

1 **Beyond  $p$ -values:**

2 **Rethinking Statistical Frameworks for Addressing the Replication Crisis**<sup>1</sup>

3  
 4 Fernando Marmolejo-Ramos<sup>#\*</sup>; Jose D. Perezgonzalez<sup>§</sup>; Raydonal Ospina<sup>^</sup>; Freddy  
 5 Hernandez-Barajas<sup>+</sup>; Mauricio Castillo<sup>//</sup>; Rafael Izbicki<sup>™</sup>; Rafael B. Stern<sup>▪</sup>; & Julian  
 6 Tejada<sup>∂</sup>

7  
 8 <sup>#</sup> College of Education, Psychology, and Social Work; Flinders University, Adelaide 5042,  
 9 SA, Australia. Email: [fernando.marmolejoramos@flinders.edu.au](mailto:fernando.marmolejoramos@flinders.edu.au) <sup>\*</sup>: corresponding author.

10 <sup>§</sup> Massey Business School, Massey University, Palmerston North 4442, New Zealand. Email:  
 11 [j.d.perezgonzalez@massey.ac.nz](mailto:j.d.perezgonzalez@massey.ac.nz)

12 <sup>^</sup> Departamento de Estatística, LInCa, Universidade Federal da Bahia, Cidade Universitária,  
 13 Salvador/BA, 40170–110, Brazil. Email: [raydonalmartinez@ufba.br](mailto:raydonalmartinez@ufba.br)

14 <sup>+</sup> Departamento de Estadística; Universidad Nacional de Colombia, sede Medellín; Colombia.  
 15 Email: [fhernanb@unal.edu.co](mailto:fhernanb@unal.edu.co)

16 <sup>//</sup> Center for Basic Research in Psychology (CIBPsi), Faculty of Psychology, Universidad de  
 17 la República; Montevideo 11200, Uruguay. Email: [castillomega@gmail.com](mailto:castillomega@gmail.com)

18 <sup>™</sup> Department of Statistics, Federal University of São Carlos (UFSCar); Brazil. Email:  
 19 [rafaelizbicki@gmail.com](mailto:rafaelizbicki@gmail.com)

20 <sup>▪</sup> Department of Statistics, University of São Paulo; Brazil. Email: [rbstern@gmail.com](mailto:rbstern@gmail.com)

21 <sup>∂</sup> Department of Psychology, Federal University of Sergipe, Brazil. Email:

22 [jtejada@academico.ufs.br](mailto:jtejada@academico.ufs.br)

23

---

<sup>1</sup> Contributed chapter for handbook “Research Handbook on the Replication Crisis” by David Trafimow (to appear in 2025 or thereabouts).

- 24       • *CRedit author statement = FM-R*: Conceptualization, Investigation, Methodology, Visualization,  
25       Writing- Original draft, Writing- Reviewing and Editing, Supervision. *J.D.P, R.O., F.H-B., M. C., R.*  
26       *I., R.S., and J.T.*: Visualization, Writing- Reviewing and Editing, Investigation, Methodology.
- 27       • *Conflict of Interests* = The authors declare that they have no known competing financial interests or  
28       personal relationships that could have appeared to influence the work reported in this paper.
- 29       • *Data availability statement* = all data and R codes are available within this manuscript
- 30       • *Acknowledgements* = none to disclose

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74

## Abstract

The replication crisis across scientific disciplines has prompted critical examination of statistical practices underpinning empirical research. This chapter analyzes how significance testing contributes to replication challenges through the lens of the TASI (Theoretical, Auxiliary, Statistical, Inferential) model. We examine three distinct contexts—inferential knowledge, region of acceptance testing, and statistical learning—to identify limitations in conventional significance testing and propose alternative frameworks. We argue that most problems with significance testing stem from researchers' tendencies to confirm rather than falsify hypotheses, regardless of statistical approach. We introduce REACT (Region of Acceptance Testing) as a stronger alternative to  $p$ -values, offering a structured decision-making process that integrates effect sizes and confidence intervals while explicitly recognizing when data is insufficient for definitive conclusions. Additionally, we propose Generalized Additive Models for Location, Scale, and Shape (GAMLSS) as a comprehensive statistical learning framework that transcends traditional hypothesis testing, focusing instead on descriptive, explanatory, and predictive modeling. The chapter ends with a concise exploration of the connection between data science and artificial intelligence.

*Keywords* = TASI model; REACT methodology; hypothesis testing; replication crisis; statistical learning.



100 *A recap on what significance testing is*

101       Significance testing is a statistical method used to determine whether a property  
102 assumed to be present in a population is compatible with what is observed in a sample  
103 (Kaplan, 2004). In other words, a significance test is conducted to determine whether the data  
104 observed in the sample is consistent with the hypothesis that is being made about a  
105 population. Significance tests are widely used in many scientific fields to assess the evidence  
106 against the null hypothesis by calculating the probability of obtaining the observed result, or a  
107 more extreme result, if the null hypothesis were true (Hernandez-Sampieri et al., 2014;  
108 Trafimow, 2019).

109       This probability is represented by the  $p$ -value, which is compared to a predetermined  
110 threshold, usually .05 (5%; although a lower threshold of .005 has also been proposed by  
111 Benjamin et al., 2018) in Social sciences, to determine whether the result is statistically  
112 significant or not. If the  $p$ -value is less than the predetermined threshold, it is typically  
113 concluded that the observed result is statistically significant, and the null hypothesis is  
114 rejected in favour of the alternative hypothesis. On the other hand, if the  $p$ -value is greater  
115 than the threshold, it is concluded that the observed result is not statistically significant, and  
116 the null hypothesis cannot be rejected (Wasserstein & Lazar, 2016).

117       In the following section, we will examine how questionable research practices (QRPs)  
118 and poor statistical practices in significance testing impact the replicability and  
119 generalizability of research findings.

120

121 *Questionable research practices impact on Transparency, Replicability and*  
122 *Generalizability.*

123       Transparency, replicability, and generalizability are closely tied to the process of  
124 validating scientific knowledge because they are essential characteristics that ensure the

125 verifiability and reliability of scientific claims (Bunge, 2018). As will be developed further,  
126 these three methodological criteria are embedded in the execution and publication of a  
127 scientific study and are interrelated in such a way that difficulties in one can directly or  
128 indirectly impact the others.

129 In recent years, growing concerns about research integrity have emerged, particularly in  
130 light of several high-profile cases of scientific misconduct, fraudulent research, and  
131 Questionable research practices (QRP) (Catanzaro, 2023; John et al., 2012; Meho & Akl,  
132 2025; O'Grady, 2024; Open science collaboration, 2015; Yarkoni, 2022). The former  
133 intentionally aims to circumvent the proper design or development of scientific research,  
134 which leads to a discussion about the ethics of scientific practice and the functioning of the  
135 scientific publishing system; whereas QRP are not necessarily carried out with that intention.  
136 On the other hand, QRP may occur due to omissions in the research design or constraints  
137 during the research process, such as limited access to a statistically appropriate sample or  
138 interruptions in the data collection process (John et al., 2012). Some of the key issues that  
139 contribute to these problems include a lack of transparency in data processing, weak  
140 experimental designs, low statistical power and poor statistical practice (Open science  
141 collaboration, 2015; Poldrack et al., 2017).

142 Some of the QRP are related to not following the typical process of designing a  
143 research and its implementation process, which constitutes a lack of transparency and could  
144 impact the replicability or generalizability of results. The QRP described here can be  
145 addressed by following the methodological steps outlined in the research design. However,  
146 the second and third QRP described could be addressed through statistical strategies. The first  
147 one, Hypothesizing After the Results are Known (HARKing) refers to the practice of  
148 researchers developing hypotheses after they have already analysed their data (Kerr, 1998).  
149 This practice increases the risk of false positives (Type error I) because hypotheses are not

150 based on a priori knowledge, but rather on the observed data. Another consequence is a  
151 reduction in the replicability of research findings because the hypothesis is not formulated  
152 based on previous studies (Kerr, 1998; Rubin, 2022).

153 In this way, Sasaki and Yamada (2023) introduced a QRP defined as Sample-size  
154 Planning After the Results are Known (SPARKing). This might be related to concerns about  
155 the study's sample size and that is justified by researchers in a posteriori fashion. As a result,  
156 the experimental design is less transparent and the difficulty observed in relation to the  
157 robustness of the results obtained is overcome with a justification created for that purpose. To  
158 determine sample size in a non-arbitrary manner, statistical power analysis software, such as  
159 G\*Power (Faul et al., 2007) or simulation studies (Kumle & Draschkow, 2021) can be  
160 employed. The next section will further explore good practices in conducting simulation  
161 studies.

162 Another widely reported QRP is the P-hacking, which refers to the selective reporting  
163 or manipulation of statistical analyses to achieve statistically significant results (Head et al.,  
164 2015). This can include running multiple analyses until a significant result is obtained, or  
165 selectively excluding data points or variables to achieve a desired outcome (Andrade, 2021).  
166 In the same line as previously reported QRP, this practice decreases the transparency of the  
167 experimental design and can produce false positives due to looking for patterns or  
168 relationships that are statistically significant, without any a priori hypotheses or a clear  
169 theoretical rationale for testing a specific hypothesis. However, Multiverse analysis (Steege  
170 et al., 2016) serves as an effective methodological strategy to mitigate p-hacking by  
171 standardizing criteria for data exclusion, variable transformation, and group selection. This  
172 approach prevents biased data selection driven by  $p$ -value attainment through the systematic  
173 visualization of all possible data processing decisions and their impact on outcomes. By  
174 doing so, it precludes the concealment of variability in findings, enabling researchers to

175 evaluate the robustness of results and detect whether outcomes shift depending on analytical  
176 decisions.

177         These QRPs, along with the concepts explored in this chapter, must be contextualized  
178 within the frameworks of statistical significance testing and studies grounded in inductive  
179 reasoning, where conclusions are derived from observations supported by inferential  
180 statistical methods (Bunge, 2018; Hernandez-Sampieri et al., 2014). It is critical to emphasize  
181 that certain research designs incorporate flexible components—such as sample selection or  
182 hypothesis formulation—which may be adjusted or omitted during the study’s execution.  
183 This consideration is not trivial, as it also necessitates critical reflection on methodological  
184 aspects prior to design—for instance, the balance between internal and ecological validity—  
185 to delineate the boundaries of the knowledge generated and ensure the generalizability of  
186 inferences.

187         In this context, one promising approach to promoting methodological rigor and  
188 addressing some of the challenges that lead some researchers to conduct QRPs, particularly in  
189 areas like sample size determination and evaluating analytical choices, lies in the application  
190 of simulation studies, which will be explored in the subsequent section.

191

### 192 ***Transparency and Replicability in Simulation Studies***

193         Simulation studies involve creating a model or simulation of a system, and then using  
194 that model to generate data that can be analyzed using statistical methods. Simulation studies  
195 are conducted to ascertain an appropriate design for a particular investigation. As previously  
196 mentioned, simulation studies can be employed for sample size planning or to assess the  
197 robustness of statistical methods under varying conditions. Additionally, these studies are  
198 valuable for comparing different statistical analysis, enabling researchers to identify the most

199 appropriate approach for their specific research context (Friedrich & Friede, 2023; Morris et  
200 al., 2019).

201 Transparency and replicability are just as important in simulation studies for estimating  
202 sample sizes, resample data from existing datasets, or serve as standalone studie to determine  
203 the properties of a statistical procedure. To be considered valid, it must be possible for other  
204 researchers to replicate the study using the same model and methods. This requires careful  
205 documentation of the simulation model and parameters, as well as clear reporting of the  
206 statistical methods used to analyze the data (Luijken et al., 2024). Simulation models are  
207 often highly specific to a particular system or set of conditions, and it can be difficult to  
208 generalize the findings to other contexts or populations. To address this issue, researchers can  
209 use sensitivity analyses to examine how changes in the model affect the results. They can also  
210 collaborate with other independent researcher groups to compare the results and produce  
211 more generalizable findings (Pawel et al., 2023).

212

### 213 *Open Science and Collaboration as a Framework for Enhancing Transparency in* 214 *Psychological Studies and AI Research*

215 Contemporary scientific practice relies on collaboration, a necessity driven not only by  
216 the inherently collective nature of scientific inquiry but also by the multistage complexity of  
217 modern research. High-quality investigations often require specialized expertise at distinct  
218 phases of the process: statisticians ensure methodological rigor, computational scientists  
219 design algorithms, psychometricians validate measurement tools, and domain experts  
220 contextualize findings. This interdependence underscores the need for coordinated efforts  
221 across disciplines to address the multifaceted demands of rigorous research.

222 This is critical even in research contexts that do not involve the collection of human  
223 data, such as computational simulation studies or the development of large-scale AI models,

224 where reproducibility remains foundational to validate models and their results. In these  
225 domains, collaborative efforts are indispensable for addressing persistent challenges in  
226 computational reproducibility and technological innovation, ensuring that models and  
227 algorithms can be independently verified, generalized across contexts, and iteratively  
228 improved by the broader scientific community (Gundersen et al., 2018; Stodden et al., 2014).  
229 Moreover, such collaboration inherently promotes transparency and replicability in IA  
230 research, as it necessitates the open sharing of research data, training datasets, code,  
231 methodologies, and findings with both the scientific community and society at large (Haibe-  
232 Kains et al., 2020; Prieto, 2022; Toribio-Flórez et al., 2021).

233 In this way, open science is an approach to research that emphasizes transparency,  
234 collaboration, and accessibility to scientific knowledge. In addition to traditional research  
235 methods, Open science also involves using new technologies and tools, such as open data  
236 repositories (i.e. Open Science Framework; Foster & Deardorff, 2017), collaborative research  
237 platforms (i.e. Overleaf), platforms to store and share code (i.e. GitHub), repositories of  
238 trained machine learning (ML) models (i.e. Hugging face, TensorFlow Hub) and pre-  
239 registration of study protocols (Pownall et al., 2021). Consequently, these technological  
240 infrastructures require workflows that facilitate collaborative efforts. While advanced tools  
241 have emerged to support joint efforts, their effective implementation depends on parallel  
242 workflows that standardize processes, streamline communication, and embed collaboration  
243 into the research design itself. Such frameworks ensure that teams can efficiently leverage  
244 technological resources while adhering to shared protocols for transparency and  
245 reproducibility.

246 As AI models become more complex and powerful, it is important to ensure that they  
247 are transparent, explainable, and ethical (Haibe-Kains et al., 2020; Kapoor & Narayanan,  
248 2023). Open science involves sharing data, software, and research materials to help ensure AI

249 systems are developed in a responsible manner, and that their potential risks and benefits are  
250 fully understood. For example, it would allow the identification and evaluation of errors such  
251 as variable repetition; the presence of low-quality data (i.e., images) within the dataset that  
252 could negatively influence the integrity and accuracy of the obtained results (Roberts et al.,  
253 2021). Also, it is necessary to consider the possibility that the material used contains personal  
254 information, sensitive data, or explicit content, which could have privacy and ethical  
255 implications.

