

General linear models (not Generalized linear model)

Linear Model	Common name
$Y = \mu + X$	Simple linear regression
$Y = \mu + A_1$	One-factorial (one-way) ANOVA
$Y = \mu + A_1 + A_2 + A_1 \times A_2$	Two-factorial (two-way) ANOVA
$Y = \mu + A_1 + X (+A_1 \times X)$	Analysis of Covariance (ANCOVA)
$Y = \mu + X_1 + X_2 + X_3$	Multiple regression
$Y = \mu + A_1 + g + A_1 \times g$	Mixed model ANOVA
$Y_1 + Y_2 = \mu + A_1 + A_2 + A_1 \times A_2$	Multivariate ANOVA (MANOVA)

Y (response) is a continuous variable
 X (predictor) is a continuous variable
 A represents categorical predictors (factors)
 g represents groups of data (more on this later)
 (+A₁ × X) - step 1 on an ANCOVA, but not in the final analysis
 Multiple factors A₁ + A₂ + etc (and their interactions)

1

Understanding and dealing with heterogeneity

 Intermediary steps before going fully mixed.....

 model

2

Let's start with a problem

Seasonal patterns of investment in reproductive and somatic tissues in the squid *Loligo forbesi*

Jennifer M. Smith^{1,2}, Graham J. Pierce¹, Alain F. Zuur² and Peter R. Boyle¹

¹ Department of Zoology, School of Biological Sciences, University of Aberdeen, Tillydrone Avenue, Aberdeen AB24 2TZ, UK
² Highland Statistics Ltd., 6 Laverock Road, Newburgh, Aberdeenshire, AB41 6FN, UK

Goal: study seasonal variation (patterns) in reproductive and somatic tissues (mating is aseasonal).

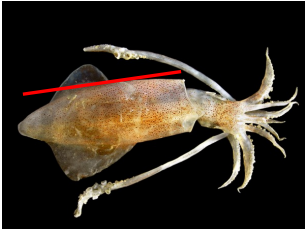
In which month there is more investment (relative to individual size, i.e., DML) in reproduction?

testis weight (mg)
 dorsal mantle length (DML; mm)

month 1
 month 2

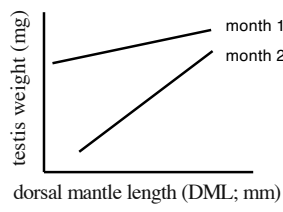
Aquat. Living Resour. 18, 341–351 (2005)
 © IEEP Sciences, IPRIMER, IRD 2005
 DOI: 10.1051/aqr/2005183
 www.edpsciences.org/aqr

3



Goal: study seasonal patterns in reproductive and somatic tissues.

In which month there is more investment (relative to individual size DML) in reproduction?



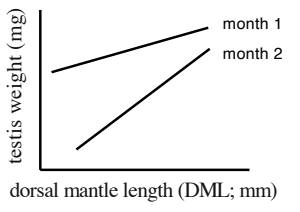
4

Data structure

Specimen	MONTH	DML	Testisweight
1017	2	136	0.006
1034	9	144	0.008
1070	12	108	0.008
1070	11	130	0.011
1019	8	121	0.012
1002	10	117	0.012
1001	5	133	0.013
1013	7	105	0.015
1002	7	109	0.017
1006	7	97	0.017
1020	9	144	0.022
1002	6	141	0.023
1039	9	125	0.024
1038	9	140	0.026
1012	12	128	0.027
1037	9	142	0.036
1001	6	139	0.036
1027	7	145	0.043
1003	7	181	0.05

768 individuals

Goal: study seasonal patterns in reproductive and somatic tissues.



5

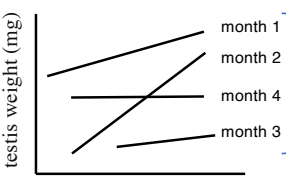
Goal: study seasonal patterns in reproductive and somatic tissues.

