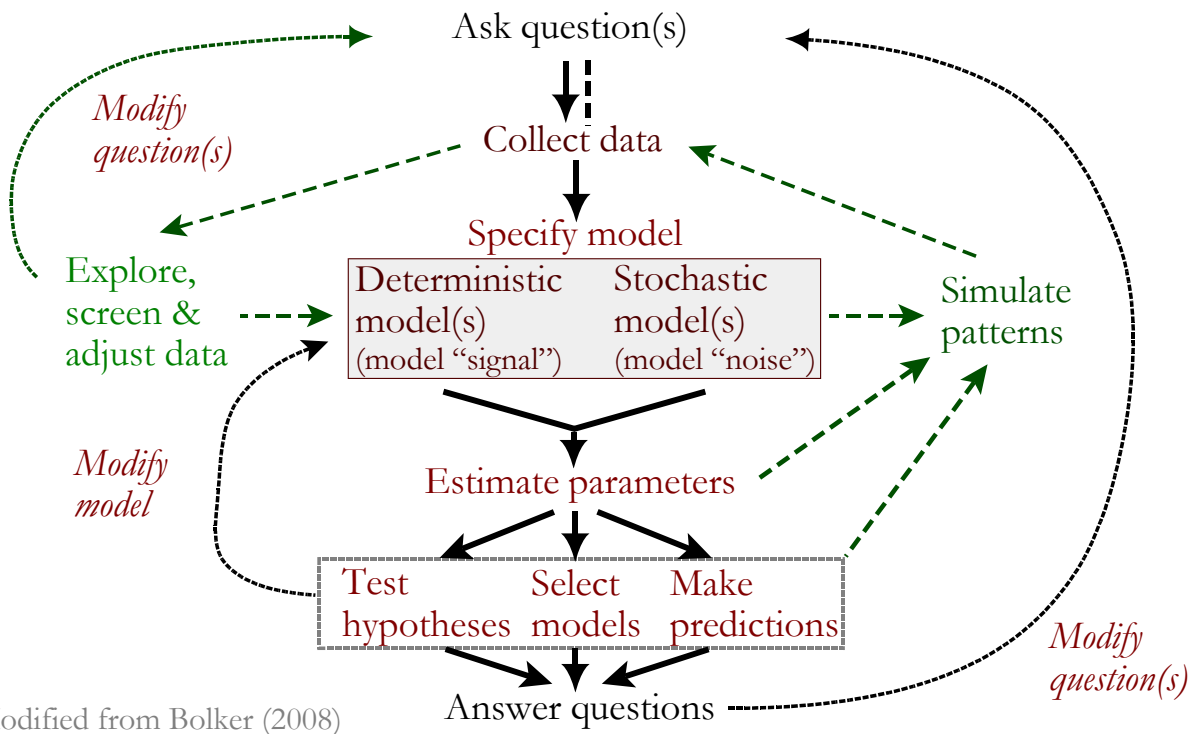


Role of Statistics... *the modeling process*



1

Landscape of Statistical Methods...

The basic statistical model:

$$Y = \text{deterministic part} + \text{stochastic part}$$



■ Univariate

■ Multivariate

■ Linear

■ Nonlinear

■ Smoothed

■ Distribution

■ Heterogeneity

■ Autocorrelation

■ Multiple levels

■ Random noise

2

Landscape of Statistical Methods...

The basic statistical model:

$$Y = \text{deterministic part} + \text{stochastic part}$$



- Univariate
- Multivariate

- Linear
- Nonlinear
- Smoothed

- Distribution
- Heterogeneity
- Autocorrelation
- Multiple levels
- Random noise

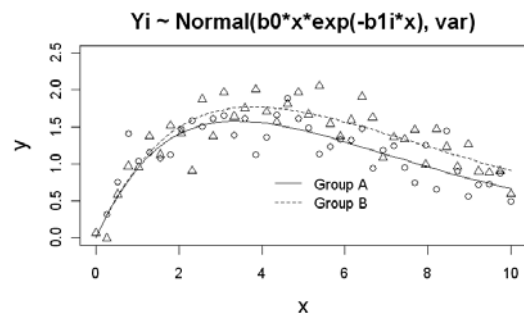
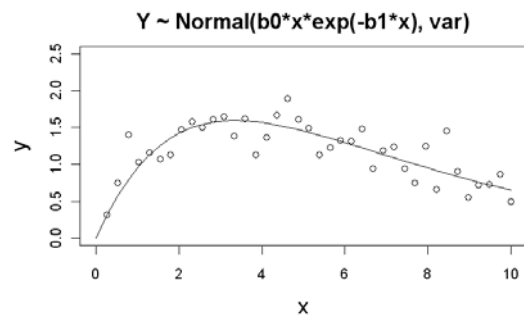
3

Landscape of Statistical Methods...

Dealing with nonlinearity

Nonlinear least squares models (NLS)

- Relax the requirement of linearity (in the parameters) but keep the requirements of independent (and normal and constant) errors (to compute Likelihoods)
- Method: *numerical least squares*



4

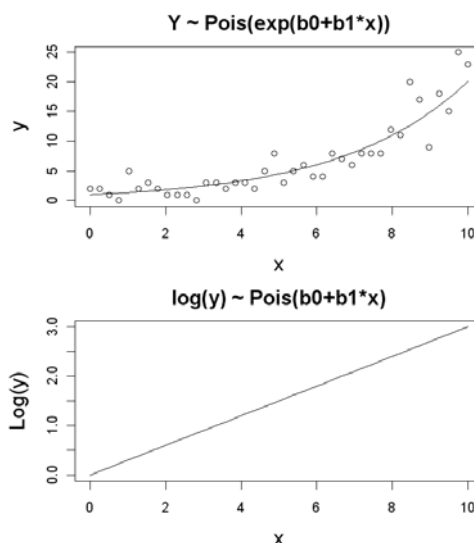
Landscape of Statistical Methods...

Dealing with nonlinearity

Generalized linear models (GLMs)

- Models that have a particular kind of *nonlinearity* and particular kinds of *nonnormally* distributed (but still independent and constant) errors
- Can fit any nonlinear relationship that has a *linearizing transformation* (link function)
- Method: *iteratively reweighted least squares* (avoids distortions in expected variance that linearizing transformation would induce)

Link $\left\{ \begin{array}{l} y = e^x \\ x = \log(y) \end{array} \right.$



5

Landscape of Statistical Methods...

