Master of Cognitive Science

Data Science Course

# Linear model II

Professor: Moran Steven

Lecturer: Maiolini Marco

Lecture 9:  27/April/2022

# Outline

- Part 1: Book report & discussion (15 minutes)

- Part 2: Mixed linear models (45 minutes)

*Break (15 minutes)*

- Part 3: Practical

# On the previous lecture...

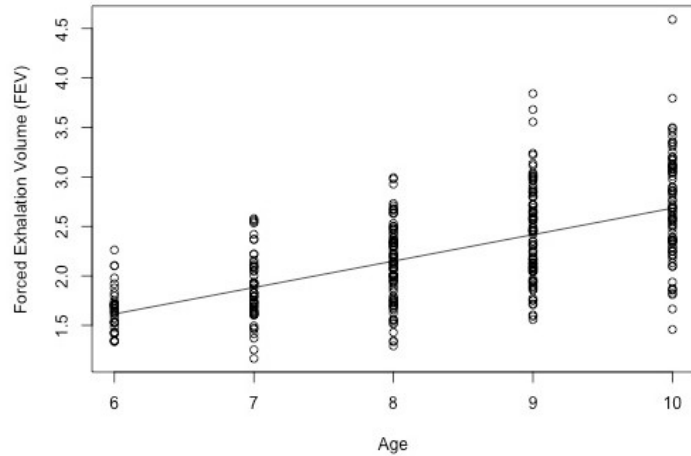| Model | Variables | Distribution | R code |
|---|---|---|---|
| Linear Regression | $Y = b_0 + b_1x$ | Normal | *lm(formula, data)* |
| General Linear Models | $Y = b_0 + b_1x_1 + b_2x_2 + ...$ | Normal | *lm(formula, data)* |
| Generalized Linear Models (GLM) | $Y = b_0 + b_1x_1 + b_2x_2 + ...$ | Any | *glm(formula, family, data)* |

# Linear Models (LMs)

In most models that we seen (LM or GLM) we were interested to quantifying the exact effect of each explanatory variable.

# Linear Models (LMs)

In most models that we seen (LM or GLM) we were interested to quantifying the exact effect of each explanatory variable.
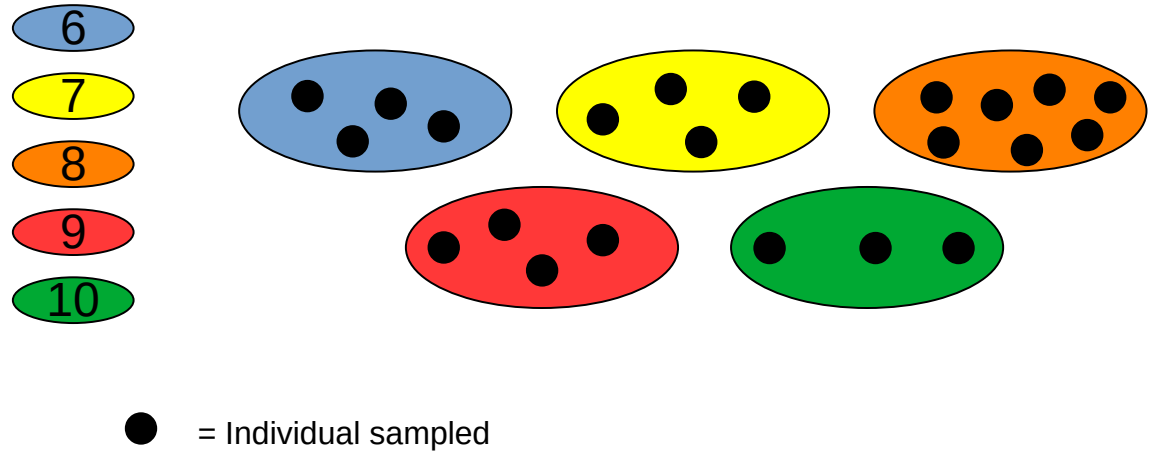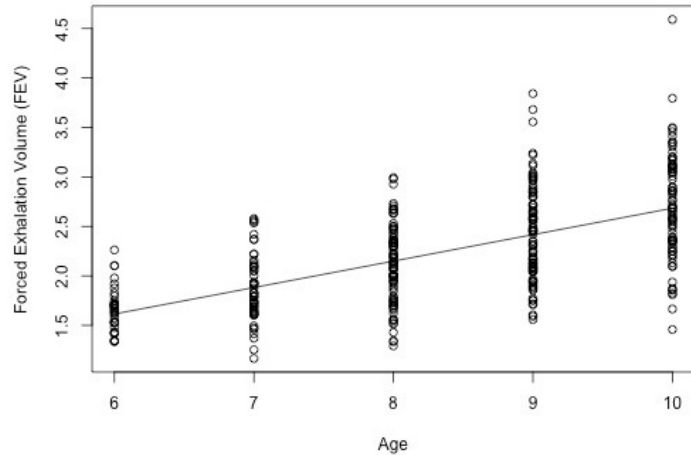
Differences between age of FEV

# Linear Models (LMs)

In most models that we seen (LM or GLM) we were interested to quantifying the exact effect of each explanatory variable.