Model of interest

$$\text{TestisWeight} = \text{constant} + \beta_1 \text{DML} + \beta_2 \text{Month} + \beta_3 (\text{DML} \times \text{Month}) + e$$

$e \sim N(0, \sigma^2)$

continuous variable continuous variable categorical variable (factor)



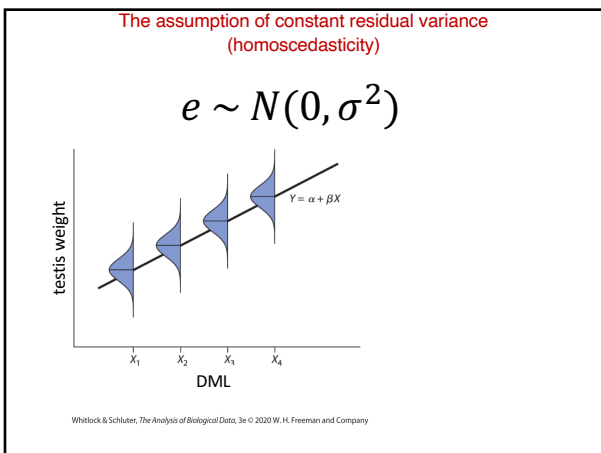
seasonal variation (environmental drivers)?

What component of the model test for the variation in slopes across months?

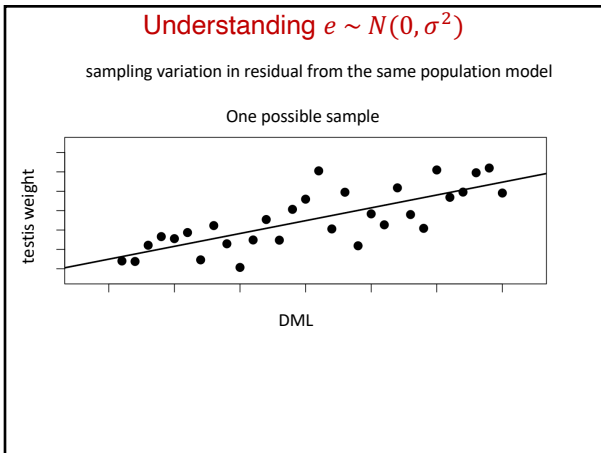
dorsal mantle length (DML; mm)
(proxy for somatic tissue)

