect are the same for all species

## Multispecies modeling

---

- 1) Is the same effect for all species realistic?
- 2) Is the same (average) probability of occurrence for all species realistic?

## Multispecies modeling

---

- 1) Is the same effect for all species realistic?
- 2) Is the same (average) probability of occurrence for all species realistic?
- 3) We usually assume that species have their own preferred environmental conditions
- 4) Some species might still like similar conditions; there is a common component
- 5) We can separate this out with GLMMs or with a “sum-to-zero” contrast

## Swiss birds: species-specific effects

---

```
model2 <- glm(occ~slp*Species,  
                data = datalong, family="binomial")
```

- ▶ One intercept per species
- ▶ One slp effect per species
- ▶ But all are relative to the first species

## Swiss birds: species-specific effects

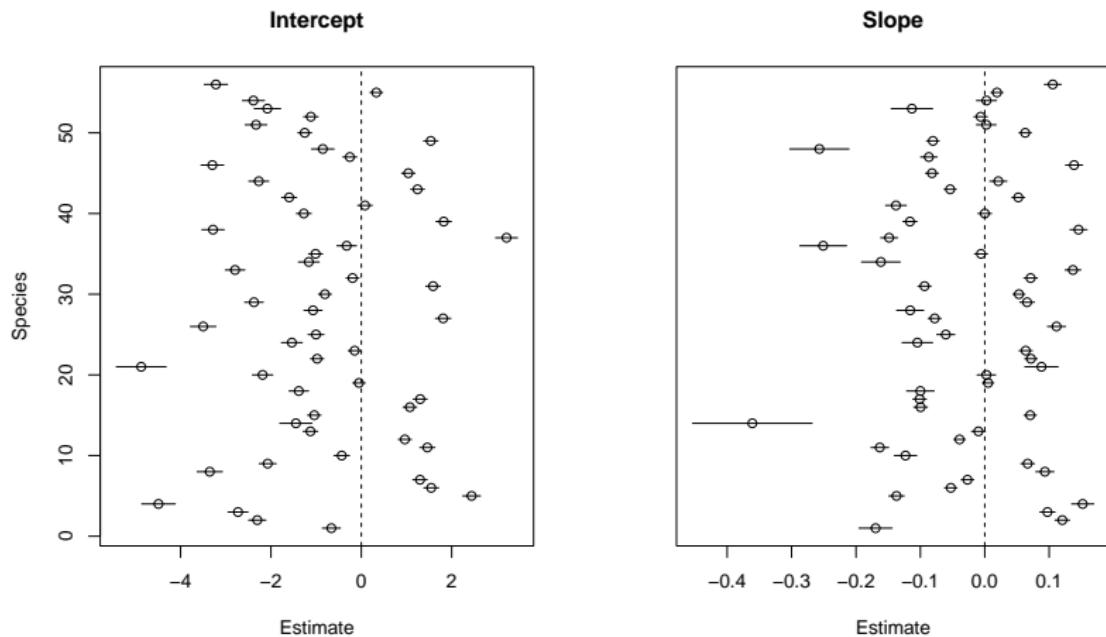
---

```
model3 <- glm(occ ~ 0 + Species + slp:Species,  
               data = datalong, family="binomial")
```

The same model, but a bit easier to interpret

- ▶ One intercept per species
- ▶ One slp effect per species
- ▶ Not relative to each other (prevents post-hoc processing of tests and CI)

## Swiss birds: results



## Interpreting Binomial GLM coefficients

---

- Below one we are more likely to not observe the species

## Interpreting Binomial GLM coefficients

---

- ▶ Below one we are more likely to not observe the species
- ▶ The likelihood of observing the species 1 at flat ground is  $\exp(-0.663)/(1+\exp(-0.663)) = 0.34$  (given by the intercept)

