

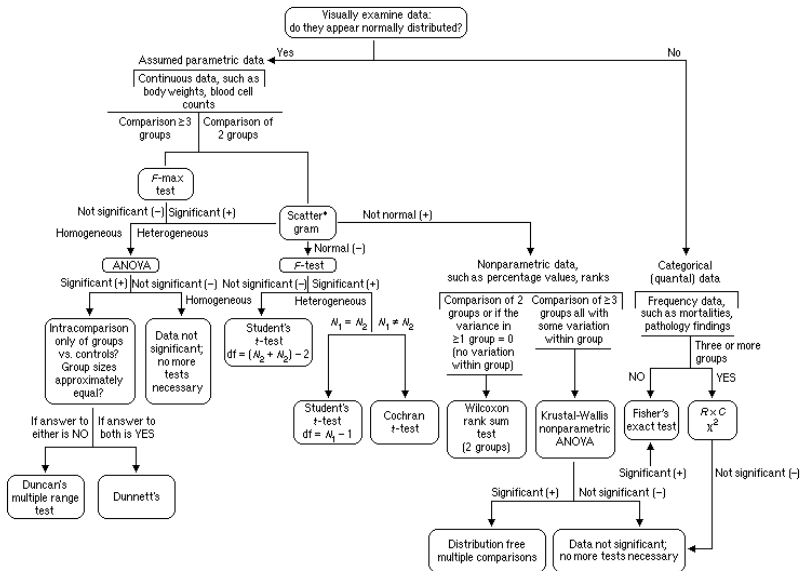
# GLM as a unified framework for data analysis

---

Francisco Rodríguez-Sánchez

<https://frodriguezsanchez.net>

# How I was taught statistics



# So many questions

- **Why** should we really use analysis Y over Z?

# So many questions

- **Why** should we really use analysis Y over Z?
- What if my data are **not Normal**?

# So many questions

- **Why** should we really use analysis Y over Z?
- What if my data are **not Normal**?
- What if they are **not independent**?

# So many questions

- **Why** should we really use analysis Y over Z?
- What if my data are **not Normal**?
- What if they are **not independent**?
- Why am I getting **different p-values** with different tests?

# So many questions

- **Why** should we really use analysis Y over Z?
- What if my data are **not Normal**?
- What if they are **not independent**?
- Why am I getting **different p-values** with different tests?
- What even is a **p-value**?

# So many questions

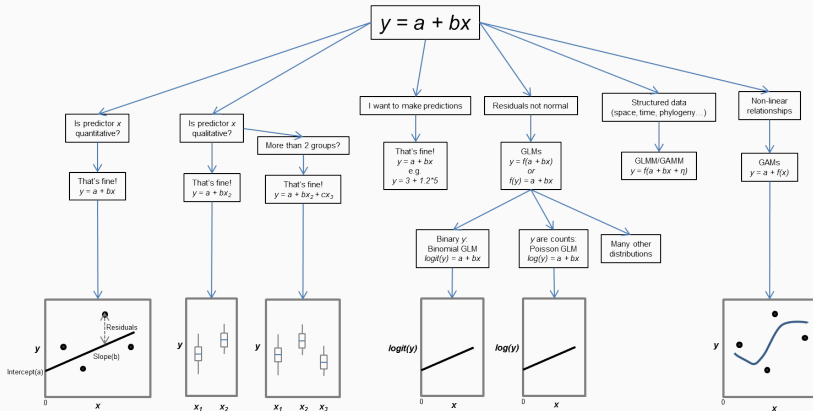
- **Why** should we really use analysis Y over Z?
- What if my data are **not Normal**?
- What if they are **not independent**?
- Why am I getting **different p-values** with different tests?
- What even is a **p-value**?
- How can I take **different factors** into account?