déjà vu

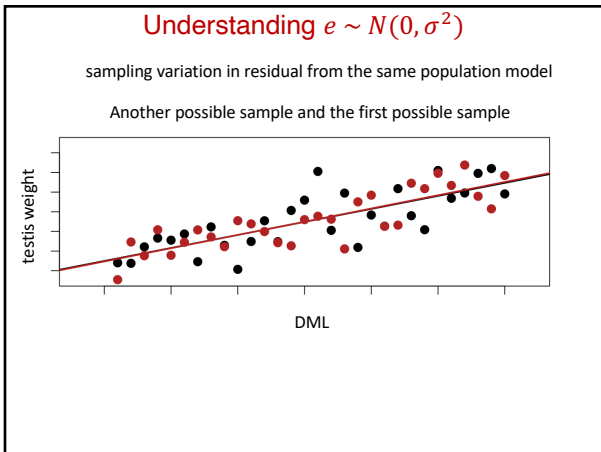
6



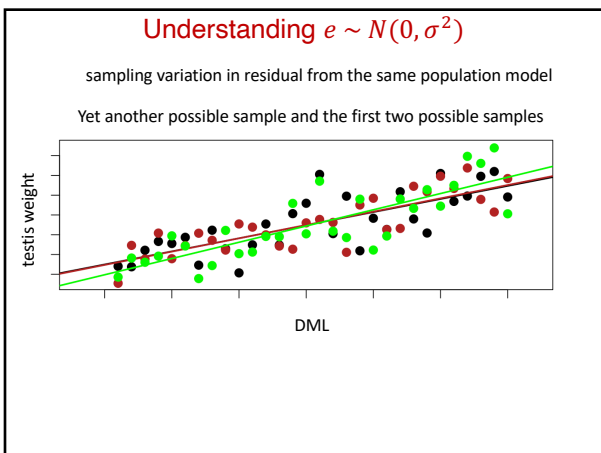
7



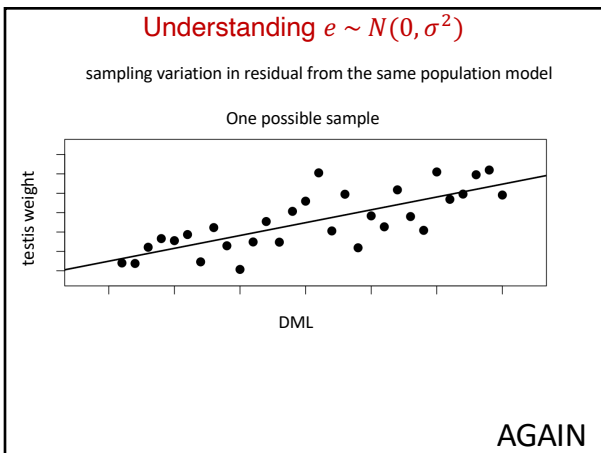
8



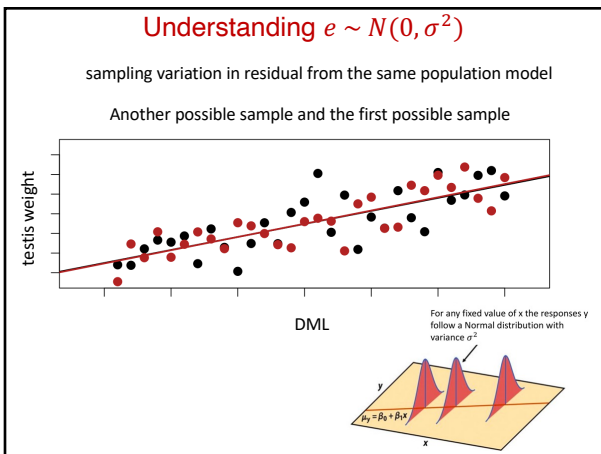
9



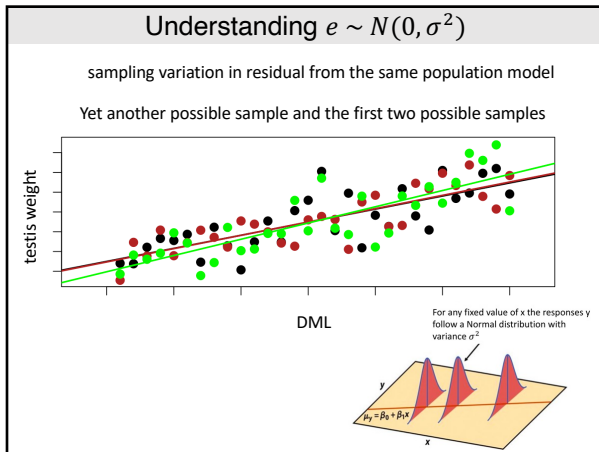
10



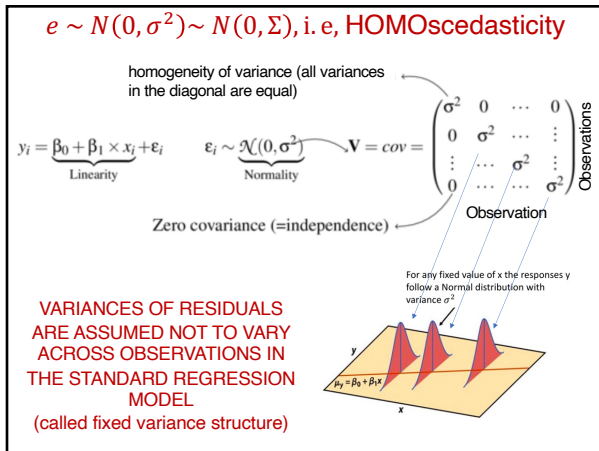
11



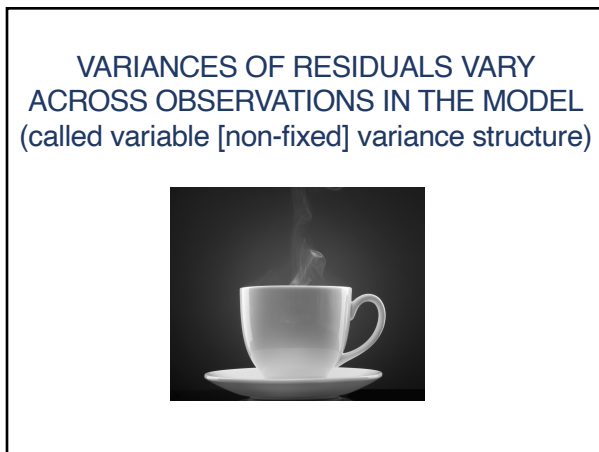
12



13



14



15

$e \neq N(0, \sigma^2) \neq N(0, \Sigma)$, i. e, HETEROscedasticity

Heteroscedasticity (variances in the diagonal are not equal)

