

Statistical research at CBD

At the interface between biology and mathematical sciences

Stefanie Muff

February 16 2022



(<https://quotefancy.com/>)

CBD is a real paradise for a statistician like me.

Research topics

- Quantitative genetics/evolutionary biology
- Movement ecology/habitat selection studies
- Interpretation and use of statistical methods
 - p -values
 - Variable importance in regression models
- Measurement error modeling
- Bayesian statistics

Research topics

- Quantitative genetics/evolutionary biology
- Movement ecology/habitat selection studies
- Interpretation and use of statistical methods
 - p -values
 - Variable importance in regression models
- Measurement error modeling
- Bayesian statistics

→ *All these topics are linked with each other.* How?

Example 1: Heritability is a measurement error problem

- A phenotype (P) can be additively decomposed into a genetic (G) and an environmental (E) component: $P = G + E$.
- The phenotypic variance σ_p^2 can be decomposed accordingly:
 $\sigma_p^2 = \sigma_g^2 + \sigma_e^2$.
- The aim is to estimate σ_g^2 , and in particular *heritability*

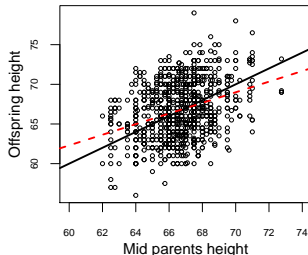
$$h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_e^2} .$$

Example 1: Heritability is a measurement error problem

- A phenotype (P) can be additively decomposed into a genetic (G) and an environmental (E) component: $P = G + E$.
- The phenotypic variance σ_p^2 can be decomposed accordingly:
 $\sigma_p^2 = \sigma_g^2 + \sigma_e^2$.
- The aim is to estimate σ_g^2 , and in particular *heritability*

$$h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_e^2} .$$

- We can estimate h^2 by *mid-parent regression*:



- In the measurement error world, assume we observe $w_i = x_i + u_i$, with $u_i \sim \mathbf{N}(0, \sigma_u^2)$, where x_i is the true variable.
- The aim is to find the slope β_x

$$y_i = \beta_0 + \beta_x x_i + \varepsilon_i ,$$

but *in practice* we regress

$$y_i = \beta_0 + \beta_x^* w_i + \varepsilon_i .$$

- In the measurement error world, assume we observe $w_i = x_i + u_i$, with $u_i \sim \mathcal{N}(0, \sigma_u^2)$, where x_i is the true variable.
- The aim is to find the slope β_x

$$y_i = \beta_0 + \beta_x x_i + \varepsilon_i ,$$

but *in practice* we regress

$$y_i = \beta_0 + \beta_x^* w_i + \varepsilon_i .$$

- Then it is known that we are actually estimating

$$\beta_x^* = \underbrace{\frac{\sigma_x^2}{\sigma_x^2 + \sigma_u^2}}_{=\lambda} \cdot \beta_x .$$

- In the measurement error world, assume we observe $w_i = x_i + u_i$, with $u_i \sim \mathcal{N}(0, \sigma_u^2)$, where x_i is the true variable.
- The aim is to find the slope β_x

$$y_i = \beta_0 + \beta_x x_i + \varepsilon_i ,$$

but *in practice* we regress

$$y_i = \beta_0 + \beta_x^* w_i + \varepsilon_i .$$

- Then it is known that we are actually estimating

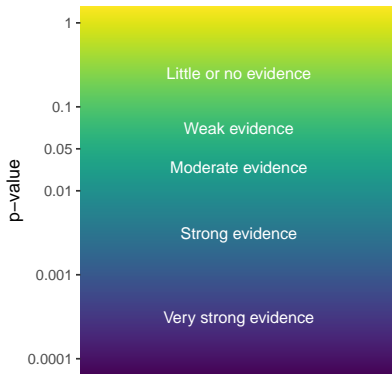
$$\beta_x^* = \underbrace{\frac{\sigma_x^2}{\sigma_x^2 + \sigma_u^2}}_{=\lambda} \cdot \beta_x .$$

- For a slope of $\beta_x = 1$, the attenuation factor $\lambda = h^2$.