256 Besides, transparency in AI systems has gained new dimensions with recent  
257 technological shifts (Larsson & Heintz, 2020). First, the emergence of reasoning models  
258 enables it to articulate its decision-making steps in natural language during problem-solving.  
259 This innovation allows users to trace the model's logical pathways (e.g., "I will first analyze  
260 the premise, then compare historical precedents..."), overcoming the "black box" reasoning  
261 nature of previous Large Language Models (LLM) and fostering trust through explainable  
262 intermediate steps that are useful to the common user. Secondly, DeepSeek's disruptive  
263 release of a free public API redefines openness standards in LLM (Liu et al., 2024). By  
264 providing unrestricted access to its architecture, DeepSeek disrupts the industry's status quo,  
265 advancing a shift toward open science principles in AI research and development. This dual  
266 transparency creates unprecedented accountability: while procedural explanations empower  
267 end-users to audit AI reasoning, open API availability lets researchers scrutinize model  
268 behavior at scale and develop or adapt bespoke systems for specific applications (Bommasani  
269 et al., 2021, Sapkota et al., 2025).

270

### 271 ***Many-Lab Approaches: Enhancing Generalizability Through Collaborative Research***

272 One of the challenges in generalizing the results of studies involving human data lies in  
273 the potential influence of cultural, economic, educational, and political variables on the

274 treatment variable. In this way, the concept of WEIRD societies (Western, Educated,  
275 Industrialized, Rich, and Democratic) has become increasingly important in recent years, as  
276 researchers have recognized that many research studies are based on samples that are not  
277 representative of the global population (Alves et al., 2022; Henrich et al., 2010). This can  
278 lead to difficulties in the generalization of the findings in fields where the variables to be  
279 investigated depend on sociocultural variables, for example in moral psychology or education  
280 studies (Alves et al., 2022; Yarkoni, 2022). These studies lack representativeness of national  
281 populations due to non-probabilistic sampling frameworks or reliance on convenience  
282 samples drawn from university student populations. In contexts marked by significant  
283 socioeconomic disparities or stratified access to education—common in low- and middle-  
284 income countries—such methodological limitations critically undermine external validity,  
285 rendering findings inapplicable to broader demographic groups (Henrich et al., 2010).

286         One way to address this issue is through the use of Big Team Science, multi-lab/many-  
287 lab approaches. These ones involve conducting large-scale, collaborative studies across  
288 multiple labs and locations, using diverse samples from different populations and cultures.  
289 This can help to ensure that research findings are more generalizable, robust and applicable to  
290 a wider range of populations (Baumgartner et al., 2023; Yarkoni, 2022).

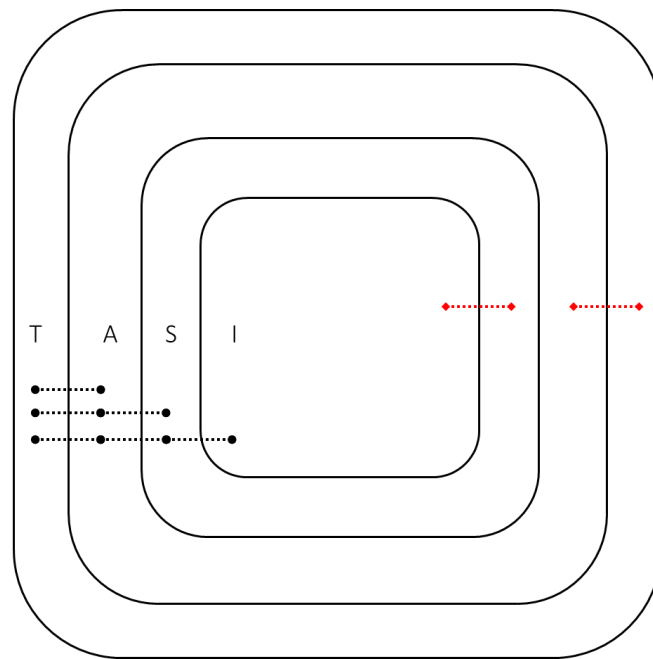
291 Many-lab approaches also help to increase the transparency and replicability of research, as  
292 multiple labs can independently test the same hypotheses and methods. This can help to  
293 identify inconsistencies or errors in the research, and improve the overall quality of the  
294 findings (i.e.: Coles et al., 2022; Morey et al., 2021; Wagenmakers et al., 2016). To  
295 successfully implement Big Team Science and multilab approaches, it is imperative to  
296 establish clear governance frameworks and robust collaboration agreements from the outset,  
297 ensuring that leadership roles, decision-making processes, and conflict resolution  
298 mechanisms are explicitly defined. As previously discussed, the Open Science Framework,

299 encompassing both technological infrastructure and conceptual workflows, provides a  
300 scaffold for collaborative research by standardizing shared protocols. Moreover, fostering an  
301 inclusive, interdisciplinary team through targeted recruitment and equitable authorship  
302 policies (Allen et al., 2014), exemplified by structured systems like the Contributor Roles  
303 Taxonomy (CRediT), not only integrates diverse perspectives but also reinforces  
304 methodological rigor and ultimately enables the production of robust, generalizable scientific  
305 insights.

306         In the next sections, we refer to the TASI model (Trafimow, 2019) as a framework to  
307 categorize and evaluate the different types of assumptions researchers make in testing  
308 theoretical predictions to ensure the validity and reliability of their findings. The TASI model  
309 stands for Theoretical, Auxiliary, Statistical, and Inferential assumptions. This model helps in  
310 understanding and evaluating the robustness of research findings by considering the  
311 following four categories of assumptions: *Theoretical Assumptions*: assumptions related to  
312 the underlying theory being tested; *Auxiliary Assumptions*: assumptions necessary to bridge  
313 the gap between theoretical concepts and actual observations; *Statistical Assumptions*:  
314 assumptions related to the statistical methods and models used in the analysis; and *Inferential*  
315 *Assumptions*: assumptions needed to make statistical inferences about populations from  
316 samples (see Figure 1).

317

318



319

320 *Figure 1.* Representation of the components in the TASI model and their relatedness. T = theoretical  
 321 assumptions, A = auxiliary assumptions, S = statistical assumptions, I = inferential assumptions. The black  
 322 round dotted lines with oval arrows represent three variants of the TASI model; TA, TAS, and the full TASI  
 323 model. The red round dotted lines with diamond arrows indicate potential ad-hoc strong bonds (i.e. tight bonds  
 324 between T and A and S and I).

325

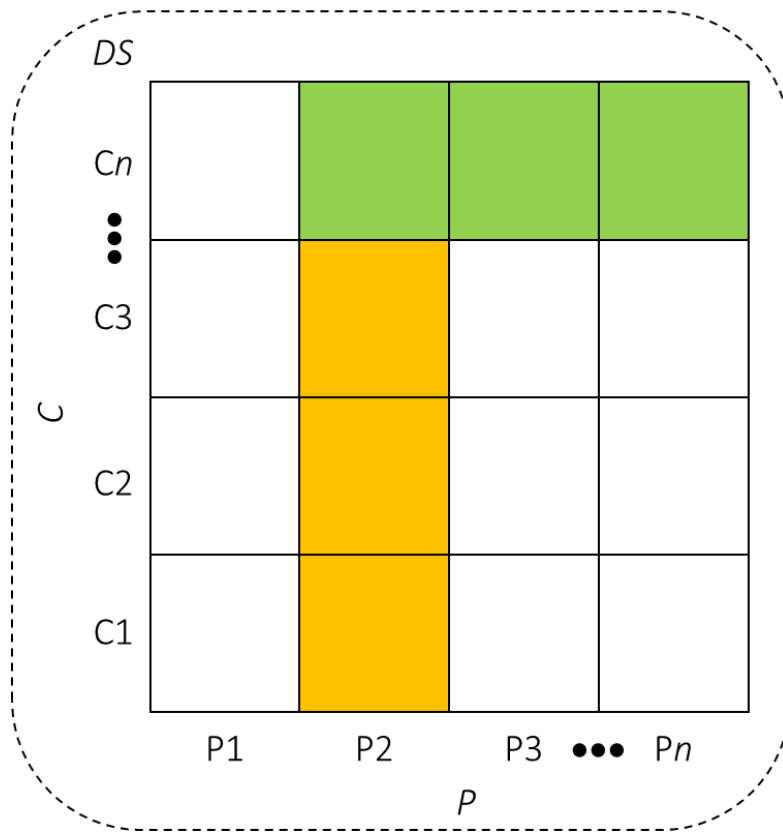
326 The simplest version of the TASI model is the TA model (see Figure 1). This model  
 327 only includes theoretical assumptions and auxiliary assumptions and is used when researchers  
 328 focus on connecting theory to observable phenomena. For example, studying working  
 329 memory processing using cognitive load theory (*T*) plus assumptions about brain activity  
 330 patterns and behavioral responses (*A*) allows researchers to connect abstract theoretical  
 331 concepts about memory and neural processing to concrete, observable measurements without  
 332 requiring complex statistical or inferential assumptions. The TAS model adds statistical  
 333 assumptions to the TA model and is used when researchers need to analyze data but do not  
 334 need to make population inferences. However, it includes choices about what statistics to use  
 335 and how to analyze data. Although more complex than TA, it is simpler than a full TASI  
 336 model. For example, studying working memory processing using cognitive load theory (*T*),



337 plus assumptions about brain activity patterns and behavioral responses (*A*), plus statistical  
338 assumptions about how to analyze brain activation patterns, including choices about  
339 averaging BOLD signals across brain regions, normalizing response times, and using  
340 parametric tests for comparing conditions (*S*) allows researchers to analyze relationships  
341 between memory load and brain activation using appropriate statistical techniques, while  
342 avoiding the additional complexity of population-level inferential assumptions. The full TASI  
343 model, as described by Trafimow (2019), is challenging because it necessitates that all  
344 inferential assumptions about populations and sampling be correct for *p*-values to be  
345 required. This full model may never be perfectly correct in real research, as it is rare for all  
346 assumptions to be entirely accurate.

347       As previously noted, the TA model explicitly links the *T* and *A* components within the  
348 TASI framework. Additionally, there is an ad hoc connection between the *S* and *I*  
349 components, given that both appear to share methodological ties by dealing with the  
350 quantitative analysis of research data. Statistical assumptions guide how data should be  
351 analyzed and determine which statistical tools are appropriate, while inferential assumptions  
352 determine how these analyses can be utilized to draw broader conclusions. These components  
353 are frequently considered together when planning data analysis strategies, and changes in  
354 statistical approaches often necessitate corresponding adjustments in inferential assumptions  
355 (see Figure 1). For example, in statistical learning, the connection between statistical and  
356 inferential components becomes evident when considering a ML model using cross-  
357 validation for performance estimation (cross-validation typically works by splitting the data  
358 into *k* subsets (folds), training the model on *k*-1 folds, testing it on the remaining fold,  
359 repeating this process so each fold serves as the test set once, and averaging the results). The  
360 *S* assumptions involve the stability of data distribution across training and test sets, the  
361 consistency of feature relationships across different data splits, and the appropriateness of

362 chosen performance metrics and validation methods. These *S* assumptions are inherently  
363 linked to their *I* counterparts: assumptions about how representative the training data is of the  
364 target population, how well the model will generalize to unseen data, and whether test set  
365 results truly indicate real-world performance. This interconnection manifests clearly in  
366 practice - if researchers' statistical assumption about data distribution stability is violated,  
367 their inferential assumption about generalization becomes compromised. Similarly, their  
368 choice of cross-validation method directly shapes what they can infer about model  
369 generalization, while their selection of performance metrics determines what conclusions  
370 they can draw about real-world applications. How researchers handle feature relationships  
371 statistically impacts their ability to make meaningful predictions about new data,  
372 demonstrating how statistical choices and inferential conclusions are tightly coupled in  
373 statistical modelling (see Figure 2).

Case A = Model 1<sub>SI</sub>  
 Case B = Model 1<sub>SI</sub> ••• Model *n*<sub>SI</sub>



 Generalisation across populations for a specific context/task  
 Generalisation across contexts/tasks for a specific population

374

375 *Figure 2.* Relationship between statistical and inferential assumptions considering target populations and  
 376 research contexts. DS = design space or space of potential experiments requiring a combination of target  
 377 populations (P) and contexts or tasks (C). In case A, generalizations across populations (green) or contexts  
 378 (orange) result from a specific SI model. In case B, generalizations across populations (green) or contexts  
 379 (orange) result from more than one SI model.

380

381 Next, the TASI model is revisited in the context of the inferential knowledge (including  
 382 significance testing, severity testing, and Bayesian approaches), region of acceptance testing,  
 383 and statistical learning.



409 comprise the first premise of the syllogism, with theory (or particular hypotheses, HA) acting  
 410 as the *Antecedent* of the premise (see Figure 3). The theory serves to set the expected  
 411 outcomes (O) to be observed, thus setting the *Consequent* for the premise too. In a nutshell,  
 412 TA links a theory (or particular hypothesis) with the theoretical results expected from such  
 413 theory. The primary concern for this first premise is its soundness. Not just the theory ought  
 414 to be sound (auxiliary assumptions play a key role here), but the link between theoretical  
 415 *Antecedent* and expected *Consequent* needs to be sound too. If the theory is true, then certain  
 416 outcomes shall be expected. If we encounter outcomes that contradict those expectations, we  
 417 shall deduce the theory is false.

418

If TA   HA	Then, O
	Not O
Thus, no TA   HA	

419

420 *Figure 3.* The *Modus Tollens* inferential process: Theoretical and auxiliary assumptions (TA) help make the  
 421 main premise sound, and the conclusion reliable.

422

### 423 ***Reliable methods for testing***

424 The second premise of *Modus Tollens* is where the testing of the first premise occurs. Here is  
 425 where we find the *Statistical and Inferential assumptions* (SI) of the TASI heuristic, as  
 426 procedures, data collection, and data analyses also need to be sound (i.e., valid, and reliable—  
 427 see Figure 4).

If TA   HA	Then, O
------------	---------

	Not O (given SI)
Thus, no TA   HA	

428

429 *Figure 4. The Modus Tollens inferential process: Statistical and inferential assumptions (SI) help make the*  
 430 *minor premise sound, and determine its reliability when denying the consequent of the main premise*

431

432 The main problem with the second premise is that it can proceed in two ways. The correct  
 433 path is towards falsifying the first premise (e.g., HA: All swans are white. O: A black swan  
 434 is observed. Therefore, Conclusion: The hypothesis is false). The incorrect path is towards  
 435 confirming the first premise (e.g., O: Another white swan is observed. Conclusion: the  
 436 hypothesis is proved, or corroborated). The former path is the formal deduction allowed by  
 437 *Modus Tollens*, thus a logical conclusion. The latter path is a formal fallacy known as  
 438 *Affirming the Consequent*, thus leading to a formally illogical conclusion.

