


# Finite sample corrections for average equivalence testing

Younes Boulaguiem<sup>1</sup>  | Julie Quartier<sup>2,3</sup> | Maria Lapteva<sup>2,3</sup> |  
Yogeshvar N. Kalia<sup>2,3</sup> | Maria-Pia Victoria-Feser<sup>1</sup> |  
Stéphane Guerrier<sup>1,2,3</sup> | Dominique-Laurent Couturier<sup>4,5</sup>

<sup>1</sup>Geneva School of Economics and Management, University of Geneva, Switzerland

<sup>2</sup>School of Pharmaceutical Sciences, University of Geneva, Switzerland

<sup>3</sup>Institute of Pharmaceutical Sciences of Western Switzerland, University of Geneva, Switzerland

<sup>4</sup>Medical Research Council Biostatistics Unit, University of Cambridge, England

<sup>5</sup>Cancer Research UK – Cambridge Institute, University of Cambridge, England

## Correspondence

Dominique-Laurent Couturier, Medical Research Council Biostatistics Unit, University of Cambridge, Cambridge, UK.  
Email:  
[dominique.couturier@mrc-bsu.cam.ac.uk](mailto:dominique.couturier@mrc-bsu.cam.ac.uk)

## Funding information

Cancer Research UK, Grant/Award Number: C9545/A29580; Innosuisse - Schweizerische Agentur für Innovationsförderung, Grant/Award Numbers: 37308.1 IP-ENG, 53622.1 IP-ENG; Schweizerischer Nationalfonds zur Förderung der Wissenschaftlichen Forschung, Grant/Award Numbers: 176843, 182684, 211007; UK Medical Research Council, Grant/Award Number: MC\_UU\_00002/14

Average (bio)equivalence tests are used to assess if a parameter, like the mean difference in treatment response between two conditions for example, lies within a given equivalence interval, hence allowing to conclude that the conditions have “equivalent” means. The *two one-sided tests* (TOST) procedure, consisting in testing whether the target parameter is respectively significantly greater and lower than some pre-defined lower and upper equivalence limits, is typically used in this context, usually by checking whether the confidence interval for the target parameter lies within these limits. This intuitive and visual procedure is however known to be conservative, especially in the case of highly variable drugs, where it shows a rapid power loss, often reaching zero, hence making it impossible to conclude for equivalence when it is actually true. Here, we propose a finite sample correction of the TOST procedure, the  $\alpha$ -TOST, which consists in a correction of the significance level of the TOST allowing to guarantee a test size (or type-I error rate) of  $\alpha$ . This new procedure essentially corresponds to a finite sample and variability correction of the TOST procedure. We show that this procedure is uniformly more powerful than the TOST, easy to compute, and that its operating characteristics outperform the ones of its competitors. A case study about econazole nitrate deposition in porcine skin is used to illustrate the benefits of the proposed method and its advantages compared to other available procedures.

## KEYWORDS

bioequivalence, interval inclusion principle, scaled average bioequivalence, similarity test, two one-sided test

## 1 | INTRODUCTION

Equivalence tests, also known as similarity or parity tests, have gained significant attention during the past two decades. They originated from the field of pharmacokinetics,<sup>1,2</sup> where they are called bioequivalence tests and have numerous applications in both research and production.<sup>3</sup> They find their most common application in the manufacturing of generic medicinal drugs, where, by proving that the generic version has a similar bioavailability to its well-studied brand-name

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Authors. *Statistics in Medicine* published by John Wiley & Sons Ltd.

counterpart, the manufacturer can considerably shorten the approval process for the generic drug.<sup>4</sup> Moreover, equivalence tests have attracted growing interest in other domains and for other types of purposes, such as in production when, for example, the mode of administration is altered or when the production site is changed,<sup>5</sup> or in the social and behavioral sciences for the evaluation of replication results and corroborating risky predictions.<sup>6</sup> Very recent literature reflects the expanding use of equivalence tests across a growing range of domains. Examples include the investigation of the equivalence of virtual reality imaging measurements by feature,<sup>7</sup> of cardiovascular responses to stimuli by sex,<sup>8</sup> of children neurodevelopment,<sup>9</sup> chemotherapy efficacy and safety by treatment,<sup>10</sup> of post-stroke functional connectivity patterns by patient group,<sup>11</sup> of risk-taking choices by moral type,<sup>12</sup> and of 2020 US presidential election turnout by political advertising condition.<sup>13</sup> Review articles have also appeared, for example, in food sciences,<sup>14</sup> in psychology,<sup>15</sup> in sport sciences,<sup>16</sup> and in pharmaceutical sciences.<sup>17</sup>

Equivalence testing implies defining an equivalence region within which the parameter of interest, such as the difference between outcome means measured under two conditions, would lie, for these conditions to be considered equivalent. Indeed, when comparing two treatments, for example, differences in therapeutic effects that belong to the equivalence region would typically be considered as negligible or irrelevant. This is different from standard equality-of-means hypothesis tests in which the null and the alternative hypothesis are interchanged and the null hypothesis states that both means are equal rather than equivalent.

Formally, a canonical form for the average equivalence problem consists of two independent random variables  $\hat{\theta}$  and  $\hat{\sigma}_v$  having the distributions

$$\hat{\theta} \sim \mathcal{N}(\theta, \sigma_v^2) \quad \text{and} \quad \frac{v\hat{\sigma}_v^2}{\sigma_v^2} \sim \chi_v^2, \quad (1)$$

where  $\theta$  and  $\sigma_v^2$  respectively denote the target equivalence parameter and its variance, depending on the number of the degrees of freedom  $v$  which is a function of the sample size and total number of parameters. This setting is very general. It covers cases where, for example, the bioequivalence parameter corresponds to the difference in means of the (logarithm of pharmacokinetic) responses between two experimental conditions, or to an element of the parameter vector of a (generalized) linear mixed effect model, like the difference between the slopes of two conditions in a longitudinal study. The hypotheses of interest are given by

$$H_0 : \theta \notin \Theta_1 \quad \text{vs} \quad H_1 : \theta \in \Theta_1 := (\delta_L, \delta_U), \quad (2)$$

where  $\delta_L$  and  $\delta_U$  are known constants. Without loss of generality, it can be assumed that the equivalence limits are symmetrical around zero. In this case,  $c := \delta_U = -\delta_L$  so that  $\Theta_1 = (-c, c)$ . Equivalence is typically investigated via the *two one-sided tests* (TOST) procedure,<sup>18</sup> consisting in testing whether the target parameter is respectively significantly greater than  $-c$  and lower than  $c$ , with the test size, or type-I error rate, controlled at the significance (or nominal) level  $\alpha$ , usually chosen as 5%.<sup>19</sup> More precisely, the TOST is level- $\alpha$ , meaning that its size is smaller or equal to  $\alpha$  (see (4) in Section 2.1, for its formal definition). The most common way of assessing equivalence is to use the *interval inclusion principle* (IIP) and check whether the  $100(1 - 2\alpha)\%$  confidence interval (CI) for the target parameter falls within the equivalence margins  $(-c, c)$ .<sup>3,20</sup> This strategy has been shown to lead to the same test decision as the TOST procedure if the CI is equi-tailed.<sup>21,22</sup>

However, it is well known that this procedure can be conservative as the size of the TOST can be considerably lower than the (specified) significance level  $\alpha$ . This induces a drop in power and therefore to a lower probability of detecting truly equivalent mean effects as such. This problem is particularly noticeable in cases where  $\sigma_v$  is relatively large. Such situations may occur, for example, when the sample size is determined using an underestimated standard deviation value obtained from a prior experiment, or with studies involving highly variable drugs and in which the sample size that would be needed to achieve reasonable values for  $\sigma_v$  is unrealistic. For that purpose, Anderson and Hauck<sup>23</sup> proposed a test that has greater power than the TOST for situations where  $\sigma_v$  is relatively large. This test, referred to here as the AH-test, can be liberal and therefore does not control the size.<sup>22</sup> In some cases, it can also lead to the equivalence declaration (ie, acceptance of equivalence through the rejection of the null hypothesis in (2) at the  $\alpha$  level) when  $\theta$ , the target parameter of interest, falls outside the equivalence interval.<sup>18</sup> Brown et al<sup>24</sup> constructed an unbiased test that is uniformly more powerful than the TOST, however, it is computationally intensive and its rejection region may exhibit rather irregular shapes in some cases.<sup>22</sup> Berger et al<sup>22</sup> therefore proposed a smoothed version. These tests cannot be assessed using the IIP and the last two are difficult to interpret due to the use of polar coordinates.<sup>25</sup>

In the specific context of average bioequivalence testing in replicated crossover designs<sup>26</sup> for highly variable drugs, that is, for cases with relatively large  $\sigma_v$ , regulatory authorities have recommended an alternative approach based on the linear scaling of the bioequivalence limits according to the value of the standard deviation within the reference group, called Scaled Average BioEquivalence (SABE),<sup>27</sup> also referred to as Average BioEquivalence with Expanding limits (ABEL) in some references,<sup>26,28,29</sup> with the constraint that  $\hat{\theta}$  lies within the bioequivalence margins  $(-c, c)$ . These recommendations were issued by the European Medicines Agency (EMA) and the US Food and Drug Administration (FDA).<sup>20,30</sup> The amount of expansion is limited by the authorities, and several recent publications have shown that the size of the SABE can be larger than the significance level  $\alpha$ <sup>20,26,31-33</sup> and therefore have proposed different ways to correct for it.<sup>28,34-36</sup> These corrections ensure that the size is smaller than or equal to  $\alpha$  and lead to acceptance regions that change more smoothly with  $\sigma_v$ .

In this article, as an alternative to previous methods, we propose a finite sample correction of the TOST procedure that simply consists in a correction of the TOST's significance level to guarantee a size- $\alpha$  test when  $\sigma_v$  is known. This correction is design-agnostic and can be used with parallel or (replicated) crossover designs, for example. The corrected significance level  $\alpha^*$  is straightforward to compute and allows to define  $100(1 - 2\alpha^*)\%$  CIs used in the classical TOST. Hence, the  $\alpha$ -TOST essentially corresponds to a finite sample continuous variability correction of the TOST procedure, that leads to an increased probability of declaring equivalence when it is true for large values of  $\sigma_v$  while maintaining a size of exactly  $\alpha$  when  $\sigma_v$  is known. Indeed, the  $\alpha$ -TOST is shown to be uniformly more powerful than the TOST and, for small to moderate values of  $\sigma_v$ , to be nearly equivalent to the TOST with a comparable power as  $\alpha^* \approx \alpha$  in such cases. Since, in practice,  $\sigma_v$  needs to be estimated from the data, a straightforward estimator for  $\alpha^*$  is also proposed. It is shown, through an extensive simulation study considering the canonical form defined in (1) and therefore valid in a wide range of settings, that the estimator remains level- $\alpha$  and its size stays close to  $\alpha$ . Our simulation study also considers a version of the TOST that adjusts the equivalence limits  $\delta_L$  and  $\delta_U$  instead of the level, to guarantee a size- $\alpha$  test and therefore referred to as the  $\delta$ -TOST. Our results show that the  $\alpha$ -TOST is both more powerful and accurate than the standard TOST and  $\delta$ -TOST, indicating that, when looking for a design-agnostic correction valid in general settings, a correction on the level ( $\alpha$ -TOST) leads to better operating characteristics. A comparison of the performance of these methods to the corrected SABE, that consists in an adjustment on both the equivalence bounds and the level, is presented in Appendix E in a simple paired setting. More adequate and extensive comparisons, considering the different adjustments proposed by regulatory agencies and other authors, including variants such as the corrected SABE, are needed in the specific case of average equivalence testing with replicated crossover designs and are left for further research.

The article is organized as follows. The  $\alpha$ -TOST is presented in Section 2 from a suitable formulation of the TOST. Its statistical properties as well as a simple algorithm to compute  $\alpha^*$  are also provided. In Section 3, an extensive simulation study is used to compare the empirical performances of the  $\alpha$ -TOST,  $\delta$ -TOST and standard TOST. In Section 4, we consider a case study for which we apply the TOST and the  $\alpha$ -TOST, as well as other available methods, in order to showcase the advantages of our proposed design-agnostic approach. Finally, Section 5 discusses some potential extensions.

## 2 | EQUIVALENCE TESTING

In this section, we present the methodology for deriving a corrected statistical equivalence test. We first present the TOST and its properties. We then define the  $\alpha$ -TOST procedure through a natural correction of the TOST, derive its statistical properties, propose an iterative procedure to compute the corrected level  $\alpha^*$  and show that this procedure is exponentially fast. We also show that the  $\alpha$ -TOST is uniformly more powerful than the TOST.

### 2.1 | The TOST procedure

For testing the hypotheses in (2), the TOST uses the two following test statistics:

$$Z_L := \frac{\hat{\theta} + c}{\hat{\sigma}_v} \quad \text{and} \quad Z_U := \frac{\hat{\theta} - c}{\hat{\sigma}_v},$$

where  $Z_L$  tests for  $H_{01} : \theta \leq -c$  vs  $H_{11} : \theta > -c$ , and  $Z_U$  tests for  $H_{02} : \theta \geq c$  vs  $H_{12} : \theta < c$ . At a significance level  $\alpha$ , the TOST therefore rejects  $H_0 := H_{01} \cup H_{02}$  (ie,  $\theta \notin \Theta_1$ ) in favor of  $H_1 := H_{11} \cap H_{12}$  (ie,  $\theta \in \Theta_1$ ) if both tests simultaneously reject their marginal null hypotheses, that is, if

$$Z_L \geq t_{\alpha, \nu} \quad \text{and} \quad Z_U \leq -t_{\alpha, \nu},$$

where  $t_{\alpha, \nu}$  denotes the upper  $\alpha$  quantile of a  $t$ -distribution with  $\nu$  degrees of freedom. The corresponding rejection region of the TOST is given by

$$C_1 := \left\{ \hat{\theta} \in \mathbb{R}, \hat{\sigma}_\nu \in \mathbb{R}_+ \mid c \geq |\hat{\theta}| + t_{\alpha, \nu} \hat{\sigma}_\nu \right\}. \quad (3)$$

Consequently, equivalence cannot be declared with the TOST for all  $\hat{\sigma}_\nu > \hat{\sigma}_{\max} := c/t_{\alpha, \nu}$ , even for  $\hat{\theta} = 0$  (see also Figure 6 of Section 4).