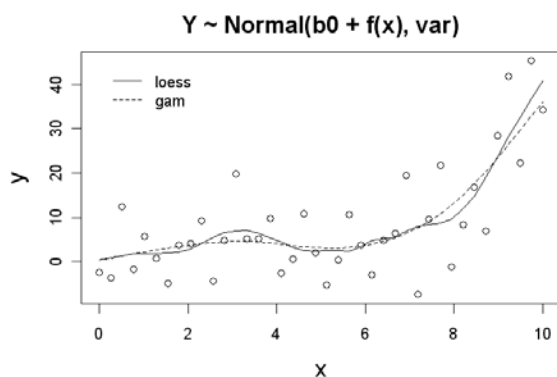
Dealing with nonlinearity

Generalized additive models (GAMs)

- Fit a *smoothing curve* through the data but keep the requirements of independent and constant *parametric* errors; hybrid parametric-nonparametric model
- Note, purely data-driven; phenomenological
- Many different types of smoothers

$$Y \sim \text{Normal}(b_0 + b_1 x, \sigma^2)$$

$$Y \sim \text{Normal}(b_0 + f(x), \sigma^2)$$



6

Landscape of Statistical Methods...

The basic statistical model:

$$Y = \text{deterministic part} + \text{stochastic part}$$



- Univariate
- Multivariate
- Linear
- Nonlinear
- Smoothed
- Distribution
- Heterogeneity
- Autocorrelation
- Multiple levels
- Random noise

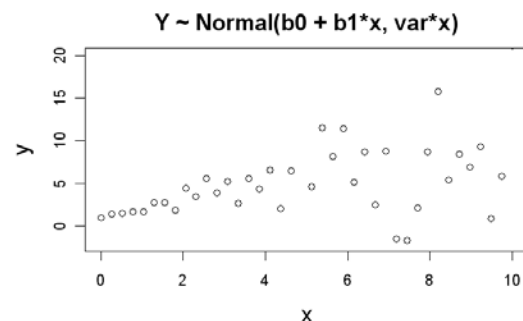
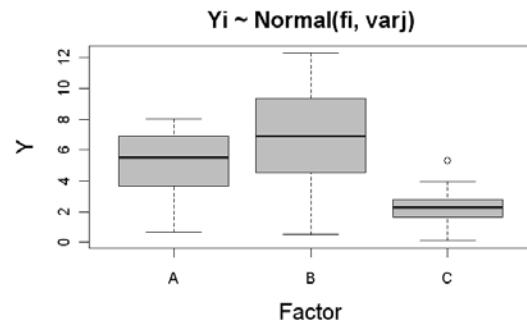
7

Landscape of Statistical Methods...

Dealing with heterogeneity

Generalized (non)linear least squares (GLS/GNLS)

- Y is continuous
- All observed values are independent and normally distributed, but with a nonconstant variance (heterogeneity)
- Method: (restricted) maximum likelihood (REML/ML)

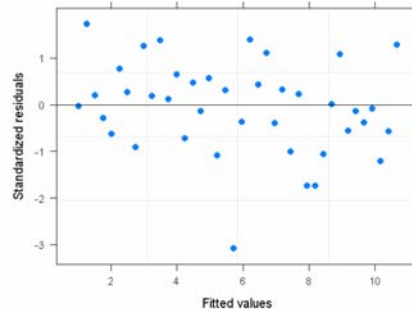
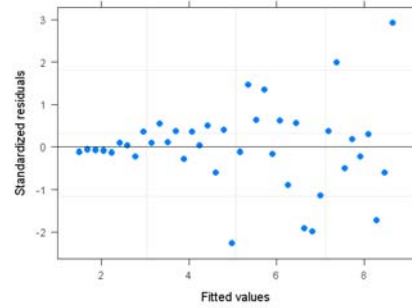
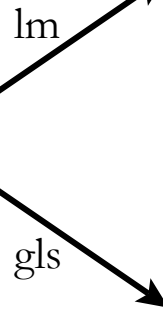
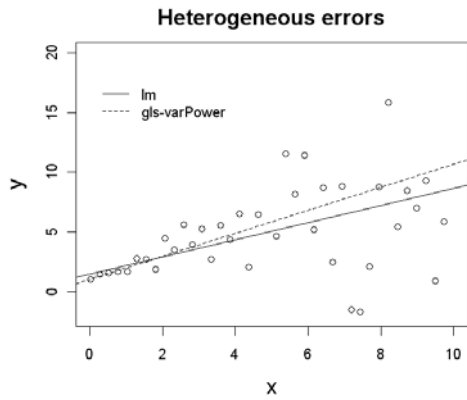


8

Landscape of Statistical Methods...

Dealing with heterogeneity

Generalized (non)linear least squares (GLS/GNLS)



	Model	df	AIC	BIC	logLik	Test	L.Ratio	p-value
fit.lm	1	3	233.5692	238.4819	-113.78458			
fit.gls2	2	4	195.3466	201.8969	-93.67328	1 vs 2	40.2226	<.0001

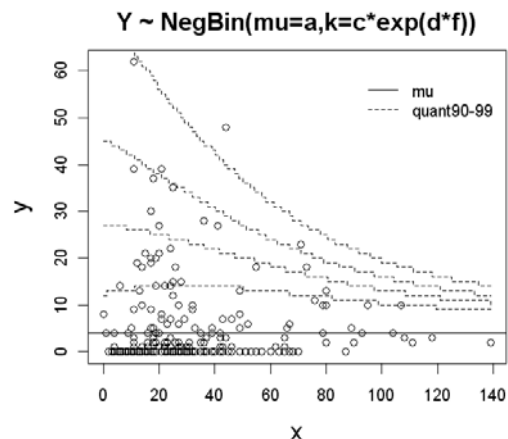
Landscape of Statistical Methods...

Dealing with heterogeneity

Customized linear/nonlinear models with *nonnormal* errors:

- In this example, the observed values of Y are *independent* (counts), distributed *negative binomial* with mean (μ) equal to a constant and the overdispersion parameter k (affecting the variance) varying exponentially as a function of X

from Bolker (2008)



	AIC	df
1	1139.333	3
2	1113.235	3

Likelihood Ratio Tests							
Model 1: fit.mu, [negbinNLL]: a+b+k							
Model 2: fit.size, [negbinNLL]: a+c+d							
	Tot	Df	Deviance	Chisq	Df	Pr(>Chisq)	
1	3	3	1133.3				
2	3	3	1107.2	26.098	0	< 2.2e-16	

Landscape of Statistical Methods...