439 Tests of significance were born at a time when there were no stringent protocols for  
 440 experimentation nor any other statistical technology down the *Modus Tollens* path. Before  
 441 those tests, both descriptive statistics and Bayesian technology were used to confirm  
 442 hypotheses, thus rendering the path into a fallacious conclusion. Fisher (1925, 1935; also  
 443 Neyman, 1967) was the first statistician providing a basis for severity testing via  
 444 experimental control and *p*-values, as a tool for falsifying hypotheses. The *p*-value has been  
 445 so much confused since then that it may seem to offer little value today, but it was created  
 446 when there was nothing else—contemplating alternative hypotheses and statistical power  
 447 would need of Neyman and Pearson's constructs (e.g., 1928, 1942), effect sizes and sample  
 448 size calculations for power would need of Cohen's contributions (e.g., 1988), and confidence  
 449 intervals would also be derived from Neyman's proposal (1935).

450           The primary concern with  $p$ -values and, thus, significance tests, is that they are misused  
451 not to falsify hypotheses but to confirm them, thus going down the path of the confirmation  
452 fallacy. It is this bias towards confirming hypotheses that seems to underlie the myriad of  
453 breaches of statistical and inferential assumptions (SI) described earlier in the chapter.  
454 Therefore, most problems with significance testing are born out of attempts by researchers—  
455 intentionally or unwittingly—to support or confirm their own theories and hypotheses.  
456 Furthermore, such problems are not exclusive to significance tests. It will not really matter if  
457 we substitute effect sizes for  $p$ -values, for example. As long as we are intended to confirm  
458 rather than falsify, then the misuse of statistical and methodological assumptions are not  
459 going to disappear.

460

#### 461 ***Refutation, Severity, In/Out Model Inferences***

462 Among approaches exhorting the *Modus Tollens* inferential tree, we have Popper's  
463 falsificationism (1962), which mostly focuses on the first premise of the *Modus Tollens*  
464 syllogism. Popper's falsificationism offers a viable solution to the problem of induction: We  
465 can learn by proposing provisional hypotheses and submitting them to stringent tests for  
466 potential refutation. Science advances by proposing novel hypotheses to test, improving weak  
467 theories by trimming off failing assumptions, and rejecting bad theories that fail to pass those  
468 tests. The most recognisable motto of Popper is that a scientific theory is demarcated from a  
469 pseudoscientific one insofar the former is open to be tested to fail, and genuine attempts are  
470 made to do so.

471

472 Fisher's (1954), and Neyman-Pearson's (1928) approaches on research methods and statistics  
473 provided a workable environment for developing the second premise of the *Modus Tollens*  
474 syllogism. Pearson and, especially, Neyman, were most concerned with mathematics than

475 with philosophy, yet provided enduring tools for statistics beyond null hypotheses and  $p$ -  
476 values, such as alternative hypotheses, Type I and Type II errors, confidence intervals, and—  
477 mostly popularised by Cohen years later (e.g., 1962)—sample size, and power.

478 Mayo's approach (e.g., 2018) follows Popper's philosophy more closely, albeit it  
479 mostly focuses on the second premise of the *Modus Tollens* syllogism. The main take of  
480 Mayo's philosophy is that testing differs in their quality, from no severity (BENT, bad  
481 evidence, no test), to weak severity, to strong severity. Plus, not just the one test is enough:  
482 Multiple severe tests in agreement are needed to better substantiate whatever statistical  
483 inference is derived from them. Working with the disjunctive-hypotheses strategy inherited  
484 from Fisher's and Neyman-Pearson's approaches, Mayo seeks a reliable method to support a  
485 hypothesis via falsifying its opposite: We have evidence for a hypothesis just to the extent it  
486 survives a stringent scrutiny with severe tests designed to falsify it. (Mayo's corroboration  
487 argument based on a principle of coincidence among tests, however, is eerily close to flip-  
488 flopping towards the fallacy of *affirming the consequent*—also Schurz, 2019, similarly  
489 criticising Popper on his corroboration argument. However, it is technically possible because  
490 the disjunctive-hypotheses strategy used by frequentists aligns with that Sherlock Holmes's  
491 intuition (Doyle, 1927) that when you have eliminated all which is impossible—a.k.a., the null  
492 hypothesis—then whatever remains, however improbable,—a.k.a., the untested alternative—  
493 must be the truth.

494 Taleb's (2007) work may be substantiating another worthy approach, also targeting the  
495 second premise of the *Modus Tollens* syllogism. Most inferential attempts have been on  
496 rejecting those hypotheses that fail to pass null hypothesis tests (whether severely or not).  
497 Taleb (2001), most concerned with not being fooled by our own inferences, provides a path  
498 to also contemplate the hypotheses that are rejected. Namely, what is not rejected remains in  
499 the corpus of current knowledge and, thus, in our awareness, so they will not surprise us. And

500 yet, it is what is rejected (the highly improbable yet highly consequential) that could come  
501 back to bite us in the end—something particularly relevant in volatile fields such as economics  
502 or safety management.

503 Finally, the Bayesian approach (e.g., Kruschke, 2015) is worth considering too.  
504 Although Bayesian philosophy is eminently enumerative and confirmatory, thus down the  
505 path of the formal fallacy of *affirming the consequent*, there is no reason why it cannot be  
506 flipped to play a falsificationist role. Indeed, rather than focusing on the most supported  
507 hypothesis as the most credible, we can test hypotheses severely and discard the least  
508 supported ones as least credible. One way of doing so is to use priors to make Bayesian tests  
509 more severe (e.g., instead of assuming priors of 50/50 when we know nothing about the  
510 hypotheses, we could ‘bias’ the priors 95/5 in favour of whatever ‘main hypothesis’ we are  
511 using). If posterior probabilities still rule against the most supported hypothesis, such results  
512 serve to falsify it in the same manner than a frequentist approach may do, and severely so.

513 The advantage of Bayesian methods lies in their ability to move beyond the singular  
514 perspective of frequentist statistics. For instance, in Jeffrey's approach (1961), analyzing the  
515 same data under both normal and Cauchy distributions is typical. However, comprehensive  
516 Bayesian methods enable the simultaneous consideration of various models, as well as  
517 model-free Markov Chain Monte Carlo simulations (e.g., Metropolis et al., 1953). A second  
518 benefit is that Bayesian approaches are more interested in the hypotheses or model  
519 themselves [ $p(H|D)$ ] and not just on the location of the data in a single distribution [ $p(D|H)$ ].  
520 Thus, using a Bayesian approach in a falsificationist manner is not only possible, but also  
521 accords well with the logical path of *Modus Tollens*, thus, rendering formally logical  
522 conclusions.

523 In sum, lack of reproducibility and lack of replication are not negative in themselves, as  
524 they signal a failure of theories and hypotheses to hold. The problem is when reproducibility

525 and replication cannot be trusted because theories become impervious to being falsified (thus,  
526 becoming unsound, even pseudo-scientific), and methods are biased towards confirming  
527 those theories (thus, becoming invalid and unreliable). Good inferential work benefits from  
528 two complementary mindsets: (a) serious attempts at putting our theories and hypotheses to  
529 severe tests that may falsify them, and (b) an overall eagerness not to be fooled by our own  
530 biases towards confirmation. The two mindsets need not be restricted to just methods and  
531 statistics, but should permeate theories and hypotheses as well. Poor attempts should be  
532 considered BENT, opening us to be fooled by our own decisions and procedures. Weak  
533 probing already works in the right direction and gives us good leeway to learn from our  
534 research work (e.g., screening out auxiliary assumptions whose role seems to be that of  
535 saving our theories, pilot studies, exploratory data analysis...). Strong probing is the  
536 cornerstone of highly sophisticated inferential work (e.g., open science, reproducibility,  
537 replication, meta-analysis...). Tests of significance are appropriate tools in the inferential  
538 toolbox. However, they are so much misunderstood and misapplied that today's researcher  
539 might be less fooled if she were to substitute alternative tools that came afterwards:  
540 confidence intervals (the other side of a test of significance, but centred on the sample data  
541 rather than a theoretical null hypothesis, e.g., Perezgonzalez, 2015), effect sizes (to probe  
542 practical importance, not just statistical significance), replication and meta-analysis (to probe  
543 via a principle of coincidence), etc. Bayesian approaches, if done with severity in mind, can  
544 be repurposed to play a falsificationist role and may prove to be quite valuable both for  
545 sensible inferences about hypotheses themselves and, more importantly, for substantiating the  
546 need for further replication in order to maximise such inferences beyond reasonable doubt.

547

548

## 549                   **REACT: A New Perspective on Hypothesis Testing**

550   Trafimow (2019) argues that  $p$ -values offer no real benefit to scientific inference. His primary  
551   critique, framed within the TASI model, is that  $p$ -values conflate effect size and sample size,  
552   obscuring their contributions to statistical conclusions. Concretely, he argues that because no  
553   TASI model is perfectly correct, the usefulness of  $p$ -values in determining whether a model is  
554   "close enough" to reality is questionable. With large sample sizes, any TASI model will  
555   eventually be rejected, even if the truth is still in the ballpark of TASI. Similarly, small  
556   samples may fail to detect even substantial effects due to high variability. Trafimow thus  
557   advocates for reporting effect sizes and sample sizes separately rather than merging them into  
558   a single measure like the  $p$ -value.

559           In this section, we review REACT (Izbicki et al., in press) and demonstrate that it  
560   provides a stronger alternative to  $p$ -values under the TASI framework. REACT extends  
561   equivalence testing (Schuirmann, 1987) and three-way hypothesis testing (Berg, 2004;  
562   Esteves et al., 2016; Coscrato et al. 2019), integrating effect sizes and confidence sets into a  
563   structured decision-making process. Unlike NHST, which forces binary conclusions of  
564   "reject" or "fail to reject" the null hypothesis, REACT introduces an agnostic decision,  
565   explicitly recognizing when the available data is insufficient to support a definitive  
566   conclusion (for an R implementation visit <https://github.com/Monoxido45/REACT>).

567

### 568   ***The REACT Method***

569   Let  $\theta$  be a parameter of interest, such as the difference in mean systolic blood pressure  
570   reduction between two antihypertensive drugs or the difference in anxiety reduction scores  
571   between two therapeutic approaches. REACT operates by defining a null hypothesis not as a  
572   single point but as a region of practical equivalence (also called the pragmatic region; Esteves  
573   et al., 2019), denoted by  $\Theta_0$ . For instance,  $\Theta_0$  may have the shape

574  $\Theta_0 = \{\theta: |\theta| < \Delta\}$ ,

575 where  $\Delta$  is a threshold chosen based on domain knowledge, representing the smallest effect  
 576 size that would be considered practically meaningful. By asking the scientist to specify  $\Delta$ ,  
 577 REACT invites them to define a plausible and relevant theoretical range for the hypothesis.  
 578 In the TASI framework, this means that REACT requires the scientist to specify a more  
 579 plausible theory.

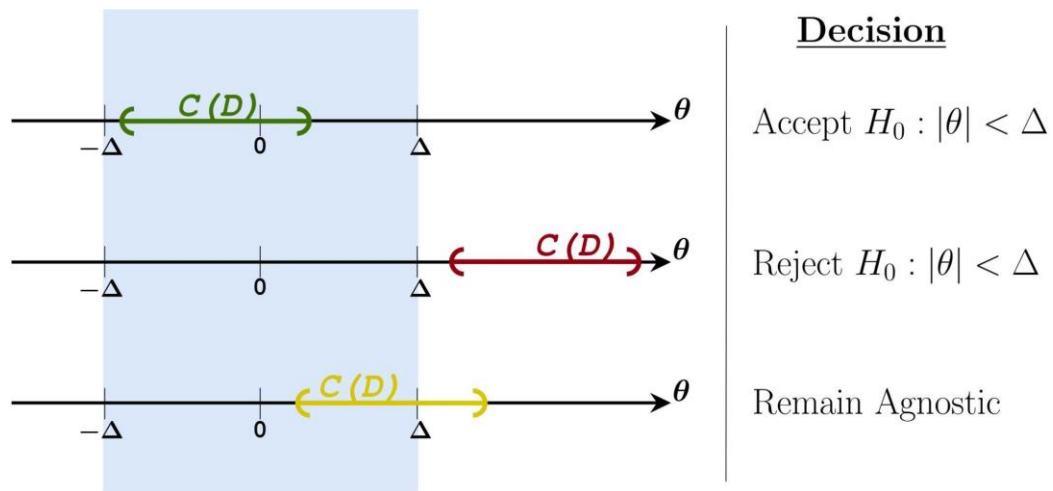
580 In many settings, it is easier to specify  $\theta$  in terms of effect sizes. For example, it could  
 581 be Cohen's  $d$ , which measures the standardized difference between two group means, with  $\Delta$   
 582  $= 0.2$  often considered a small effect. In a medical application,  $\theta$  might also be the distance  
 583 from the hazard ratio and 1, and  $\Delta = 0.1$  could define the smallest increase in risk deemed  
 584 clinically relevant. In economics,  $\theta$  might be a percentage change in income, with  $\Delta = 2\%$   
 585 setting the threshold for a meaningful policy effect.

586 Once  $\Theta_0$  is defined, the method then consists of the following steps, illustrated in  
 587 Figure 5:

- 588 1. **Construct a Confidence Region:** Given a dataset  $\mathcal{D}$ , construct a confidence set  $C(\mathcal{D})$   
 589 for the parameter  $\theta$ . This set may be frequentist or Bayesian.
- 590 2. **Compare to the Equivalence Region:** Compare the confidence set to the pre-  
 591 specified equivalence range  $\Theta_0$ , which defines a meaningful difference from the null.

592 One of the following conclusions is drawn:

- 593 ○ Accept  $H_0: \theta \in \Theta_0$  if  $C(\mathcal{D}) \subseteq \Theta_0$ , meaning all plausible values fall within the  
 594 equivalence region.
- 595 ○ Reject  $H_0: \theta \in \Theta_0$  if  $C(\mathcal{D}) \cap \Theta_0 = \emptyset$ , indicating that all plausible values lie outside  
 596 the equivalence region.
- 597 ○ Remain agnostic otherwise, acknowledging that the data does not provide  
 598 sufficient resolution to make a strong conclusion.



599

600 *Figure 5.* Illustration of the REACT procedure. A confidence set for  $C(D)$  is build. If the entire set lies within  
 601 the null hypothesis  $H_0: \theta \in \Theta_0$ , the hypothesis is accepted. If it lies entirely outside, it is rejected. Otherwise, the  
 602 procedure remains inconclusive.