$y_i = \beta_0 + \beta_1 \times x_i + \epsilon_i$ $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Linearity Normality

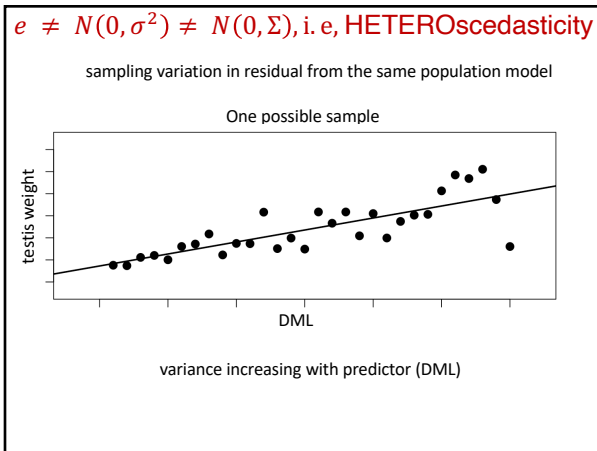
$V = cov = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & \vdots \\ \vdots & \dots & \sigma_3^2 & \vdots \\ 0 & \dots & \dots & \sigma_4^2 \end{pmatrix}$

Observations

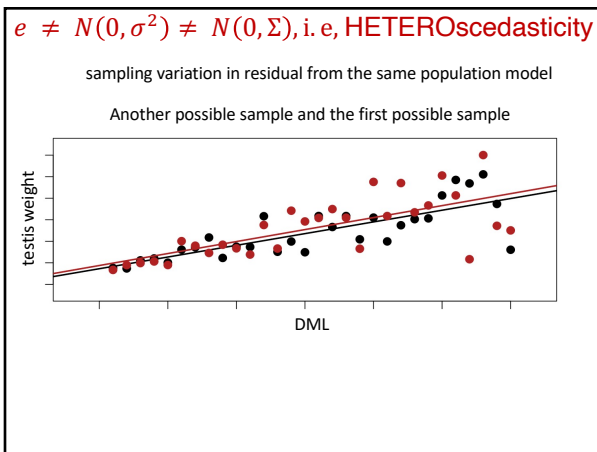
Zero covariance (=independence)

VARIANCES OF RESIDUALS VARY ACROSS OBSERVATIONS IN THE MODEL (called variable [non-fixed] variance structure)