The basic statistical model:

$$Y = \text{deterministic part} + \text{stochastic part}$$



- Univariate
- Multivariate
- Linear
- Nonlinear
- Smoothed
- Distribution
- Heterogeneity
- Autocorrelation
- Multiple levels
- Random noise

Landscape of Statistical Methods...

Dealing with (auto-)correlated errors

1. Spatial auto-correlation

Linear model

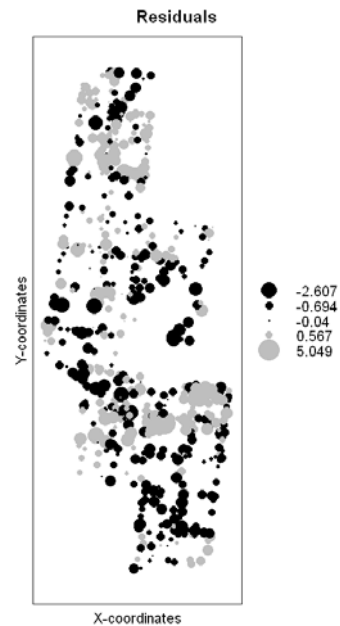
$$BFI_i = b_0 + b_1 \cdot Wetness_i + \varepsilon_i$$

$$\varepsilon_i \sim Normal(0, \sigma^2)$$

$$cov(\varepsilon_i, \varepsilon_j) = \begin{cases} \sigma^2, & i = j \\ 0, & \text{else} \end{cases}$$

Constant variance assumption

Independence assumption



from Zuur et al (2009)

Landscape of Statistical Methods...

Dealing with (auto-)correlated errors

1. Spatial auto-correlation

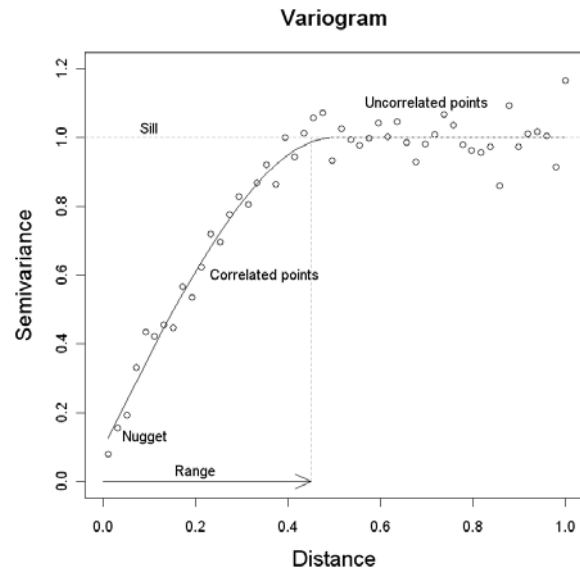
(Semi-)Variogram method:

Semivariance:

$$\gamma(x_i, x_j) = \frac{1}{2} E[(Z(x_i) - Z(x_j))^2]$$

Experimental variogram:

$$\hat{\gamma}(h) = \frac{1}{2N(h)} \sum_{i=1}^{N(h)} [z(x_i + h) - z(x_i)]^2$$



from Zuur et al (2009)

13

Landscape of Statistical Methods...

Dealing with (auto-)correlated errors

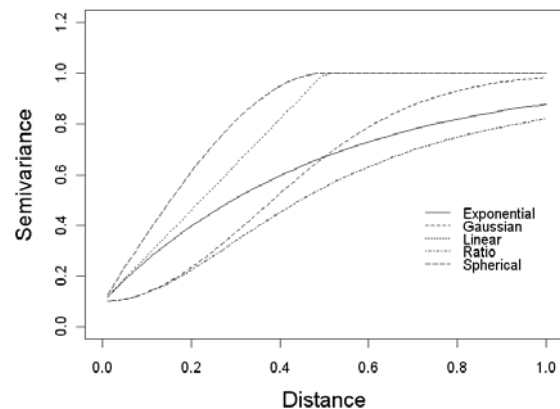
1. Spatial auto-correlation

Theoretical variograms:

- *Exponential* correlation (corExp)
- *Gaussian* correlation (corGaus)
- *Linear* correlation (corLin)
- *Rational quadratic* correlation (corRatio)
- *Spherical* correlation (corSpher)

$$\text{cor}(\varepsilon_i, \varepsilon_j) = \begin{cases} 1, & i = j \\ h(\varepsilon_i, \varepsilon_j, \rho), & \text{else} \end{cases}$$

Theoretical variograms



14

Landscape of Statistical Methods...

Dealing with (auto-)correlated errors

1. Spatial auto-correlation

Experimental variogram:

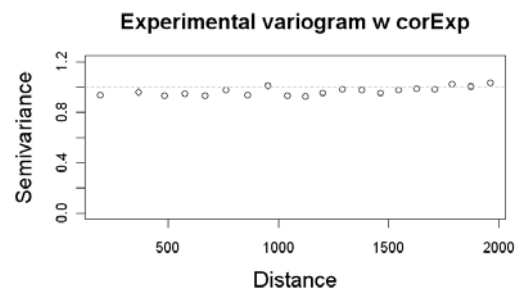
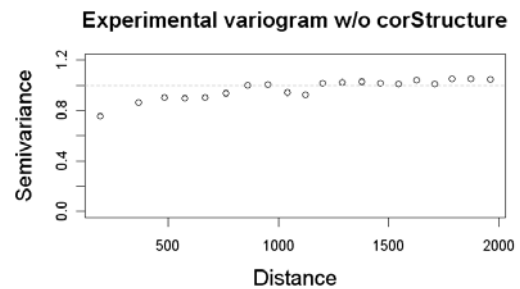
```
df      AIC
B.none  3 2844.541
B.spher 5 2737.012
B.lin   5 2848.510
B.ratio 5 2732.930
B.gaus  5 2736.292
B.exp   5 2732.224
```

$$\text{cor}(\varepsilon_i, \varepsilon_j) = \begin{cases} 1, & i = j \\ 1 - e^{-\frac{d}{r}}, & \text{else} \end{cases}$$

d = distance

r = estimated range

```
Model df      AIC      BIC    logLik  Test  L.Ratio p-value
B.none  1  3 2844.541 2857.365 -1419.271
B_exp   2  5 2732.224 2753.597 -1361.112 1 vs 2 116.3175 <.0001
```



from Zuur et al (2009)

15

Landscape of Statistical Methods...