Differences between age of FEV



The samples of each level are selected randomly, but the levels (Age) of X are not randomly chosen

# Linear Models (LMs)

In most models that we seen (LM or GLM) we were interested to quantifying the exact effect of each explanatory variable.

Differences between sex in tool use acquisition in Bonobos

Boose et al., 2013; "Sex differences in tool use acquisition in bonobos (pan paniscus)", American Journal of Primatology.
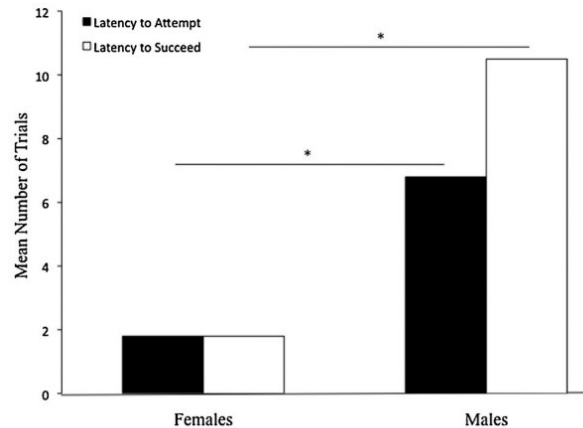
# Linear Models (LMs)

In most models that we seen (LM or GLM) we were interested to quantifying the exact effect of each explanatory variable.

Differences between sex in tool use acquisition in Bonobos



Boose et al., 2013; "Sex differences in tool use acquisition in bonobos (pan paniscus)", American Journal of Primatology.

# Linear Models (LMs)

In most models that we seen (LM or GLM) we were interested to quantifying the exact effect of each explanatory variable.
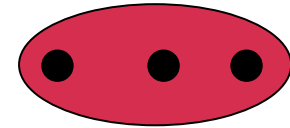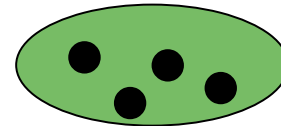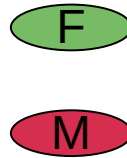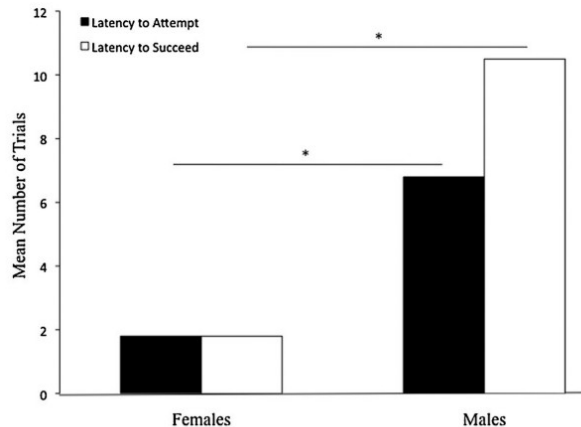
Differences between sex in tool use acquisition in Bonobos



= Individual sampled

The samples of each level are selected randomly, but the levels (Sex) of X are not randomly chosen

Boose et al., 2013; "Sex differences in tool use acquisition in bonobos (pan paniscus)", American Journal of Primatology.

# Linear Models (LMs)

The levels of our study (Age, Sex) are called <u>FIXED EFFECT.</u>

- Each level of <u>Fixed Factor</u> is not derived randomly among infinite possibilities

- In my experiment, I took all the possible levels of <u>Fixed Factor</u>

- If I change the levels of <u>Fixed Factor</u>, I have to change the Hypothesis and repeat the experiment

# Linear Mixed Models (GLMMs)

Most of Biological and Behavioural data involves RANDOM EFFECTS, whose purpose is instead to quantify the variation among units.

# Linear Mixed Models (GLMMs)

Most of Biological and Behavioural data involves RANDOM EFFECTS, whose purpose is instead to quantify the variation among units.

Differences between age of FAV, recorded in 2 schools.

# Linear Mixed Models (GLMMs)

Most of Biological and Behavioural data involves RANDOM EFFECTS, whose purpose is instead to quantify the variation among units.

Differences between age of FAV, recorded in 2 schools.