## Interpreting Binomial GLM coefficients

---

- ▶ Below one we are more likely to not observe the species
- ▶ The likelihood of observing the species 1 at flat ground is  $\exp(-0.663)/(1+\exp(-0.663)) = 0.34$  (given by the intercept)
  - ▶ At 0.5 there is equal probability to observe and not

## Interpreting Binomial GLM coefficients

---

- ▶ Below one we are more likely to not observe the species
- ▶ The likelihood of observing the species 1 at flat ground is  $\exp(-0.663)/(1+\exp(-0.663)) = 0.34$  (given by the intercept)
  - ▶ At 0.5 there is equal probability to observe and not
- ▶ Odds ratio for the first species at slp 0 is  $\exp(-0.663) = 0.52:1$

## Interpreting Binomial GLM coefficients

---

- ▶ Below one we are more likely to not observe the species
- ▶ The likelihood of observing the species 1 at flat ground is  $\exp(-0.663)/(1+\exp(-0.663)) = 0.34$  (given by the intercept)
  - ▶ At 0.5 there is equal probability to observe and not
- ▶ Odds ratio for the first species at slp 0 is  $\exp(-0.663) = 0.52:1$ 
  - ▶ We are two times more likely to not observe the species on flat ground, than to observe it

## Interpreting Binomial GLM coefficients

---

- ▶ Below one we are more likely to not observe the species
- ▶ The likelihood of observing the species 1 at flat ground is  $\exp(-0.663)/(1+\exp(-0.663)) = 0.34$  (given by the intercept)
  - ▶ At 0.5 there is equal probability to observe and not
- ▶ Odds ratio for the first species at slp 0 is  $\exp(-0.663) = 0.52:1$ 
  - ▶ We are two times more likely to not observe the species on flat ground, than to observe it
- ▶ This decreases by  $\exp(-0.17)$  for every unit of slp $0.52*\exp(-0.17) = 0.52*0.84 = 0.44$

## Interpreting Binomial GLM coefficients

---

- ▶ Below one we are more likely to not observe the species
- ▶ The likelihood of observing the species 1 at flat ground is  $\exp(-0.663)/(1+\exp(-0.663)) = 0.34$  (given by the intercept)
  - ▶ At 0.5 there is equal probability to observe and not
- ▶ Odds ratio for the first species at slp 0 is  $\exp(-0.663) = 0.52:1$ 
  - ▶ We are two times more likely to not observe the species on flat ground, than to observe it
- ▶ This decreases by  $\exp(-0.17)$  for every unit of slp $0.52*\exp(-0.17) = 0.52*0.84 = 0.44$

## Contrasts

There are other “contrast” treatments in R than “dummy”

- ▶ We can instead use “sum-to-zero” contrasts
  - ▶ If the sum is zero, the mean must be too
  - ▶ The coefficient of the last species is set to the negative sum

```
(contr <- contr.sum(levels(datalong$Species)))
```

```
##                                     [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]
## Falco_subbuteo                  1    0    0    0    0    0    0    0    0
## Anthus_trivialis                0    1    0    0    0    0    0    0    0
## Phylloscopus_bonelli              0    0    1    0    0    0    0    0    0
## Tetrao_tetrix                   0    0    0    1    0    0    0    0    0
## Parus_caeruleus                 0    0    0    0    1    0    0    0    0
## Dendrocopos_major                 0    0    0    0    0    1    0    0    0
## Garrulus_glandarius              0    0    0    0    0    0    0    1    0
## Carduelis_spinus                  0    0    0    0    0    0    0    0    1
## Loxia_curvirostra                 0    0    0    0    0    0    0    0    1
## Phylloscopus_trochilus            0    0    0    0    0    0    0    0    0
## Certhia_brachydactyla             0    0    0    0    0    0    0    0    0
## Sylvia_borin                      0    0    0    0    0    0    0    0    0
## Phoenicurus_phoenicurus            0    0    0    0    0    0    0    0    0
```

## Swiss birds: species-specific responses with common effect

```
model4 <- glm(occ~0+slp+Species + slp:Species, data = datalong, family = "binomial",
               contrasts = list(Species = contr))
coef(model4)[1]
```

```
##          slp
## -0.02703675
```

- ▶ One intercept per species
- ▶ One slp effect that is the same for all species (the mean of effects)
- ▶ One slp effect per species, relative to the common effect

The benefit: the average effect gets a statistical test.