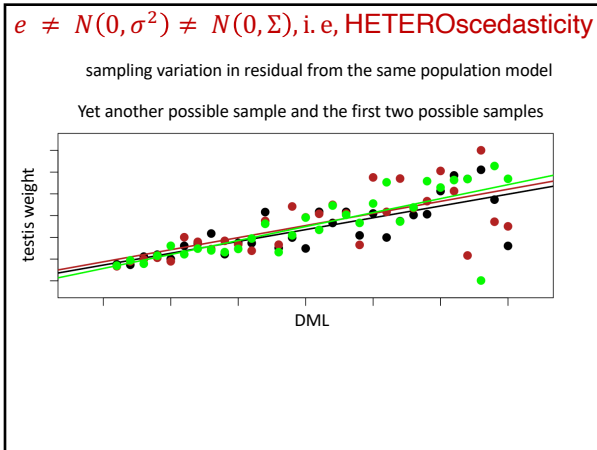
16



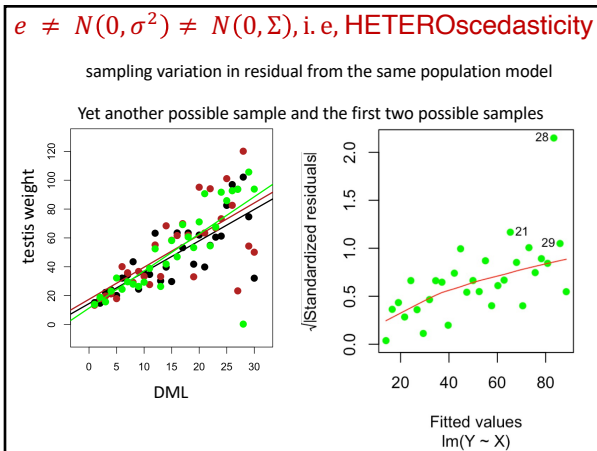
17



18



19



20

$e \neq N(0, \sigma^2) \neq N(0, \Sigma)$, i. e, HETEROscedasticity

How was variance heterogeneity generated in these examples?

```

21
22 n=30
23 X = 1:n
24 e = rnorm(n, 0, X)
25 Y = constant + slopeX * X + e
    
```

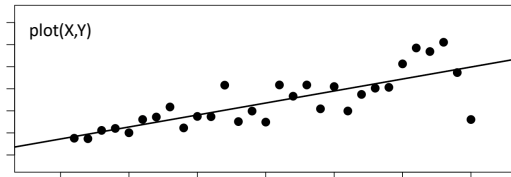
21

$e \neq N(0, \sigma^2) \neq N(0, \Sigma)$, i.e, HETEROscedasticity

How was variance heterogeneity generated in these examples?

```

21
22 n=30
23 X = 1:n
24 e = rnorm(n,0,X)
25 Y = constant + slopeX * X + e
    
```

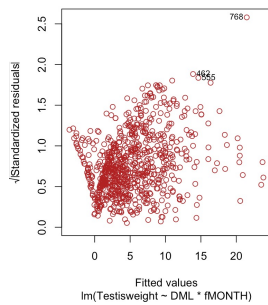


22

Goal: study seasonal patterns in reproductive and somatic tissues

Going back to the model of interest

$$\text{TestisWeight} = \text{constant} + \beta_1 \text{DML} + \beta_2 \text{Month} + \beta_3 (\text{DML} \times \text{Month}) + e$$



Residuals are highly heteroscedastic

```

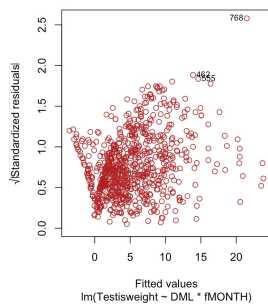
> bptest(M1)
studentized Breusch-Pagan test
data: M1
BP = 160.08, df = 23, p-value < 2.2e-16
    
```

23

Goal: study seasonal patterns in reproductive and somatic tissues.

Going back to the model of interest

$$\text{TestisWeight} = \text{constant} + \beta_1 \text{DML} + \beta_2 \text{Month} + \beta_3 (\text{DML} \times \text{Month}) + e$$



What are the origins (or proxies) of change in residual variance?

```

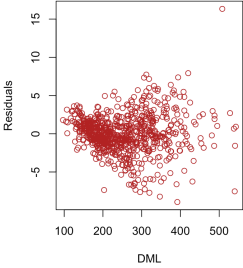
> bptest(M1)
studentized Breusch-Pagan test
data: M1
BP = 160.08, df = 23, p-value < 2.2e-16
    
```

24

Goal: study seasonal patterns in reproductive and somatic tissues.

Variance changes as a function of DML

TestisWeight = constant + β_1 DML + β_2 Month + β_3 (DML \times Month) + e



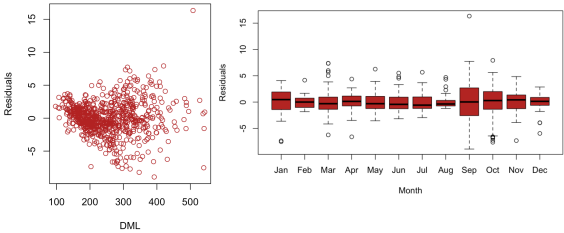
What are the origins (or proxies) of change in residual variance?

25

Goal: study seasonal patterns in reproductive and somatic tissues.