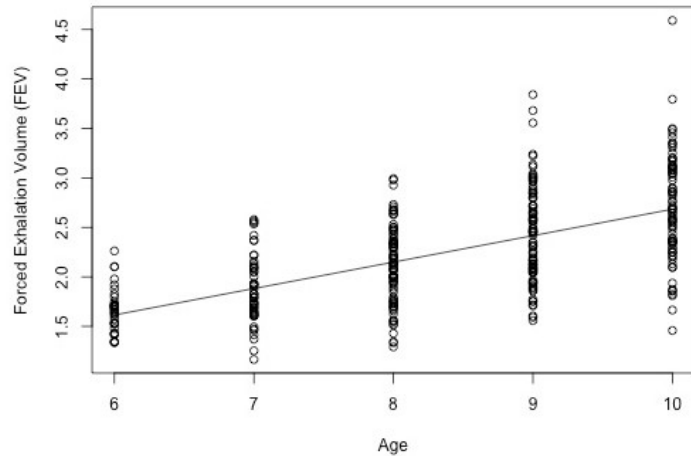The samples of each level are selected randomly, as the Schools

# Linear Mixed Models (GLMMs)

Most of Biological and Behavioural data involves RANDOM EFFECTS, whose purpose is instead to quantify the variation among units.

Differences between sex in tool use acquisition in Bonobos, in 3 different groups

Boose et al., 2013; "Sex differences in tool use acquisition in bonobos (pan paniscus)", American Journal of Primatology.

# Linear Mixed Models (GLMMs)

Most of Biological and Behavioural data involves RANDOM EFFECTS, whose purpose is instead to quantify the variation among units.
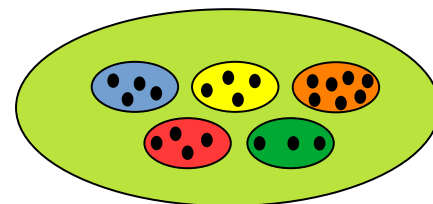
Differences between sex in tool use acquisition in Bonobos, in 3 different groups



Boose et al., 2013; "Sex differences in tool use acquisition in bonobos (pan paniscus)", American Journal of Primatology.

# Linear Mixed Models (GLMMs)

Most of Biological and Behavioural data involves RANDOM EFFECTS, whose purpose is instead to quantify the variation among units.
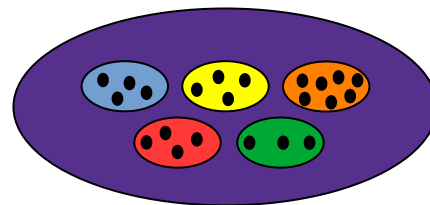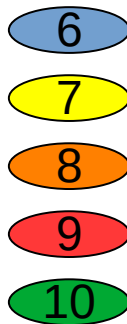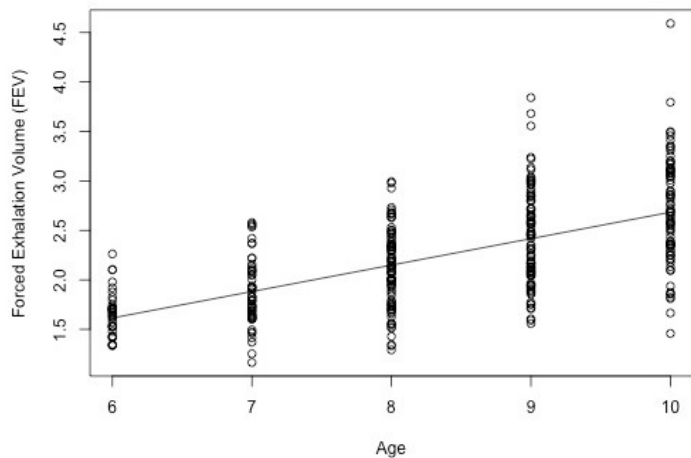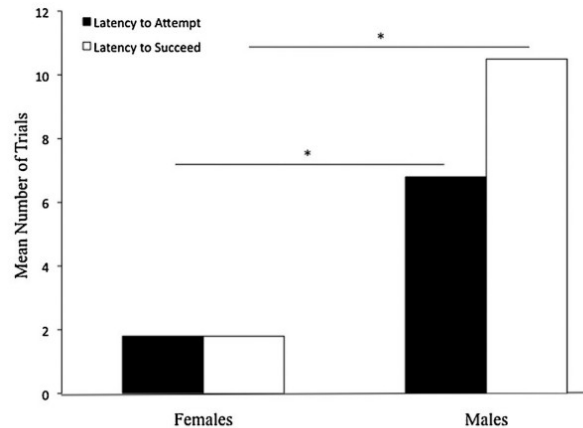
Differences between sex in tool use acquisition in Bonobos, in 3 different groups



<u>The samples of each level are selected randomly, as the groups</u>

Boose et al., 2013; "Sex differences in tool use acquisition in bonobos (pan paniscus)", American Journal of Primatology.