Dealing with (auto-)correlated errors

Other methods for dealing with spatial and temporal auto-correlation:

- Autocovariate models
- Spatial eigenvector mapping (SEVM)
- Generalized least squares (GLS)
 - Conditional autoregressive models (CAR)
 - Simultaneous autoregressive models (SAR)
 - Generalized linear mixed models (GLMM)
- Generalized estimating equations (GEE)
- Wavelet-revised model (WRM)



16

Landscape of Statistical Methods...

The basic statistical model:

$$Y = \underbrace{\text{deterministic part}} + \underbrace{\text{stochastic part}}$$



- Univariate
- Multivariate
- Linear
- Nonlinear
- Smoothed
- Distribution
- Heterogeneity
- Autocorrelation
- Multiple levels
- Random noise

17

Landscape of Statistical Methods...

Dealing with multiple levels of error

- *Zero-inflated data* – models for data with too many zeros involving a mixture of two distributions
- *Nested data* – models for nested or blocked data, which divide observations into discrete groups according to their spatial or temporal locations, or other characteristics
- *Observation-process models* – models with process error and measurement (observation) error in the same model



18

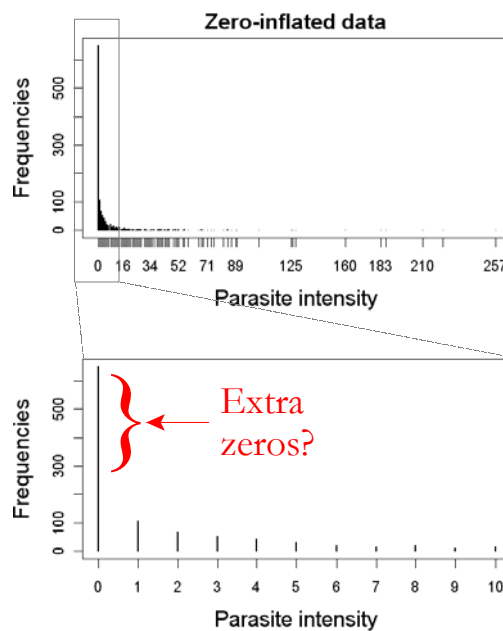
Landscape of Statistical Methods...

Dealing with multiple levels of error

1. Zero-inflated data

Why so many zeros?

- *True zeros* (positive zeros or true negatives) – structural errors, absent because the habitat is not suitable
- *False zeros* (false negatives) – due either to study design (in the wrong place or at the wrong time), survey method (low detectability), or observer error



from Zuur et al (2009)

Landscape of Statistical Methods...

Dealing with multiple levels of error

1. Zero-inflated data

Where do the zeros come from?

Zero-altered models

I am not here, but the habitat is good!

You didn't see me!

You thought I was a spotted owl



I am not here because the habitat is not good

>0 owls



Here we are!

Zero-inflated models

I am not here, but the habitat is good!

You didn't see me!



You thought I was a spotted owl

Count process

Zero mass

0 owls



I am not here because the habitat is not good

>0 owls



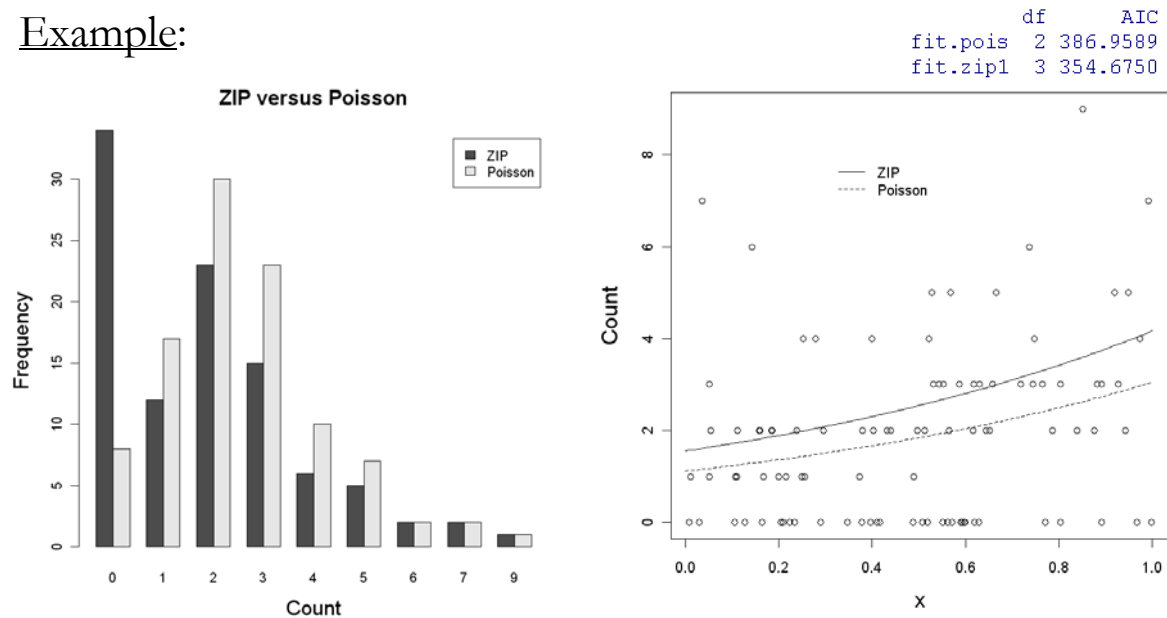
Here we are!

Landscape of Statistical Methods...