603 REACT presents several advantages over traditional NHST methods (see Izbicki et al. 2025  
 604 for details):

- 605 1. **Identifying Evidence of Absence vs. Absence of Evidence:** REACT explicitly  
 606 distinguishes between cases where the null hypothesis is strongly supported and cases  
 607 where the data is inconclusive. This is in contrast to standard hypothesis tests, in  
 608 which a non-rejection of the null hypothesis can be interpreted as acceptance or  
 609 failure to reject it (Edwards et al. 1963, Neyman 1976).
- 610 2. **Logical Coherence:** Unlike traditional NHST, which can lead to paradoxical  
 611 conclusions, REACT ensures that conclusions remain logically consistent across  
 612 multiple hypothesis tests. For instance, concluding that medication A is better than B  
 613 and that B is better than C, entails the conclusion that A is better than C.
- 614 3. **Application to Multiple Hypotheses:** REACT naturally extends to testing multiple  
 615 null hypotheses simultaneously while controlling family-wise errors. This approach  
 616 eliminates the need for ad hoc corrections such as Bonferroni adjustments.

617 4. **Consistency:** REACT prevents automatic rejections in large samples and ensures that  
 618 small samples do not obscure meaningful effects. When the region of practical  
 619 equivalence holds, the test will accept it as the sample size grows; otherwise, it will  
 620 eventually lead to rejection. Additionally, for small samples with wide confidence  
 621 intervals, the test remains agnostic to the null hypothesis. Thus, this approach directly  
 622 addresses the limitations of  $p$ -values highlighted by Trafimow (2019), particularly  
 623 their restricted informativeness in hypothesis testing.

624

#### 625 *Examples illustrating how to REACT*

626 To illustrate REACT in practice, we apply it to assess whether specific variables significantly  
 627 impact certain outcomes. We consider three examples: the association between daily red  
 628 meat consumption and colon cancer (Di Maso et al., 2013), the relationship between sending  
 629 SMS while driving and the risk of motor vehicle accidents (Kogani et al., 2020), and the link  
 630 between smoking and lung cancer (Matos et al., 1998).

631 Let  $\theta$  represent the odds ratio between the outcome and the variable of interest. We  
 632 define three pragmatic hypotheses in the form  $\Theta_0 = \{\theta: \theta > \Delta\}$ . In the first hypothesis, we set  
 633  $\Delta = 1.0$ , indicating the presence of at least some effect. For the second hypothesis, we use  
 634  $\Delta = 1.5$ , often considered to be a small effect size (Cohen, 1988). For the third hypothesis, we  
 635 set  $\Delta = 3.5$ , representing a medium effect size. Table 1 summarizes the results, which we  
 636 describe in the sequence.

637

638 *Red Meat Consumption and Esophagus Cancer.* The estimated odds ratio is 1.46, with  
 639 a 95% confidence interval of (1.23, 1.72). For  $H_0: \theta > 1.0$ , the entire confidence interval lies  
 640 above 1, so REACT accepts that red meat consumption has some effect on colon cancer risk.  
 641 However, for  $H_0: \theta > 1.5$ , the confidence interval includes values below and above 1.5, making

642 REACT agnostic about the presence of a small effect; more data would be needed to reach a  
 643 conclusion. For  $H_0:\theta>3.5$ , the confidence interval lies entirely below 3.5, so REACT rejects  
 644 the hypothesis of a medium effect.

645 *Texting While Driving and Motor Vehicle Accidents.* The estimated odds ratio is 2.3,  
 646 with a 95% confidence interval of (1.2, 4.4). For  $H_0:\theta>1.0$ , the confidence interval is entirely  
 647 above 1, so REACT accepts the presence of at least some effect. For  $H_0:\theta>1.5$ , the confidence  
 648 interval also remains above 1.5, so REACT accepts that texting while driving has at least a  
 649 small effect on the risk of motor vehicle accidents. However, for  $H_0:\theta>3.5$ , the confidence  
 650 interval includes values below 3.5, so REACT remains agnostic about whether the effect  
 651 reaches a medium magnitude, suggesting that more data are necessary to draw a definitive  
 652 conclusion.

653 *Smoking and Lung Cancer.* The estimated odds ratio is 8.5, with a 95% confidence  
 654 interval of (4.3, 16.7). For  $H_0:\theta>1.0$ , the confidence interval lies entirely above 1, so REACT  
 655 accepts that smoking has some effect on lung cancer risk. Similarly, for  $H_0:\theta>1.5$  and  
 656  $H_0:\theta>3.5$ , the confidence interval remains above the respective thresholds, indicating that  
 657 REACT accepts that the effect size for smoking on lung cancer is at least medium.

658

659 *Table 1.* REACT results for three problems and hypotheses

660

	Confidence Interval for Odds Ratio $\theta$	$H_0:\theta>1.0$ (at least some effect)	$H_0:\theta>1.5$ (at least a small effect)	$H_0:\theta>3.5$ (at least a medium effect)
Meat vs Esophagus Cancer	(1.23; 1.72)	Accept	Agnostic	Reject
Sending SMS vs Accident	(2.3; 4.4)	Accept	Accept	Agnostic
Smoking vs Lung Cancer	(4.3; 16.7)	Accept	Accept	Accept

661

662

663 These examples illustrate how REACT provides more nuanced conclusions by incorporating  
664 confidence intervals and predefined effect size thresholds. This approach enables researchers  
665 to distinguish between different effect magnitudes and the strength of the supporting  
666 evidence. In contrast, relying solely on  $p$ -values would only test the null hypothesis that the  
667 odds ratio equals one, without offering insight into the practical relevance of the observed  
668 effects. Indeed, in all three datasets, such a test would reject the null hypothesis.

669

### 670 *Limitations, Assumptions and Conclusions*

671 REACT assumes the ASI components of TASI are approximately correct. Thus, both  
672 rejection and acceptance of the theory are contingent on these components being good  
673 approximations. However, contrary to standard point null theories, the ASI assumptions can  
674 be good approximations in practice. Moreover, one can use nonparametric methods within  
675 REACT, which rely on weaker assumptions (Lassance et al., 2025).

676

677 REACT also depends on the choice of the threshold  $\Delta$ . The value of  $\Delta$  represents the  
678 smallest effect size deemed practically meaningful, and its specification often involves  
679 subjective decisions based on domain knowledge. If the chosen  $\Delta$  is too small, the method  
680 might frequently reject hypotheses, while a large  $\Delta$  could lead to excessive acceptance.  
681 Several strategies can be used for selecting  $\Delta$ , including usage of prior literature, clinical  
682 relevance, or heuristic rules like the minimal clinically important difference. To assess the  
683 robustness of conclusions to the choice of  $\Delta$ , researchers can perform a sensitivity analysis by  
684 examining the results under a range of plausible thresholds.

684

685 In sum, the REACT framework not only mitigates the shortcomings of  $p$ -values  
686 identified within Trafimow's TASI model but also offers a more nuanced perspective on  
687 hypothesis testing. By allowing for agnostic conclusions, REACT aligns more closely with  
688 the scientific principle of acknowledging uncertainty rather than forcing binary decisions.

688 Additionally, REACT provides flexibility by accommodating both frequentist and Bayesian  
689 approaches when constructing confidence regions. This adaptability ensures that researchers  
690 can choose the method best suited to their assumptions, making the inference process more  
691 flexible and context-specific.

692

### 693 **Statistical learning instead of hypothesis testing via GAMLSS**

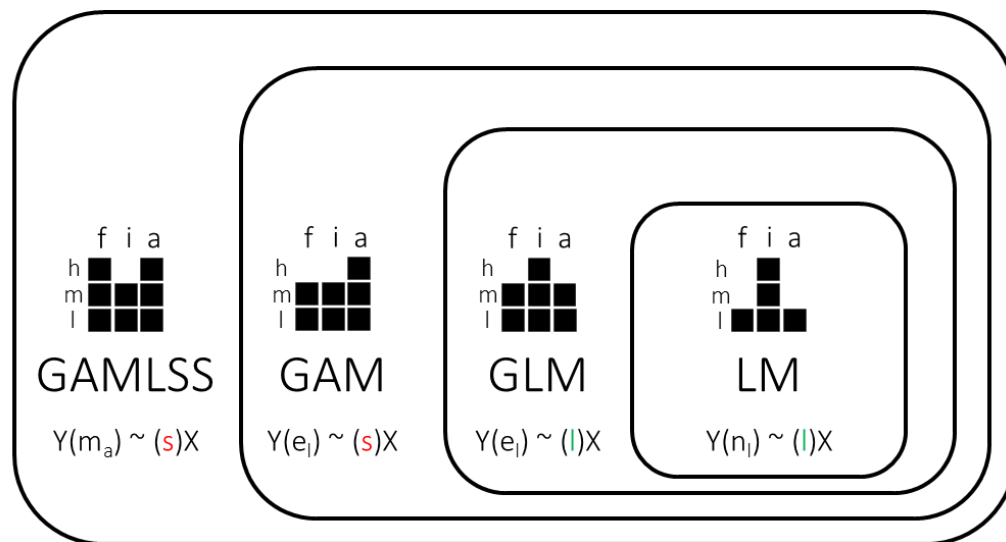
694 The SI sub-model is crucial to the TASI framework as it determines research validity and  
695 reliability. *S* assumptions govern statistical methods (distribution, independence, variance),  
696 while *I* assumptions address inferential reasoning (generalizability, causality, significance  
697 interpretation). These components form an inseparable bond—statistical integrity enables  
698 sound inference, and both must be satisfied for valid findings. Statistical violations inevitably  
699 compromise inferential quality. Together, they create a fundamental partnership that ensures  
700 research robustness and replicability, forming the backbone of scientific inquiry rather than  
701 merely serving as technical requirements.

702 Hypothesis testing serves as the traditional framework within the SI model, where a  
703 null hypothesis (assuming no effect or difference) is contrasted with an alternative  
704 hypothesis. This process involves assessing the probability of false rejection and results in a  
705 binary decision: either rejecting or failing to reject the null hypothesis. In contrast, we  
706 advocate for statistical learning as an alternative approach that entirely circumvents  
707 hypothesis testing, instead emphasizing descriptive, explanatory, and predictive statistical  
708 models. The proposed framework for this approach is Generalised Additive Models for  
709 Location, Scale, and Shape (GAMLSS) (Rigby & Stasinopoulos, 2005).

710 GAMLSS is a statistical learning framework that expands on traditional regression  
711 models by modeling all parameters of the response variable's distribution, rather than  
712 focusing solely on the mean. It extends linear models (LM) by supporting non-normal

713 distributions from the exponential family and employing link functions. Additionally,  
 714 GAMLSS goes beyond generalized linear models (GLM) by integrating non-linear smooth  
 715 functions of predictors. It further advances generalized additive models (GAM) by modeling  
 716 all parameters of the distribution (not just location) and accommodating a wider variety of  
 717 distributions (Stasinopoulos, Rigby, & de Bastiani, 2018) (see Figure 6).

718



719

720 *Figure 6.* Illustration of the relationship among four regression approaches. LM = linear model (a.k.a.,  
 721 ordinary least squares or general linear model); GLM = generalised linear models; GAM = generalised  
 722 additive models; GAMLSS = generalised additive models for location, scale and shape. In the models ' $Y(\bullet) \sim$   
 723 ' $(\bullet)X$ ', Y stands for the dependent (continuous or discrete) variable and X is a design matrix (i.e. set of  
 724 covariates or explanatory variables). The (fixed-effects) explanatory variables in X are categorical and/or  
 725 numerical (note that only GAMLSS includes smoothers for both kinds of covariates). Random effects are  
 726 modellable in all methods but only GAM and GAMLSS can also model random effects via smoothers (e.g.  
 727 temporal and spatial effects). Modelling of the dependent variable  $Y(\bullet)$ : n = Normal distribution only; e =  
 728 family of exponential distributions only; m = several family of distributions (note that  $m \supset e \supset n$ ). The  
 729 subscript by the dependent variable's distribution indicates that while LM, GLM, and GAM can assess the  
 730 effects of covariates on the dependent variable's location 'l' only (e.g. the mean in the case of LM),  
 731 GAMLSS can do so for all, 'a', the dependent variable's parameters (i.e. location, scale, skewness, and

732 kurtosis). Modelling of the covariates in  $X$ :  $s$  = smoothers,  $l$  = linear (note that  $s \supset l$ ). The black squares  
733 represent each regression approach's levels of flexibility ( $f$ ), interpretability ( $i$ ) and accuracy ( $a$ ) such that  
734 these can be high ( $h$ ), medium ( $m$ ) or low ( $l$ ) (these levels are based on the authors' experience and are  
735 inspired in Figure 2.7 in James et al (2021) and Figure 12 in Barredo Arieta et al (2020). Note all these  
736 metrics need to be assessed empirically). Symbol ' $\supset$ ' to be understood as superset.

737

738 As a statistical learning framework, GAMLSS enables the understanding, analysis, and  
739 extraction of patterns from data to support predictions or inferences. It connects traditional  
740 statistical methods (e.g., beta weights in parameters) with modern ML approaches (e.g., using  
741 neural networks to model numeric covariates) (Stasinopoulos et al., 2017; Rigby,  
742 Stasinopoulos, Heller, & de Bastiani, 2020; Stasinopoulos et al., 2024). Furthermore, as a  
743 regression framework, GAMLSS quantifies relationships while accounting for natural  
744 variation and random error. When applied correctly, with careful consideration of its  
745 underlying assumptions, GAMLSS models can produce reliable statistical inferences.  
746 Additionally, GAMLSS emphasizes interpretability by separating data into deterministic and  
747 random components, making complex relationships more understandable—ensuring that  
748 users can grasp how the model operates internally and how it transforms inputs into outputs.

749 The following section positions GAMLSS within the broader context of ML and related  
750 concepts. Additionally, it shows how GAMLSS can generate descriptive, explanatory, and  
751 predictive models, which can serve as a novel substitute for conventional hypothesis testing.