# Linear Mixed Models (GLMMs)

The levels that quantify our variance are called <u>RANDOM EFFECT.</u>

- Each level of <u>Random Effect</u> is collected randomly among infinite possibilities

- In my experiment, I could not take all the possible levels of <u>Random Effect</u>

- If I change the levels of <u>Random Effect</u>, I have not to change the Hypothesis.

# Linear Mixed Models (LMMs)

Example

You want to measure how much it is the food intake for a species in two habitat types (Forest and Plantation), measured in three different sites.

| food_intake | intake_rate | habitat | site |
|---|---|---|---|
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteC |
| 9 | 0.9 | Plantation | siteC |
| 18 | 1.8 | Forest | siteC |
| 17 | 1.7 | Forest | siteC |
| 3 | 0.3 | Forest | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteB |
| 21 | 2.1 | Forest | siteB |
| 20 | 2 | Forest | siteB |
| 4 | 0.4 | Forest | siteC |
| 3 | 0.3 | Plantation | siteC |
| 15 | 1.5 | Forest | siteB |
| 9 | 0.9 | Forest | siteB |
| 13 | 1.3 | Forest | siteB |
| 0 | 0 | Plantation | siteC |
| 2 | 0.2 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 2 | 0.2 | Plantation | siteB |
| 5 | 0.5 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 0 | 0 | Plantation | siteB |

# Linear Mixed Models (LMMs)

Example

You want to measure how much it is the food intake for a species in two habitat types
(Forest and Plantation), measured in three different sites.

Fixed Factor ⟶ Habitat

| food_intake | intake_rate | habitat | site |
|---|---|---|---|
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteC |
| 9 | 0.9 | Plantation | siteC |
| 18 | 1.8 | Forest | siteC |
| 17 | 1.7 | Forest | siteC |
| 3 | 0.3 | Forest | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteB |
| 21 | 2.1 | Forest | siteB |
| 20 | 2 | Forest | siteB |
| 4 | 0.4 | Forest | siteC |
| 3 | 0.3 | Plantation | siteC |
| 15 | 1.5 | Forest | siteB |
| 9 | 0.9 | Forest | siteB |
| 13 | 1.3 | Forest | siteB |
| 0 | 0 | Plantation | siteC |
| 2 | 0.2 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 2 | 0.2 | Plantation | siteB |
| 5 | 0.5 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 0 | 0 | Plantation | siteB |

# Linear Mixed Models (LMMs)

Example

You want to measure how much it is the food intake for a species in two habitat types (Forest and Plantation), measured in three different sites.

Fixed Factor ──────────▶ Habitat

Random Effect ──────────▶ Sites

| food_intake | intake_rate | habitat | site |
|---|---|---|---|
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteC |
| 9 | 0.9 | Plantation | siteC |
| 18 | 1.8 | Forest | siteC |
| 17 | 1.7 | Forest | siteC |
| 3 | 0.3 | Forest | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteB |
| 21 | 2.1 | Forest | siteB |
| 20 | 2 | Forest | siteB |
| 4 | 0.4 | Forest | siteC |
| 3 | 0.3 | Plantation | siteC |
| 15 | 1.5 | Forest | siteB |
| 9 | 0.9 | Forest | siteB |
| 13 | 1.3 | Forest | siteB |
| 0 | 0 | Plantation | siteC |
| 2 | 0.2 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 2 | 0.2 | Plantation | siteB |
| 5 | 0.5 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 0 | 0 | Plantation | siteB |

# Linear Mixed Models (LMMs)

Example

You want to measure how much it is the food intake for a species in two habitat types
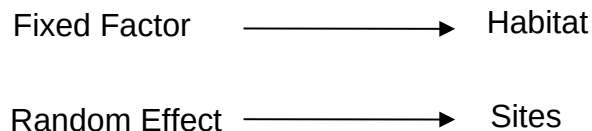(Forest and Plantation), measured in three different sites.

Fixed Factor ⟶ Habitat

Random Effect ⟶ Sites

| food_intake | intake_ra te | habitat | site |
|---|---|---|---|
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteC |
| 9 | 0.9 | Plantation | siteC |
| 18 | 1.8 | Forest | siteC |
| 17 | 1.7 | Forest | siteC |
| 3 | 0.3 | Forest | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteB |
| 21 | 2.1 | Forest | siteB |
| 20 | 2 | Forest | siteB |
| 4 | 0.4 | Forest | siteC |
| 3 | 0.3 | Plantation | siteC |
| 15 | 1.5 | Forest | siteB |
| 9 | 0.9 | Forest | siteB |
| 13 | 1.3 | Forest | siteB |
| 0 | 0 | Plantation | siteC |
| 2 | 0.2 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 2 | 0.2 | Plantation | siteB |
| 5 | 0.5 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 0 | 0 | Plantation | siteB |

In R:

*lmer(Response variable ~ Fixed Factor + (1|Random), data = dataset )*

# Linear Mixed Models (LMMs)

Example

You want to measure how much it is the food intake for a species in two habitat types
(Forest and Plantation), measured in three different sites.