By design corresponds to the result from our previous model:

```
mean(coef(model3)[-c(1:ncol(y))]) = -0.0270367
```

## Swiss birds

---

The three models have the same number of parameters, but are just differently parameterized. So, their log-likelihoods are the same:

```
c(logLik(model2), logLik(model3), logLik(model4))
```

```
## [1] -66008.4 -66008.4 -66008.4
```

## Swiss birds: species-specific responses with common effect

---

```
##  
## Call:  
## glm(formula = occ ~ 0 + slp + Species + slp:Species, family = "binomial",  
##       data = datalong, contrasts = list(Species = contr))  
##  
## Coefficients:  
##  
##                Estimate Std. Error z value Pr(>|z|)  
## slp                 -0.0270367  0.0013931 -19.408 < 2e-16 ***  
## SpeciesFalco_subbuteo      -0.6629549  0.1024622  -6.470 9.78e-11 ***  
## SpeciesAnthus_trivialis     -2.3028109  0.0972711 -23.674 < 2e-16 ***  
## SpeciesPhylloscopus_bonelli    -2.7295986  0.1127350 -24.213 < 2e-16 ***  
## SpeciesTetrao_tetrix          -4.4890112  0.1896457 -23.671 < 2e-16 ***  
## SpeciesParus_caeruleus        2.4406078  0.1003095  24.331 < 2e-16 ***  
## SpeciesDendrocopos_major        1.5471247  0.0838875  18.443 < 2e-16 ***  
## SpeciesGarrulus_glandarius      1.2992201  0.0813667  15.967 < 2e-16 ***  
## SpeciesCarduelis_spinus         -3.3544392  0.1426728 -23.511 < 2e-16 ***  
## SpeciesLoxia_curvirostra       -2.0729215  0.0955422 -21.696 < 2e-16 ***  
## SpeciesPhylloscopus_trochilus    -0.4337227  0.0880349 -4.927 8.36e-07 ***  
## SpeciesCerthia_brachyactyla      1.4614529  0.0852668  17.140 < 2e-16 ***  
## SpeciesSylvia_borin             0.9639614  0.0750798  12.839 < 2e-16 ***  
## SpeciesPhoenicurus_phoenicurus     -1.1252670  0.0827956 -13.591 < 2e-16 ***  
## SpeciesHippolais_icterina        -1.4481047  0.1822145 -7.947 1.91e-15 ***  
## SpeciesPyrrhula_pyrrhula          -1.0409062  0.0757068 -13.749 < 2e-16 ***  
## SpeciesEmberiza_citrinella        1.0753753  0.0771653  13.936 < 2e-16 ***  
## SpeciesMuscicapa_striata          1.3028822  0.0793269  16.424 < 2e-16 ***  
## SpeciesPicus_canus              -1.3840458  0.1104130 -12.535 < 2e-16 ***  
## SpeciesPicus_viridis              -0.0517754  0.0699471 -0.740 0.459174  
## SpeciesAccipiter_gentilis           -2.1845813  0.1151459 -18.972 < 2e-16 ***
```