Dealing with multiple levels of error

1. Zero-inflated data

Example:



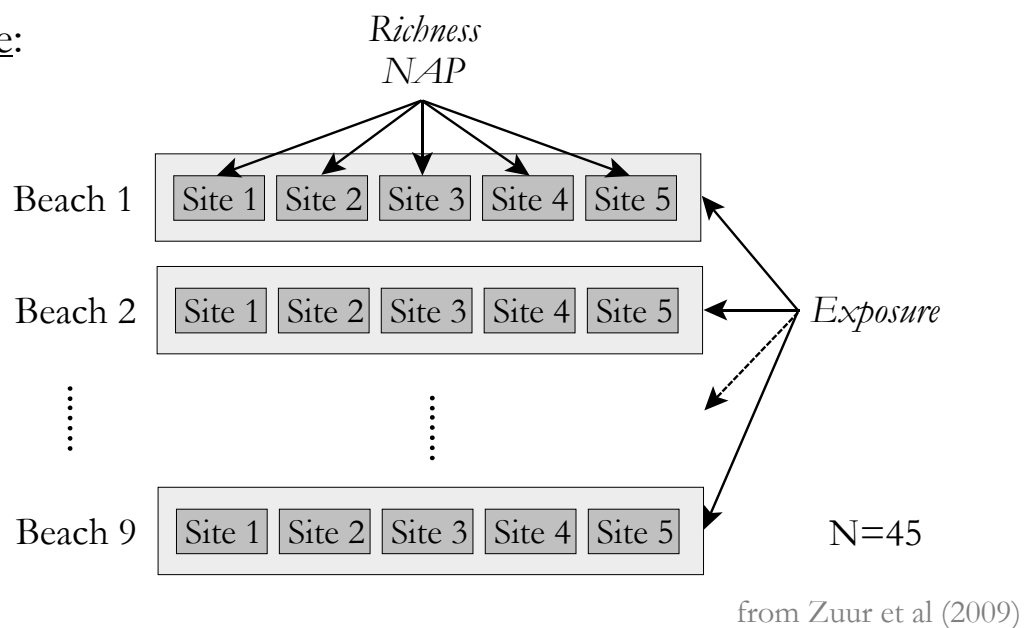
21

Landscape of Statistical Methods...

Dealing with multiple levels of error

2. Nested data

Example:



22

Landscape of Statistical Methods...

Dealing with multiple levels of error

2. Nested data

(a) Single-level model:

- Ignore nested structure

$$R_{ij} = \alpha + \beta_1 \cdot NAP_{ij} + \beta_2 Exposure_i + \varepsilon_{ij}$$

$$\varepsilon_{ij} \sim Normal(0, \sigma^2)$$

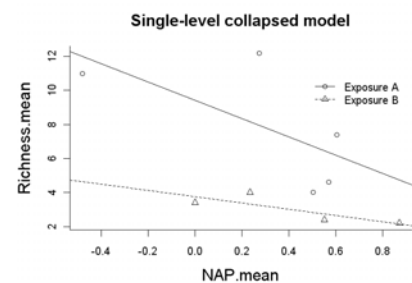
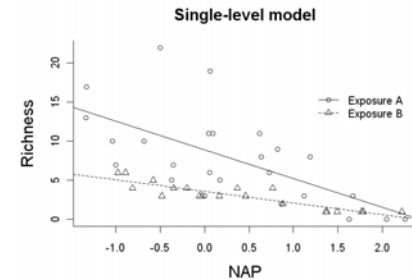
$$i = 1, \dots, 9$$

$$j = 1, \dots, 5$$

- Collapse to group level

$$\bar{R}_i = \alpha + \beta_1 \cdot \bar{NAP}_i + \beta_2 Exposure_i + \varepsilon_i$$

Neither approach is ideal



23

Landscape of Statistical Methods...

Dealing with multiple levels of error

2. Nested data

(b) Two-stage method:

- Stage 1 – separate regression for each beach

$$R_{ij} = \alpha + \beta_i \cdot NAP_{ij} + \varepsilon_{ij}$$

$$j = 1, \dots, 5$$

- Stage 2 – model estimated regression coefficients as a function of group-level covariates (exposure)

$$\hat{\beta}_i = \eta + \tau \cdot Exposure_i + b_i$$

$$i = 1, \dots, 9$$

Stage 1 result:

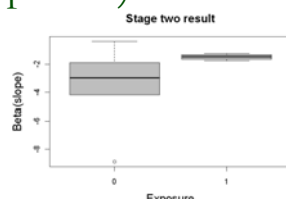
beach	betas
1	-0.3718279
2	-4.1752712
3	-1.7553529
4	-1.2485766
5	-8.9001779
6	-1.3885120
7	-1.5176126
8	-1.8930665
9	-2.9675304

Stage 2 result:

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -3.662      1.099   -3.332  0.0126 *
fExposure91   2.184      1.649    1.325  0.2268
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.458 on 7 degrees of freedom
Multiple R-squared: 0.2005,    Adjusted R-squared: 0.08625
F-statistic: 1.755 on 1 and 7 DF,  p-value: 0.2268
    