752

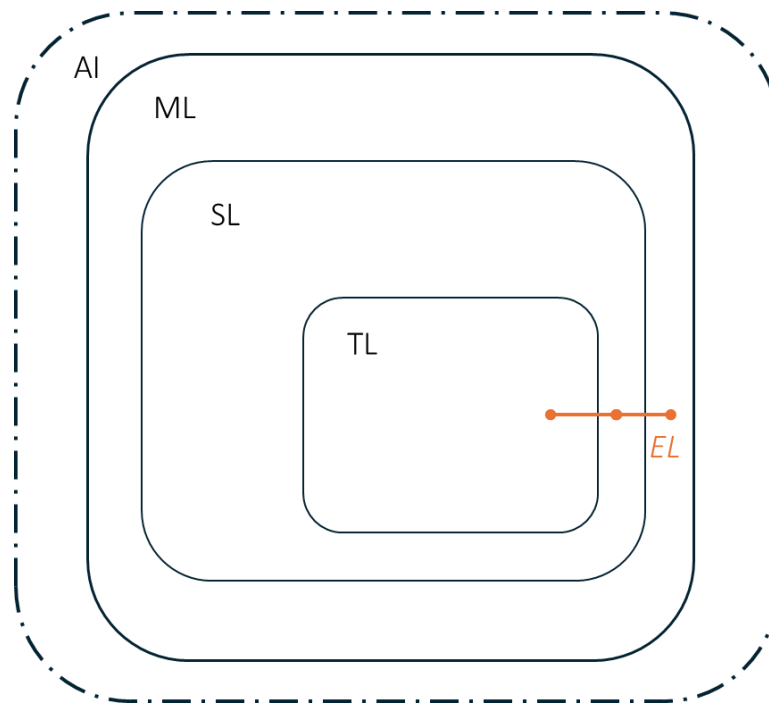
### 753 ***Description, explanation, and prediction models via GAMLSS***

754 ML provides the computational tools that enable artificial intelligence (AI) systems to  
755 learn from data and improve through experience. These ML algorithms, from simple linear  
756 models to complex deep neural networks, derive their theoretical foundations from statistical

757 principles (Friedrich et al., 2022; Min et al., 2024). AI systems built using these ML  
758 techniques rely fundamentally on statistical concepts like probability theory, hypothesis  
759 testing, and experimental design to ensure their reliability and validity (Friedrich et al., 2022).  
760 The synergy between these fields is particularly evident in areas like explainable AI (XAI),  
761 where statistical interpretability methods help make AI decision-making transparent (Min et  
762 al., 2024). Rather than separate disciplines, ML, AI and statistics form an interconnected  
763 framework - with ML providing the algorithmic implementation, AI delivering the intelligent  
764 applications, and statistics supplying the mathematical rigor (Friedrich et al., 2022).

765 ML is a subfield of AI, and within ML lies statistical learning (SL), a branch grounded  
766 in statistics that prioritizes model interpretability and uncertainty quantification (James et al.,  
767 2021). GAMLSS exemplify this approach. A specialized method within SL is targeted  
768 learning (TL), which focuses on estimating causal or statistical parameters with minimal bias  
769 (Coyle et al., 2023; Van der Laan & Starman, 2014). Notably, recent research suggests that  
770 GAMLSS models can be regularized to enhance causal inference (Marmolejo-Ramos et al.,  
771 2023). Additionally, ensemble learning (EL) is a versatile technique applicable across SL,  
772 TL, and broader ML. Inspired by the "wisdom of crowds," EL combines multiple models to  
773 improve predictive accuracy beyond what any single model can achieve (Sagi & Rokach,  
774 2018). In ML, EL is exemplified by random forest, a bootstrap aggregation (bagging) method  
775 that combines multiple decision trees to enhance predictive performance (Breiman, 2001b).  
776 Similarly, in classical statistical hypothesis testing, Stouffer's z-score method represents a  
777 form of EL by aggregating  $p$ -values from independent tests (Stouffer et al., 1949; Hoang &  
778 Dickhaus, 2022). The type of data used determines whether ML methods are categorized as  
779 supervised, semi-supervised, or unsupervised learning algorithms (see Figure 7).

780



781

782 *Figure 7.* Relationship among artificial intelligence (AI), machine learning (ML), statistical learning (SL),  
 783 targeted learning (TL), and ensemble learning (EL). Supervised learning algorithms work with labeled data  
 784 (input-output pairs) to perform prediction or classification (as in SL), estimate causal parameters (as in TL), or  
 785 implement methods like decision trees (in EL). Unsupervised learning algorithms analyze unlabeled data to  
 786 discover patterns, though this approach is less common in SL (which emphasizes inference over exploration)  
 787 and TL (which typically requires labeled data for causal inference), but appears in EL through clustering  
 788 ensembles or consensus clustering (Vega-Pons & Ruiz-Shulcloper, 2011). Semi-supervised learning algorithms  
 789 utilize small amounts of labeled data alongside larger unlabeled datasets for label propagation or self-training,  
 790 with specific applications in TL for causal inference with missing data, in EL through ensemble semi-supervised  
 791 classifiers (e.g., Zhao & Liu, 2021), and in emerging SL approaches such as semi-supervised GAMs (Culp,  
 792 2011) and semi-supervised GLMs (Tu et al., 2024).

793

794 SL encompasses a variety of methods and data types, but its core objective remains  
 795 consistent: to build descriptive, explanatory, and predictive models that balance accuracy  
 796 with interpretability. Descriptive models provide compact summaries of data, serving as a  
 797 foundation for further analysis without delving into causality or prediction. Explanatory  
 798 models, in contrast, focus on testing causal hypotheses, relying on theoretical assumptions to

799 uncover underlying mechanisms while minimizing bias. Predictive models, meanwhile,  
800 prioritize forecasting future outcomes, often employing data-driven techniques that favor  
801 accuracy over theoretical clarity (Shmueli, 2010).

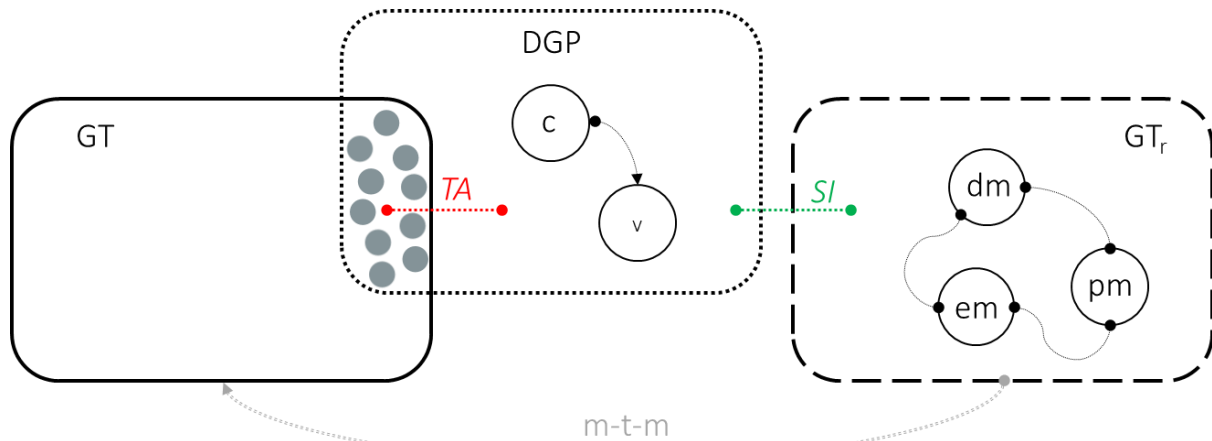
802         The choice of approach in SL hinges on how researchers conceptualize the data-  
803 generating process (DGP)—the underlying system that produces the observed data (Breiman,  
804 2001a). In some cases, the DGP is treated as stochastic, meaning it follows a probabilistic  
805 structure that can be modeled explicitly. Here, the emphasis is on statistical inference, where  
806 researchers estimate parameters, test hypotheses, and validate findings through measures like  
807 statistical significance. This data-driven approach aligns closely with explanatory modeling,  
808 as it seeks to uncover the true mechanisms behind the data while maintaining interpretability.  
809 However, when dealing with complex explanatory models where such inherent  
810 interpretability is elusive, achieving reliable understanding necessitates an interactive analysis  
811 that juxtaposes multiple explanatory methods to counter the misleading potential of any  
812 single perspective (Baniecki et al., 2023).

813         In other cases, the DGP is seen as too complex or unknown—effectively a "black box."  
814 Rather than trying to decipher the underlying process, researchers focus on finding  
815 algorithms that can reliably predict outcomes (Y) from inputs (X), regardless of theoretical  
816 fidelity. This algorithm-driven approach aligns with predictive modeling, where the primary  
817 goal is maximizing accuracy, even if it means sacrificing interpretability. Certain class of  
818 techniques in machine learning thrive in this domain, with models validated based on their  
819 performance on unseen data rather than their adherence to theoretical expectations.

820         We argue, though, while descriptive models are often overlooked in formal theory-  
821 building, they play a crucial role in exploratory analysis and can inform explanatory research  
822 by refining measurement tools or suggesting new theoretical directions. However,  
823 explanatory models are not without limitations—they face parameterization uncertainty

824 (ambiguity in model or parameter selection) and description uncertainty (gaps in theoretical  
 825 understanding), both of which can affect their reliability (Draper, 1995) (see Figure 8).

826



827

828 *Figure 8.* Human-data lifecycle. There exists a ground truth (GT), theorized to be generated by an underlying  
 829 system (the data generation process; DGP). The dotted region denotes the unknown overlap between GT and  
 830 DGP (i.e.  $P(\text{GT} \cap \text{DGP})$  is indeterminate). The TA component of the TASI model seeks  
 831 to link GT and DGP, where the DGP governs the formation of latent constructs (c)  
 832 and their manifest variables (v). The SI component of the TASI model  
 833 operationalizes GT into an empirical representation ( $\text{GT}_r$ ) via descriptive (dm),  
 834 explanatory (em), and predictive (pm) models. However,  $\text{GT}_r$  is subject to distortion by systematic error (bias)  
 835 and random error (variance), formalized as  $\text{GT}_r = \text{GT} + \text{bias} + \text{error}$ . Critically, the alignment between  $\text{GT}_r$  and  
 836 GT remains unquantified (i.e.,  $P(\text{GT}_r \cap \text{GT})$  is unknown; but see below). All statistical  
 837 models (i.e. dm, em, and pm) can be used to modify the GT through behavior  
 838 modification techniques. The "model-then-modify" (m-t-m) strategy uses  
 839 statistical models ( $\text{GT}_r$ ) to actively influence and adjust real-world behaviors (GT), aiming to make  
 840 those behaviors match the model's descriptions, explanations, or predictions. This establishes a cause-and-effect  
 841 relationship where the model representation directly impacts the actual system ( $\text{GT}_r \rightarrow \text{GT}$ ). Thus,  $P(\text{GT}_r \cap$   
 842  $\text{GT})$  is not merely unknown but potentially manipulated through intentional  
 843 intervention. This introduces a significant complication for the TASI framework,  
 844 as the statistical model (empirical representation) does not just approximate

845 the GT but can directly reshape it. This reshaping carries the risk that the  
846 actual system (GT) might diverge from the original DGP. While this divergence  
847 might misleadingly appear to improve predictive accuracy, it could actually be  
848 distorting the underlying reality it seeks to capture. (Note: This concept draws  
849 from Shmueli & Tafti, 2023, although their 'predict-then-modify' approach  
850 specifically targeted predictive models).

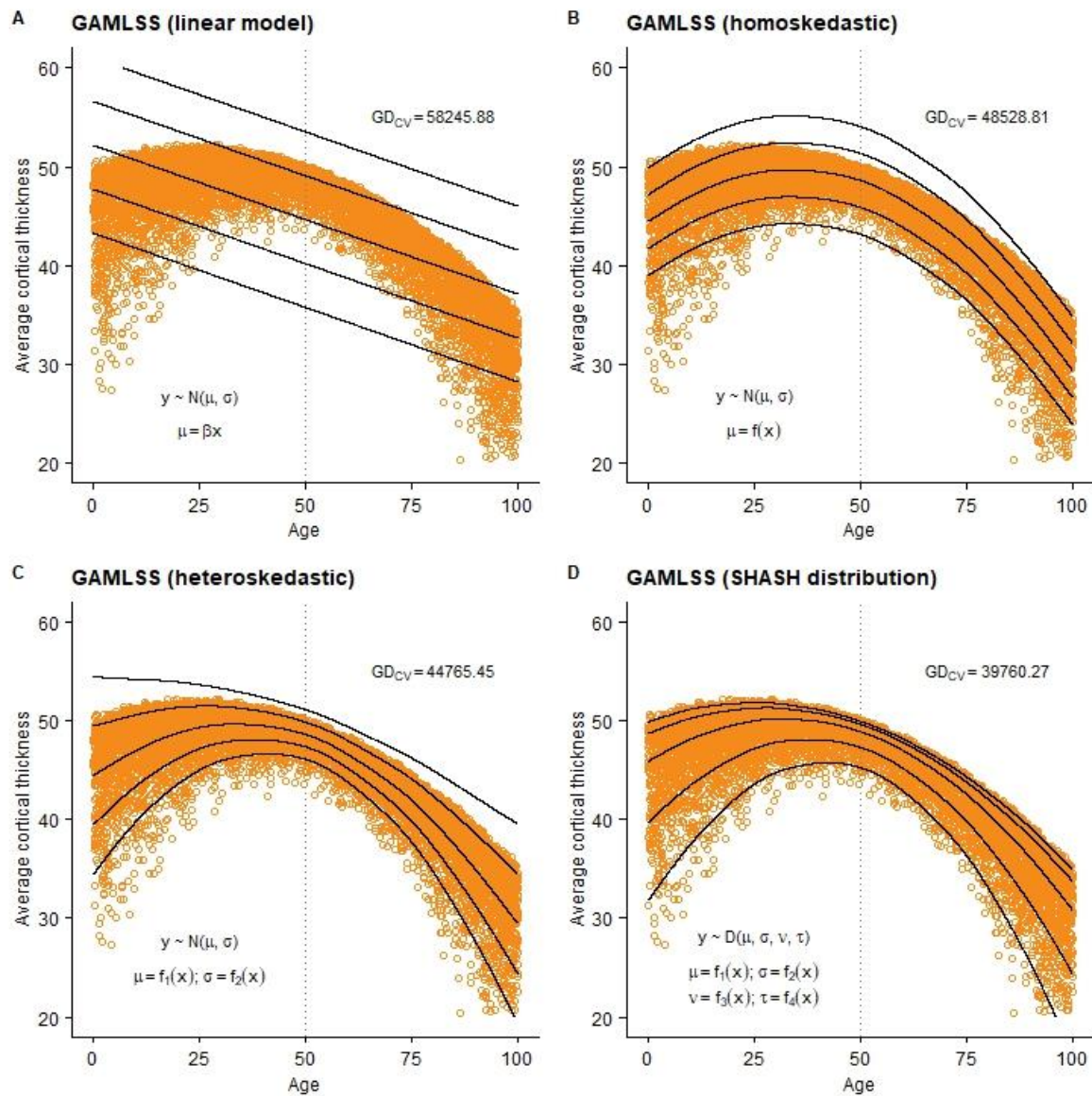
851

### 852 ***GAMLSS in action***

853 *Modelling cortical thickness.* It is well-established that aging is linked to widespread cortical  
854 thinning, especially in the frontal and temporal regions, which tend to be more susceptible  
855 than other areas (Salat et al., 2004; Piccolo et al., 2016). The relationship between average  
856 cortical thickness and age reveals a distinct pattern: a steep decline in MRI-based thickness  
857 estimates during childhood and adolescence, followed by a mild, continuous thinning from  
858 early adulthood onward, and an accelerated thinning trend in older adulthood (Vidal-Piñeiro  
859 et al., 2020). This non-linear association is illustrated in Figure 9.