With respect to the TOST's size, given  $\alpha$ ,  $\theta$ ,  $\sigma_\nu$ ,  $\nu$ , and  $c$ , the probability of declaring equivalence can be expressed as follows<sup>37</sup>:

$$\begin{aligned} p(\alpha, \theta, \sigma_\nu, \nu, c) &:= \Pr\left(Z_L \geq t_{\alpha, \nu} \text{ and } Z_U \leq -t_{\alpha, \nu} \mid \alpha, \theta, \sigma_\nu, \nu, c\right) \\ &= Q_\nu\left(-t_{\alpha, \nu}, \frac{\theta - c}{\sigma_\nu}, \lambda\right) - Q_\nu\left(t_{\alpha, \nu}, \frac{\theta + c}{\sigma_\nu}, \lambda\right), \end{aligned} \quad (4)$$

where  $\lambda := \frac{c\sqrt{\nu}}{\sigma_\nu t_{\alpha, \nu}}$  and  $Q_\nu(t, y, z)$  corresponds to a special case of Owen's Q-function<sup>38</sup> defined as

$$Q_\nu(t, y, z) := \frac{\sqrt{2\pi}}{\Gamma\left(\frac{1}{2}\nu\right) 2^{\frac{1}{2}(\nu-2)}} \int_0^z \Phi\left(\frac{tx}{\sqrt{\nu}} - y\right) x^{\nu-1} \varphi(x) dx,$$

where  $\varphi(x)$  and  $\Phi(x)$  denote the probability and cumulative distribution functions of a standard normal distribution, respectively.

Then, for given values of  $\alpha$ ,  $\sigma_\nu$ , and  $\nu$ , the TOST's size is defined as the supremum of (4),<sup>39</sup> and is given by

$$\omega(\alpha, c, \sigma_\nu, \nu) := \sup_{\theta \notin \Theta_1} p(\alpha, \theta, \sigma_\nu, \nu, c) = p(\alpha, c, \sigma_\nu, \nu, c) = Q_\nu(-t_{\alpha, \nu}, 0, \lambda) - Q_\nu\left(t_{\alpha, \nu}, \frac{2c}{\sigma_\nu}, \lambda\right). \quad (5)$$

We can then deduce that the TOST is level- $\alpha$ , by noting that, for  $\sigma_\nu > 0$ , we have

$$\begin{aligned} \omega(\alpha, c, \sigma_\nu, \nu) &< Q_\nu(-t_{\alpha, \nu}, 0, \lambda) = \frac{\sqrt{2\pi}}{\Gamma\left(\frac{1}{2}\nu\right) 2^{\frac{1}{2}(\nu-2)}} \int_0^\lambda \Phi\left(-\frac{t_{\alpha, \nu}x}{\sqrt{\nu}}\right) x^{\nu-1} \varphi(x) dx \\ &< \frac{\sqrt{2\pi}}{\Gamma\left(\frac{1}{2}\nu\right) 2^{\frac{1}{2}(\nu-2)}} \int_0^\infty \Phi\left(-\frac{t_{\alpha, \nu}x}{\sqrt{\nu}}\right) x^{\nu-1} \varphi(x) dx = \Pr(T_\nu \leq -t_{\alpha, \nu}) = \alpha, \end{aligned} \quad (6)$$

where  $T_\nu$  denotes a random variable following a  $t$ -distribution with  $\nu$  degrees of freedom, so that

$$\lim_{\sigma_\nu \rightarrow 0} \omega(\alpha, c, \sigma_\nu, \nu) = \alpha.$$

Thus, while the TOST is indeed level- $\alpha$ , it actually never achieves a size of  $\alpha$ , except in the theoretical case of  $\sigma_\nu = 0$ , as already highlighted by several authors.<sup>36</sup> When  $\hat{\sigma}_\nu$  is small, the difference between the size and  $\alpha$  is marginal, but as  $\hat{\sigma}_\nu$  approaches or exceeds  $c/t_{\alpha, \nu}$ , this difference increases, leading to a high probability of the TOST failing to detect equivalence when it exists. As a solution to this issue, we suggest an alternative approach, the  $\alpha$ -TOST, that corrects the size of the TOST for a large range of values of  $\hat{\sigma}_\nu$  and still allows to assess equivalence by means of confidence intervals, as depicted in Figure 5 of Section 4.

## 2.2 | The $\alpha$ -TOST

A corrected version of the TOST can theoretically be constructed by adjusting the significance level and using  $\alpha^*$  instead of  $\alpha$  in the standard TOST procedure, where

$$\alpha^* := \alpha^*(\sigma_v) = \underset{\gamma \in [\alpha, 0.5]}{\operatorname{argzero}} [\omega(\gamma, c, \sigma_v, \nu) - \alpha], \quad (7)$$

with  $\omega(\gamma, c, \sigma_v, \nu)$  defined in (5). The dependence of  $\alpha^*$  on  $\alpha$  and  $\nu$  is omitted from the notation as these quantities are known. A similar type of correction was also used to amend the significance level of the SABE procedure by Labes and Schütz<sup>28</sup> and Ocaña and Muñoz<sup>35</sup> (see also Palmes et al<sup>40</sup> for power adjustment). However, in these cases, the corrected significance level was reduced (instead of increased like in (7)) so that the size does not exceed the significance level of  $\alpha$ . The aim of these corrections is therefore not the same as the one proposed here. Furthermore, the size of the  $\alpha$ -TOST is guaranteed to be exactly  $\alpha$  when  $\sigma_v$  is known, which is not the case for these competing methods.

In Appendix A, we demonstrate that the existence of  $\alpha^*$  relies on a simple condition that is satisfied in most settings of practical importance. In particular, this requirement can be translated into a maximal value for the estimated standard error  $\hat{\sigma}_v$ , that is  $\hat{\sigma}_v < \frac{2c}{\Phi^{-1}(\alpha+0.5)}$ . Moreover, since,  $\alpha^*(\sigma_v)$  is a population size quantity as it depends on the unknown quantity  $\sigma_v$ , a natural estimator for its sample value is given by

$$\hat{\alpha}^* := \alpha^*(\hat{\sigma}_v) = \underset{\gamma \in [\alpha, 0.5]}{\operatorname{argzero}} [\omega(\gamma, c, \hat{\sigma}_v, \nu) - \alpha]. \quad (8)$$

Hence, in practice, based on the (estimated) corrected significance level  $\hat{\alpha}^*$ , the  $\alpha$ -TOST procedure rejects the non-equivalence null hypothesis in favor of the equivalence one at the significance level  $\alpha$ , if  $Z_L > t_{1-\hat{\alpha}^*, \nu}$  and  $Z_U < -t_{1-\hat{\alpha}^*, \nu}$ . In Appendix B, we study the asymptotic properties of  $\hat{\alpha}^*$  and show that  $\hat{\alpha}^* = \alpha^* + o_p(\nu^{-1})$ . Informally, this result implies that the uncertainty associated to  $\hat{\alpha}^*$  is (asymptotically) negligible compared to the uncertainty associated to  $\hat{\theta}$  and  $\hat{\sigma}_v$  as these terms have slower convergence rates in that  $\hat{\theta} = \theta + \mathcal{O}_p(\nu^{-1/2})$  and  $\hat{\sigma}_v = \sigma_v + \mathcal{O}_p(\nu^{-1})$ . This result also suggests that the  $\alpha$ -TOST procedures based on  $\alpha^*$  or on  $\hat{\alpha}^*$  are expected to provide very similar finite sample performances.

In Section 3, we consider an extensive Monte Carlo simulation study to compare the empirical performances of different methods when  $\sigma_v$  needs to be estimated. For the  $10^4$  simulation settings we considered (ie, 100 values for  $\sigma_v$  and 100 for  $\nu$  covering most combinations of interest, see Simulation 2 in Table 1), we find that the empirical size of the  $\alpha$ -TOST is generally closer to the nominal level  $\alpha$  in comparison to the other methods (see Figure 2 in Section 3). We also find that in less than 1% of the settings, the  $\alpha$ -TOST procedure can be slightly liberal, with a maximal empirical size of 0.05311 (see Figure D3 in Appendix D). However, this behavior can mostly be explained by the randomness associated to our large scale simulation.

The corrected significance level  $\hat{\alpha}^*$  can easily be computed using the following iterative approach. At iteration  $k$ , with  $k \in \mathbb{N}$ , we define

$$\hat{\alpha}^{*(k+1)} = \alpha + \hat{\alpha}^{*(k)} - \omega(\hat{\alpha}^{*(k)}, c, \hat{\sigma}_v, \nu), \quad (9)$$

with  $\omega(\alpha, c, \sigma_v, \nu)$  given in (5) and where the procedure is initialized at  $\hat{\alpha}^{*(0)} = \alpha$ . This simple iterative approach converges exponentially fast to  $\hat{\alpha}^*$  as it can be shown that

$$|\hat{\alpha}^{*(k+1)} - \hat{\alpha}^*| < \frac{1}{2} \exp(-bk),$$

for some positive constant  $b$  (see Appendix C for more details).

Finally, since the conclusion of  $\alpha$ -TOST considers an interval computed using a smaller value than  $t_{\alpha, \nu}$  compared to the TOST, the  $\alpha$ -TOST rejection interval is necessarily larger than its TOST counterpart as  $\hat{\alpha}^* > \alpha$ . This implies that the  $\alpha$ -TOST is uniformly more powerful than the TOST, and explains cases like the one encountered in the porcine skin case study presented in Section 4, in which equivalence is declared using the  $\alpha$ -TOST but not with the TOST (which has an empirical power of zero given  $\hat{\sigma}_v$ ).

**TABLE 1** Parameter values used in each simulation, where  $c$  denotes the tolerance limit,  $\nu$  the number of degrees of freedom,  $\theta$  the target parameter and  $\sigma_\nu$  its standard deviation,  $\alpha$  the target significance level and  $B$  the number of Monte Carlo samples per simulation.

	Simulation			
	1	2	3	4
$c$	log(1.25) $\approx$ 0.2231			
$\nu$	15, 30, 45	5, 6, 7, ..., 100, 250, 500, 750, 1000 (100 values)		45
$\sigma_\nu$	0.08, 0.12, 0.16	100 evenly spaced values between 0.01 and 0.3		0.08, 0.12, 0.16
$\theta$	30 evenly spaced values between 0 and 0.26	$c$	0	30 evenly spaced values between 0 and 0.26
$\alpha$	0.05			
$B$	$10^5$			
Design	General (canonical form)			Paired
Methods	TOST, $\alpha$ -TOST and $\delta$ -TOST			TOST, $\alpha$ -TOST, $\delta$ -TOST, SABE and cSABE

Appendix E

### 3 | SIMULATION STUDY

In this section, we conduct an extensive Monte Carlo simulation study with parameters settings per simulation reported in Table 1. Simulations 1 to 3, performed under the canonical form defined in (1) and therefore valid in a wide range of settings, assess the empirical performances of the  $\alpha$ -TOST and compare them to the ones of the standard TOST and  $\delta$ -TOST methods, where the latter, defined below, considers a correction on the equivalence limits rather than on the level to reach a size of  $\alpha$ . Simulation 4, presented in Appendix E, investigates the empirical performances of these methods with the ones of the design-specific SABE and corrected SABE, where the latter consists in an adjustment on both the level and the equivalence limits. In that simulation, we consider a paired design setting that is closely related to the example considered in our case study and that allows us to estimate the within-subject variability of the reference treatment required by SABE-like methods. All simulations consider a target significance level of 5%, a value of  $c$  equal to log(1.25) and  $10^5$  Monte Carlo samples per configuration.

Formally, the  $\delta$ -TOST is defined as follows

$$\delta^* := \delta^*(\sigma_\nu) = \underset{\delta \in [c, \infty)}{\operatorname{argzero}} [\omega(\alpha, \delta, \sigma_\nu, \nu) - \alpha]. \quad (10)$$

Using the same arguments as in Appendix A, we can easily demonstrate that a unique solution always exists, regardless of the value of the standard error for the  $\delta$ -TOST. However, an exponentially fast iterative algorithm cannot be used to find the solution for this method. This highlights an important practical advantage of the  $\alpha$ -TOST over the  $\delta$ -TOST and, to a larger extent, over the corrected SABE, which relies on both Monte Carlo integration and numerical optimization procedures to define its correction (see Appendix E).

In our simulations, the empirical performances of the TOST,  $\alpha$ -TOST and  $\delta$ -TOST are defined using the following steps:

1. **Simulation:** for a given Monte Carlo sample  $b = 1, \dots, B$ :
  - (a) simulate a value for  $\hat{\theta}_b \sim \mathcal{N}(\theta, \sigma_\nu)$  given the values  $\theta$  and  $\sigma_\nu$  of interest,
  - (b) simulate a value for  $t = \frac{\nu \hat{\sigma}_b^2}{\sigma_\nu^2} \sim \chi_\nu^2$  and set  $\hat{\sigma}_{\nu,b} = \sqrt{t \sigma_\nu^2 / \nu}$ .
2. **Finite sample adjustments:** for a given Monte Carlo iteration  $b = 1, \dots, B$ :
  - (a)  $\alpha$ -TOST:

- (i) compute  $\hat{\alpha}_b^*$  using (the algorithm associated to) (7) with  $\hat{\sigma}_{v,b}$ ,
- (ii) compute  $t_{\hat{\alpha}_b^*,v}$ .

(b)  $\delta$ -TOST: compute  $\hat{\delta}_b^*$  using (10) with  $\hat{\sigma}_{v,b}$ .