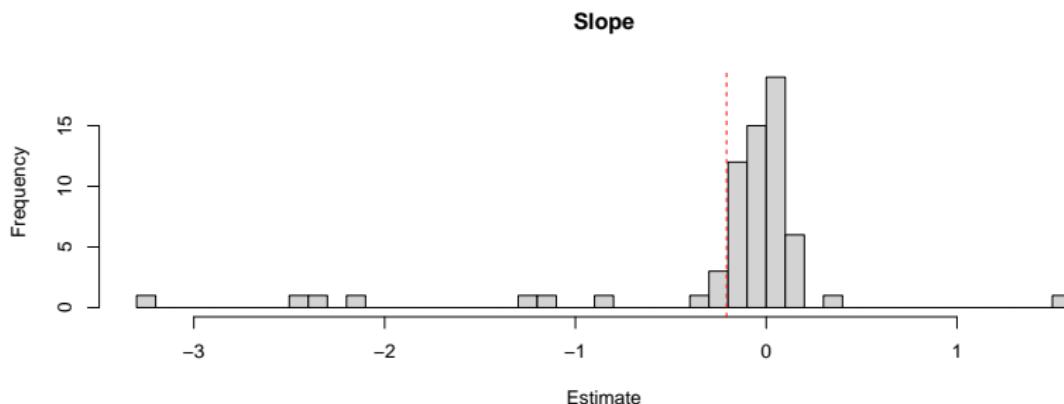
## Swiss birds: conclusions

---

We can conclude that fewer bird species occur in steep places

Most species are more negatively affected than the average

Some species are positively affected by slope, but most negatively



## Interpreting the coefficients

---

Or with predict:

```
predict(model3, newdata =  
       data.frame(Species = factor("Falco_subbuteo", levels = colnames(y)), slp = 1),  
       type = "response")  
  
##           1  
## 0.3031163
```

## Example: macroinvertebrate counts in USA desert

---

**Observation** process: count of macroinvertebrates in three “dips”

**Alternatively:** The proportion of a species in a dip

**Alternatively:** Was this species found in the dip

## Example: macroinvertebrate counts in USA desert

---

- ▶ Data by Pina and Lougheed 2022
- ▶ Counts of 14 species, in 2018 and 2019, in 14 wetlands
- ▶ Main goal: assess impacts of water quality on macroinvertebrates



# The abundance data

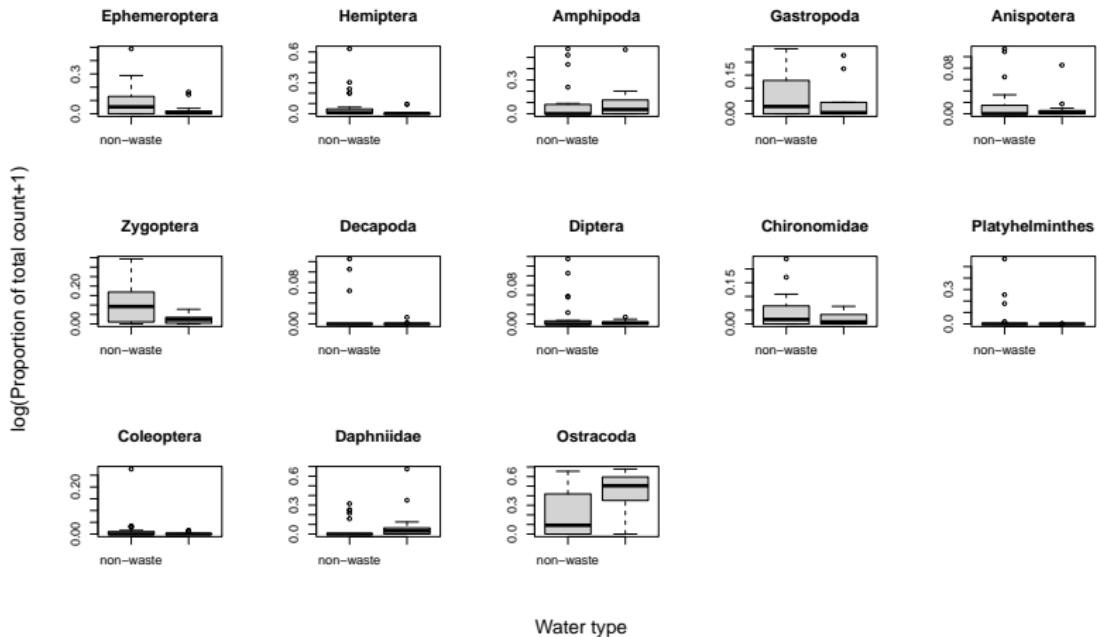
Ephemeroptera	Hemiptera	Amphipoda	Gastropoda	Anisoptera	Zygoptera	Decapoda	Diptera	Chiron
0	1	100	8	0	1	12	0	0
21	0	0	5	3	10	0	0	1
11	3	5	0	4	32	0	0	0
0	1	11	0	0	0	0	0	0
80	1	60	25	0	7	0	0	2
9	15	0	6	0	10	0	0	1
30	0	0	10	7	25	0	0	1
10	1	60	190	4	60	0	0	0
6	0	0	5	0	3	2	2	2
5	41	0	5	0	2	0	0	0
32	2	0	0	1	35	0	0	0
6	25	0	31	2	13	0	0	0
28	15	400	0	21	14	0	0	0
28	15	400	0	21	14	0	0	0
6	150	200	70	2	48	3	4	4
6	150	200	70	2	48	3	4	4
9	1	0	35	14	50	0	0	7
26	1	0	1	10	21	0	0	12
0	1	0	2	2	7	0	0	0
1	0	60	6	1	9	0	0	0
13	0	15	26	0	8	0	0	1
87	4	0	0	0	3	0	0	0
2	1	11	1	2	1	0	0	1