Variance changes as a function of DML x Month (interaction)

TestisWeight = constant + β_1 DML + β_2 Month + β_3 (DML \times Month) + e

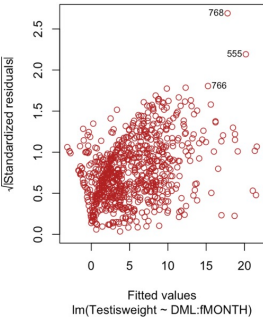


26

Goal: study seasonal patterns in reproductive and somatic tissues.

Variance changes as a function of DML x Month (interaction)

TestisWeight = constant + β_1 DML + β_2 Month + β_3 (DML \times Month) + e



27

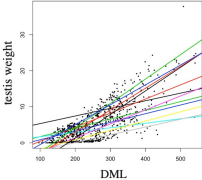
Variance changes as a function of Month

TestisWeight = constant + β_1 DML + β_2 Month + β_3 (DML \times Month) + e

$e \sim N(0, \sigma^2)$ \Rightarrow This assumption does not hold

If the DML by Month interaction is significant, we know that the slopes of DML change as a function of Month (i.e., ANCOVA).

If the slopes for DML change across months, then assuming one single slope for all the data will generate heteroscedasticity, i.e., perhaps residuals are homoscedastic but only within models per month.



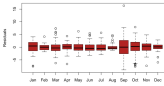
28

Variance changes as a function of Month

$e_{ij} \sim N(0, \sigma_j^2) \quad j = 1, \dots, 12$

	Specimen 1	...	Specimen 768
Specimen 1	e_{ij}	0	...
.	0	e_{ij}	...
.	e_{ij}
Specimen 768	0

Variance-covariance matrix



29

Variance changes as a function of Month

$e_{ij} \sim N(0, \sigma_j^2) \quad j = 1, \dots, 12$

How is this variance structure included in the model?

Ordinary Least Square GLS (fixed variance):


$$\beta = (X^T X)^{-1} X^T Y$$

Generalized Least Square GLS (variable variance):

$$\beta = (X^T W X)^{-1} X^T W Y$$

30

How to account for variance differences?



31

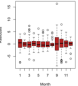
Variance changes as a function of Month
 How is this variance structure included in the model?
 Generalized Least Square GLS (variable variance):

$$\beta = (X^T W X)^{-1} X^T W Y \quad W \sim 1/f(\Sigma)$$

$\Sigma =$

Specimen 1	e_{ij}	0	...	0
⋮	0	e_{ij}	...	⋮
⋮	⋮	...	e_{ij}	⋮
Specimen 768	0	e_{ij}

Variance-covariance matrix

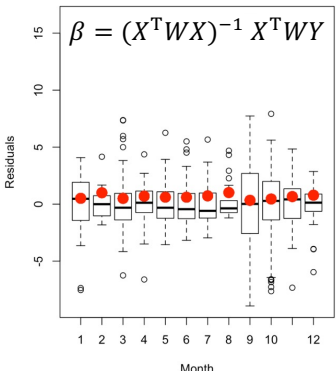


W is the reciprocal of a function of the variance-covariance matrix, but this function can take different forms (e.g., square root of residuals) or more complex structures. Using the reciprocal, specimens (within months here) with large residual will influence less the regression.

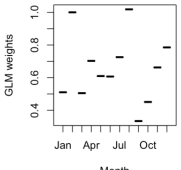
32

Variance changes as a function of Month & Weights are set inversely (reciprocal) to that variance