### 3. Empirical probability of declaring equivalence:

(a) TOST:

$$\text{empirical probability} = \frac{1}{B} \sum_{b=1}^B \eta \left( c \geq |\hat{\theta}_b| + t_{\alpha,v} \hat{\sigma}_{v,b} \right),$$

where  $\eta(\cdot)$  denote the indicator function with  $\eta(A) = 1$  if  $A$  is true and zero otherwise,

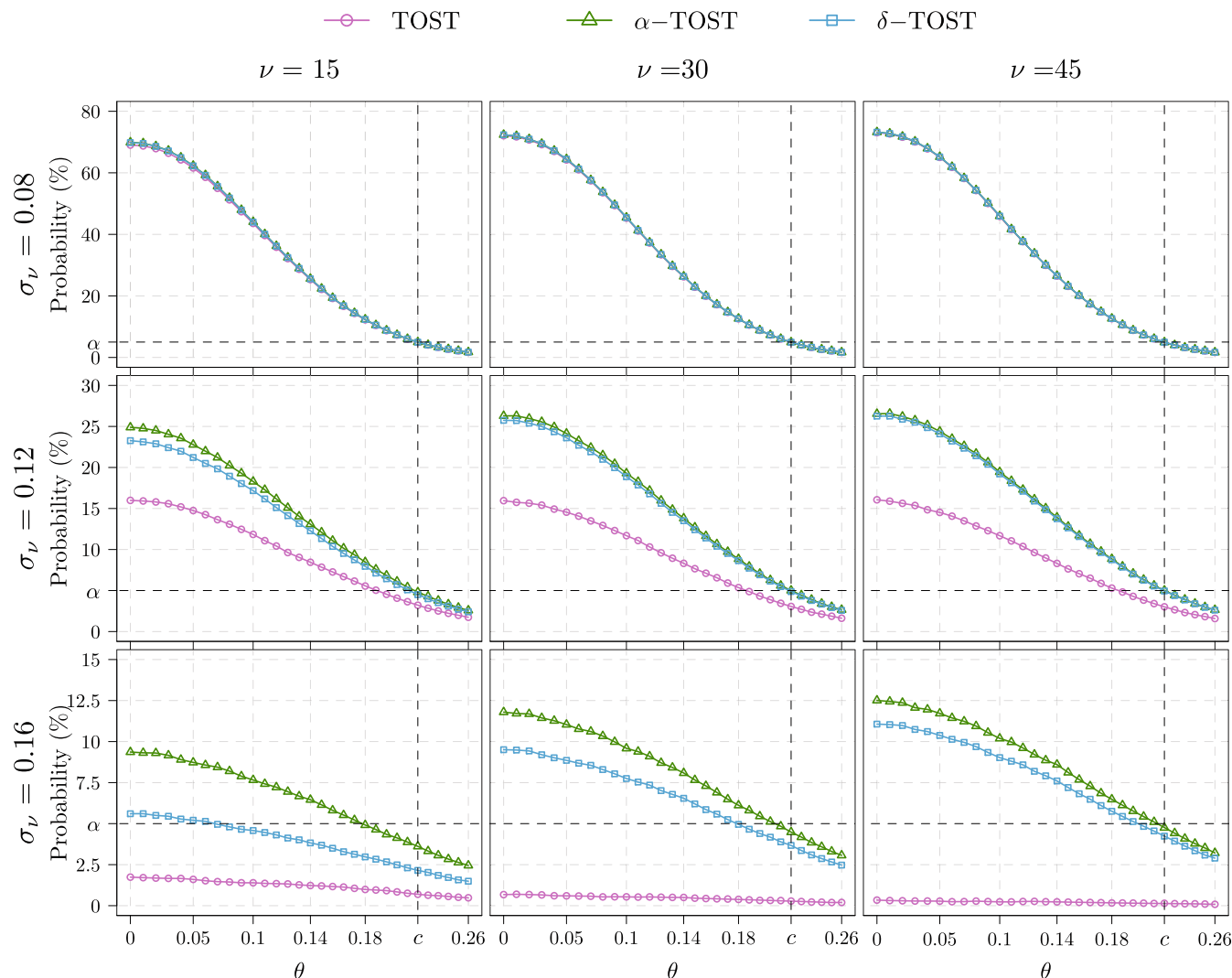
- (b)  $\alpha$ -TOST: same as the TOST in Step 3(a) but replacing  $t_{\alpha,v}$  by  $t_{\hat{\alpha}_b^*,v}$ ,
- (c)  $\delta$ -TOST: same as the TOST in Step 3(a) but replacing  $c$  by  $\hat{\delta}_b^*$ .

Simulation 1 investigates the probability of declaring equivalence for varying values of  $\theta$  allowing to study both the power and the size of each methods for combinations of selected values of  $v$  and  $\sigma_v$ . Simulation results are presented in Figure 1, which shows, for each method of interest, the empirical probability of declaring equivalence as a function of  $\theta$  for different combinations of values of  $v$  (rows) and  $\sigma_v$  (columns). For small values of  $\sigma_v$ , the empirical performance of all methods are similar. For moderate to large values of  $\sigma_v$ , we can note that the TOST is conservative, with an empirical size far smaller than the nominal level  $\alpha = 5\%$  when  $\theta = c$ , and that it quickly reaches an empirical power of 0 for large values of  $\sigma_v$ . On the other hand, the  $\alpha$ -TOST and  $\delta$ -TOST have a higher power throughout, are generally size- $\alpha$  but are a bit conservative for large values of  $\sigma_v$  and relatively small values of  $v$ . This deviation from the nominal level  $\alpha$  for the  $\alpha$ -TOST and  $\delta$ -TOST is due to the estimation error induced by using  $\hat{\sigma}_v$  instead of  $\sigma_v$  to construct an adjustment to the TOST. However, our simulation results confirm that such adjustment, either on the level or the equivalence bounds, considerably improves both the size and power in finite samples, especially with larger  $\sigma_v$  where it prevents it from becoming 0. Moreover, this simulation suggests that the  $\alpha$ -TOST outperforms the  $\delta$ -TOST, indicating that an adjustment on the level provides both a more accurate and a more powerful test than an adjustment on the equivalence bounds.

Simulations 2 and 3, respectively performed with  $\theta = \log(c)$  and  $\theta = 0$ , investigate the empirical size and power for  $10^4$  settings defined as combinations of 100 values of  $v$  and 100 values of  $\sigma_v$  chosen to cover *most* cases of practical interest. Figures D1, D2 and D3 in Appendix D respectively show the results of Simulation 2 for the standard TOST,  $\alpha$ -TOST and  $\delta$ -TOST. Each figure consists in a heatmap displaying the empirical size, computed by replacing  $\sigma_v$  and  $\theta$  by realizations of their random variables in (1) to reproduce the parameter estimation process, for all combinations of values of  $v$  and  $\sigma_v$  of interest. Figure D1 shows that the TOST is size- $\alpha$  only for relatively small values of  $\sigma_v$  (below 0.09), and that its size decreases abruptly as  $\sigma_v$  increases to reach 0. We can note that the value of  $v$  does not seem to have an important effect on the size of the TOST. In comparison, Figures D2 and D3 show that both the  $\delta$ -TOST and  $\alpha$ -TOST are size- $\alpha$  for a larger number of combinations of values for  $\sigma_v$  and  $v$  and that the probability of being size- $\alpha$  increases with  $v$  for a given value of  $\sigma_v$ . A comparison of Figures D2 and D3 shows that the  $\alpha$ -TOST is both more powerful and more accurate than the  $\delta$ -TOST overall, a conclusion in agreement with results of Simulation 1. A look at the proportion of configurations with an empirical size significantly greater than  $\alpha$ —as assessed by a two-sided binomial exact test at the 1% level performed on the results obtained on the  $10^5$  Monte Carlo samples per setting—shows that the  $\alpha$ -TOST procedure is slightly liberal in 0.9% of the configurations considered in Simulation 2, with a maximal empirical size of 0.05311, compared to 0.0528 for the TOST and  $\delta$ -TOST. This behavior can largely be attributed to the randomness inherent to our large-scale simulation.

Figure 2 summarizes the results of Simulation 2 by displaying, for each method of interest, a histogram of the empirical sizes obtained over the  $10^4$  configurations considered in the simulation. A comparison of these histograms clearly shows that the  $\alpha$ -TOST has an empirical size closer to the nominal level  $\alpha$  for a larger number of settings compared to the  $\delta$ -TOST and standard TOST, which shows a large clump-at-zero (31.5%) corresponding to configurations with a power of 0.

The heatmaps in Figures D4–D6 in Appendix D show the results of Simulation 3 by displaying the power of each method for the same configurations as considered in Simulation 2. The results show that, over our  $10^4$  configurations of interest, the  $\alpha$ -TOST method is the most powerful, followed by the  $\delta$ -TOST method and then the standard TOST. As expected, the power of all methods goes to one asymptotically as  $\sigma_v \rightarrow 0$  (see Appendix B for details). Figure 3 summarizes the results of Simulation 3 by displaying, for each pair of methods, the histogram of their differences in power for all configurations. The results show that the  $\alpha$ -TOST is overall the most powerful, followed by the  $\delta$ -TOST then by the TOST.

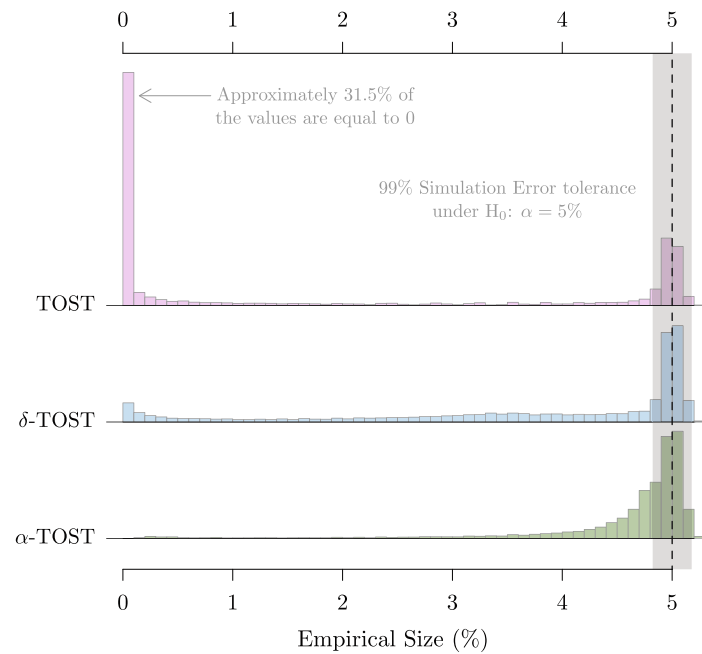


**FIGURE 1** Empirical probability of declaring equivalence (y-axis) as a function of  $\theta$  (x-axis),  $\nu$  (columns), and  $\sigma_\nu$  (rows), for the TOST (pink circles), the  $\alpha$ -TOST (green triangles), and the  $\delta$ -TOST (blue squares). Refer to the settings of Simulation 1 in Table 1 for details. In all configurations considered here,  $\alpha$ -TOST shows a similar or greater power than the TOST and  $\delta$ -TOST while remaining more accurate in terms of empirical size.

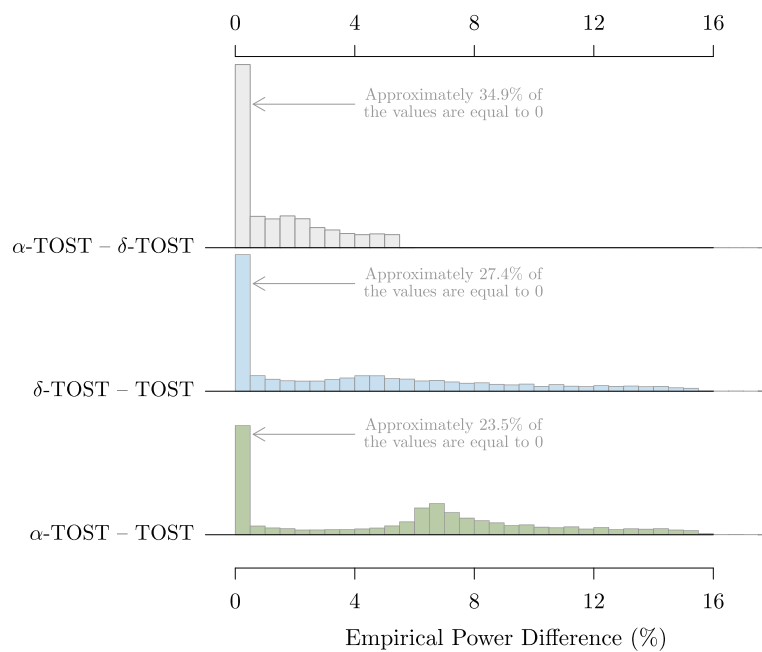
In summary, the simulation studies considered here suggests that a correction of the TOST provides more power and better accuracy in finite samples, with considerably large improvements when  $\sigma_\nu$  is large. Moreover, the  $\alpha$ -TOST appears to provide a better performance than the  $\delta$ -TOST, indicating that an adjustment on the level rather than on the equivalence bounds is preferable to enhance sample properties of equivalence tests. Results of Simulation 4, considering a paired study and additional correction methods, also suggest that adjusting the level of the TOST leads to better operating characteristics over competing methods, including the corrected SABE. More adequate and extensive simulations are needed though to compare these methods in the design-specific context required by regulatory agencies when assessing bioequivalence. Such simulations should consider the different adjustments proposed by regulatory agencies and are left for further research.

Finally, note that the idea of improving the size of the TOST is not new as Cao and Mathew<sup>41</sup> have proposed a correction based on the adjustment of the critical values defined as a non-increasing continuous function of the sample standard deviation to reduce the conservatism of the TOST. More particularly, they defined adjustment constants for specific values of  $\hat{\sigma}_\nu$  and used linear interpolation for the adjacent values. The lower panel of Figure F1 in Appendix F, compares the critical values obtained with the method of Cao and Mathew to the ones obtained by the  $\alpha$ -TOST for different values of  $\hat{\sigma}_\nu$  and  $\nu$ . We can note that, for all values of  $\nu$  considered here and for values of  $\hat{\sigma}_\nu$  above 0.1, the corrected critical values<sup>41</sup> correspond to a piecewise version of the critical values obtained with the  $\alpha$ -TOST when  $\nu$  is large. Therefore, their correction appears to be an approximation of the  $\alpha$ -TOST, evaluated asymptotically, that is, at  $\nu \rightarrow \infty$ .