# The environment data

Year	Hydro	Water_Type	Conductivity	DOC	TDN	Turbidity	Alkalinity	Total_CHL	Correc
2018	permanent	non-waste	4.060	2.846	0.306	4.40	63.050	2.231	
2018	permanent	waste	2.582	23.160	3.544	60.73	412.400	95.211	
2018	permanent	non-waste	8.563	28.120	2.450	17.40	363.708	16.915	
2018	permanent	non-waste	15.710	75.040	7.160	38.50	457.625	26.160	
2018	ephemeral	non-waste	1.029	4.012	0.386	3.78	198.042	12.657	
2018	ephemeral	non-waste	1.204	5.356	0.491	24.70	168.042	30.353	

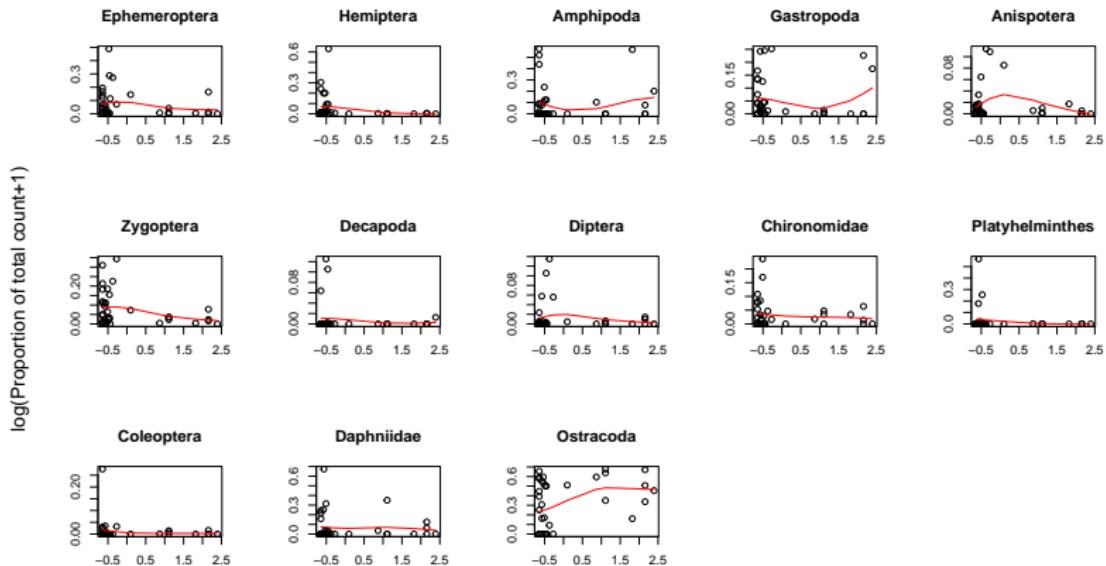
- ▶ 11 environmental variables
  - ▶ Water chemistry
  - ▶ Water type
  - ▶ Presence of hydro power
  - ▶ Permanent or temporary wetland