```



24

Landscape of Statistical Methods...

Dealing with multiple levels of error

2. Nested data

(c) Mixed-effects model:

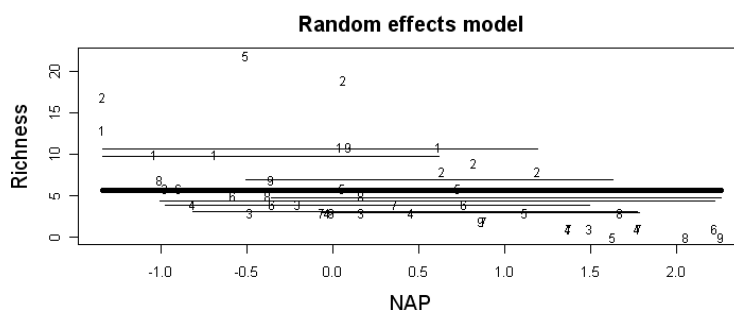
- *Random effects model* – allows intercept to vary among beaches but no slope

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

$$\mu_i = \beta_{0_{i|j}}$$

$$\beta_{0_i} \approx \text{Norm}(\mu_{\beta_0}, \sigma_{\beta_0}^2)$$

$$\mu_{\beta_0} = \gamma_0 + \gamma_1 \cdot \text{Exposure}_i$$



25

Landscape of Statistical Methods...

Dealing with multiple levels of error

2. Nested data

(c) Mixed-effects model:

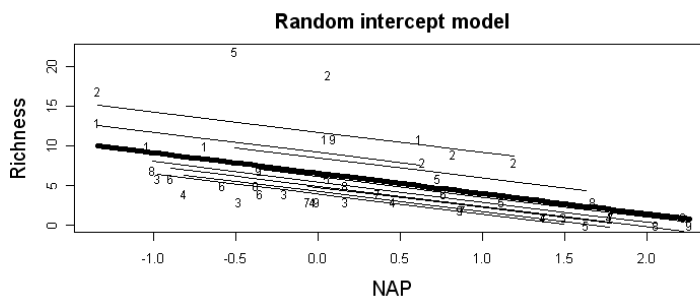
- *Random intercept model* – allows intercept to vary among beaches

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

$$\mu_i = \beta_{0_{i|j}} + \beta_{1_{i|j}} \cdot \text{NAP}_{ij}$$

$$\beta_{0_i} \approx \text{Norm}(\mu_{\beta_0}, \sigma_{\beta_0}^2)$$

$$\mu_{\beta_0} = \gamma_0 + \gamma_1 \cdot \text{Exposure}_i$$



26

Landscape of Statistical Methods...

Dealing with multiple levels of error

2. Nested data

(c) Mixed-effects model:

- *Random intercept and slope model* – allows intercept and slope to vary among beaches

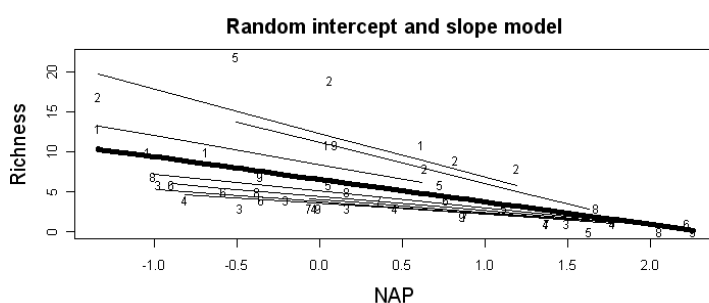
$$R_{ij} \sim \text{Normal}(\mu_i, \sigma^2)$$

$$\mu_i = \beta_{0_{ij}} + \beta_{1_{ij}} \cdot \text{NAP}_{ij}$$

$$\begin{pmatrix} \beta_{0_i} \\ \beta_{1_i} \end{pmatrix} \approx \text{Norm} \left(\begin{pmatrix} \mu_{\beta_0} \\ \mu_{\beta_1} \end{pmatrix}, \begin{pmatrix} \sigma_{\beta_0}^2, \rho \sigma_{\beta_0} \sigma_{\beta_1} \\ \rho \sigma_{\beta_0} \sigma_{\beta_1}, \sigma_{\beta_1}^2 \end{pmatrix} \right)$$

$$\mu_{\beta_0} = \gamma_0 + \gamma_1 \cdot \text{Exposure}_i$$

$$\mu_{\beta_1} = \tau_0 + \tau_1 \cdot \text{Exposure}_i$$



27

Landscape of Statistical Methods...

Dealing with multiple levels of error

3. Observation-Process models

Models that account for the ecological process and the observation process separately in a single model

“Few animals are so conspicuous that they are always detected at each survey.”
MacKenzie et al. (2002)

- When detection bias is suspected to be significant, it is necessary to account for it in the model to achieve accurate estimates of the parameters associated with the ecological process of interest



28

Landscape of Statistical Methods...

Dealing with multiple levels of error

3. Observation-Process models

Example:

Estimate occupancy rate of an invasive species of crab along a coastline in relation to the percent of the substrate covered by cobbles - a potentially important habitat covariate



site	survey.1	survey.2	survey.3	waterClarity.1	waterClarity.2	waterClarity.3	% cover cobbles
1	0	0	0	3.06	1.14	1.92	75.1
2	0	0	0	1.79	0.72	0.54	79.9
3	1	1	1	6.61	9.18	5.43	28.1
4	1	1	1	8.68	8.51	7.92	19.4
5	0	0	0	2.49	1.68	2.91	91.0
6	1	0	1	9.98	6.80	8.44	100.0
7	1	1	0	7.95	7.38	8.74	90.2
.							
.							
.							
100	0	0	0	6.59	8.41	8.31	84.6

from Royle and Dorazio (2008)

29

Landscape of Statistical Methods...

Dealing with multiple levels of error

3. Observation-Process models

Statistical model:

$$\begin{array}{l}
 \text{Process model} \left\{ \begin{array}{l} z_i \sim \text{Bern}(\psi_i) \\ \text{Logit}(\psi_i) = \beta_0 + \beta_1 \cdot \text{Cobble}_i \end{array} \right. \text{Process covariate} \\
 \text{Observation model} \left\{ \begin{array}{l} y_{ij} \sim \text{Bern}(p_{ij} \cdot z_i) \\ \text{Logit}(p_{ij}) = \alpha_0 + \alpha_1 \cdot \text{waterClarity}_{ij} \end{array} \right. \text{Observation covariate}
 \end{array}$$

Z_i = Unobserved state variable
(presence/absence at site i)

y_{ij} = Observed data (detected/not detected at site i on survey j)

α_i, β_i = Parameters to estimate

30

Landscape of Statistical Methods...

Dealing with multiple levels of error

3. Observation-Process models

Model selection:

$$(1) \quad z_i \sim \text{Bern}(\psi_i) \\ y_{ij} \sim \text{Bern}(p_{ij} \cdot z_i)$$

$$(2) \quad z_i \sim \text{Bern}(\psi_i) \\ \text{Logit}(\psi_i) = \beta_0 + \beta_1 \cdot \text{Cobble}_i \\ y_{ij} \sim \text{Bern}(p_{ij} \cdot z_i)$$

$$(3) \quad z_i \sim \text{Bern}(\psi_i) \\ y_{ij} \sim \text{Bern}(p_{ij} \cdot z_i) \\ \text{Logit}(p_{ij}) = \alpha_0 + \alpha_1 \cdot \text{waterClarity}_{ij}$$

$$(4) \quad z_i \sim \text{Bern}(\psi_i) \\ \text{Logit}(\psi_i) = \beta_0 + \beta_1 \cdot \text{Cobble}_i \\ y_{ij} \sim \text{Bern}(p_{ij} \cdot z_i) \\ \text{Logit}(p_{ij}) = \alpha_0 + \alpha_1 \cdot \text{waterClarity}_{ij}$$

Model	n	K	AIC	Δ AIC	AICwt	R-squared	AICwtCum	Ψ	S.E.(Ψ)
p(clarity), psi(cobbles)	100	4	210.67	0.000	0.679	0.310	0.679	0.418	0.088
p(clarity), psi(.)	100	3	212.16	1.497	0.321	0.281	1.000	0.420	0.070
p(.), psi(cobbles)	100	3	237.55	26.885	0.000	0.042	1.000	0.318	0.067
p(.), psi(.)	100	2	239.31	28.642	0.000	0.000	1.000	0.320	0.050

31

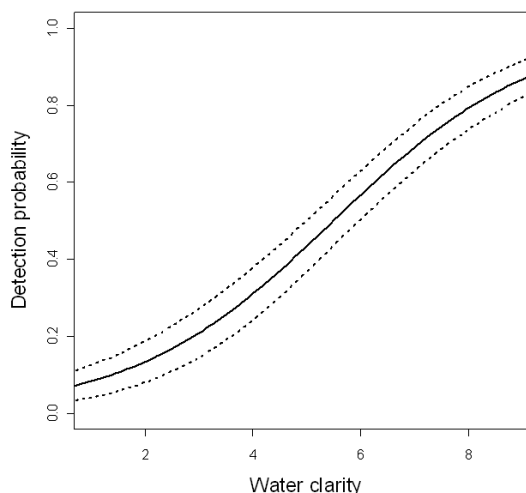
Landscape of Statistical Methods...

Dealing with multiple levels of error

3. Observation-Process models