Fixed Factor ——————→ Habitat

Random Effect ——————→ Sites

*lmer(intake_rate ~ habitat + (1|site), data = … )*

In R:

*lmer(Response variable ~ Fixed Factor + (1|Random), data = dataset )*

| food_intake | intake_rate | habitat | site |
|---|---|---|---|
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteC |
| 9 | 0.9 | Plantation | siteC |
| 18 | 1.8 | Forest | siteC |
| 17 | 1.7 | Forest | siteC |
| 3 | 0.3 | Forest | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteB |
| 21 | 2.1 | Forest | siteB |
| 20 | 2 | Forest | siteB |
| 4 | 0.4 | Forest | siteC |
| 3 | 0.3 | Plantation | siteC |
| 15 | 1.5 | Forest | siteB |
| 9 | 0.9 | Forest | siteB |
| 13 | 1.3 | Forest | siteB |
| 0 | 0 | Plantation | siteC |
| 2 | 0.2 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 2 | 0.2 | Plantation | siteB |
| 5 | 0.5 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 0 | 0 | Plantation | siteB |

# Linear Mixed Models (LMMs)

Example

You want to measure how much it is the food intake for a species in two habitat types (Forest and Plantation), measured in three different sites.

Fixed Factor $\longrightarrow$ Habitat

Random Effect $\longrightarrow$ Sites

*lmer(intake_rate ~ habitat + (1|site), data = … )*

⚠️ Remember to check if the residuals are normally distributed

*A ← residuals(lmer(intake_rate ~ habitat + (1|site), data = … ))*
*Shapiro.test(A)*

In R:

*lmer(Response variable ~ Fixed Factor + (1|Random), data = dataset )*

| food_intake | intake_rate | habitat | site |
|---|---|---|---|
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteA |
| 7 | 0.7 | Plantation | siteC |
| 9 | 0.9 | Plantation | siteC |
| 18 | 1.8 | Forest | siteC |
| 17 | 1.7 | Forest | siteC |
| 3 | 0.3 | Forest | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 2 | 0.2 | Forest | siteA |
| 2 | 0.2 | Plantation | siteA |
| 0 | 0 | Plantation | siteA |
| 0 | 0 | Forest | siteB |
| 21 | 2.1 | Forest | siteB |
| 20 | 2 | Forest | siteB |
| 4 | 0.4 | Forest | siteC |
| 3 | 0.3 | Plantation | siteC |
| 15 | 1.5 | Forest | siteB |
| 9 | 0.9 | Forest | siteB |
| 13 | 1.3 | Forest | siteB |
| 0 | 0 | Plantation | siteC |
| 2 | 0.2 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 2 | 0.2 | Plantation | siteB |
| 5 | 0.5 | Plantation | siteB |
| 1 | 0.1 | Plantation | siteB |
| 0 | 0 | Plantation | siteB |

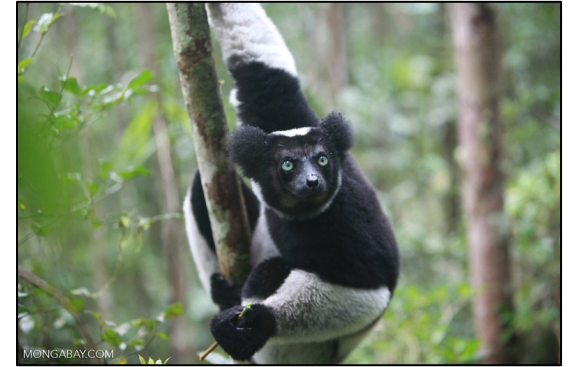# Generalized Linear Mixed Models (GLMMs)

GLMMs combine the proprieties of Linear Mixed Models (LMM) and Generalized Linear Models (GLM), which handle non normal data.