**FIGURE 2** Histograms of the empirical size (%) of the TOST (first line), the  $\delta$ -TOST (second line), and  $\alpha$ -TOST (third line), computed from the results displayed in Figures D1–D3 in Appendix D, respectively. Overall, the  $\alpha$ -TOST maintains a size of  $\alpha$  for a larger proportion of parameters' in comparison to the other methods.



**FIGURE 3** Histograms of empirical power differences (%) for each pair combinations of the TOST, the  $\delta$ -TOST, and  $\alpha$ -TOST, computed from the results displayed in Figures D1–D3 in Appendix D. Overall, the  $\alpha$ -TOST is the most powerful, followed by the  $\delta$ -TOST then by the TOST.

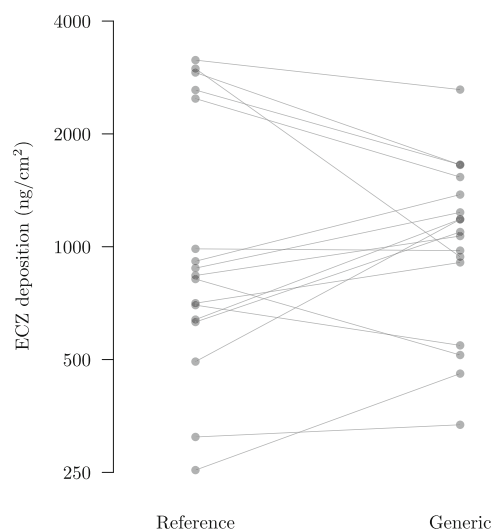
## 4 | EVALUATION OF BIOEQUIVALENCE FOR ECONAZOLE NITRATE DEPOSITION IN PORCINE SKIN

Quartier et al<sup>42</sup> studied the cutaneous bioequivalence of two topical cream products: a Reference Medicinal Product (RMP) and an approved generic containing econazole nitrate (ECZ), an antifungal medication used to treat skin infections. The evaluation of the putative bioequivalence is based on the determination of the cutaneous biodistribution profile of ECZ observed after application of the RMP and the generic product. The dataset we analyse in this section consists in 17 pairs of comparable porcine skin samples on which measurements of ECZ deposition were collected using both creams. Figure 4 presents the data, collected via a simple paired design, in which each pig delivered two skin samples respectively treated with one of the two drugs of interest. Such designs, possibly attractive for studies not involving regulators, are *stricto sensu* incompatible with the use of design-specific SABE-like corrections<sup>26</sup> and therefore interesting to showcase the advantages of our design-agnostic method. In order to assess bioequivalence of both topical treatments, the TOST and  $\alpha$ -TOST procedures, based on a paired  $t$ -test statistic considering ECZ levels on the logarithmic scale, are conducted using  $c = \delta_U = -\delta_L = \log(1.25) \approx 0.223$ . Although the way to define bioequivalence limits for topical products is still being discussed,<sup>43</sup> we believe the chosen limits to be reasonable.<sup>42</sup>

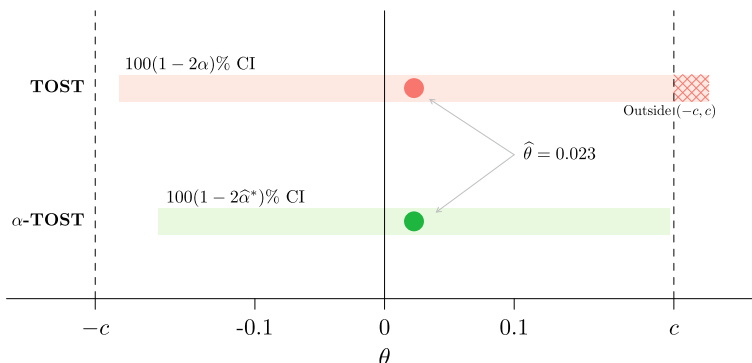
Figure 5 shows the CIs corresponding to both approaches. The  $100(1 - 2\alpha)\%$  TOST confidence interval for the mean of the paired differences in ECZ levels equals  $[-0.204, 0.250]$ , given that  $\hat{\theta} = 0.023$ ,  $\hat{\sigma}_v = 0.134$ ,  $\nu = 16$ , and  $\alpha = 5\%$ . As its upper bound exceeds the upper bioequivalence limit, the classical TOST procedure does not allow us to conclude that the topical products are (on average) equivalent. To reach a size of 5%, the  $\alpha$ -TOST procedure uses in this case a significance level of  $\hat{\alpha}^* = 7.48\%$  leading to a confidence interval of  $[-0.166, 0.211]$ . This CI being strictly embedded within the  $(-c, c)$  bioequivalence limits, the  $\alpha$ -TOST procedure allows to declare bioequivalence, hence illustrating the increase in power induced by the increased significance level considered to reach a size of 5%. Note that in this case, given  $\hat{\sigma}_v^2$  and  $\nu$ , the empirical power of the TOST is zero (regardless of  $\hat{\theta}$ ) as  $t_{0.05,16} \hat{\sigma}_v > c$ , where  $t_{\alpha,\nu}$  denotes the upper quantile  $\alpha$  of a  $t$ -distribution with  $\nu$  degrees of freedom; see Appendix E. Since the  $\alpha$ -TOST guarantees a size of  $\alpha$  (for all sample sizes), the conclusion brought in by the  $\alpha$ -TOST is more trustworthy.

To gain additional insight into the benefits conferred by our approach, we also compare the characteristics and conclusion of the  $\alpha$ -TOST to other available methods in Table 2 as well as their rejection region as a function of  $\hat{\theta}$  and  $\hat{\sigma}_v$  in Figure 6. We considered here the AH-test, the TOST,  $\alpha$ -TOST and  $\delta$ -TOST. The AH-test does not satisfy the IIP, but represents a good proxy for the other tests without this property and is relatively easy to implement. Among the level- $\alpha$  tests, the  $\alpha$ -TOST is the only one leading to bioequivalence declaration.

Figure 6 shows the combinations of values for  $\hat{\theta}$  and  $\hat{\sigma}_v$  leading to bioequivalence declaration in the setting of the porcine skin dataset, that is, with  $c = \log(1.25)$  and  $\nu = 16$ . The rejection regions of the different methods almost perfectly



**FIGURE 4** Econazole nitrate deposition levels (y-axis) measured using the reference and generic creams (x-axis) on 17 pairs of comparable porcine skin samples (lines).

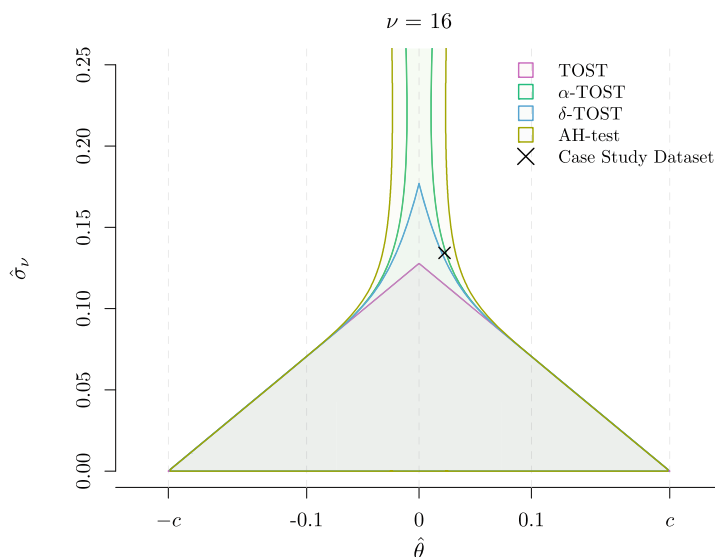


**FIGURE 5**  $100(1 - 2\alpha)\%$  and  $100(1 - 2\hat{\alpha}^*)\%$  confidence intervals of the TOST and  $\alpha$ -TOST procedures for the mean of the paired log differences in ECZ levels obtained with the reference and generic creams with  $\alpha = 5\%$  and  $\hat{\alpha}^* = 7.48\%$ . The dashed vertical lines correspond to the used lower and upper bioequivalence limits with  $c = \log(1.25)$ . Comparison of the CI of each approach to the bioequivalence limits leads to the declaration of bioequivalence for the  $\alpha$ -TOST procedure and not for the classic TOST approach due to its CI upper limit exceeding  $c$  (hatched area).

**TABLE 2** Bioequivalence declaration (yes/no) for the econazole nitrate deposition in porcine skin data using the AH-test, TOST,  $\alpha$ -TOST, and  $\delta$ -TOST.

Method	IIP	Level- $\alpha$	Size- $\alpha$	Bioequivalence declaration
AH-test	No	No	No	Yes
TOST	Yes	Yes	No	No
$\alpha$ -TOST	Yes	Yes*	Yes*	Yes
$\delta$ -TOST	Yes	Yes*	Yes*	No

Note: The estimated parameter values are  $\hat{\sigma}_v = 0.134$ ,  $\nu = 16$ ,  $\hat{\theta} = 0.023$  and  $\alpha = 5\%$ . The columns IIP, Level- $\alpha$  and Size- $\alpha$ , respectively indicate if each method satisfies the IIP, if its size is bounded by  $\alpha$  and if its size is exactly  $\alpha$ . The symbol \* specifies that the property is valid when the standard error  $\sigma_v$  is known.



**FIGURE 6** Bioequivalence test rejection regions as a function of  $\hat{\theta}$  ( $x$ -axis) and  $\hat{\sigma}_v$  ( $y$ -axis) per method considered in Table 2 (colored areas) showing combinations of values for  $\hat{\theta}$  and  $\hat{\sigma}_v$  leading to equivalence declaration in the setting of the porcine skin dataset, that is, with  $c = \log(1.25)$  and  $\nu = 16$ . The rejection regions of the different methods almost perfectly overlap for values of  $\hat{\sigma}_v$  below 0.09 and differ for larger values. Regardless of  $\hat{\sigma}_v$ , the TOST cannot declare bioequivalence for large values of  $\hat{\sigma}_v$  (greater than approximately 0.12 here) and the  $\delta$ -TOST for approximately  $\hat{\sigma}_v > 0.17$ , while the  $\alpha$ -TOST and AH-test can, with the rejection region of the  $\alpha$ -TOST embedded in the too liberal one of the AH-test. The symbol  $\times$  represents the analysed data set in the acceptance/rejection regions where  $\hat{\theta} = 0.023$  and  $\hat{\sigma}_v = 0.134$ .

overlap for values of  $\hat{\sigma}_v$  below 0.09 and differ for larger values. Regardless of  $\hat{\sigma}_v$ , the TOST cannot declare bioequivalence for large values of  $\hat{\sigma}_v$  (greater than approximately 0.12 here) and the  $\delta$ -TOST for approximately  $\hat{\sigma}_v > 0.17$ , while the  $\alpha$ -TOST and AH-test can, with the rejection region of the  $\alpha$ -TOST embedded in the too liberal one of the AH-test. In Figure 6, the symbol  $\times$  shows the values of  $\hat{\theta}$  and  $\hat{\sigma}_v$  obtained in the porcine skin dataset. These coordinates lead to the declaration of (average) bioequivalence of the two topical products with the AH-test and  $\alpha$ -TOST procedures, the later one only being size- $\alpha$ .

## 5 | DISCUSSION

The canonical framework treated in this article is given in (1) and therefore concerns differences that can be assumed to be normally distributed (in finite samples), with a known finite sample distribution for  $\hat{\sigma}_v$ . This framework covers a quite large spectrum of data settings, such as the standard two-period crossover experimental design,<sup>44</sup> and could be extended to include covariates to possibly reduce residual variance. Extensions to non-linear cases, such as for example binary outcomes,<sup>45-49</sup> would follow the same logic, but would require a specific treatment due to the nature of the responses and to the use of link functions. Such extensions also deserve some attention but are left for further research.

For sample size calculations, we could, in principle, proceed with the  $\alpha$ -TOST, for given values of  $c$ ,  $\theta$ , and  $\sigma_v$ . However, when considering high levels of power, the correction is negligible and we have  $\alpha^* \approx \alpha$  as shown in Section 3, so that the sample size can be computed using the TOST, as implemented in standard packages. The  $\alpha$ -TOST approach would then be used to assess equivalence and show its benefits when the observed value of  $\sigma_v$  is unexpectedly large compared to the one considered in the sample size calculation either due to (lack of) chance or to an underestimated value obtained from a prior experiment.

## ACKNOWLEDGEMENTS

S. Guerrier is supported by the SNSF Grants #176843 and # 211007 as well as by the Innosuisse Grants #37308.1 IP-ENG and #53622.1 IP-ENG. Y. Boulaguiem and M.-P. Victoria-Feser are partially supported by the SNSF Grant #182684. D.-L. Couturier is partially supported by the Cancer Research UK Grant C9545/A29580 and UK Medical Research Council Grant MC\_UU\_00002/14. Y. N. Kalia thanks the University of Geneva for teaching assistantships for J. Quartier and M. Lapteva, for providing financial support for the purchase of the Waters Xevo TQ-MS detector, and also thanks the Fondation Ernst and Lucie Schmidheiny and the Société Académique de Genève for providing equipment grants.

## CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

## DATA AVAILABILITY STATEMENT

The econazole nitrate deposition data as well as an implementation of the method proposed in this article are available in the cTOST R package available on CRAN.