## Visually inspect the data: categorical covariate



Is there a common effect?

## Visually inspect the data: continuous covariate



NO<sub>3</sub>

Is there a common effect?

# Wetlands: species-specific responses with common effect

```
model5 <- glm(Count~0+Species+N03+N03:Species, data = long, family = "poisson", contrasts = list(Species = con
```

```
##   SpeciesEphemeroptera      SpeciesHemiptera      SpeciesAmphipoda
##   2.67111728              2.36082751          3.86298300
##   SpeciesGastropoda        SpeciesAnisoptera      SpeciesZygoptera
##   2.55157209              1.21032834          2.70445223
##   SpeciesDecapoda          SpeciesDiptera       SpeciesChironomidae
##   -0.51635148              0.44859004          2.18230972
## SpeciesPlatyhelminthes    SpeciesColeoptera      SpeciesDaphniidae
##   0.93797899              0.46774680          4.02318285
##   SpeciesOstracoda         N03                 Species1:N03
##   5.76826016              -0.02504095          0.04674984
##   Species2:N03             Species3:N03       Species4:N03
##   -0.74506167              0.21656771          0.51506145
##   Species5:N03             Species6:N03       Species7:N03
##   0.14244009              0.12796859          -0.43765004
##   Species8:N03             Species9:N03       Species10:N03
##   0.06845114              0.34248963          -0.95659058
##   Species11:N03            Species12:N03
##   0.48870124              -0.05675712
```

## Wetlands: species-specific responses with common effect

Count data is usually overdispersed; we might want to switch to a NB.

```
coef(model6)
```

```
##      SpeciesEphemeroptera          SpeciesHemiptera          Species
##                  2.67109283          2.24282790          2.24282790
##      SpeciesGastropoda           SpeciesAnisopota          Species
##                  2.56892440          1.20450595          1.20450595
##      SpeciesDecapoda            SpeciesDiptera           SpeciesCh
##                 -0.49779235          0.44836747          0.44836747
## SpeciesPlatyhelminthes        SpeciesColeoptera          Species
##                  0.80901413          0.43073882          0.43073882
##      SpeciesOstracoda             N03
##                  0.80901413          0.43073882          0.43073882
```

## Wetlands: conclusions

---

We can conclude that NO<sub>3</sub> has, **on average**, a negative effect on our species pool (but this is not statistically significant)