- Have <u>Random Effect</u>

- Can handle data with poisson, binomial or normal distribution
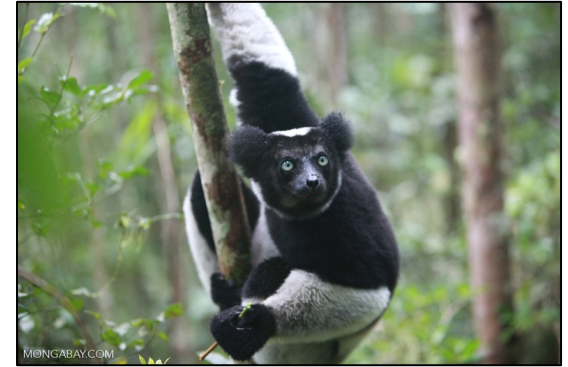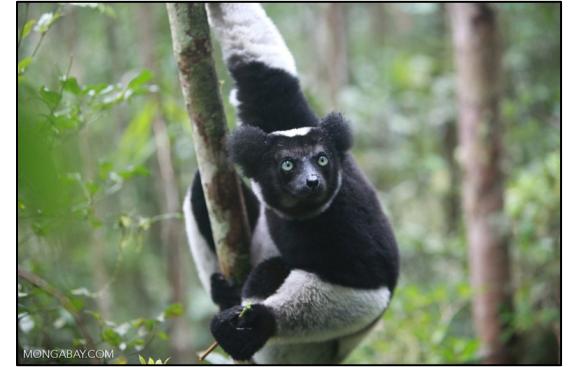
- Can handle interactions and additive effects

# Generalized Linear Mixed Models (GLMMs)

Example

You want to check if there is a rhythmic difference among context in the indris' song.
You  take the data from different individuals, groups and songs.

# Generalized Linear Mixed Models (GLMMs)

Example

You want to check if there is a rhythmic difference among context in the indris' song.
You take the data from different individuals, groups and songs.

Response ⟶ Rhythm

# Generalized Linear Mixed Models (GLMMs)

Example

You want to check if there is a rhythmic difference among context in the indris' song.
You take the data from different individuals, groups and songs.
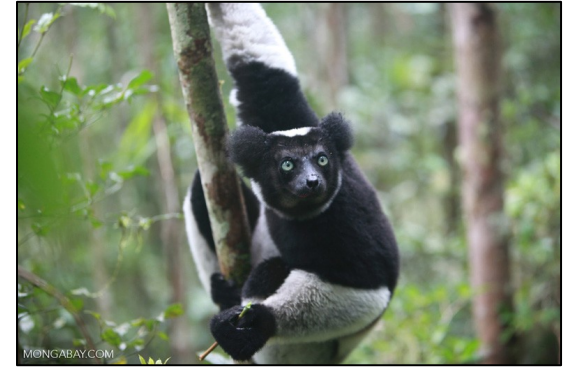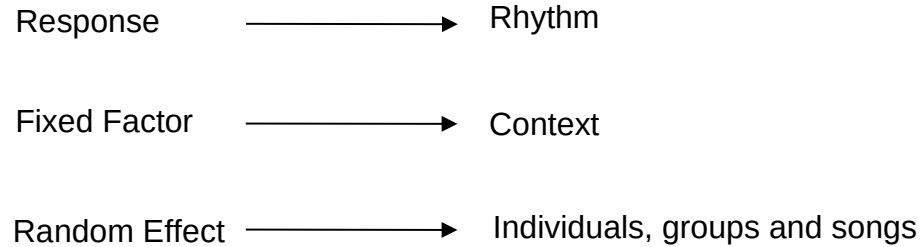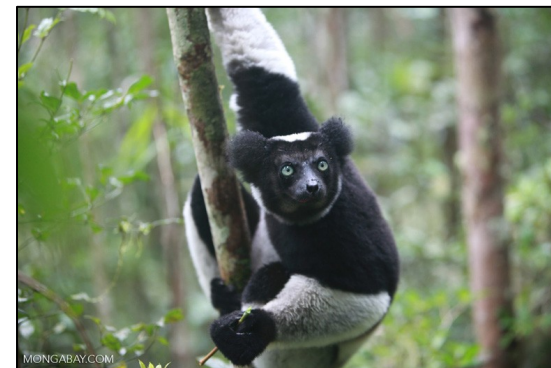
Response $\longrightarrow$ Rhythm

Fixed Factor $\longrightarrow$ Context

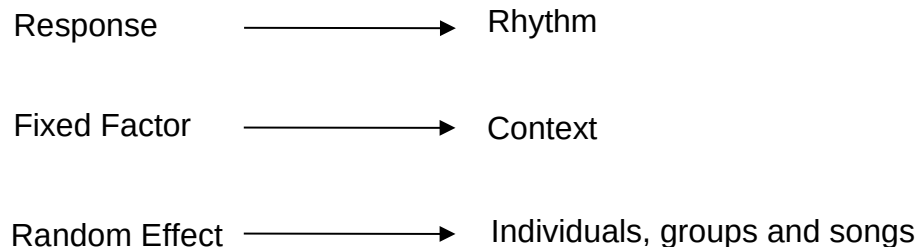# Generalized Linear Mixed Models (GLMMs)

Example

You want to check if there is a rhythmic difference among context in the indris' song.
You take the data from different individuals, groups and songs.