## ORCID

Younes Boulaguiem  <https://orcid.org/0000-0003-0795-0714>

## REFERENCES

1. Metzler C. Bioavailability—a problem in equivalence. *J Pharm Sci.* 1974;30:309-317.
2. Westlake T. Symmetrical confidence intervals for bioequivalence trials. *Biometrics.* 1976;32:741-744.
3. Pallmann P, Jaki T. Simultaneous confidence regions for multivariate bioequivalence. *Stat Med.* 2017;36(29):4585-4603.
4. Senn S. *Statistical Issues in Drug Development.* 3rd ed. Hoboken, NJ: Wiley; 2021.
5. Patterson S, Jones B. *Bioequivalence and Statistics in Clinical Pharmacology.* Boca Raton, FL: Chapman & Hall/CRC; 2006.
6. Lakens D. Equivalence tests: a practical primer for *t*-tests, correlations, and meta-analyses. *Soc Psychol Personal Sci.* 2017;8:355-362.
7. Sureshkumar H, Xu R, Erukulla N, Wadhwa A, Zhao L. “Snap on” or not? A validation on the measurement tool in a virtual reality application. *J Digit Imaging.* 2022;35(3):692-703.
8. O'Brien MW, Kimmerly DS. Is “not different” enough to conclude similar cardiovascular responses across sexes? *Am J Physiol Heart Circ Physiol.* 2022;322:H355-H358.
9. Wehrle FM, Bartal T, Adams M, et al. Similarities and differences in the neurodevelopmental outcome of children with congenital heart disease and children born very preterm at school entry. *J Pediatr.* 2022;250:29-37.e1.

10. Sansone P, Giaccari LG, Aurilio C, et al. Comparative efficacy of Tapentadol versus Tapentadol plus duloxetine in patients with chemotherapy-induced peripheral neuropathy. *Cancer*. 2022;14:4002.
11. Branscheidt M, Ejaz N, Xu J, et al. No evidence for motor-recovery-related cortical connectivity changes after stroke using resting-state fMRI. *J Neurophysiol*. 2022;127:637-650.
12. Feri F, Giannetti C, Guarnieri P. Risk-taking for others: an experiment on the role of moral discussion. *J Behav Exp Financ*. 2023;37:100735. doi:10.1016/j.jbef.2022.100735
13. Aggarwal M, Allen J, Coppock A, et al. A 2 million-person, campaign-wide field experiment shows how digital advertising affects voter turnout. *Nat Hum Behav*. 2023;7:332-341.
14. Meyners M. Equivalence tests—a review. *Food Qual Prefer*. 2012;26:231-245.
15. Lakens D, Scheel AM, Isager PM. Equivalence testing for psychological research: a tutorial. *Adv Methods Pract Psychol Sci*. 2018;1:259-269.
16. Mazzolari R, Porcelli S, Bishop DJ, Lakens D. Myths and methodologies: the use of equivalence and non-inferiority tests for interventional studies in exercise physiology and sport science. *Exp Physiol*. 2022;107:201-212.
17. Wang K, Li Y, Chen B, et al. In vitro predictive dissolution test should be developed and recommended as a bioequivalence standard for the immediate-release solid oral dosage forms of the highly variable mycophenolate Mofetil. *Mol Pharm*. 2022;19:2048-2060.
18. Schuirmann DJ. A comparison of the two one-sided tests procedure and the power approach for assessing the equivalence of average bioavailability. *J Pharmacokinet Biopharm*. 1987;15(6):657-680.
19. Berger RL. Multiparameter hypothesis testing and acceptance sampling. *Dent Tech*. 1982;24(4):295-300.
20. Muñoz J, Alcaide D, Ocaña J. Consumer's risk in the EMA and FDA regulatory approaches for bioequivalence in highly variable drugs. *Stat Med*. 2016;35:1933-1943.
21. Hsu JC, Hwang JTG, Liu HK, Ruberg SJ. Confidence intervals associated with tests for bioequivalence. *Biometrika*. 1994;81:103-114.
22. Berger RL, Hsu JC. Bioequivalence trials, intersection-union tests and equivalence confidence sets. *Stat Sci*. 1996;11:283-319.
23. Anderson S, Hauck WW. A new procedure for testing equivalence in comparative bioavailability and other clinical trials. *Commun Stat Theory Methods*. 1983;12:2663-2692.
24. Brown HJT, Munk A. An unbiased test for the bioequivalence problem. *Ann Stat*. 1997;25:2345-2367.
25. Liu JP, Chow SC. Bioequivalence trials, intersection-union tests and equivalence confidence set: comment. *Stat Sci*. 1996;11:306-312.
26. Schütz H, Labes D, Wolfsegger MJ. Critical remarks on reference-scaled average bioequivalence. *J Pharm Pharm Sci*. 2022;25:285-296.
27. Guideline on the investigation of bioequivalence. European Medicines Agency; 2010. [https://www.ema.europa.eu/en/documents/scientific-guideline/guideline-investigation-bioequivalence-rev1\\_en.pdf](https://www.ema.europa.eu/en/documents/scientific-guideline/guideline-investigation-bioequivalence-rev1_en.pdf) Accessed July 10, 2023.
28. Labes D, Schütz H. Inflation of type I error in the evaluation of scaled average bioequivalence, and a method for its control. *Pharm Res*. 2016;33:2805-2814.
29. Tothfalusi L, Endrenyi L, Arieta AG. Evaluation of bioequivalence for highly variable drugs with scaled average bioequivalence. *Clin Pharmacokinet*. 2009;48:725-743.
30. Davit BM, Chen ML, Conner DP, et al. Implementation of a reference-scaled average bioequivalence approach for highly variable generic drug products by the US Food and Drug Administration. *AAPS J*. 2012;14:915-924.
31. Wonnemann M, Frömke C, Koch A. Inflation of the type I error: investigations on regulatory recommendations for bioequivalence of highly variable drugs. *Pharm Res*. 2015;32:135-143.
32. Endrenyi L, Tothfalusi L. Bioequivalence for highly variable drugs: regulatory agreements, disagreements, and harmonization. *J Pharmacokinet Pharmacodyn*. 2019;46:117-126.
33. Molins E, Labes D, Schütz H, Cobo E, Ocaña J. An iterative method to protect the type I error rate in bioequivalence studies under two-stage adaptive 2x2 crossover designs. *Biom J*. 2021;63:122-133.
34. Tothfalusi L, Endrenyi L. An exact procedure for the evaluation of reference-scaled average bioequivalence. *AAPS J*. 2016;18:476-489.
35. Ocaña J, Muñoz J. Controlling type I error in the reference-scaled bioequivalence evaluation of highly variable drugs. *Pharm Stat*. 2019;18:583-599.
36. Deng Y, Zhou XH. Methods to control the empirical type I error rate in average bioequivalence tests for highly variable drugs. *Stat Methods Med Res*. 2020;29:1650-1667.
37. Phillips K. Power of the two one-sided tests procedure in bioequivalence. *J Pharmacokinet Biopharm*. 1990;18:137-144.
38. Owen DB. A special case of a bivariate non-central *t*-distribution. *Biometrika*. 1965;52:437-446.
39. Lehmann EL. *Testing Statistical Hypothesis*. 2nd ed. New York: Wiley; 1986.
40. Palmes C, Bluhmki T, Funke B, Bluhmki E. Asymptotic properties of the two one-sided *t*-tests—new insights and the Schuirmann-constant. *Int J Biostat*. 2022;18:19-38.
41. Cao L, Mathew T. A simple numerical approach towards improving the two one-sided test for average bioequivalence. *Biom J*. 2008;50:205-211.
42. Quartier J, Capony N, Lapteva M, Kalia YN. Cutaneous biodistribution: a high-resolution methodology to assess bioequivalence in topical skin delivery. *Pharmaceutics*. 2019;11:484.
43. Committee for Medicinal Products for Human Use. Draft guideline on quality and equivalence of topical products; 2018. [https://www.ema.europa.eu/en/documents/scientific-guideline/draft-guideline-quality-equivalence-topical-products\\_en.pdf](https://www.ema.europa.eu/en/documents/scientific-guideline/draft-guideline-quality-equivalence-topical-products_en.pdf) Accessed July 10, 2023.
44. Jones B, Kenward MG. *Design and Analysis of Cross-Over Trials*. Boca Raton, FL: CRC Press; 2014.
45. Dunnett CW, Gent M. Significance testing to establish equivalence between treatments, with special reference to data in the form of 2x2 tables. *Biometrics*. 1977;33:593-602.

46. Tu D. On the use of the ratio or the odds ratio of cure rates in establishing therapeutic equivalence of non-systemic drugs with binary clinical endpoints. *J Biopharm Stat.* 1998;8:263-282.
47. Schouten H, Kester A. A simple analysis of a simple crossover trial with a dichotomous outcome measure. *Stat Med.* 2010;29:193-198.
48. Lui KJ, Chang KC. Test non-inferiority (and equivalence) based on the odds ratio under a simple crossover trial. *Stat Med.* 2011;30:1230-1242.
49. Ostrovski V. Testing equivalence to binary generalized linear models with application to logistic regression. *Stat Probab Lett.* 2022;191:109658.
50. Rudin W. *Principles of Mathematical Analysis.* 2nd ed. New York: McGraw Hill; 1976.
51. van der Vaart AW. *Asymptotic Statistics.* Cambridge, UK: Cambridge University Press; 2000.
52. Federer H. *Geometric Measure Theory.* New York: Springer; 2014.

**How to cite this article:** Boulaguiem Y, Quartier J, Lapteva M, et al. Finite sample corrections for average equivalence testing. *Statistics in Medicine.* 2024;43(5):833-854. doi: 10.1002/sim.9993

## APPENDIX A. EXISTENCE OF $\alpha^*$

In this section, we state the conditions for  $\alpha^*$ , defined in (7), to be a singleton. Fixing  $\alpha$ ,  $c$ ,  $\sigma_v$ , and  $\nu$ , we simplify our notation so that  $\omega(\gamma) := \omega(\gamma, c, \sigma_v, \nu)$  and let

$$\mathcal{A} := \left\{ x \in [\alpha, 0.5] \mid \omega(x) > 0 \right\}.$$

The function  $\omega(\gamma)$ , defined in (5), is continuously differentiable and strictly increasing in  $\gamma$  in  $\mathcal{A}$ . From (6), we have that  $\alpha \geq \omega(\alpha)$ . Thus, it is sufficient to show that  $\alpha < \alpha_{\max} := \lim_{\alpha \rightarrow 0.5^-} \omega(\alpha)$ , where  $\alpha \rightarrow 0.5^-$  denotes the limit from below 0.5, to ensure that  $\alpha^*$  is a singleton. Let  $T_{\nu, \delta}$  denote a random variable following a non-central  $t$ -distribution with  $\nu$  degrees of freedom and noncentrality parameter  $\delta = 2c/\sigma_v$ . Then, we have

$$\begin{aligned} \alpha_{\max} &= \lim_{\gamma \rightarrow 0.5^-} \omega(\gamma) = \lim_{\gamma \rightarrow 0.5^-} \left\{ Q_{\nu} \left( -t_{\gamma, \nu}, 0, \frac{c\sqrt{\nu}}{\sigma_v t_{\gamma, \nu}} \right) - Q_{\nu} \left( t_{\gamma, \nu}, \delta, \frac{c\sqrt{\nu}}{\sigma_v t_{\gamma, \nu}} \right) \right\} \\ &= \Pr(T_{\nu, 0} \leq 0) - \Pr(T_{\nu, \delta} \leq 0) = 0.5 - \Pr(T_{\nu, \delta} \leq 0) = 0.5 - \Phi(-\delta) = \Phi(\delta) - 0.5. \end{aligned}$$

Thus, the condition  $\alpha < \alpha_{\max}$  can be expressed as follows

$$\Phi(\delta) - 0.5 > \alpha \iff \frac{2c}{\sigma_v} > \Phi^{-1}(\alpha + 0.5) \iff \sigma_v < \frac{2c}{\Phi^{-1}(\alpha + 0.5)}.$$

Therefore, the condition  $\sigma_v < \frac{2c}{\Phi^{-1}(\alpha + 0.5)}$  implies that  $\alpha < \alpha_{\max}$  and consequently that  $\alpha^*$  is a singleton. The existence of  $\hat{\alpha}^*$  follows the same argument but replacing  $\sigma_v$  by  $\hat{\sigma}_v$  in our condition.

## APPENDIX B. ASYMPTOTIC PROPERTIES

In this section, we study the convergence rates of  $\hat{\theta}$ ,  $\hat{\sigma}_v^2$ , and  $\hat{\alpha}^*$ . In particular, we show that

$$\hat{\theta} = \theta + \mathcal{O}_p(\nu^{-1/2}), \tag{B1}$$

$$\hat{\sigma}_v = \sigma_v + \mathcal{O}_p(\nu^{-1}), \tag{B2}$$

$$\hat{\alpha}^* = \alpha^* + o_p(\nu^{-1}). \tag{B3}$$

These results are based on the following standard regularity conditions. First, there exists a positive constant  $\sigma^2$  such that  $\sigma^2 := \lim_{\nu \rightarrow \infty} \nu \sigma_v^2$ . Second, the sequences

$$\left\{ \frac{\partial}{\partial \alpha} \omega(\alpha, c, \sigma_v, \nu) \right\}_{\nu \in \mathcal{M}} \quad \text{and} \quad \left\{ \frac{\partial}{\partial \sigma_v} \omega(\alpha, c, \sigma_v, \nu) \right\}_{\nu \in \mathcal{M}},$$

where  $\mathcal{M} \subseteq \mathbb{R}$ , converge uniformly in  $\alpha$  and  $\sigma_v$ . This second condition implies by Theorem 7.17 of Rudin<sup>50</sup> that,

$$\lim_{v \rightarrow \infty} \frac{\partial}{\partial \alpha} \omega(\alpha, c, \sigma_v, v) = \frac{\partial}{\partial \alpha} \lim_{v \rightarrow \infty} \omega(\alpha, c, \sigma_v, v) = \frac{\partial}{\partial \alpha} \alpha = 1.$$

Similarly,

$$\lim_{v \rightarrow \infty} \frac{\partial}{\partial \sigma_v} \omega(\alpha, c, \sigma_v, v) = 0.$$

Using these regularity conditions, we start by proving (B1). From (1), we have that  $\hat{\theta} \sim \mathcal{N}(\theta, \sigma_v^2)$  and, thus, by Markov's inequality, for any  $M > 0$ , we obtain

$$\Pr\left(\sqrt{v}|\hat{\theta} - \theta| \geq M\right) \leq \frac{v}{M^2} \sigma_v^2 = \frac{\sigma^2}{M^2} + o(1).$$

Therefore, we have  $\hat{\theta} = \theta + \mathcal{O}_p(v^{-1/2})$ , which verifies (B1).