# So many questions

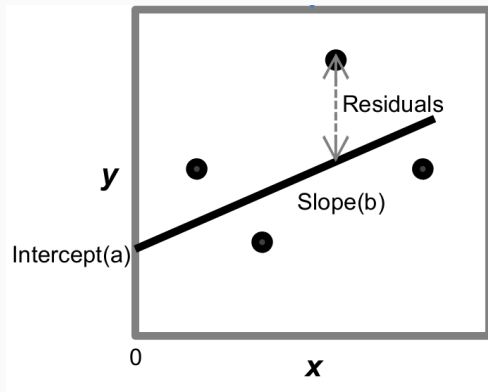
- **Why** should we really use analysis Y over Z?
- What if my data are **not Normal**?
- What if they are **not independent**?
- Why am I getting **different p-values** with different tests?
- What even is a **p-value**?
- How can I take **different factors** into account?
- Can I make **predictions**?

# A unified framework



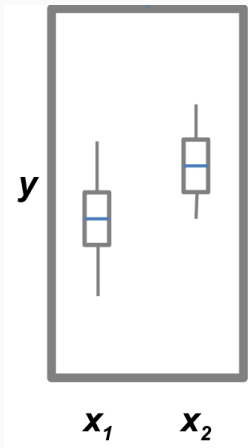
# Linear regression

$$y = a + bx$$



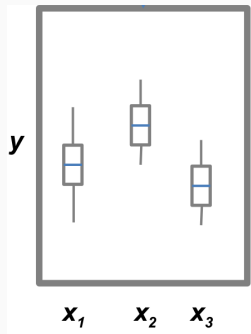
## Is predictor X qualitative?

$$y = a + bx_2$$



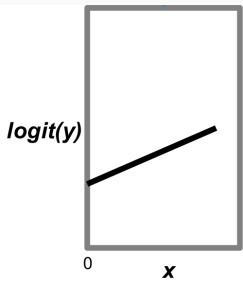
## More than 2 groups?

$$y = a + bx_2 + cx_3$$



## My data (residuals) are not Normal

$$y = f(a + bx)$$

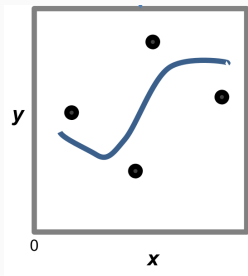


My data are structured (space, time, phylogeny)

$$y = f(a + bx + \eta)$$

# Relationships are not linear

$$y = a + f(x)$$





t-tests

ANOVA

regression

...

are special cases of GLM

With GLM we can analyse  
many different types of data  
using many predictors  
(quantitative & qualitative)

**Unified, coherent framework** for data analysis with many extensions:

- **GLMM** (mixed models): accomodate data structure & variation (space, time, phylogeny)

**Unified, coherent framework** for data analysis with many extensions:

- **GLMM** (mixed models): accomodate data structure & variation (space, time, phylogeny)
- **GAMM** (generalised additive models): non-linear relationships

**Unified, coherent framework** for data analysis with many extensions:

- **GLMM** (mixed models): accomodate data structure & variation (space, time, phylogeny)
- **GAMM** (generalised additive models): non-linear relationships
- **Model-based multivariate** statistics

**Unified, coherent framework** for data analysis with many extensions:

- **GLMM** (mixed models): accomodate data structure & variation (space, time, phylogeny)
- **GAMM** (generalised additive models): non-linear relationships
- **Model-based multivariate** statistics
- **Bayesian** modelling

The Generalised Linear Model (GLM) is a particularly reasonable vantage point on statistical analyses, as **many tests and procedures are special cases** of the GLM. The downside of that (and any other) vantage point is that **we first have to climb it**. There are the morass of unfamiliar terminology, the scree slopes of probability and the cliffs of distributions. **The vista, however, is magnificent**. From the GLM, t-test, ANOVA and regression neatly arrange themselves into regular patterns, and we can see the paths leading towards the horizon: to time series analyses, Bayesian statistics, spatial statistics and so forth.

Dormann 2020