860

861



862

863

864 *Figure 9.* Simulated data representing the non-linear negative association between cortical thickness ( $y$  axis) and865 age ( $x$  axis) ( $n=1e^4$ ) (adapted from Dinga et al., 2021, Figure 1). A: An LM assuming a normal distribution for

866 the outcome variable, with a linear predictor for the location parameter. B: A GAM assuming a normal

867 distribution for the outcome variable, with a nonlinear predictor for the location parameter. C: A GAM

868 assuming a normal distribution for the outcome variable, with nonlinear predictors for both location and scale

869 parameters. D: A GAMLSS using the third form of the four-parameter Sinh-Arcsinh (SHASH) distribution

870 (Jones &amp; Pewsey, 2009), with nonlinear predictors for location and scale while holding skewness and kurtosis

871 constant. All models were fitted using GAMLSS (see Figure 6). The five curvy solid lines on each plot represent

872 the estimated conditional quantiles of the response variable distribution at each predictor value. From bottom to

873 top, the lines indicate: q5 curve (5th percentile, lower outer bound), q25 curve (25th percentile, first quartile),  
874 q50 curve (median, 50th percentile), q75 curve (75th percentile, third quartile), and q95 curve (95th percentile,  
875 upper outer bound). The vertical dotted grey line marks the age of 50.  $GD_{CV}$  = Global Deviance CV (Cross-  
876 Validation) ( $GD_{CV}$  is a metric that reflects a model's predictive ability, as opposed to solely how well it fits the  
877 observed data. The model with the lowest  $GD_{CV}$  value is generally preferred because it indicates better  
878 generalizability, making  $GD_{CV}$  especially useful when comparing different models).

879

880 Figure 9 shows that both the traditional LM and the GAM approaches do not fit the data well,  
881 although the GAMs are an improvement over the LM. The GAMLSS SHASH model  
882 provides the best fit overall, both descriptively and predictively, as indicated by the  $GD_{CV}$ .  
883 The five quantile curves, representing specific percentiles of the predicted distribution, can be  
884 understood as prediction intervals that illustrate the probable range of individual observations  
885 at each predictor value. Thus, for example, the LM predicts that at the age of 50, the average  
886 cortical thickness is about 45 'mock-meters', the other three models predict this thickness to  
887 be 49 'mock-meters'. The cortical thickness predictions for a 50-year-old at the 95th  
888 percentile vary across the different models. Both the LM and GAMLSS homoskedastic  
889 models predict a thickness of 54 'mock-meters'. This value falls outside the range of the  
890 observed data. In contrast, the GAMLSS heteroskedastic and GAMLSS SASHO models  
891 predict thicknesses of 51 and 50 'mock-meters', respectively. These values are within the  
892 range of the data. The GAMLSS SHASHO model is only a valid explanatory model if age is  
893 the only factor associated with cortical thickness and if cortical thickness has been measured  
894 with valid metrics given by the most suitable brain imaging tool. Otherwise, the SHASH  
895 GAMLSS model is best used as a descriptive and predictive model (the R code for this  
896 example can be found at <https://cutt.ly/mrfSoIs4>). The example above demonstrates how  
897 GAMLSS operates in a simple linear regression scenario. For a multiple regression case, and

898 to explore visual techniques such as worm plots (van Buuren & Fredriks, 2001) and bucket  
899 plots (de Bastiani et al., 2022), please visit <https://cutt.ly/mrfSy8sm>.)

900 GAMLSS is a highly flexible SL framework that enables the development of diverse  
901 models. Beyond its general utility, GAMLSS modeling has had significant real-world impact.  
902 For instance, it has been employed in constructing centile estimates for the World Health  
903 Organization Child Growth Standards (Borghi et al., 2006), which aim to establish population  
904 reference ranges across age and sex to better detect clinically atypical measurements  
905 throughout life. Additionally, GAMLSS has proven effective in modeling relationships  
906 between gray matter volume and age (Bethlehem et al., 2022). The framework has also been  
907 extended in various ways, including: BAMLSS (Umlauf, Klein, & Zeileis, 2018), a Bayesian  
908 variant of GAMLSS, NAMLSS (Thielmann et al., 2024) and semi-structured distributional  
909 regression (Rügamer, Kolb, & Klein, 2024), which integrate neural networks with GAMLSS,  
910 and hybrid approaches combining decision trees with GAMLSS (Constable et al., 2023). It is  
911 worth noting that artificial neural networks (ANNs) fundamentally operate as regression  
912 models. The estimated parameters (weights and biases) in ANNs function similarly to  
913 regression coefficients, mapping input variables ( $x$ ) to outputs ( $y$ )—precisely the goal of  
914 regression. In fact, generalized additive models (GAMs), a sub-model in GAMLSS, can be  
915 viewed as a special case of ANNs (Cheng & Titterton, 1994). This connection is explicitly  
916 leveraged in modern statistical learning tools that model tabular and multimodal data within a  
917 regression framework (see Figure 1 in Rügamer et al., 2024).

918 Regression models are inherently interpretable (as noted by Lehmann, 2008). Since  
919 GAMLSS is fundamentally a regression model, it possesses a similar level of interpretability.  
920 The primary objective of GAMLSS as an SL approach is to facilitate statistical abductive  
921 learning—a process where the data analyst draws logical inferences aimed at finding the  
922 simplest and most probable conclusion from multiple (GAMLSS) models.

## 923 Discussion

924 In this book chapter we argue that the replication crisis in science is fundamentally a testing  
925 crisis rooted in the misapplication of statistical frameworks rather than inherent flaws in  
926 significance testing itself. Through the TASI framework, we showed that researchers' bias  
927 toward confirming rather than falsifying hypotheses undermines the logical foundation of  
928 scientific inquiry, regardless of whether  $p$ -values, effect sizes, or other metrics are employed.  
929 Both REACT and GAMLSS offer compelling alternatives that address this fundamental issue  
930 by providing more nuanced frameworks for statistical analysis. REACT explicitly  
931 acknowledges uncertainty through its agnostic decision option, preventing premature  
932 conclusions when data is insufficient, while GAMLSS transcends binary hypothesis testing  
933 altogether by focusing on comprehensive modeling of distributions across multiple  
934 parameters. These approaches represent a paradigm shift from the reductive nature of  
935 significance testing toward statistical frameworks that better accommodate the complexity of  
936 real-world phenomena. Below, we offer further reflections on the connections between data  
937 science and artificial intelligence

938 As noted earlier, modern AI models—particularly LLMs—can now explain their  
939 reasoning in natural language, reducing the "black box" problem and improving transparency.  
940 A growing practice involves describing a statistical problem in plain language (e.g., "Create  
941 R code for robust linear regression with spline-modeled covariates") and having the LLM  
942 generate functional, well-commented code (see Figure 1 in Marmolejo-Ramos et al., 2023).  
943 This approach, sometimes called 'vibe coding', streamlines the coding process by translating  
944 intuitive prompts into executable scripts. Beyond one-off code generation, specialized  
945 platforms (e.g., julius.ai; and 'data interpreter', Hong et al., 2024) and ambitious projects  
946 (e.g., The Automatic Statistician, Steinruecken et al., 2019; Sakana AI's AI Scientist, Lu et  
947 al., 2024) aim to automate data analysis and even scientific research. While these tools

948 accelerate processing, they remain just that—tools. The human analyst retains full control,  
949 deciding what data to analyze, how to combine datasets, and which statistical methods (AI-  
950 suggested or otherwise) to apply. In this way, LLM-powered tools act as collaborators,  
951 enhancing—not replacing—human expertise and judgment.

952         To conclude this chapter, we emphasize that the collaboration between AI-powered  
953 data analysis tools and human analysts offers significant benefits for teaching and learning  
954 data science. LLM-based statistical tools can play a key role in education by helping to  
955 develop course materials (e.g. R shiny apps via vibe coding), deliver personalized learning  
956 experiences, and clarify complex statistical concepts (Ellis & Slade, 2023). The ideal  
957 outcome of this human-AI partnership is well-informed model and hypothesis testing, where  
958 human expertise guides AI-generated insights to ensure rigorous, meaningful analysis.

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

**References**

- 978 Allen, L., Scott, J., Brand, A., Hlava, M., & Altman, M. (2014). Publishing: Credit where  
979 credit is due. *Nature*, 508(7496), 312–313. <https://doi.org/10.1038/508312a>
- 980 Alves, M. V., Ekuni, R., Hermida, M. J., & Lisboa, J. V. (Eds.). (2022). Cognitive Sciences  
981 and Education in Non-WEIRD Populations: A Latin American Perspective. *Springer*  
982 *Nature*. DOI: <https://doi.org/10.1007/978-3-031-06908-6>
- 983 Andrade, C. (2021). HARKing, cherry-picking, p-hacking, fishing expeditions, and data  
984 dredging and mining as questionable research practices. *The Journal of Clinical*  
985 *Psychiatry*, 82(1), 25941. <https://doi.org/10.4088/JCP.20f13804>
- 986 Baker, M. (2016). 1,500 scientists lift the lid on reproducibility. *Nature*, 533(7604).  
987 <https://doi.org/10.1038/533452a>
- 988 Baniecki, H., Parzych, D., & Biecek, P. (2023). The grammar of interactive explanatory  
989 model analysis. *Data Mining and Knowledge Discovery*.  
990 <https://doi.org/10.1007/s10618-023-00924-w>
- 991 Barch, D. M., & Yarkoni, T. (2013). Introduction to the special issue on reliability and  
992 replication in cognitive and affective neuroscience research. *Cognitive, Affective &*  
993 *Behavioral Neuroscience*, 13(4), 687–689. [https://doi.org/10.3758/s13415-013-0201-](https://doi.org/10.3758/s13415-013-0201-7)  
994 [7](https://doi.org/10.3758/s13415-013-0201-7)
- 995 Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A.,  
996 Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F.  
997 (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies,

- 998 opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82-115.
- 999 <https://doi.org/10.1016/j.inffus.2019.12.012>
- 1000 Baumgartner, H., Alessandrini, N., Byers-Heinlein, K., Frank, M. C., Hamlin, K.,  
1001 Soderstrom, M., ... & Coles, N. A. (2023). How to build up big team science: A  
1002 practical guide for large-scale collaborations. *PsyArXiv*.
- 1003 <https://doi.org/10.31234/osf.io/j7mt4>
- 1004 Benjamin DJ, Berger JO, Johnson VE et al (2018) Redefine statistical significance. *Nat Hum*  
1005 *Behav* 2:6–10 <https://doi.org/10.1038/s41562-017-0189-z>
- 1006 Berg, N., 2004: No-decision classification: an alternative to testing for statistical  
1007 significance. *the Journal of socio-Economics*, 33, no. 5, 631–65.
- 1008 <https://doi.org/10.1016/j.socec.2004.09.036>
- 1009 Bethlehem, R. A. I., Seidlitz, J., White, S. R., Vogel, J. W., Anderson, K. M., Adamson, C.,  
1010 Adler, S., Alexopoulos, G. S., Anagnostou, E., Areces-Gonzalez, A., Astle, D. E.,  
1011 Auyeung, B., Ayub, M., Bae, J., Ball, G., Baron-Cohen, S., Beare, R., Bedford, S. A.,  
1012 Benegal, V., ... Alexander-Bloch, A. F. (2022). Brain charts for the human lifespan.  
1013 *Nature*, 604(7906), 525–533. <https://doi.org/10.1038/s41586-022-04554-y>
- 1014 Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... & Liang, P.  
1015 (2021). On the opportunities and risks of foundation models. *arXiv preprint*.
- 1016 <https://doi.org/10.48550/arXiv.2108.0725>
- 1017 Borghi, E., de Onis, M., Garza, C., Van den Broeck, J., Frongillo, E.A., Grummer-Strawn, L.,  
1018 Van Buuren, S., Pan, H., Molinari, L., Martorell, R., Onyango, A.W., Martines, J.C.,  
1019 & (2006). Construction of the World Health Organization child growth standards:

- 1020 selection of methods for attained growth curves. *Statistics in Medicine*, 25, 247-265.
- 1021 <https://doi.org/10.1002/sim.2227>
- 1022 Breiman, L. (2001a). Statistical modeling: The two cultures. *Statistical Science*, 16(3), 199-
- 1023 231. <https://doi.org/10.1214/ss/1009213726>
- 1024 Breiman, L. (2001b). Random forests. *Machine Learning*, 45(1), 5–32.
- 1025 <https://doi.org/10.1023/A:1010933404324>
- 1026 Bunge, M. (2018). *La ciencia: su método y su filosofía* (Vol. 1). Laetoli.
- 1027 Catanzaro, M. (2023, November 27). *Saudi universities lose highly cited researchers after*
- 1028 *payment schemes raise ethics concerns*. *Science*.
- 1029 <https://doi.org/10.1126/science.zhs1429>
- 1030 Cheng, B., & Titterington, D. M. (1994). Neural networks: A review from a statistical
- 1031 perspective. *Statistical Science*, 9(1), 2–54. <https://doi.org/10.1214/ss/117701063>
- 1032 Cohen, J. (1962). The statistical power of abnormal-social psychological research: A
- 1033 review. *Journal of Abnormal and Social Psychology*, 65(3), 145-153.
- 1034 <https://doi.org/10.1037/h0045186>
- 1035 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.).
- 1036 Psychology Press. <https://doi.org/10.4324/9780203771587>
- 1037 Coles, N. A., March, D. S., Marmolejo-Ramos, F., Larsen, J. T., Arinze, N. C., Ndukaihe, I.
- 1038 L., ... & Liuzza, M. T. (2022). A multi-lab test of the facial feedback hypothesis by
- 1039 the many smiles collaboration. *Nature Human Behaviour*, 1-12.
- 1040 <https://doi.org/10.1038/s41562-022-01458-9>