$$\beta = (X^T W X)^{-1} X^T W Y$$

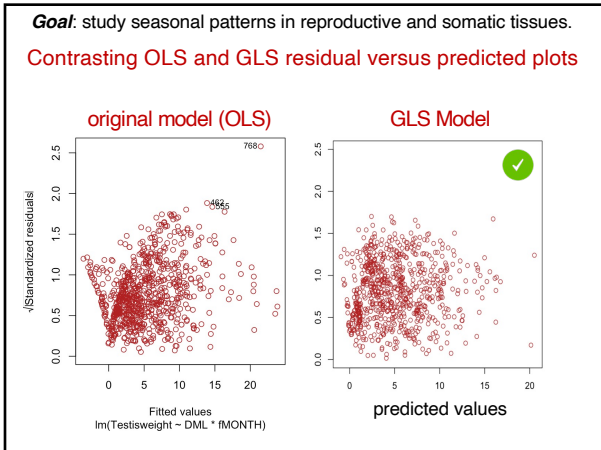


Weights

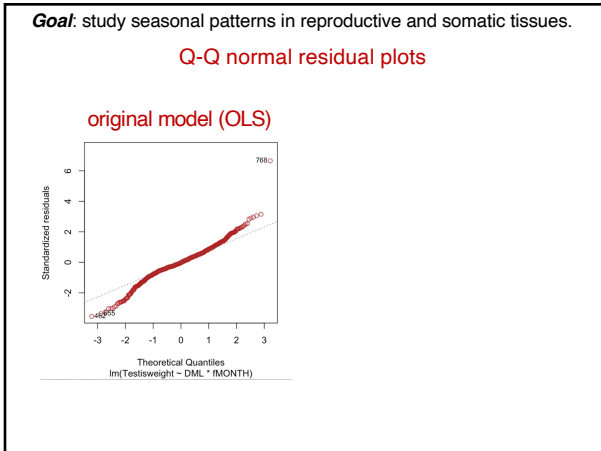


The weight of each individual is reciprocal to the residual variance of the month in which it was sampled.

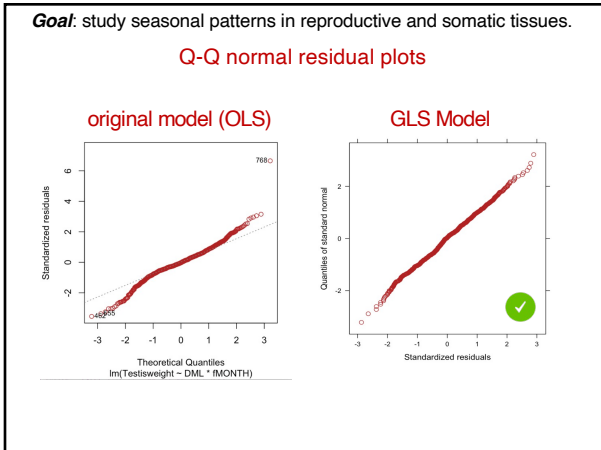
33



34



35



36

Seasonal patterns of investment in reproductive and somatic tissues in the squid *Loligo forbesi*

Jennifer M. Smith^{1,2}, Graham J. Pierce¹, Alain F. Zuur² and Peter R. Boyle¹

¹ Department of Zoology, School of Biological Sciences, University of Aberdeen, Tillydrone Avenue, Aberdeen AB24 2TZ, UK
² Highland Statistics Ltd., 6 Laverock Road, Newburgh, Aberdeenshire, AB41 6FN, UK

Goal: study seasonal patterns in reproductive and somatic tissues.

In which month there is more investment (proportionally to amount of somatic tissues) in reproduction?

testis weight (mg)

dorsal mantle length (DML; mm)
(proxy for somatic tissue)

month 1
month 2

Aquat. Living Resour. 18, 341-351 (2005)
 © EEP Sciences, IPREMER, IRD 2005
 DOI: 10.1051/aqr/2005183
 www.edpsciences.org/aqr

37

Goal: study seasonal patterns in reproductive and somatic tissues.

ANOVA results for GLS model

```
> anova(M.gls)
Denom. DF: 744
```

	numDF	F-value	p-value
(Intercept)	1	3615.591	<.0001
DML	1	1648.534	<.0001
fMONTH	11	76.560	<.0001
DML : fMONTH	11	28.592	<.0001

TestisWeight = constant + β_1 DML + β_2 Month + β_3 (DML \times Month) + e

38

Interaction between dorsal mantle length (DML) and month indicating clear differences in reproductive investment among months (seasons)

testis weight

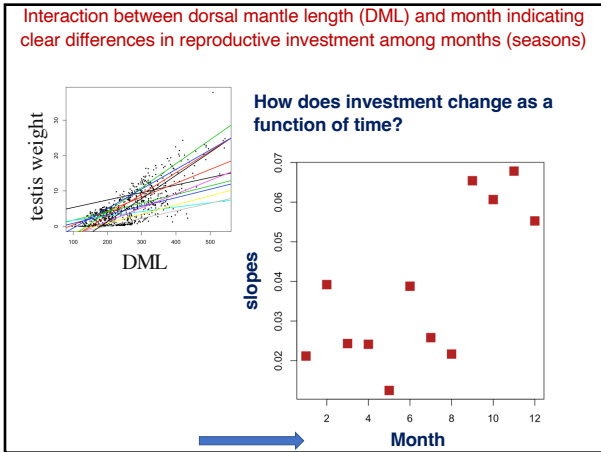
DML

In which month there is more investment in reproduction?

```
> anova(M.gls)
Denom. DF: 744
```