Response ⟶ Rhythm

Fixed Factor ⟶ Context

Random Effect ⟶ Individuals, groups and songs

# Generalized Linear Mixed Models (GLMMs)

Example

You want to check if there is a rhythmic difference among context in the indris' song.
You take the data from different individuals, groups and songs.
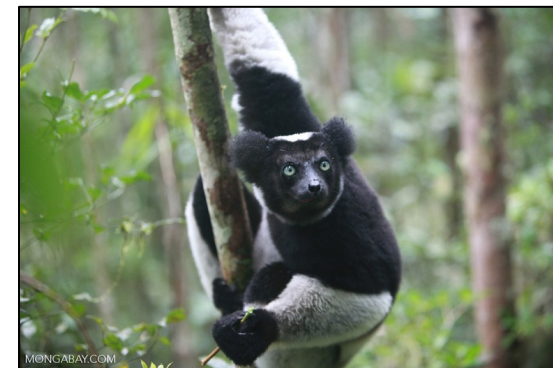
Response $\longrightarrow$ Rhythm

Fixed Factor $\longrightarrow$ Context

Random Effect $\longrightarrow$ Individuals, groups and songs



In R:

*lmer(Response variable ~ Fixed Factor + (1|RE1) + (1|RE2), family = …, data = dataset)*

# Generalized Linear Mixed Models (GLMMs)

Example

You want to check if there is a rhythmic difference among context in the indris' song.
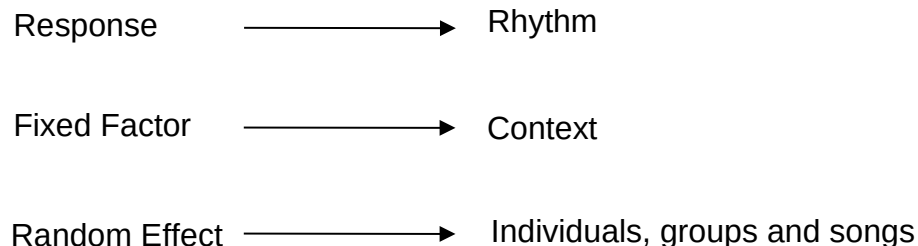You take the data from different individuals, groups and songs.

Response     →     Rhythm

Fixed Factor     →     Context

Random Effect     →     Individuals, groups and songs



*lmer(rhythm ~ context + (1|individual) + (1|group) + (1|song), family = …,*
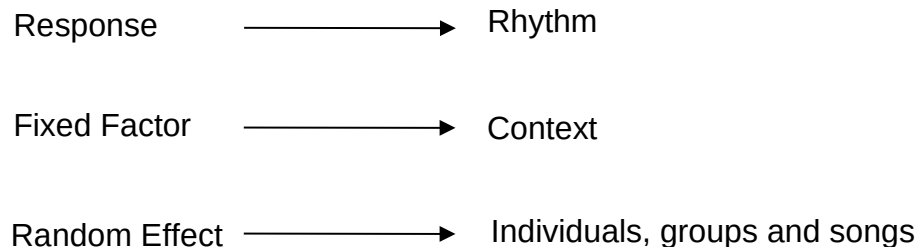*data = … )*

In R:

*lmer(Response variable ~ Fixed Factor + (1|RE1) + (1|RE2), family = …, data = dataset)*

# Generalized Linear Mixed Models (GLMMs)
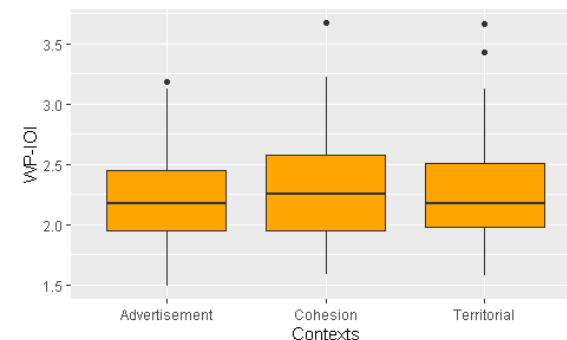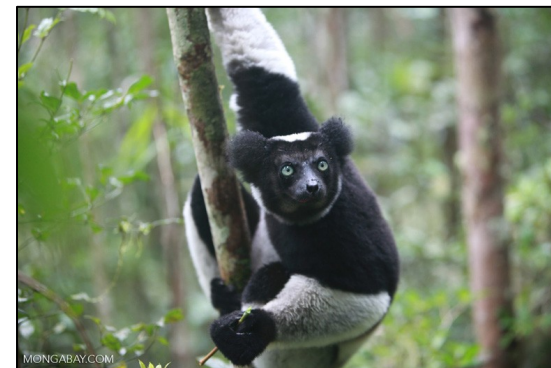
Example

You want to check if there is a rhythmic difference among context in the indris' song. You take the data from different individuals, groups and songs.