- 1041 Constable, P. A., Loh, L., Prem-Senthil, M., & Marmolejo-Ramos, F. (2023). Visual search  
1042 and childhood vision impairment: A GAMLSS-oriented multiverse analysis approach.  
1043 *Attention, Perception, & Psychophysics*, 85(3), 968–977.  
1044 <https://doi.org/10.3758/s13414-023-02670-z>
- 1045 Coscrato, V., R. Izbicki, and R. B. Stern, 2019: Agnostic tests can control the type I  
1046 and type II errors simultaneously. *Brazilian Journal of Probability and Statistics* 34.2  
1047 (2020): 230-250. <https://doi.org/10.1214/19-BJPS431>
- 1048 Coyle, J. R., Hejazi, N. S., Malenica, I., Phillips, R. V., Arnold, B. F., Mertens, A., Benjamin-  
1049 Chung, J., Cai, W., Dayal, S., Colford, J. M., Hubbard, A. E., & van der Laan, M. J.  
1050 (2023). Targeted learning. In *Wiley StatsRef: Statistics Reference Online*. John Wiley  
1051 & Sons, Ltd. <https://doi.org/10.1002/9781118445112.stat08414>
- 1052 Culp, M. (2011). On Propagated Scoring for Semisupervised Additive Models. *Journal of the*  
1053 *American Statistical Association*, 106(493), 248–259.  
1054 <https://doi.org/10.1198/jasa.2011.tm09316>
- 1055 De Bastiani, F., Stasinopoulos, D. M., Rigby, R. A., Heller, G. Z., & Silva, L. A. (2022).  
1056 Bucket plot: A visual tool for skewness and kurtosis comparisons. *Brazilian Journal*  
1057 *of Probability and Statistics*, 36(3), 421-440. <https://doi.org/10.1214/22-BJPS533>
- 1058 Di Maso, M., Talamini, R., Bosetti, C., Montella, M., Zucchetto, A., Libra, M., Negri,  
1059 E., Levi, F., La Vecchia, C., Franceschi, S., Serraino, D., & Polesel, J. (2013). Red  
1060 meat and cancer risk in a network of case-control studies focusing on cooking  
1061 practices. *Annals of oncology : official journal of the European Society for Medical*  
1062 *Oncology*, 24(12), 3107–3112. <https://doi.org/10.1093/annonc/mdt392>

- 1063 Dinga, R., Fraza, C. J., Bayer, J. M. M., Kia, S. M., Beckmann, C. F., & Marquand, A. F.  
1064 (2021). Normative modeling of neuroimaging data using generalized additive models  
1065 of location scale and shape. *bioRxiv*. <https://doi.org/10.1101/2021.06.14.448106>
- 1066 Doyle, A. C. (1927). *The case-book of Sherlock Holmes*. John Murray.
- 1067 Draper, D. (1995). Assessment and propagation of model uncertainty (with discussion).  
1068 *Journal of the Royal Statistical Society, Series B*, 57, 45-97.  
1069 <https://doi.org/10.1111/j.2517-6161.1995.tb02015.x>
- 1070 Edwards, W., Lindman, H., & Savage, L. J. (1963). Bayesian statistical inference for  
1071 psychological research. *Psychological Review*, 70(3), 193–242.  
1072 <https://doi.org/10.1037/h0044139>
- 1073 Ellis, A. R., & Slade, E. (2023). A New Era of Learning: Considerations for ChatGPT as a  
1074 Tool to Enhance Statistics and Data Science Education. *Journal of Statistics and Data  
1075 Science Education*, 31(2), 128-133. <https://doi.org/10.1080/26939169.2023.2223609>
- 1076 Esteves, L. G., Izbicki, R., Stern, J.M., Stern, R.B. (2019). Pragmatic hypotheses in  
1077 the evolution of science. *Entropy*. 2019 Sep 11;21(9):883.  
1078 <https://doi.org/10.3390/e21090883>
- 1079 Esteves, L. G., R. Izbicki, J. M. Stern, and R. B. Stern (2016): The logical consistency  
1080 of simultaneous agnostic hypothesis tests. *Entropy*, 18, no. 7, 256.  
1081 <https://doi.org/10.3390/e18070256>
- 1082 Fanelli, D. (2018). Opinion: Is science really facing a reproducibility crisis, and do we  
1083 need it to? *PNAS*, 115, 2628-2631. <https://doi.org/10.1073/pnas.1708272114>

- 1084 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G\*Power 3: A flexible  
1085 statistical power analysis program for the social, behavioral, and biomedical sciences.  
1086 Behavior Research Methods, 39(2), 175–191. <https://doi.org/10.3758/BF03193146>
- 1087 Fisher, R. A. (1925). *Statistical methods for research workers*. Oliver and Boyd.
- 1088 Fisher, R. A. (1935). *The design of experiments*. Oliver and Boyd.
- 1089 Foster, E. D., & Deardorff, A. (2017). Open Science Framework (OSF). *Journal of the*  
1090 *Medical Library Association, JMLA*, 105(2), 203–206.  
1091 <https://doi.org/10.5195/jmla.2017.88>
- 1092 Frias-Navarro, D., Pascual-Llobell, J., Pascual-Soler, M., Perezgonzalez, J., &  
1093 Berrios-Riquelme, J. (2020). Replication crisis or an opportunity to improve scientific  
1094 production? *European Journal of Education*, 55, 618-631.  
1095 <https://doi.org/10.1111/ejed.12417>
- 1096 Friedrich, S., & Friede, T. (2023). On the role of benchmarking data sets and simulations in  
1097 method comparison studies. *Biometrical Journal*, 00, e2200212.  
1098 <https://doi.org/10.1002/bimj.202200212>
- 1099 Friedrich, S., Antes, G., Behr, S., Binder, H., Brannath, W., Dumpert, F., Ickstadt, K.,  
1100 Kestler, H. A., Lederer, J., Leitgöb, H., Pauly, M., Steland, A., Wilhelm, A., & Friede,  
1101 T. (2022). Is there a role for statistics in artificial intelligence? *Advances in Data*  
1102 *Analysis and Classification*, 16(3), 823-846. [https://doi.org/10.1007/s11634-021-](https://doi.org/10.1007/s11634-021-00455-6)  
1103 [00455-6](https://doi.org/10.1007/s11634-021-00455-6)
- 1104 Grimes, D. R. (2019). *The irrational ape*. Simon & Schuster.

- 1105 Gundersen, O. E., Gil, Y., & Aha, D. W. (2018). On reproducible AI: Towards reproducible  
1106 research, open science, and digital scholarship in AI publications. *AI Magazine*, 39(3),  
1107 56-68. <https://doi.org/10.1609/aimag.v39i3.2816>
- 1108 Haibe-Kains, B., Adam, G.A., Hosny, A. et al. Transparency and reproducibility in artificial  
1109 intelligence. *Nature* 586, E14–E16 (2020). <https://doi.org/10.1038/s41586-020-2766-y>
- 1110 Head, M. L., Holman, L., Lanfear, R., Kahn, A. T., & Jennions, M. D. (2015). The extent and  
1111 consequences of p-hacking in science. *PLoS Biology*, 13(3), e1002106.  
1112 <https://doi.org/10.1371/journal.pbio.1002106>
- 1113 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world? *The*  
1114 *Behavioral and Brain Sciences*, 33(2-3), 61–135.  
1115 <https://doi.org/10.1017/S0140525X0999152X>
- 1116 Hensel, P. G. (2021). Reproducibility and replicability crisis: How management compares to  
1117 psychology and economics—A systematic review of literature. *European Management*  
1118 *Journal*, 39(5), 577-594. <https://doi.org/10.1016/j.emj.2021.01.002>
- 1119 Hernández Sampieri, R., Fernández Collado C., & Baptista Lucio, P. (2014). *Metodología de*  
1120 *la Investigación*. México: McGraw-Hill.
- 1121 Hoang, A.-T., & Dickhaus, T. (2022). Combining independent p-values in replicability  
1122 analysis: A comparative study. *Journal of Statistical Computation and Simulation*, 92  
1123 (10), 2184–2204. <https://doi.org/10.1080/00949655.2021.2022678>
- 1124 Hong, S., Lin, Y., Liu, B., Liu, B., Wu, B., Zhang, C., Wei, C., Li, D., Chen, J., Zhang, J.,  
1125 Wang, J., Zhang, L., Zhang, L., Yang, M., Zhuge, M., Guo, T., Zhou, T., Tao, W.,  
1126 Tang, X., Lu, X., Zheng, X., Liang, X., Fei, Y., Cheng, Y., Gou, Z., Xu, Z., & Wu, C.

- 1127 (2024). Data Interpreter: An LLM Agent For Data Science. arXiv.  
1128 <https://arxiv.org/pdf/2402.18679>
- 1129 Ioannidis, J. P. (2005). Why most published research findings are false. *PLoS Medicine*, 2(8).  
1130 <https://doi.org/10.1371/journal.pmed.0020124>
- 1131
- 1132 Izbicki, R., Cabezas, L. M. C., Colugnatti, F. A. B., Lassance, R. F. L., de Souza, A. A. L., &  
1133 Stern, R. B. (in press). REACT to NHST: Sensible conclusions for meaningful  
1134 hypotheses. *The Quantitative Methods for Psychology*.
- 1135 REACT to NHST: Sensible conclusions for meaningful hypotheses
- 1136 James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). *An Introduction to Statistical*  
1137 *Learning: with Applications in R (2nd ed.)*. Springer.
- 1138 Jeffreys, H. (1961). *Theory of probability (3rd ed.)*. Clarendon Press.
- 1139 John, L. K., Loewenstein, G., & Prelec, D. (2012). Measuring the prevalence of questionable  
1140 research practices with incentives for truth-telling. *Psychological Science*, 23, 524–  
1141 532. <https://doi.org/10.1177/0956797611430953>
- 1142 Jones, M. C., & Pewsey, A. (2009). Sinh-arcsinh distributions. *Biometrika*, 96 (4), 761–780,  
1143 <https://doi.org/10.1093/biomet/asp053>
- 1144 Kaplan, D. (2004). *The Sage handbook of quantitative methodology for the social sciences*.  
1145 Sage.
- 1146 Kapoor, S., & Narayanan, A. (2023). Leakage and the reproducibility crisis in machine-  
1147 learning-based science. *Patterns*, 4(9). <https://doi.org/10.1016/j.patter.2023.100804>

- 1148 Kerr, N. L. (1998). HARKing: Hypothesizing After the Results are Known. *Personality and*  
1149 *Social Psychology Review*, 2(3), 196–217.  
1150 [https://doi.org/10.1207/s15327957pspr0203\\_4](https://doi.org/10.1207/s15327957pspr0203_4)
- 1151 Kogani, M., Almasi, S. A., Ansari-Mogaddam, A., Dalvand, S., Okati-Aliabad, H.,  
1152 Tabatabaee, S. M., & Almasi, S. Z. (2020). Relationship between using cell phone  
1153 and the risk of accident with motor vehicles: An analytical cross-sectional study.  
1154 *Chinese journal of traumatology = Zhonghua chuang shang za zhi*, 23(6), 319–323.  
1155 <https://doi.org/10.1016/j.cjtee.2020.08.002>
- 1156 Kruschke, J. K. (2015). *Doing Bayesian data analysis. A tutorial with R, JAGS, and*  
1157 *Stan (2nd ed.)*. Academic Press.
- 1158 Kumle, L., Vö, M. L.-H., & Draschkow, D. (2021). Estimating power in (generalized)  
1159 linear mixed models: An open introduction and tutorial in R. *Behavior Research*  
1160 *Methods*, 53(6), 2528-2543. <https://doi.org/10.3758/s13428-021-01546-0>
- 1161 Larsson, S., & Heintz, F. (2020). Transparency in artificial intelligence. *Internet*  
1162 *policy review*, 9(2), 1-16. <https://doi.org/10.14763/2020.2.1469>
- 1163 Lassance RF, Izbicki R, Stern RB. (2025) Adding imprecision to hypotheses: A  
1164 Bayesian framework for testing practical significance in nonparametric settings.  
1165 *International Journal of Approximate Reasoning*. 2025 Mar 1;178:109332.  
1166 <https://doi.org/10.1016/j.ijar.2024.109332>
- 1167 Lehmann, E. L. (2008). On the history and use of some standard statistical models. *IMS*  
1168 *Collections, Probability and Statistics: Essays in Honor of David A. Freeman*, 2, 114-  
1169 126. <http://doi.org/10.1214/1939403070000000419>

- 1170 Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., ... & Piao, Y. (2024). Deepseek-  
1171 v3 technical report. arXiv preprint. <https://arxiv.org/abs/2412.19437>
- 1172 Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., & Ha, D. (2024). The AI Scientist:  
1173 Towards Fully Automated Open-Ended Scientific Discovery. arXiv.  
1174 <https://arxiv.org/abs/2408.06292>
- 1175 Luijken K., Lohmann A., Alter U., Claramunt Gonzalez J., Clouth F. J., Fossum J. L.,  
1176 Hesen L., Huizing A. H. J., Ketelaar J., Montoya A. K., Nab L., Nijman R. C. C.,  
1177 Penning de Vries B. B. L., Tibbe T. D., Wang Y. A. and Groenwold R. H. H. (2024)  
1178 *Replicability of simulation studies for the investigation of statistical methods: the*  
1179 *RepliSims project*R. Soc. Open Sci.11231003. <http://doi.org/10.1098/rsos.231003>
- 1180 Marmolejo-Ramos, F., Simon, T., & Abadia, R. (2023, August 28). The power of large  
1181 language models to augment human learning. *Digital Dexterity*.  
1182 [https://digitaldexterity.edublogs.org/2023/08/28/the-power-of-large-language-models-](https://digitaldexterity.edublogs.org/2023/08/28/the-power-of-large-language-models-to-augment-human-learning/)  
1183 [to-augment-human-learning/](https://digitaldexterity.edublogs.org/2023/08/28/the-power-of-large-language-models-to-augment-human-learning/)
- 1184 Marmolejo-Ramos, F., Tejo, M., Brabec, M., Kuzilek, J., Joksimovic, S., Kovanovic, V.,  
1185 González, J., Kneib, T., Bühlmann, P., Kook, L., Briseño-Sánchez, G., & Ospina, R.  
1186 (2023). Distributional regression modelling via GAMLSS. An overview through a  
1187 data set from learning analytics. *WIRES Data Mining and Knowledge Discovery*, 13  
1188 (1), e1479. <https://doi.org/10.1002/widm.1479>
- 1189 Matos, E., Vilensky, M., Boffetta, P., & Kogevinas, M. (1998). Lung cancer and  
1190 smoking: a case-control study in Buenos Aires, Argentina. *Lung cancer (Amsterdam,*  
1191 *Netherlands)*, 21(3), 155–163. [https://doi.org/10.1016/s0169-5002\(98\)00055-5](https://doi.org/10.1016/s0169-5002(98)00055-5)

- 1192 Mayo, D. G. (2018). *Statistical inference as severe testing*. Cambridge University  
1193 Press.
- 1194 Meho, L. I., & Akl, E. A. (2025). Using bibliometrics to detect questionable  
1195 authorship and affiliation practices and their impact on global research metrics: A  
1196 case study of 14 universities. *Quantitative Science Studies*, 1-36. DOI:  
1197 [https://doi.org/10.1162/qss\\_a\\_00339](https://doi.org/10.1162/qss_a_00339)
- 1198 Metropolis, N. Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E.  
1199 (1953). Equations of state calculations by fast computing machines. *Journal of*  
1200 *Chemical Physics*, 21, 1087-1091. <https://doi.org/10.1063/1.1699114>
- 1201 Min, J., Song, X., Zheng, S., King, C. B., Deng, X., & Hong, Y. (2024). *Applied statistics in*  
1202 *the era of artificial intelligence: A review and vision*. <https://arxiv.org/abs/2412.10331>
- 1203 Morey, R. D., Kaschak, M. P., Díez-Álamo, A. M., Glenberg, A. M., Zwaan, R. A., Lakens,  
1204 D., ... & Ziv-Crispel, N. (2022). A pre-registered, multi-lab non-replication of the  
1205 action-sentence compatibility effect (ACE). *Psychonomic Bulletin & Review*, 29(2),  
1206 613-626. <https://doi.org/10.3758/s13423-021-01927-8>
- 1207 Morris, T. P., White, I. R., & Crowther, M. J. (2019). Using simulation studies to evaluate  
1208 statistical methods. *Statistics in Medicine*, 38(11), 2074–2102.  
1209 <https://doi.org/10.1002/sim.8086>
- 1210 Neyman, J., and Pearson, E. S. (1928). On the use and interpretation of certain test  
1211 criteria for purposes of statistical inference: part I. *Biometrika*, 20A, 175–240.  
1212 <https://doi.org/10.2307/2331945>