Example 2: The interpretation of p -values is relevant in most studies

Opinion

Rewriting results sections in the language of evidence



Huge echo! We received three letters, people wrote three blog posts and the Twitter community had very controversial opinions.

Rewriting results sections in the language of evidence

Citation Data: Trends in Ecology & Evolution, ISSN: 0169-5347, Vol: 37, Issue: 3, Page: 203-210
Publication Year: 2022

2 Citations | 232 Captures | 2 Mentions | 783 Social Media

Metric Options: Counts 1 Year 3 Year

This review has 708 Twitter interactions across 3 URLs.
It has received 123 tweets and 585 retweets.

Sylvaine Giakoumi
@Sylvaine_G

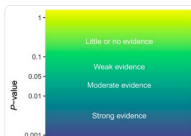
Beyond "binary statistics", great paper!

Gerard Ricardo @GerardRicardo01
How to interpret p-values without dichotomising to 'statistically significant'. #pvalue cell.com/trends/ecology...



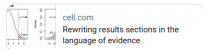
Gerard Ricardo
@GerardRicardo01

How to interpret p-values without dichotomising to 'statistically significant'. #pvalue cell.com/trends/ecology...



KRCEcological
@EcologicalKrc

Rewriting results sections in the language of evidence: Trends in Ecology & Evolution



8:46 PM · Feb 9, 2022

Reply Copy link

[Explore what's happening on Twitter](#)



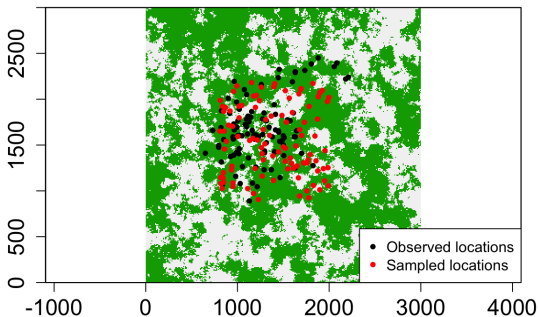
Daniel Lakens
@lakens

I've rarely seen worse advice than trying to interpret p-values directly as measures of evidence. This does not work because of Lindley's paradox, and this graph is misleading (a $p = 0.04$ can be evidence *against* a hypothesis). If you want to talk evidence, use likelihoods.

[Tweet übersetzen](#)

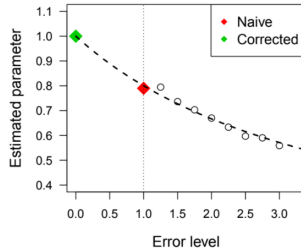
Example 3: Movement ecology

Aim: Accounting for GPS error in movement ecology studies.

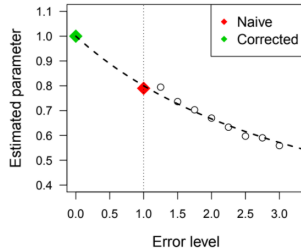


Complex problem, because GPS error propagates into covariate error in the regression of interest.

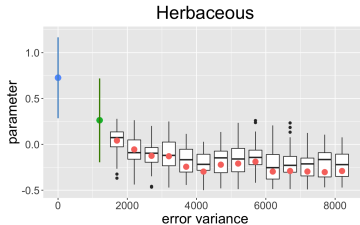
- Idea: Simulation Extrapolation (SIMEX):



- Idea: Simulation Extrapolation (SIMEX):



- Generalize this to GPS error.¹



¹Master thesis by Clara Panchaud

Back to quantitative genetics

- Remember the $P = G + E$ decomposition
- In statistical terms:

$$y_i = \mu + g_i + \varepsilon_i ,$$

where

$$\begin{aligned} (g_1, \dots, g_n)^\top &\sim \mathbf{N}(\mathbf{0}, \sigma_g^2 \cdot \mathbf{G}) , \\ \varepsilon_i &\sim \mathbf{N}(\mathbf{0}, \sigma_e^2) . \end{aligned}$$

Back to quantitative genetics

- Remember the $P = G + E$ decomposition
- In statistical terms:

$$y_i = \mu + g_i + \varepsilon_i ,$$

where

$$\begin{aligned}(g_1, \dots, g_n)^\top &\sim \mathbf{N}(\mathbf{0}, \sigma_g^2 \cdot \mathbf{G}) , \\ \varepsilon_i &\sim \mathbf{N}(\mathbf{0}, \sigma_e^2) .\end{aligned}$$