Next, we study the convergence rate of  $\hat{\sigma}_v$ . Let  $Y$  denote a random variable following a  $\chi$  distribution with  $v$  degrees of freedom, that is,  $Y \sim \chi_v$ . We have

$$\mathbb{E}[Y] = \frac{\sqrt{2}\Gamma\{(v+1)/2\}}{\Gamma(v/2)} = \sqrt{v} \left\{ 1 - \frac{1}{4(v+1)} + \mathcal{O}(v^{-2}) \right\},$$

where  $\Gamma(\cdot)$  is the gamma function, and where the second equality can be obtained using Stirling's approximation for the Gamma function. Moreover, we have

$$\text{var}(Y) = v - \mathbb{E}^2[Y] = v - v \left\{ 1 - \frac{1}{4(v+1)} + \mathcal{O}(v^{-2}) \right\}^2 = v - v \left\{ 1 - \frac{1}{2(v+1)} + \mathcal{O}(v^{-2}) \right\} = \frac{v}{2(v+1)} + \mathcal{O}(v^{-1}).$$

From (1), we have  $v\hat{\sigma}_v^2 \sigma_v^{-2} \sim \chi_v^2$  using Markov's inequality, for any  $M > 0$ , we obtain

$$\begin{aligned} \Pr(v|\hat{\sigma}_v - \sigma_v| \geq M) &\leq \frac{v^2 \mathbb{E}\left[(\hat{\sigma}_v - \sigma_v)^2\right]}{M^2} = \frac{v^2}{M^2} \left\{ \text{var}(\hat{\sigma}_v) + (\mathbb{E}[\hat{\sigma}_v] - \sigma_v)^2 \right\} \\ &= \frac{v^2}{M^2} \left\{ \frac{\sigma_v^2}{v} \text{var}(Y) + \left( \sqrt{\frac{\sigma_v^2}{v}} \mathbb{E}[Y] - \sigma_v \right)^2 \right\} \\ &= \frac{v^2}{M^2} \left( \frac{\sigma_v^2}{v} \left\{ \frac{v}{2(v+1)} + \mathcal{O}(v^{-1}) \right\} + \left[ \sigma_v \left\{ 1 - \frac{1}{4(v+1)} + \mathcal{O}(v^{-2}) \right\} - \sigma_v \right]^2 \right) \\ &= \frac{v^2}{M^2} \left\{ \frac{\sigma_v^2}{2(v+1)} + \frac{\sigma_v^2}{16(v+1)^2} + \mathcal{O}(v^{-3}) \right\} = \frac{v^2 \sigma_v^2}{2(v+1)M^2} + \mathcal{O}(v^{-1}) \\ &= \frac{v\{\sigma^2 + o(1)\}}{2(v+1)M^2} + \mathcal{O}(v^{-1}) = \frac{\sigma^2}{2M^2} + o(1). \end{aligned}$$

Thus, we have  $\hat{\sigma}_v = \sigma_v + \mathcal{O}_p(v^{-1})$ , which verifies (B2).

Finally, we consider the convergence rate of  $\hat{\alpha}^*$ . Using (B2) and the continuity in  $\sigma_v$  of  $\omega(\alpha, c, \sigma_v, v)$ , we have by the continuous mapping theorem that  $\omega(\alpha, c, \hat{\sigma}_v, v) = \omega(\alpha, c, \sigma_v, v) + o_p(1)$ , which by Lemma 5.10 of Van der Vaart<sup>51</sup> implies that

$$\hat{\alpha}^* = \alpha^* + o_p(1). \quad (\text{B4})$$

Next, there exists a point  $(\tilde{\alpha}, \tilde{\sigma}_v)$  on the line segment from  $(\alpha^*, \sigma_v)$  to  $(\hat{\alpha}^*, \hat{\sigma}_v)$  such that

$$\omega(\hat{\alpha}^*, c, \hat{\sigma}_v, v) - \omega(\alpha^*, c, \sigma_v, v) = \frac{\partial}{\partial \alpha} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, v) (\hat{\alpha}^* - \alpha^*) + \frac{\partial}{\partial \sigma_v} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, v) (\hat{\sigma}_v - \sigma_v),$$

where

$$\frac{\partial}{\partial \alpha} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu) = \frac{\partial}{\partial x} \omega(x, c, \tilde{\sigma}_v, \nu) \Big|_{x=\tilde{\alpha}} \quad \text{and} \quad \frac{\partial}{\partial \sigma_v} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu) = \frac{\partial}{\partial y} \omega(\tilde{\alpha}, c, y, \nu) \Big|_{y=\tilde{\sigma}_v}.$$

From (7) and (8), we have  $\omega(\alpha^*, c, \sigma_v, \nu) = \alpha$  and  $\omega(\hat{\alpha}^*, c, \hat{\sigma}_v, \nu) = \alpha$ , implying that

$$\frac{\partial}{\partial \alpha} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu) (\hat{\alpha}^* - \alpha^*) + \frac{\partial}{\partial \sigma_v} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu) (\hat{\sigma}_v - \sigma_v) = 0.$$

Since  $\frac{\partial}{\partial \alpha} \omega(\gamma, c, \kappa, \nu)$  converges uniformly, for all  $\varepsilon > 0$ , there exists  $N_\nu > 0$  such that for all  $\gamma \in [\alpha, 0.5)$ , for all  $\kappa \in \mathbb{R}_+$  and for all  $\nu \geq N_\nu$ , we have that  $|\frac{\partial}{\partial \alpha} \omega(\gamma, c, \kappa, \nu) - 1| \leq \varepsilon$ . Thus, we have

$$\lim_{\nu \rightarrow \infty} \frac{\partial}{\partial \alpha} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu) = 1 \quad \text{and} \quad \lim_{\nu \rightarrow \infty} \frac{\partial}{\partial \sigma_v} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu) = 0.$$

Consequently, we obtain

$$\lim_{\nu \rightarrow \infty} \frac{\frac{\partial}{\partial \sigma_v} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu)}{\frac{\partial}{\partial \alpha} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu)} = 0.$$

For sufficiently large  $\nu$ , we have

$$|\hat{\alpha}^* - \alpha^*| = \left| \frac{\frac{\partial}{\partial \sigma_v} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu)}{\frac{\partial}{\partial \alpha} \omega(\tilde{\alpha}, c, \tilde{\sigma}_v, \nu)} \right| |\hat{\sigma}_v - \sigma_v| = o(|\hat{\sigma}_v - \sigma_v|) = o_p(\nu^{-1}),$$

since, from (B2), we have  $|\hat{\sigma}_v - \sigma_v| = \mathcal{O}_p(\nu^{-1})$ . Therefore, we obtain  $\hat{\alpha}^* = \alpha^* + o_p(\nu^{-1})$ , which verifies (B3) and concludes the proof.

### APPENDIX C. CONVERGENCE RATE OF THE ITERATIVE APPROACH FOR $\alpha^*$

Using the notation of Appendix A and for  $\gamma \in \mathcal{A}$ , we have that  $\omega(\gamma)$  is continuously differentiable and such that  $0 < \dot{\omega}(\gamma) < 2$ , where

$$\dot{\omega}(\gamma) := \frac{\partial}{\partial x} \omega(x) \Big|_{x=\gamma}.$$

Next, we define

$$T(\gamma) := \alpha + \gamma - \omega(\gamma).$$

For all  $\alpha_1, \alpha_2 \in \mathcal{A}$ , we have the mean value theorem stating that

$$T(\alpha_1) - T(\alpha_2) = \alpha_1 - \alpha_2 - \omega(\alpha_1) + \omega(\alpha_2) = \alpha_1 - \alpha_2 + \dot{\omega}(\alpha_3)(\alpha_2 - \alpha_1),$$

where  $\alpha_3 = \tau \alpha_1 + (1 - \tau) \alpha_2$  for some  $\tau \in [0, 1]$ . Thus, we obtain

$$\left| T(\alpha_1) - T(\alpha_2) \right| = \left| \{1 - \dot{\omega}(\alpha_3)\} (\alpha_1 - \alpha_2) \right| = \left| 1 - \dot{\omega}(\alpha_3) \right| \left| \alpha_1 - \alpha_2 \right| < \left| \alpha_1 - \alpha_2 \right|.$$

Then, using Kirszbraun theorem,<sup>52</sup> we can extend the function  $T(\gamma)$  with respect to  $\gamma \in \mathcal{A}$  to a contraction map from  $\mathbb{R}$  to  $\mathbb{R}$ . Thus, Banach fixed point theorem ensures that  $T(\alpha^{*(k)})$  converges as  $k \rightarrow \infty$ . We then define the limit of the sequence  $\{\alpha^{*(k+1)}\}_{k \in \mathbb{N}}$  as  $\alpha^*$ , which is the unique fixed point of the function  $T(\gamma)$ . Indeed, we have

$$\alpha^* = T(\alpha^*) = \alpha + \alpha^* - \omega(\alpha^*).$$



By rearranging terms, we have

$$\alpha^* = \underset{\gamma \in \mathcal{A}}{\operatorname{argzero}} \omega(\gamma) - \alpha = \underset{\gamma \in [\alpha, 0.5]}{\operatorname{argzero}} \omega(\gamma) - \alpha,$$

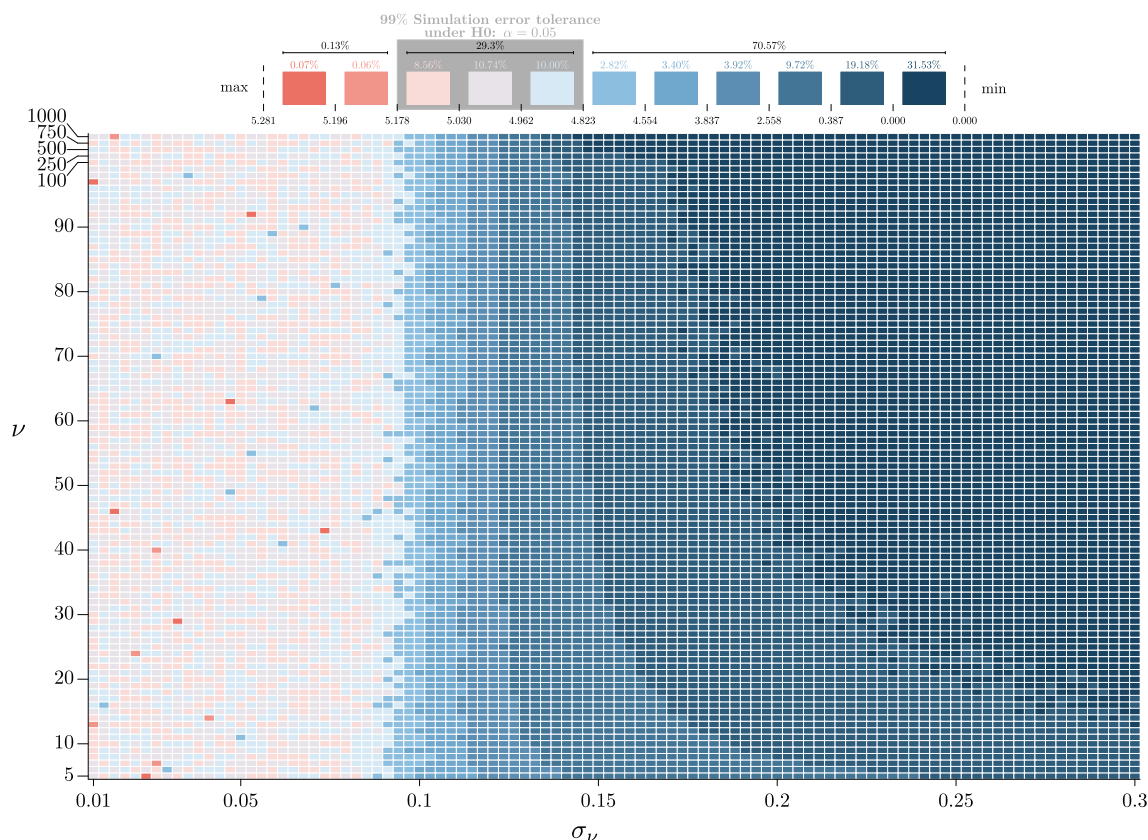
concluding the convergence of the sequence  $\{\alpha^{*(k+1)}\}_{k \in \mathbb{N}}$ . As a result, there exists some  $0 < \epsilon < 1$  such that for  $k \in \mathbb{N}$  we have

$$|\alpha^{*(k+1)} - \alpha^*| < \epsilon^k |\alpha^* - \alpha| < \frac{1}{2} \exp(-bk),$$

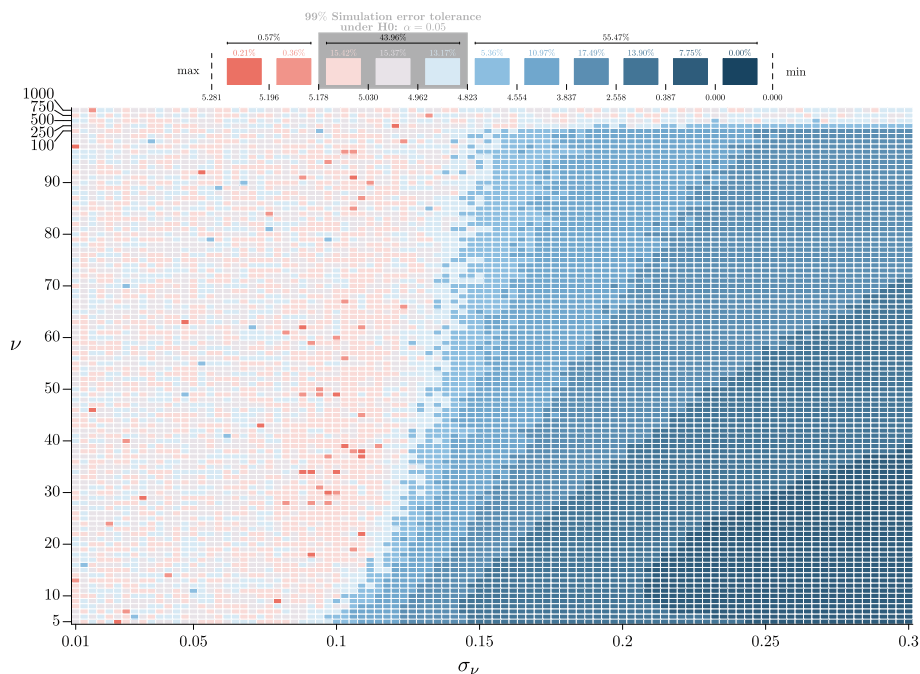
for some positive constant  $b$ .