Response ⟶ Rhythm

Fixed Factor ⟶ Context

Random Effect ⟶ Individuals, groups and songs

*lmer(rhythm ~ context + (1|individual) + (1|group) + (1|song), family = …, data = … )*

In R:

*lmer(Response variable ~ Fixed Factor + (1|RE1) + (1|RE2), family = …, data = dataset)*

# Model Selection

**Ockham razor**

# Model Selection

**Ockham's razor** $\longrightarrow$ *"Pluralitas non est ponenda sine neccesitate"*

Between two theories that have same predictions, <u>the simpler the better.</u>

# Model Selection

**Ockham's razor** ⟶ *"Pluralitas non est ponenda sine neccesitate"*

Between two theories that have same predictions, <u>the simpler the better.</u>

When comparing different models we want to find out which model explain better my data regardless the **Individual independent variables** in the model.

It doesn't mean that there is a right and a wrong model!

# Model Selection

**Ockham's razor** ⟶ *"Pluralitas non est ponenda sine neccesitate"*

Between two theories that have same predictions, <u>the simpler the better.</u>

When comparing different models we want to find out which model explain better my data regardless the **Individual independent variables** in the model.

⊘ It doesn't mean that there is a right and a wrong model!

Akaike's information criteria (AIC):

> AIC estimates model complexity. It works estimating the expected performance of model's predictions, for that scope it use observed data and hypothetical sample generated by the same model.
>
> The best model show the smallest value; a difference within **4 - 7** units indicate less support, a difference over **10** indicate that the worse model can be omitted

# Model Selection

**Ockham's razor** ⟶ *"Pluralitas non est ponenda sine neccesitate"*

Between two theories that have same predictions, <u>the simpler the better.</u>

When comparing different models we want to find out which model explain better my data regardless the **Individual independent variables** in the model.

⚠ It doesn't mean that there is a right and a wrong model!

Akaike's information criteria (AIC):

> AIC estimates model complexity. It works estimating the expected performance of model's predictions, for that scope it use observed data and hypothetical sample generated by the same model.
>
> The best model show the smallest value; a difference within **4 - 7** units indicate less support, a difference over **10** indicate that the worse model can be omitted

In R:
*AIC(model)*

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

*Model1 ← lmer(water_infection ~ river_tract + (1|id_rivers), data = ... )*

*Summary(Model1)*

*AIC(Model1)*

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

*Model1 ← lmer(water_infection ~ river_tract + (1|id_rivers), data = ... )*

*Summary(Model1)*

*AIC(Model1)*

AIC BIC logLik deviance df. Resid

4838 4920 -2410 4820 900

**AIC = 4838**

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

*Model1 ← lmer(water_infection ~ river_tract + (1|id_rivers), data = ... )*

*Summary(Model1)*

*AIC(Model1)*

AIC BIC logLik deviance df. Resid

4838 4920 -2410 4820 900

**AIC = 4838**

*Model2 ← lmer(water_infection ~ river_tract + slope + (1|id_rivers), data = ... )*

*Summary(Model2)*

*AIC(Model2)*

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

*Model1 ← lmer(water_infection ~ river_tract + (1|id_rivers), data = ... )*

*Summary(Model1)*

*AIC(Model1)*

AIC BIC logLik deviance df. Resid

4838 4920 -2410 4820 900

**AIC = 4838**

*Model2 ← lmer(water_infection ~ river_tract + slope + (1|id_rivers), data = ... )*

*Summary(Model2)*

*AIC(Model2)*

AIC BIC logLik deviance df. Resid

4810 4911 -2287 4780 899

**AIC = 4810**

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

*Model1 ← lmer(water_infection ~ river_tract + (1|id_rivers), data = ... )*

*Summary(Model1)*

*AIC(Model1)*

AIC BIC logLik deviance df. Resid

4838 4920 -2410 4820 900

**AIC = 4838**

*Model2 ← lmer(water_infection ~ river_tract + slope + (1|id_rivers), data = ... )*

*Summary(Model2)*

*AIC(Model2)*

4838 – 4810 = 28

AIC BIC logLik deviance df. Resid

4810 4911 -2287 4780 899

**AIC = 4810**

# Model Selection

Example

You want to check if there is a difference in water infection difference depending on the river tract sampled.
You sample different rivers.

*Model1 ← lmer(water_infection ~ river_tract + (1|id_rivers), data = ... )*

*Summary(Model1)*

*AIC(Model1)*

AIC BIC logLik deviance df. Resid

4838 4920 -2410 4820 900

**AIC = 4838**

*Model2 ← lmer(water_infection ~ river_tract + slope + (1|id_rivers), data = ... )*

*Summary(Model2)*

*AIC(Model2)*

AIC BIC logLik deviance df. Resid

4810 4911 -2287 4780 899

4838 – 4810 = <u>28</u>

**AIC = 4810**

Model2 is better