“Animal model”

Back to quantitative genetics

- Remember the $P = G + E$ decomposition
- In statistical terms:

$$y_i = \mu + g_i + \varepsilon_i ,$$

where

$$\begin{aligned}(g_1, \dots, g_n)^\top &\sim \mathbf{N}(\mathbf{0}, \sigma_g^2 \cdot \mathbf{G}) , \\ \varepsilon_i &\sim \mathbf{N}(\mathbf{0}, \sigma_e^2) .\end{aligned}$$

“Animal model”

- In wild system, information about the relatedness among individuals (\mathbf{G}) is historically derived from pedigrees.

Quantitative genetics with genomic data

- Relatedness information is more and more derived from genomic data → A paradise for statisticians!
- Problem: The $n \times n$ relatedness matrix \mathbf{G} is dense. Computation becomes prohibitive for many individuals.

²PhD project by Janne Hetle

³Meuwissen et al. (2001)

Quantitative genetics with genomic data

- Relatedness information is more and more derived from genomic data → A paradise for statisticians!
- Problem: The $n \times n$ relatedness matrix \mathbf{G} is dense. Computation becomes prohibitive for many individuals.
- Idea: *Marker-based regression*² (equivalent to animal model³):

$$y_i = \mu + \underbrace{\sum_j \underbrace{(SNP_{ij} \cdot u_j)}_{g_{ij}}}_{g_i} + \varepsilon_i .$$

²PhD project by Janne Hetle

³Meuwissen et al. (2001)

Quantitative genetics with genomic data

- Relatedness information is more and more derived from genomic data → A paradise for statisticians!
- Problem: The $n \times n$ relatedness matrix \mathbf{G} is dense. Computation becomes prohibitive for many individuals.
- Idea: *Marker-based regression*² (equivalent to animal model³):

$$y_i = \mu + \underbrace{\sum_j \underbrace{(SNP_{ij} \cdot u_j)}_{g_{ij}}}_{g_i} + \varepsilon_i .$$

- Problem: Too many unknowns, $p \gg n$.

²PhD project by Janne Hetle

³Meuwissen et al. (2001)

We need statistical methods that make us ready for the upcoming challenges with genomic data from wild study systems!

We need statistical methods that make us ready for the upcoming challenges with genomic data from wild study systems!

Ideas (Jannes PhD):

- Use singular value decomposition (SVD) for dimension reduction⁴.
- Use a Bayesian approach to fit a full model with additional fixed and random effects (crucial in wild systems).
- Use a ridge-based shrinkage prior to ensure efficiency.

⁴Ødegaard et al. (2018)

We need statistical methods that make us ready for the upcoming challenges with genomic data from wild study systems!

Ideas (Jannes PhD):

- Use singular value decomposition (SVD) for dimension reduction⁴.
- Use a Bayesian approach to fit a full model with additional fixed and random effects (crucial in wild systems).
- Use a ridge-based shrinkage prior to ensure efficiency.

→ In brief, we are combining several statistical concepts (SVD, Bayes, shrinkage, mixed modeling,...) to hopefully develop an efficient and useful method for evolutionary biology.

⁴Ødegaard et al. (2018)

Summary and Outlook

- Almost unlimited possibilities for statisticians at CBD.
- Several Master and PhD students with exciting projects (Clara Panchaud, Vebjørn Rekkebo, Kenneth Aase, Janne Hetle, Emma Skarstein. . .)
- More ideas than resources.
- (Too?) many collaborations.

A final impression

An email from a CBD postdoc last week (10.2.22):

I am a biologist with very superficial training in math and stats, but I like to dabble in making models and thinking analytically. I always feel like I am reinventing the wheel when I do... Or rather I am making a terrible square wheel that doesn't quite do the job as well as it could.

Which is not a very good feeling!

So my question is this: what can we do about it? Obviously people with your skills are in short supply in biology departments.

I don't know if there would be a way to organize little discussion sessions where you and perhaps other experts can give feedback to people like me that would present their ideas.

That could spur new collaborations, or at least spare us some embarrassment.