#### APPENDIX D. EMPIRICAL SIZE AND POWER COMPARISONS FOR THE TOST, $\alpha$ -TOST AND $\Delta$ -TOST

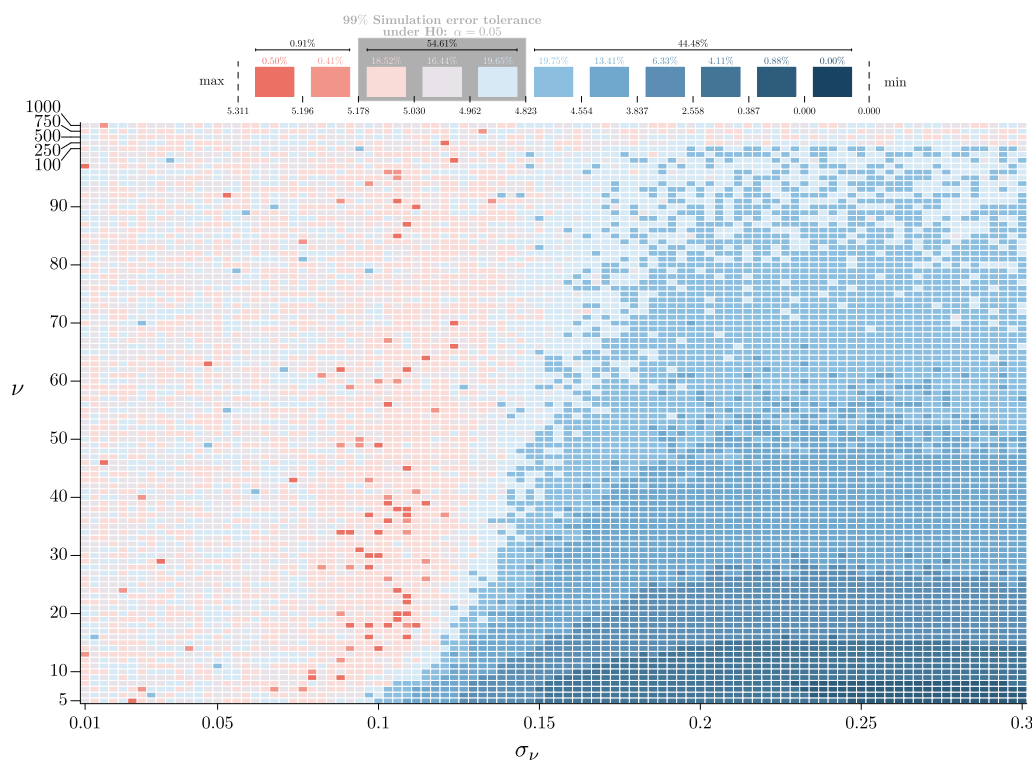
In this section, we perform an extensive simulation study, for the evaluation of the empirical size and power of the  $\alpha$ -TOST, compared to the TOST and  $\delta$ -TOST, by varying the values of  $\nu$  and  $\sigma_\nu$  over a large grid. The power (ie, when  $\theta = 0$ ) and the size (ie, when  $\theta = c$ ) are computed by replacing both  $\sigma_\nu$  and  $\theta$  by realizations of their corresponding random variables in (1), that is, reproducing the case of parameter estimation. The simulation settings we consider are given in Simulations 2 and 3 of Table 1 for the size and power respectively.



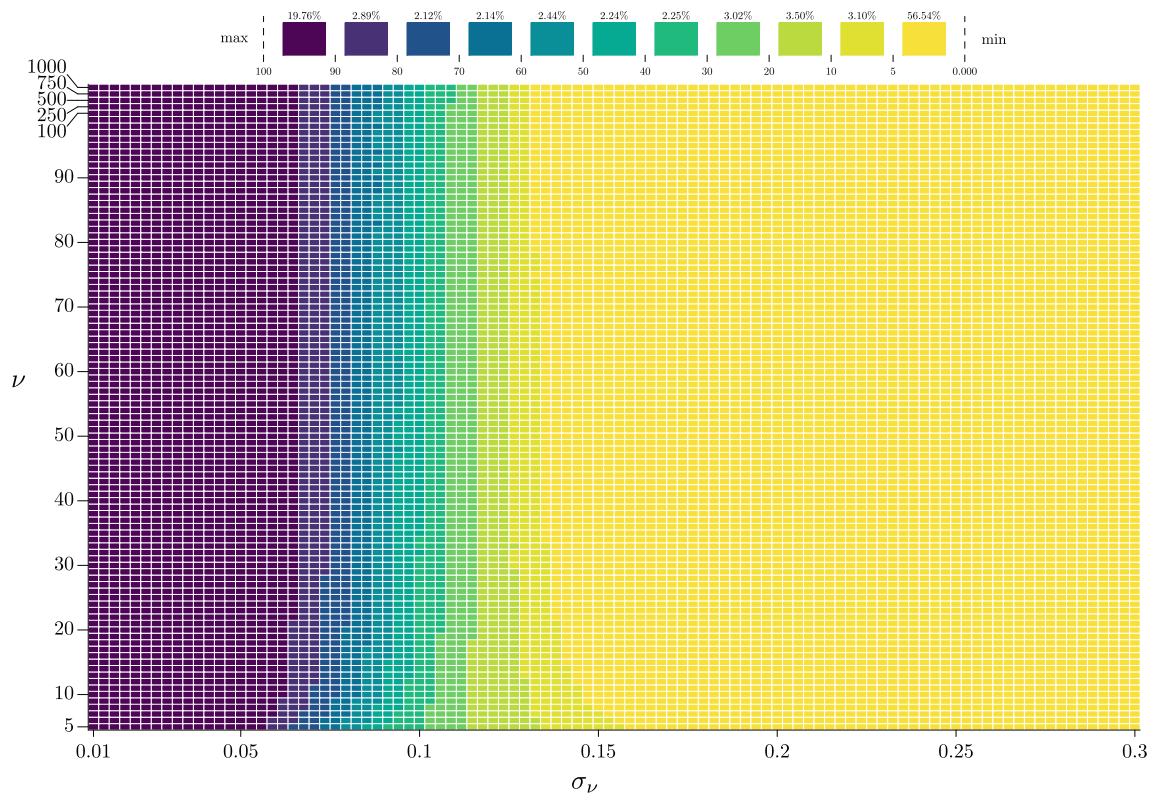
**FIGURE D1** Heatmap representing the empirical size in % (color gradient) for the TOST, computed using the setting of Simulation 2 in Table 1, as a function of  $\sigma_\nu$  (x-axis) and  $\nu$  (y-axis). The lighter colors highlighted in the top legend correspond to the  $\alpha = 5\%$  nominal level, up to a simulation error assessed by a two-sided binomial exact test at the 1% level performed on the results obtained on the  $10^5$  Monte Carlo samples per setting.



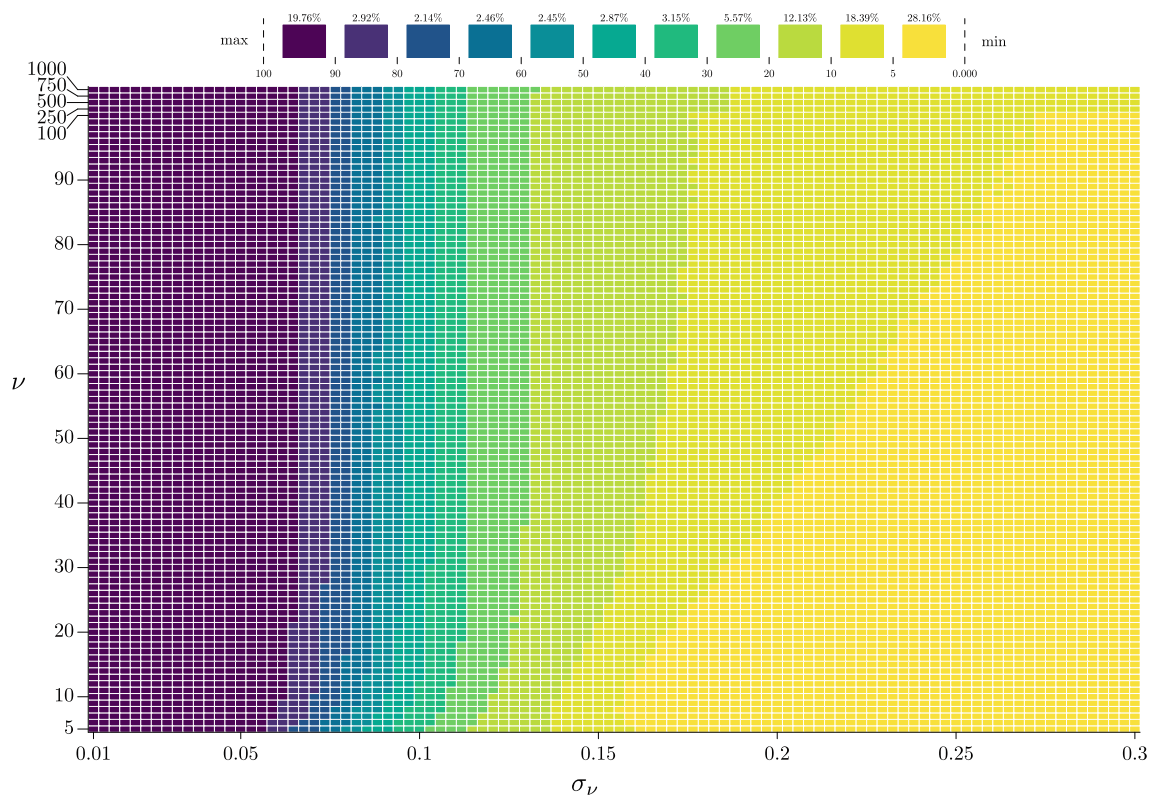
**FIGURE D2** Heatmap representing the empirical size in % (color gradient) for the  $\delta$ -TOST, computed using the setting of Simulation 2 in Table 1, as a function of  $\sigma_\nu$  ( $x$ -axis) and  $\nu$  ( $y$ -axis). The lighter colors highlighted in the top legend correspond to the  $\alpha = 5\%$  nominal level, up to a simulation error assessed by a two-sided binomial exact test at the 1% level performed on the results obtained on the  $10^5$  Monte Carlo samples per setting.



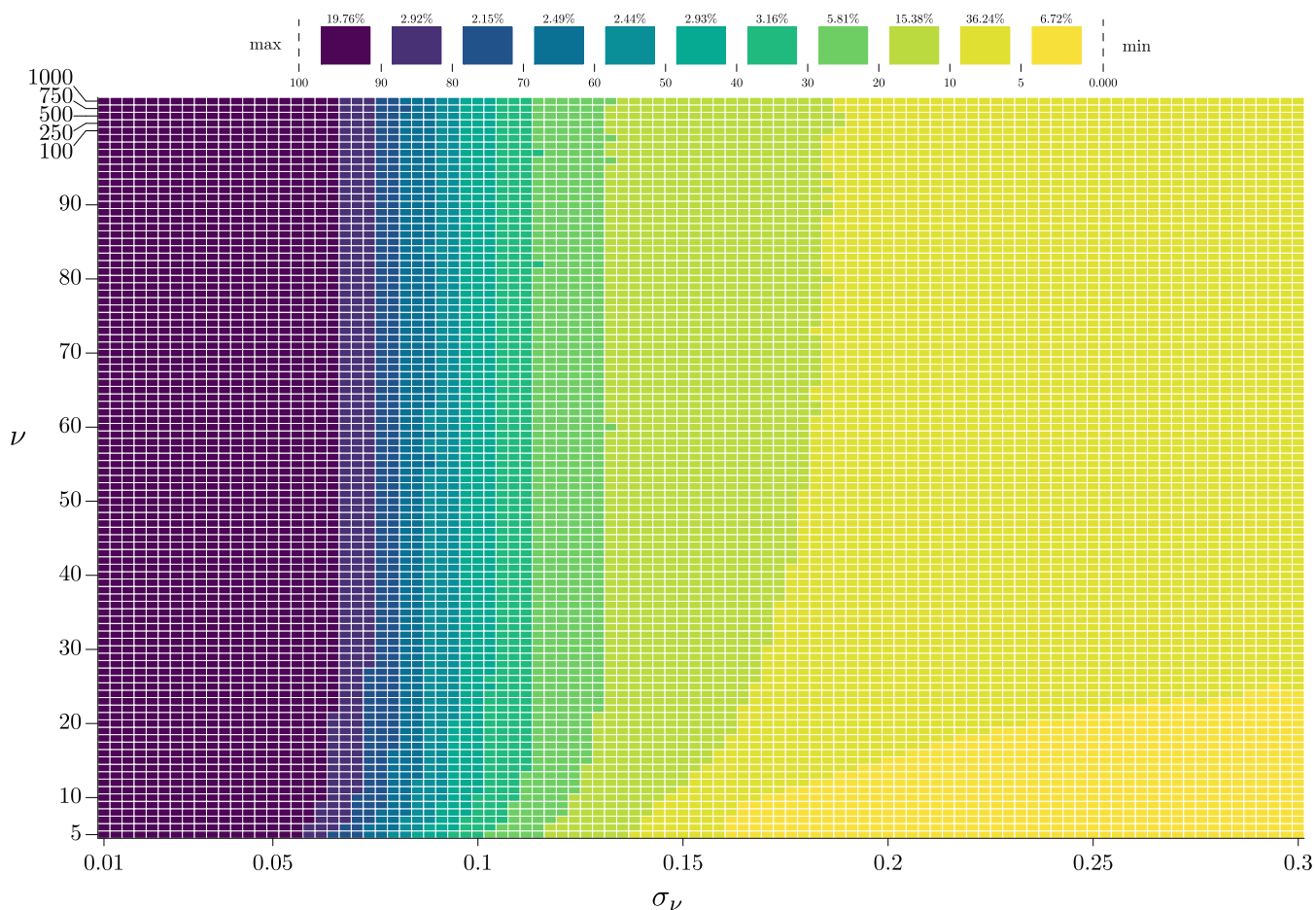
**FIGURE D3** Heatmap representing the empirical size in % (color gradient) for the  $\alpha$ -TOST, computed using the setting of Simulation 2 in Table 1, as a function of  $\sigma_\nu$  ( $x$ -axis) and  $\nu$  ( $y$ -axis). The lighter colors highlighted in the top legend correspond to the  $\alpha = 5\%$  nominal level, up to a simulation error assessed by a two-sided binomial exact test at the 1% level performed on the results obtained on the  $10^5$  Monte Carlo samples per setting.