Some species are more negatively affected than the average, some more positive

Some species are positively affected by NO<sub>3</sub>, but most negatively

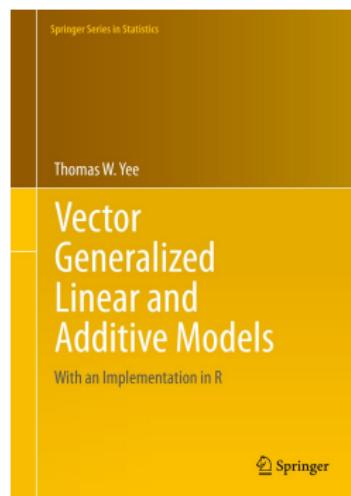
# Interpreting the coefficients

```
## SpeciesEphemeroptera      SpeciesHemiptera      SpeciesAmphipoda
##          2.67109283          2.24282790          3.85799261
## SpeciesGastropoda        SpeciesAnisopota      SpeciesZygoptera
##          2.56892440          1.20450595          2.70360747
## SpeciesDecapoda          SpeciesDiptera       SpeciesChironomidae
##          -0.49779235          0.44836747          2.15746463
## SpeciesPlatyhelminthes   SpeciesColeoptera     SpeciesDaphniidae
##          0.80901413          0.43073882          4.02257593
## SpeciesOstracoda          N03                Species1:N03
##          5.75226893          -0.03750298          0.06151128
## Species2:N03              Species3:N03       Species4:N03
##          -1.10550200          0.27812794          0.46946826
## Species5:N03              Species6:N03       Species7:N03
##          0.25306729          0.15672356          -0.32428547
## Species8:N03              Species9:N03       Species10:N03
##          0.09136078          0.48958482          -1.33007265
## Species11:N03             Species12:N03
##          0.62804597          -0.06027316
```

- ▶ Negative means a decrease in the response and positive increase
- ▶ More specifically here: the coefficient is multiplicative decrease in  $\exp(\text{intercept})$  for a unit change in N03
- ▶ E.g., for "Ephemeroptera":  $\exp(2.243) * \exp(-0.038 + 0.062) = 9.4 * 1.025$

# Vector GLMs

- ▶ One GLM per species
- ▶ Each gets their own dispersion parameter
- ▶ Slightly more flexible than what we have done so far



## Methods in Ecology and Evolution

*Methods in Ecology and Evolution* 2012, 3, 471–474



doi: 10.1111/j.2041-210X.2012.00190.x

**mvabund – an R package for model-based analysis of multivariate abundance data**

Yi Wang<sup>1,2</sup>, Ulrike Naumann<sup>1</sup>, Stephen T. Wright<sup>1</sup>, and David I. Warton<sup>1,3\*</sup>

## Fitting vector GLMs

---

A few software implementations exist:

- ▶ The VGAM R-package
- ▶ The glmmTMB R-package
- ▶ The gllvm R-package

Clearly, we will use the last one.

## VGLM Likelihood

---

- ▶ We use MLE for estimation
- ▶ With a distribution in the “exponential family” (for fixed  $\phi$ )

All GLMs have the likelihood:

$$\mathcal{L}(y_{ij}; \Theta) = \exp\left\{\frac{y_{ij} \eta_{ij} - b(\eta_{ij})}{a(\phi_j)} + c(y_{ij}, \phi_j)\right\} \quad (4)$$

So, now we have  $\phi_j$  instead of  $\phi$

# gllvm

Received: 7 May 2019 | Accepted: 5 September 2019

DOI: 10.1111/2041-210X.13303

**APPLICATION**

Methods in Ecology and Evolution 

## gllvm: Fast analysis of multivariate abundance data with generalized linear latent variable models in R

Jenni Niku<sup>1</sup>  | Francis K. C. Hui<sup>2</sup> | Sara Taskinen<sup>1</sup> | David I. Warton<sup>3</sup> 

- ▶ Originally published in 2019 by Niku et al. I “joined in” shortly after
- ▶ For model-based multivariate analysis of community ecological data
- ▶ Models are fitted in C++ (Kristensen et al. 2015)
- ▶ Can fit many different models: VGLM(M), JSIM, and ordination

## Downsides

---

- ▶ VGLM defaults to 1 dispersion parameter per species
- ▶ VGLM assumes 1 parameter per species per covariate
- ▶ This does not tend to work very well for real (sparse) community data
- ▶ VGLM assumes independence of species
- ▶ Does not include random effects (pseudoreplication, autocorrelation)

## Summary

---

- ▶ GLMs are fun, but not usually suitable for multispecies data
- ▶ VGLMs; fitting one model per species gives more flexibility
- ▶ This facilitates adding components that are shared across species
- ▶ Which is especially helpful when working with random effects

So far we have assumed that species do not influence each other