	numDF	F-value	p-value
(Intercept)	1	3615.591	<.0001
DML	1	1648.534	<.0001
fMONTH	11	76.560	<.0001
DML : fMONTH	11	28.592	<.0001

39



40

Important points

There many reasons and ways in which residual variance can change and the types of function (e.g., square root or more complex functions or structures).

We can apply different structures and pick the one that best fit the data (next lecture).

GLS per se is not a mixed model as we will discuss this issue later in details! But they are really important and key to understand variance heterogeneity; and are often used in mixed-models.

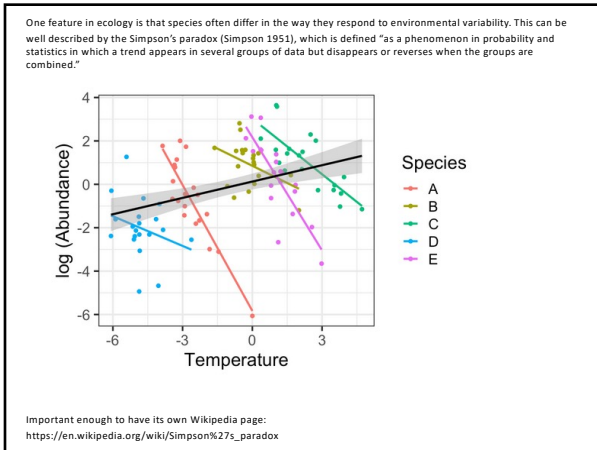
41

Next: a quick look into the general goals of a mixed model using Simpson's paradox.

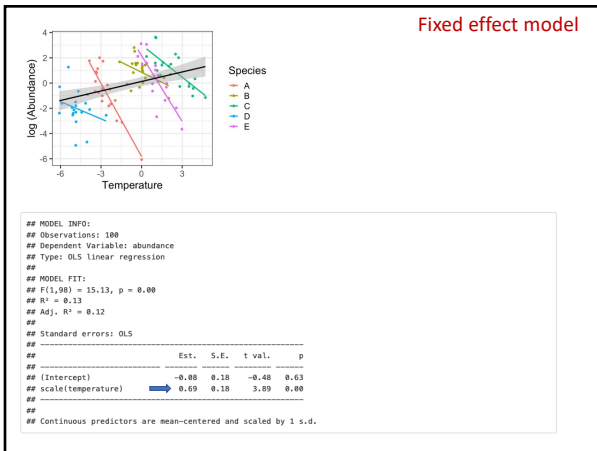
"A phenomenon in probability and statistics in which a trend appears in several groups of data but disappears or reverses when the groups are combined."

Important enough to have its own Wikipedia page:
https://en.wikipedia.org/wiki/Simpson%27s_paradox

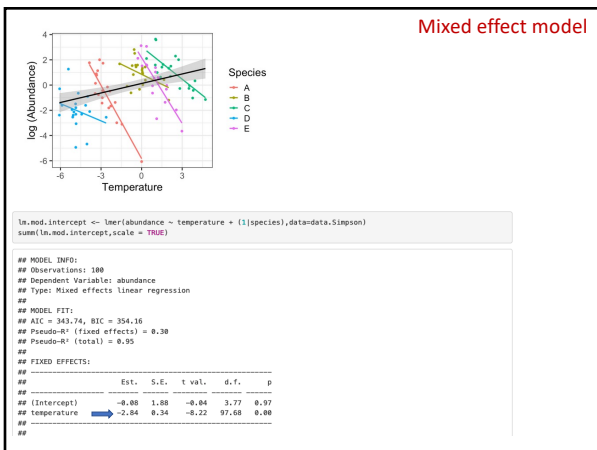
42



43



44



45