**FIGURE D4** Heatmap representing the empirical power in % (color gradient) for the TOST, computed using the setting of Simulation 3 in Table 1, as a function of  $\sigma_\nu$  (x-axis) and  $\nu$  (y-axis).



**FIGURE D5** Heatmap representing the empirical power in % (color gradient) for the  $\delta$ -TOST, computed using the setting of Simulation 3 in Table 1, as a function of  $\sigma_\nu$  (x-axis) and  $\nu$  (y-axis).



**FIGURE D6** Heatmap representing the empirical power in % (color gradient) for the  $\alpha$ -TOST, computed using the setting of Simulation 3 in Table 1, as a function of  $\sigma_\nu$  (x-axis) and  $\nu$  (y-axis).

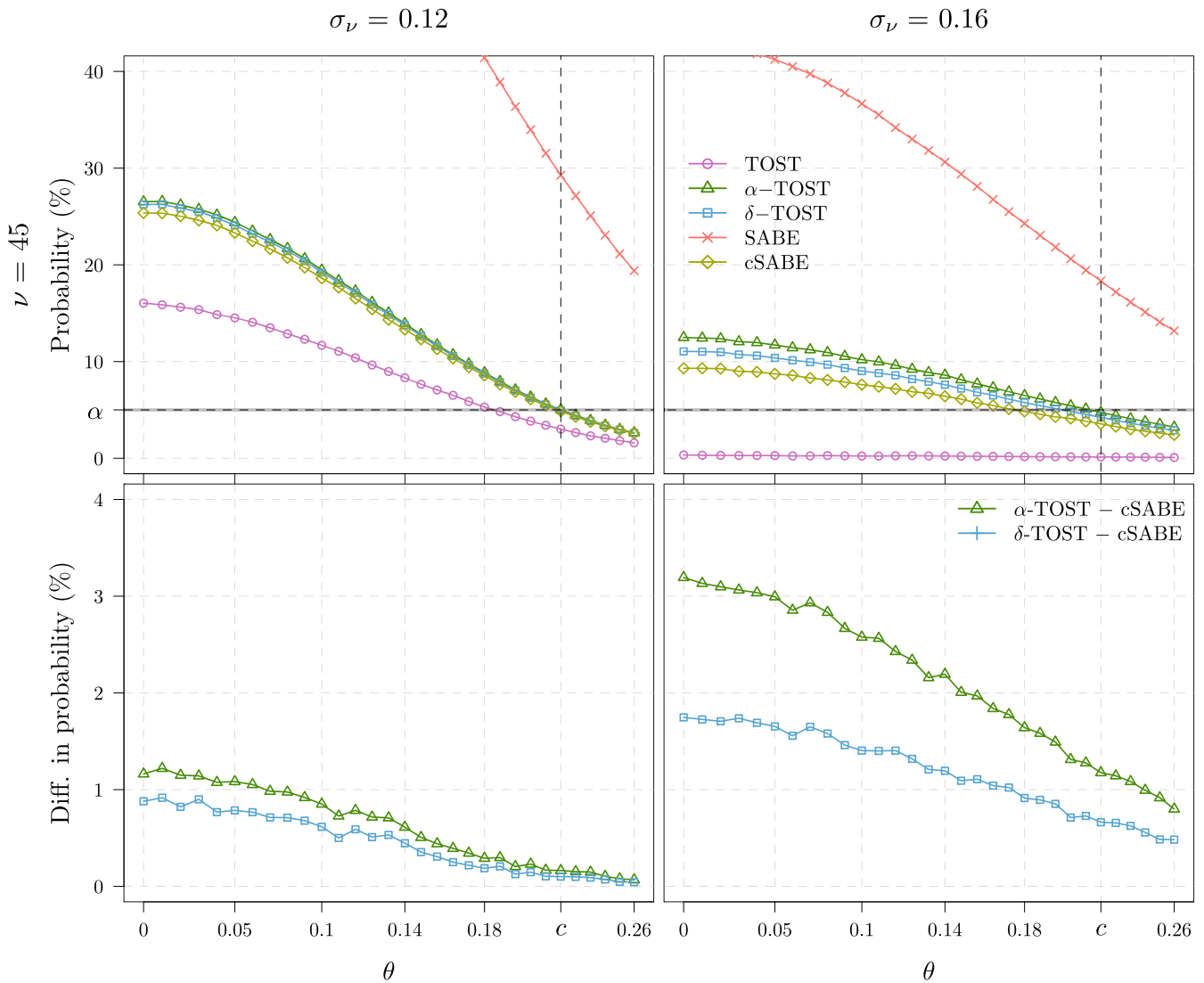
## APPENDIX E. EMPIRICAL COMPARISONS WITH THE CORRECTED SABE

In this section, we present a simulation study to compare the power and level of the TOST,  $\alpha$ -TOST,  $\delta$ -TOST, the EMA implementation of the SABE and its corrected version proposed by Labes and Schütz,<sup>28</sup> which they call the iteratively adjusted  $\alpha$  of the Average BioEquivalence Expanding Limits (ABEL). The aim is to compare testing procedures that either correct for the level ( $\alpha$ -TOST), for the equivalence limits ( $\delta$ -TOST) or both (the corrected SABE) in a paired setting closely related to the case study presented in Section 4 and allowing to estimate the within-subject variability of the reference treatment required by SABE-like methods. We should stress that the SABE and its corrected version are usually used with replicated cross-over designs<sup>26</sup> and that their use in a simple paired design can be viewed as a relaxation of the constraints imposed by regulatory authorities that still allows to validly investigate their finite sample properties. The corrected SABE method is implemented as defined by the EMA guidelines, by computing the size through Monte Carlo integration ( $10^5$  simulations), and applying a correction on the level to match the original one,  $\alpha$ , only when the test is liberal given the adjusted (bio)equivalence limits. Note that, as the SABE procedure also requires the estimated  $\hat{\theta}$  to lie inside the *original* equivalence bounds to declare bioequivalence, the size can only be computed using Monte Carlo simulations. The model we consider is given by

$$X_{ij} = \theta_j + u_i + \epsilon_{ij}, \quad u_i \sim \mathcal{N}(0, \sigma_1^2), \quad \epsilon_{ij} \sim \mathcal{N}(0, \sigma_2^2),$$

with  $i = 1, \dots, n$  and  $j = 1, 2$ . By taking the paired differences we obtain

$$D_i = X_{i,1} - X_{i,2} = \theta + \epsilon_{i,1} - \epsilon_{i,2} = \theta + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, 2\sigma_2^2).$$



**FIGURE E1** First row: Empirical probability of declaring bioequivalence (y-axis) computed using the setting of Simulation 4, as a function of  $\theta$  (x-axis) and  $\sigma_\nu$  (columns), with  $\nu = 45$ , for the TOST (pink circles), the  $\alpha$ -TOST (green triangles), the  $\delta$ -TOST (blue squares), the SABE (red crosses), and the corrected SABE (light-green diamonds). The tight gray area stands for a 99% simulation error tolerance interval of (4.84, 5.16) corresponding to  $\alpha = 5\%$  and  $B = 10^5$  Monte Carlo samples. Second row: The difference of the empirical probabilities between the  $\alpha$ -TOST and the cSABE (green triangles), and between the  $\delta$ -TOST and the cSABE (blue triangles). Empirically, the TOST is quite conservative while the SABE is very liberal. In terms of power, the  $\alpha$ -TOST uniformly dominates the other two methods and the  $\delta$ -TOST uniformly dominates the cSABE.

Thus, we have  $\hat{\theta} = \bar{D} \sim \mathcal{N}(\theta, 2\sigma_2^2/n)$ . By defining  $\sigma_\nu^2 = 2\sigma_2^2/n$  and  $\nu = n - 1$ , we have  $\hat{\theta} \sim \mathcal{N}(\theta, \sigma_\nu^2)$  and  $\frac{\nu\hat{\sigma}_\nu^2}{\sigma_\nu^2} \sim \chi_\nu^2$ . In this setting, the coefficient of variation (CV) is  $CV = \sqrt{\exp(\sigma_2^2) - 1} = \sqrt{\exp\left(\frac{n\sigma_\nu^2}{2}\right) - 1}$ , and we consider the following estimator

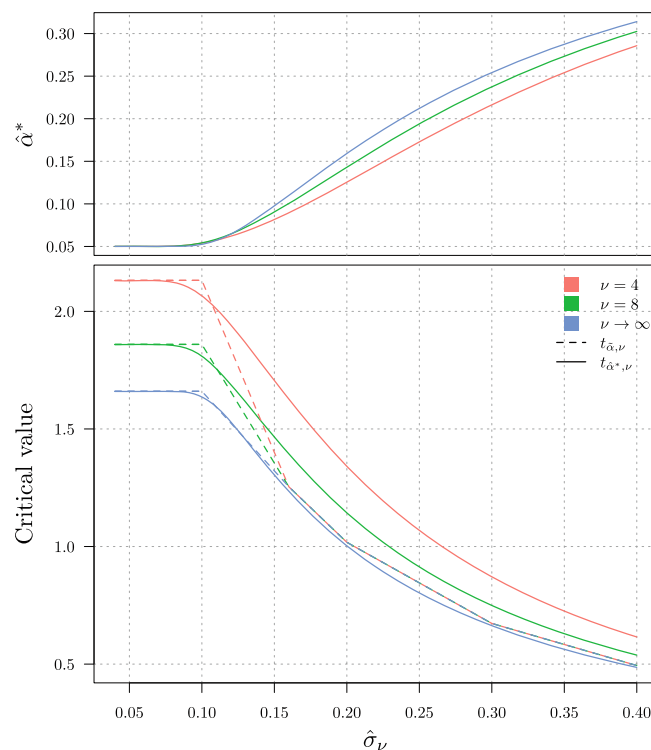
$$\widehat{CV} = \sqrt{\exp\left(\frac{n\hat{\sigma}_\nu^2}{2}\right) - 1}.$$

The parameters and settings considered in this Simulation are reported in Table 1 under Simulation 4. The first row of Figure E1 shows the empirical probabilities of declaring equivalence (y-axes) of each method (colored lines) as a function of  $\theta$  (x-axes) for two values  $\sigma_\nu$  (columns) when  $\nu = 45$ . Comparing the empirical size (obtained when  $c = \theta$ ) of the different

methods, we can note that the TOST is quite conservative while the SABE is very liberal for the considered values of  $\sigma_\nu$ . The  $\alpha$ -TOST,  $\delta$ -TOST and cSABE are size- $\alpha$  when  $\sigma_\nu = 0.12$  with their respective empirical size lying inside the 99% simulation error tolerance interval of (4.84, 5.16), corresponding to  $\alpha = 5\%$  and  $B = 10^5$  Monte Carlo simulations. On the other hand, for a value of  $\sigma_\nu = 0.16$ , none of these methods is size- $\alpha$  as all empirical sizes lie below the simulation error tolerance interval (with estimated values of approximately 0.0475, 0.0424, and 0.0357 for the  $\alpha$ -TOST, the  $\delta$ -TOST and the cSABE, respectively). The second row of the Figure E1 shows the difference in the probabilities of declaring equivalence of the  $\alpha$ - and  $\delta$ -TOST methods compared to the cSABE (y-axes) as a function of  $\theta$  (x-axes) for two values of  $\sigma_\nu$  (columns) when  $\nu = 45$ . Empirically, we can note that, in terms of power, the  $\alpha$ -TOST uniformly dominates the other two methods and the  $\delta$ -TOST uniformly dominates the cSABE. This again suggests that an adjustment on the level of the TOST is the most effective way to improve the finite sample properties of equivalence testing.

## APPENDIX F. COMPARISON OF THE $\alpha$ -TOST WITH CAO AND MATHEW'S METHOD

In Figure F1, the critical values for different values of  $\nu$ , obtained by Cao and Mathew<sup>41</sup> (Table 1) and the ones obtained using the  $\alpha$ -TOST, are compared. One can see that the method of Cao and Mathew<sup>41</sup> appears to be an approximation of the  $\alpha$ -TOST, evaluated asymptotically, that is, at  $\nu \rightarrow \infty$ .



**FIGURE F1** Upper panel: Values of  $\hat{\alpha}^*$  of the  $\alpha$ -TOST (y-axis) as a function of  $\hat{\sigma}_\nu$  (x-axis) for different values of  $\nu$  (colored lines). Lower panel: Comparison of the critical values (y-axis) obtained by the method of Cao and Mathew<sup>41</sup> (dashed lines showing  $t_{\hat{\alpha},\nu}$ ) and of the  $\alpha$ -TOST (solid lines showing  $t_{\hat{\alpha}^*,\nu}$ ) as a function of  $\hat{\sigma}_\nu$  (x-axis) for different values of  $\nu$  (colored lines). Note that for all values of  $\nu$  considered here and for values of  $\hat{\sigma}_\nu$  above 0.1, the critical values of Cao and Mathew<sup>41</sup> correspond to a piecewise version of the critical values of the  $\alpha$ -TOST obtained when  $\nu$  is large. Therefore, their correction appears to be an approximation of the  $\alpha$ -TOST, evaluated asymptotically, that is, at  $\nu \rightarrow \infty$ .