- 1213 Neyman, J. (1935). On the problem of confidence intervals. *The Annals of*  
1214 *Mathematical Statistics*, 6, 111–116. <https://doi.org/10.1214/aoms/1177732585>
- 1215 Neyman, J. (1942). Basic ideas and some recent results of the theory of testing  
1216 statistical hypotheses. *Journal of the Royal Statistical Society*, 105, 292–327.  
1217 <https://doi.org/10.2307/2980436>
- 1218 Neyman, J. (1967). R. A. Fisher (1890-1962): an appreciation. *Science*, 156, 1456–  
1219 1460. <https://doi.org/10.1126/science.156.3781.1456>
- 1220 Neyman, J. (1976). Tests of statistical hypotheses and their use in studies of natural  
1221 phenomena. *Communications in Statistics - Theory and Methods*, 5(8), 737–751.  
1222 <https://doi.org/10.1080/03610927608827392>
- 1223 Nosek, B. A., Spies, J. R., & Motyl, M. (2012). Scientific utopia: II. Restructuring  
1224 incentives and practices to promote truth over publishability. *Perspectives on*  
1225 *Psychological Science*, 7(6), 615-631. <https://doi.org/10.1177/1745691612459058>
- 1226 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.  
1227 *Science*, 349(6251), <https://doi.org/10.1126/science.aac4716>
- 1228 O’Grady, C. (2024, October 21). Springer Nature retracts 75 papers connected to Spanish  
1229 university head. *Science*. <https://doi.org/10.1126/science.zvmx7id>
- 1230 Pawel, S., Kook, L., & Reeve, K. (2023). Pitfalls and potentials in simulation studies:  
1231 Questionable research practices in comparative simulation studies allow for spurious  
1232 claims of superiority of any method. *Biometrical Journal*, 2200091.  
1233 <https://doi.org/10.1002/bimj.202200091>

- 1234 Perezgonzalez, J. D. (2015). Confidence intervals and tests are two sides of the same research  
1235 question, *Frontiers in Psychology* 28(6). <https://doi.org/10.3389/fpsyg.2015.00034>
- 1236 Piccolo, L. R., Merz, E. C., He, X., Sowell, E. R., Noble, K. G., & Pediatric Imaging,  
1237 Neurocognition, Genetics Study. (2016). Age-related differences in cortical thickness  
1238 vary by socioeconomic status. *PLOS ONE*, 11(9), e0162511.  
1239 <https://doi.org/10.1371/journal.pone.0162511>
- 1240 Poldrack, R. A., Baker, C. I., Durnez, J., Gorgolewski, K. J., Matthews, P. M., Munafò, M.  
1241 R., ... & Yarkoni, T. (2017). Scanning the horizon: towards transparent and  
1242 reproducible neuroimaging research. *Nature Reviews Neuroscience*, 18(2), 115-126.  
1243 <https://doi.org/10.1038/nrn.2016.167>
- 1244 Popper, k. (1962). *Conjectures and refutations: The growth of scientific knowledge*.  
1245 Basic Books.
- 1246 Pownall M, Pennington CR, Norris E, et al. Evaluating the Pedagogical Effectiveness of  
1247 Study Preregistration in the Undergraduate Dissertation. *Advances in Methods and*  
1248 *Practices in Psychological Science*. 2023;6(4).  
1249 <https://doi.org/10.1177/25152459231202724>
- 1250 Prieto, D. (2022). Ciencia Abierta: desafíos y oportunidades para Uruguay y el Sur Global.  
1251 *Information*, 27(1), 253-283. <https://doi.org/10.35643/Info.27.1.5>
- 1252 Rigby, R. A., & Stasinopoulos, D. M., Heller, G., & de Bastiani, F. (2020). *Distributions for*  
1253 *Modeling Location, Scale, and Shape*. Taylor & Francis.
- 1254 Rigby, R. A., & Stasinopoulos, D. M. (2005). Generalized Additive Models for Location,  
1255 Scale and Shape. *Applied Statistics*, 54, 507-554.

- 1256 Roberts, M., Driggs, D., Thorpe, M., Gilbey, J., Yeung, M., Ursprung, S., ... & Schönlieb, C.  
1257 B. (2021). Common pitfalls and recommendations for using machine learning to  
1258 detect and prognosticate for COVID-19 using chest radiographs and CT scans. *Nature*  
1259 *Machine Intelligence*, 3(3), 199-217. <https://doi.org/10.1038/s42256-021-00307-0>
- 1260 Rubin, M. (2022). The costs of HARKing. *British Journal for the Philosophy of Science*,  
1261 73(2), 535-560. <https://doi.org/10.1093/bjps/axz050>
- 1262 Rugamer, D., Kolb, C., & Klein, N. (2024). Semi-structured distributional regression. *The*  
1263 *American Statistician*, 78 (1), 88-99. <https://doi.org/10.1080/00031305.2022.2164054>
- 1264 Sagi, O., & Rokach, L. (2018). Ensemble learning: A survey. *WIREs Data Mining and*  
1265 *Knowledge Discovery*, 8(4), e1249. <https://doi.org/10.1002/widm.1249>
- 1266 Salat, D. H., Buckner, R. L., Snyder, A. Z., Greve, D. N., Desikan, R. S. R., Busa, E., Morris,  
1267 J. C., Dale, A. M., & Fischl, B. (2004). Thinning of the cerebral cortex in aging.  
1268 *Cerebral Cortex*, 14(7), 721–730. <https://doi.org/10.1093/cercor/bhh032>
- 1269 Sapkota, R., Raza, S., & Karkee, M. (2025). Comprehensive analysis of transparency and  
1270 accessibility of chatgpt, deepseek, and other sota large language models. arXiv  
1271 preprint. <https://arxiv.org/abs/2502.18505>
- 1272 Sasaki, K. & Yamada, Y. (2023). SPARKing: Sample-size planning after the results are  
1273 known. *Front. Hum. Neurosci.* 17:912338.  
1274 <https://doi.org/10.3389/fnhum.2023.912338>
- 1275 Schuirmann, D. J., 1987: A comparison of the two one-sided tests procedure and the power  
1276 approach for assessing the equivalence of average bioavailability. *Journal of*

- 1277 pharmacokinetics and biopharmaceutics, 15, 657–680.
- 1278 <https://doi.org/10.1007/bf01068419>
- 1279 Schurz, G. (2019). *Hume's problem solved: The optimality of meta-induction*. The MIT  
1280 Press.
- 1281 Shmueli, G., & Tafti, A. (2023). How to "improve" prediction using behavior modification.  
1282 *International Journal of Forecasting*, 39(2), 541-555.
- 1283 <https://doi.org/10.1016/j.ijforecast.2022.07.008>
- 1284 Shmueli, G. (2010). To explain or to predict? *Statistical Science*, 25(3), 289–310.
- 1285 <https://doi.org/10.1214/10-STS330>
- 1286 Stasinopoulos, D. M., Rigby, R. A., Heller, G. Z., Voudouris, V., & De Bastiani, F. (2017).  
1287 *Flexible Regression and Smoothing: Using GAMLSS in R*. Chapman and Hall/CRC.
- 1288 Stasinopoulos, M. D., Kneib, T., Klein, N., Mayr, A., & Heller, G. Z. (2024). *Generalized*  
1289 *Additive Models for Location, Scale and Shape: A Distributional Regression*  
1290 *Approach, with Applications*. Cambridge University Press.
- 1291 Stasinopoulos, M. D., Rigby, R. A., & de Bastiani, F. (2018). GAMLSS: A distributional  
1292 regression approach. *Statistical Modelling*, 18 (3-4), 248-273.
- 1293 <https://doi.org/10.1177/1471082X18759144>
- 1294 Steegen, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing Transparency  
1295 Through a Multiverse Analysis. *Perspectives on Psychological Science*, 11(5), 702-  
1296 712. <https://doi.org/10.1177/1745691616658637>
- 1297 Steinruecken, C., Smith, E., Janz, D., Lloyd, J., & Ghahramani, Z. (2019). The Automatic  
1298 Statistician. In Hutter, F., Kotthoff, L., & Vanschoren, J. (Eds.), *Automated Machine*

- 1299            *Learning. The Springer Series on Challenges in Machine Learning* (pp. 161-173).  
1300            Springer. [https://doi.org/10.1007/978-3-030-05318-5\\_9](https://doi.org/10.1007/978-3-030-05318-5_9)
- 1301    Stodden, V., Leisch, F., & Peng, R. D. (Eds.). (2014). *Implementing reproducible research*.  
1302            CRC Press.
- 1303    Stouffer, S.A., Suchman, E.A., DeVinney, L.C., Star, S.A., & Williams, R.M. (1949). *The*  
1304            *American Soldier: Adjustment during Army Life, Vol. 1*. Princeton University Press,  
1305            Princeton, NJ.
- 1306            Taleb, N. N. (2001). *Fooled by randomness*. Random House.
- 1307            Taleb, N. N. (2007). *The black swan. The impact of the highly improbable*. Random  
1308            House.
- 1309    Thielmann, A. F., Kruse, R.-M., Kneib, T., & Säfken, B. (2024). Neural additive models for  
1310            location scale and shape: A framework for interpretable neural regression beyond the  
1311            mean. In S. Dasgupta, S. Mandt, and Y. Li (Eds.), *Proceedings of the 27th*  
1312            *International Conference on Artificial Intelligence and Statistics (Vol. 238)*. PMLR.  
1313            <https://proceedings.mlr.press/v238/frederik-thielmann24a.html>
- 1314    Toribio-Flórez, D., Anneser, L., de Oliveira-Lopes, F., Pallandt, M., Tunn, I., Windel, H., &  
1315            Max Planck PhDnet Open Science Group (2021). Where Do Early Career Researchers  
1316            Stand on Open Science Practices? A Survey Within the Max Planck Society.  
1317            *Frontiers in Research Metrics and Analytics*, 5, 586992.  
1318            <https://doi.org/10.3389/frma.2020.586992>
- 1319            Trafimow, D. (2019). A taxonomy of model assumptions on which ‘p’ is based and  
1320            implications for added benefit in the sciences. *International Journal of Social*

- 1321 *Research Methodology*, 22(6), 571-583.
- 1322 <https://doi.org/10.1080/13645579.2019.1610592>
- 1323 Tu, J., Liu, W., & Mao, X. (2024). Distributed Estimation on Semi-Supervised Generalized  
1324 Linear Model. *Journal of Machine Learning Research*, 25, 1-41.
- 1325 <http://jmlr.org/papers/v25/22-0670.html>
- 1326 Umlauf, N., Klein, N., & Zeileis, A. (2018). BAMLSS: Bayesian Additive Models for  
1327 Location, Scale and Shape (and Beyond). *Journal of Computational and Graphical*  
1328 *Statistics*, 27(3), 612–627. <https://doi.org/10.1080/10618600.2017.1407325>
- 1329 van Buuren, S., & Fredriks, M. (2001). Worm plot: simple diagnostic device for modelling  
1330 growth reference curves. *Statistics in Medicine*, 20, 1259-1277.
- 1331 <https://doi.org/10.1002/sim.746>
- 1332 Van der Laan, M. J., & Starmans, R. J. C. M. (2014). Entering the era of data science:  
1333 Targeted learning and the integration of statistics and computational data analysis.  
1334 *Advances in Statistics*, 502678. <https://doi.org/10.1155/2014/502678>
- 1335 Vega-Pons, S., & Ruiz-Shulcloper, J. (2011). A survey of clustering ensemble algorithms.  
1336 *International Journal of Pattern Recognition and Artificial Intelligence*, 25(3), 337-  
1337 372. <https://doi.org/10.1142/S0218001411008683>
- 1338 Vidal-Piñeiro, D., Parker, N., Shin, J., et al. (2020). Cellular correlates of cortical thinning  
1339 throughout the lifespan. *Scientific Reports*, 10, 21803. [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-020-78471-3)  
1340 [020-78471-3](https://doi.org/10.1038/s41598-020-78471-3)
- 1341 Wagenmakers, E.-J., Beek, T., Dijkhoff, L., Gronau, Q. F., Acosta, A., Adams, R. B.,  
1342 Albohn, D. N., Allard, E. S., Benning, S. D., Blouin-Hudon, E.-M., Bulnes, L. C.,

- 1343 Caldwell, T. L., Calin-Jageman, R. J., Capaldi, C. A., Carfagno, N. S., Chasten, K. T.,  
1344 Cleeremans, A., Connell, L., DeCicco, J. M., ... Zwaan, R. A. (2016). Registered  
1345 Replication Report: Strack, Martin, & Stepper (1988). *Perspectives on Psychological*  
1346 *Science*, 11(6), 917–928. <https://doi.org/10.1177/1745691616674458>
- 1347 Wasserstein, R. & Lazar, N. (2016) The ASA Statement on p-Values: Context, Process, and  
1348 Purpose, *The American Statistician*, 70(2), 129-133.  
1349 <https://doi.org/10.1080/00031305.2016.1154108>
- 1350 Westlin, C., Theriault, J. E., Katsumi, Y., Nieto-Castanon, A., Kucyi, A., Ruf, S. F., Brown,  
1351 S. M., Pavel, M., Erdogmus, D., Brooks, D. H., Quigley, K. S., Whitfield-Gabrieli, S.,  
1352 & Barrett, L. F. (2023). Improving the study of brain-behavior relationships by  
1353 revisiting basic assumptions. *Trends in Cognitive Sciences*, 27(3), 246–257.  
1354 <https://doi.org/10.1016/j.tics.2022.12.015>
- 1355 Yarkoni, T. (2022) The generalizability crisis. *Behavioral and Brain Sciences* 45, e1: 1–78.  
1356 <https://doi.org/10.1017/S0140525X20001685>
- 1357 Zhao, J., & Liu, N. (2021). A safe semi-supervised classification algorithm using multiple  
1358 classifiers ensemble. *Neural Processing Letters*, 53(4), 2603–2616.  
1359 <https://doi.org/10.1007/s11063-020-10191-1>