Detectability function:

- Estimating the detectability function can be useful in the design of future studies



32

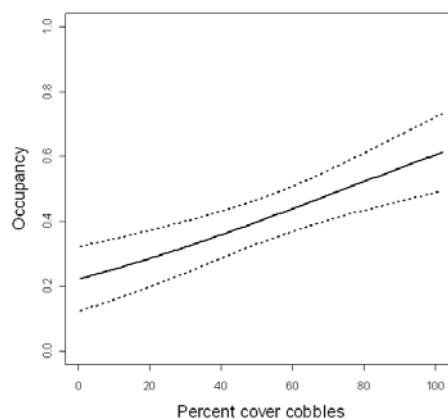
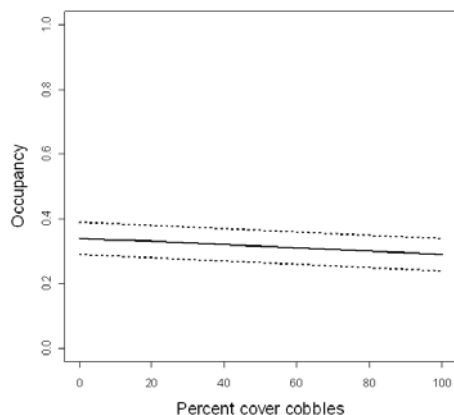
Landscape of Statistical Methods...

Dealing with multiple levels of error

3. Observation-Process models

Occupancy function:

- Estimates of occupancy can be significantly biased and even misleading if detectability is not taken into account



33

Landscape of Statistical Methods...

The basic statistical model:

$$Y = \text{deterministic part} + \text{stochastic part}$$



■ Univariate

■ Multivariate

■ Linear

■ Nonlinear

■ Smoothed

■ Distribution

■ Heterogeneity

■ Autocorrelation

■ Multiple levels

■ Random noise

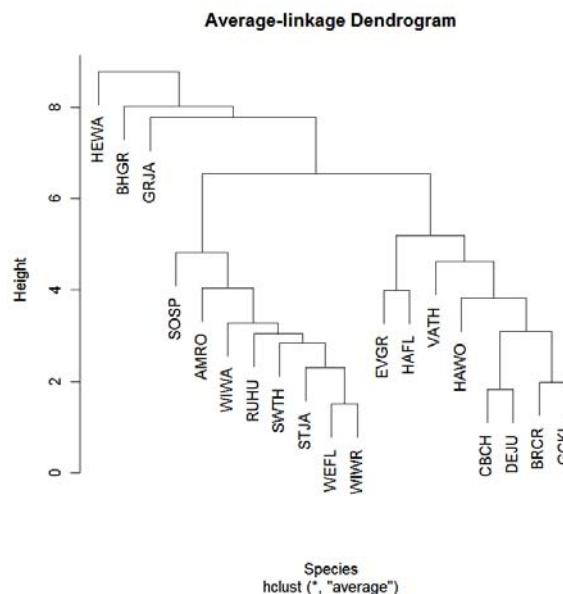
34

Landscape of Statistical Methods...

Multivariate methods

Finding groups (Cluster analysis)

- Large family of techniques with similar goals; operating on data sets for which *pre-specified, well-defined groups do "not" exist*; characteristics of the data are used to assign entities into artificial groups



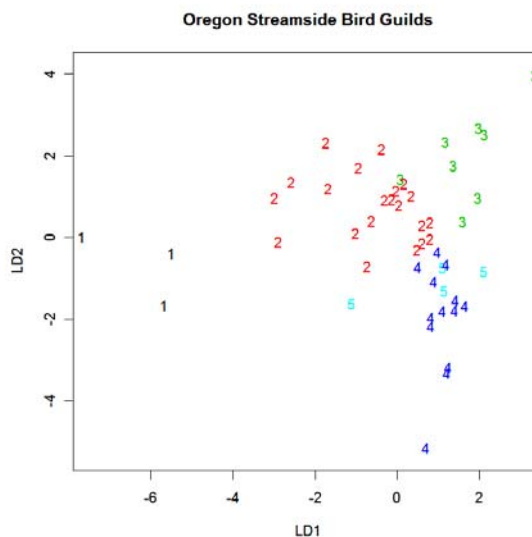
35

Landscape of Statistical Methods...

Multivariate methods

Testing/describing differences among groups (e.g., DA, ISA, mCART, MRPP, MANTEL)

- Large family of different methods for testing and/or describing differences among *pre-specified, well-defined groups* based on a set of discriminating variables



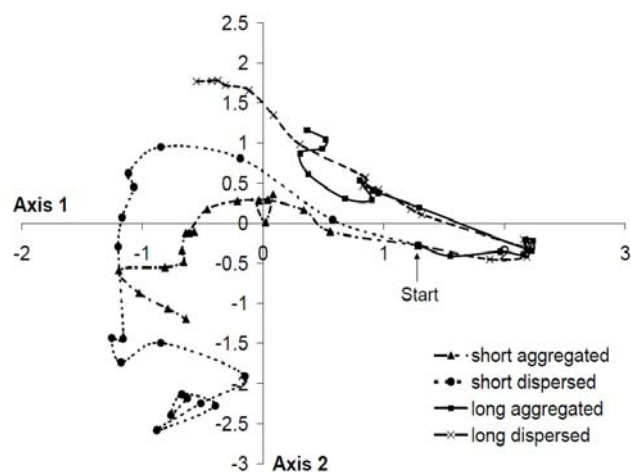
36

Landscape of Statistical Methods...

Multivariate methods

Unconstrained ordination (e.g., PCA, CA, NMDS)

- A family of different methods for organizing sampling entities (e.g., species, sites, observations, etc.) along continuous gradients based on a set of interdependent variables



37

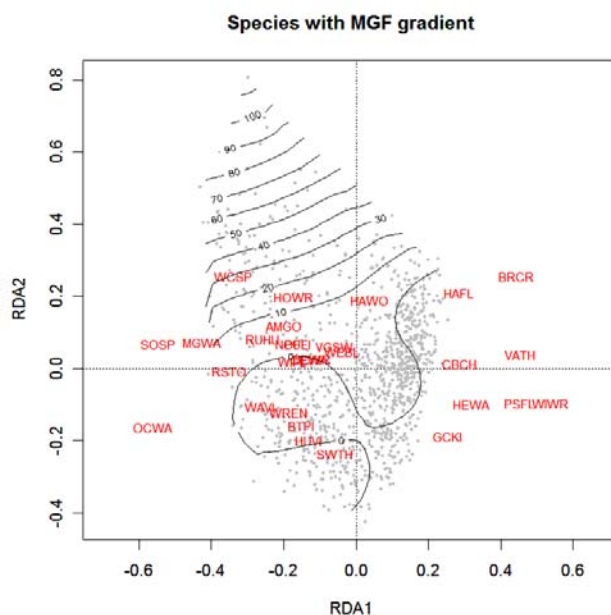
Landscape of Statistical Methods...

Multivariate methods

Constrained ordination (e.g., RDA, CCA, CAPS)

- A family of different methods for extending unconstrained ordination in which the solution is constrained to be expressed by ancillary variables

Triplot (samples, species, environment)

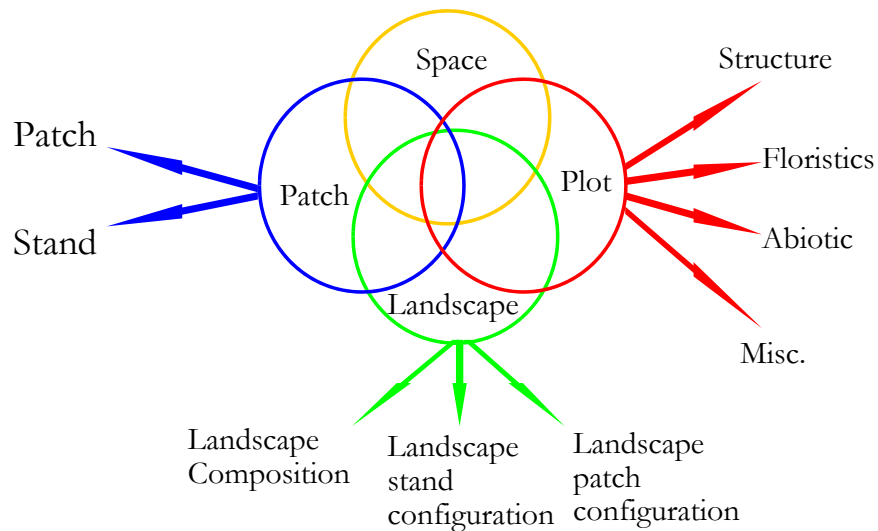


38

Landscape of Statistical Methods...

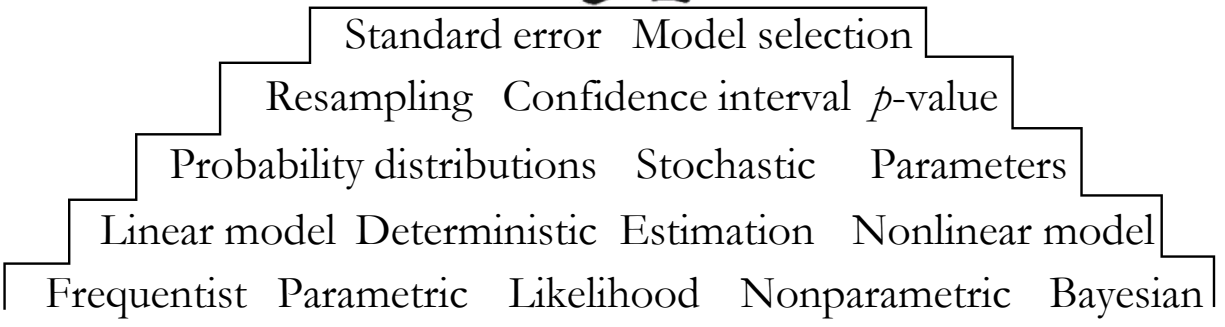
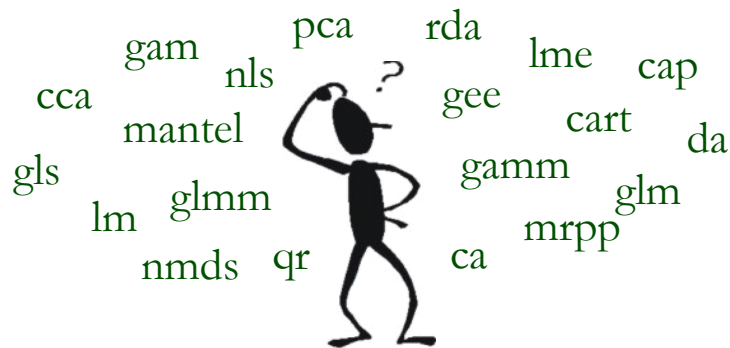
Multivariate methods

Ordination variance partitioning



39

Foundation for Understanding Statistical Methods...



40

Foundation for Understanding Statistical Methods...



Statistical foundations