

Letter to the Editor on:  
“Comparing Groups of Time Dependent Data Using  
Locally Weighted Scatterplot Smoothing  
Alpha-Adjusted Serial T-tests” by Niiler (2020)

Dominik Liebl\*  
*University Bonn*

1 **1. Introduction**

2 The main contribution in Niiler (2020) is a new  $\alpha$  level correction  
3 method for multiple testing in a functional data context. A secondary  
4 contribution is the proposal to use the nonparametric LOESS smoothing  
5 method of Cleveland & Devlin (1988) to overcome problems with irregu-  
6 larly sampled functional data. Unfortunately, both contributions generally  
7 lead to invalid statistical inferences and I will explain these issues in the  
8 following.

9 In Section 2, I briefly explain the  $\alpha$  level correction method of Niiler  
10 (2020) and introduce necessary notation. In Section 3, I explain the sta-  
11 tistical shortcomings in Niiler (2020) and demonstrate that the  $\alpha$  level  
12 correction of Niiler (2020) leads to invalid inferences. A short conclusion  
13 is given in Section 4.

14 **2. The  $\alpha$  level correction method of Niiler (2020)**

Niiler (2020) considers the statistical problem of testing for differ-  
ences between mean functions of two groups ( $a$  and  $b$ ) using two inde-  
pendent samples of biomechanical functional curve data. Exemplary data  
are shown, for instance, in Fig 1 in Niiler (2020). The testing is done using  
a series of  $M$  many two-sample test statistics

$$t(g_j) = \frac{\hat{\mu}_a(g_j) - \hat{\mu}_b(g_j)}{\sqrt{\frac{\hat{\sigma}_a^2(g_j)}{N_a} + \frac{\hat{\sigma}_b^2(g_j)}{N_b}}}, \quad j = 1, \dots, M, \quad (1)$$

15 where each test statistic  $t(g_j)$  conducts a statistical hypothesis test specific  
16 to a grid point  $g_j$  with, for instance,  $0\% \leq g_1 < \dots < g_j < \dots < g_M \leq$   
17  $100\%$  when using a standardized time (gait cycle) domain  $[0\%, 100\%]$ . The  
18 estimates  $\hat{\mu}_a(g_j)$ ,  $\hat{\mu}_b(g_j)$ ,  $\hat{\sigma}_a^2(g_j)$ , and  $\hat{\sigma}_b^2(g_j)$  denote the mean and variance  
19 estimates of groups  $a$  and  $b$  at grid point  $g_j$ , and  $N_a$  and  $N_b$  denote the  
20 samples sizes. If one uses the classic sample mean and variance estimates,  
21 the test statistic in (1) becomes the classic Welch’s  $t$ -test. Niiler (2020),

---

\*Corresponding author: Dominik Liebl, University Bonn, Institute of Finance and  
Statistics, Adenauerallee 24-26, 53113 Bonn (Germany), dliebl@uni-bonn.de. This  
research was supported by the Hausdorff Center of Mathematics (HCM), a Cluster of  
Excellence at the University financed by the German Research Foundation (DFG).

22 however, suggests using the mean and variance estimates as computed by  
 23 the R function `loess` for local polynomial regression (Cleveland & Devlin,  
 24 1988; R Core Team, 2021).

25 The  $j$ th test statistic,  $t(g_j)$ , tests the null hypothesis,  $H_0$ , of equal  
 26 means against the two-sided alternative,  $H_1$ :

$$27 \quad \begin{aligned} H_0(g_j): & \quad \mu_a(g_j) = \mu_b(g_j) \\ H_1(g_j): & \quad \mu_a(g_j) \neq \mu_b(g_j) \end{aligned}$$

28 However, one is generally not interested in the test decision at a single  
 29 grid point,  $g_j$ , but one uses the whole family of  $M$  many test statistics  
 30  $\{t(g_1), \dots, t(g_M)\}$  to find regions in  $[0\%, 100\%]$  over which the mean func-  
 31 tions are statistically different from each other. This leads to a severe  
 32 multiple testing problem since the number of tests,  $M$ , can be arbitrarily  
 33 large.

Statistical multiple testing procedures must control the *family-wise*  
 type I error rate. I.e., the probability of observing type I errors in at least  
 one of the tests  $\{t(g_1), \dots, t(g_M)\}$  must be bounded from above by the  
 pre-chosen significance level  $\alpha$ ,

$$P_{H_0}(\text{reject } H_0(g_j) \text{ for at least one } j \in \{1, \dots, M\}) \leq \alpha \quad (2)$$

34 with, for instance,  $\alpha = 0.05$ .

35 To control the family-wise type I error rate, one could use, for instance,  
 36 the classic Bonferroni correction, where each point-wise null hypothesis  
 37  $H_0(g_j)$  is tested at the reduced significance level of  $\alpha' = \alpha/M$ . However,  
 38 for biomechanical curve data, Bonferroni corrections can result in unnec-  
 39 essarily conservative (low power) testing procedures since biomechanical  
 40 curve data are typically relatively smooth and, therefore, the point-wise  
 41 test statistics  $\{t(g_1), \dots, t(g_M)\}$  are typically strongly correlated with each  
 42 other. These correlations are ignored by the Bonferroni correction. Similar  
 43 issues arise with other standard  $\alpha$  level corrections such as, for instance,  
 44 the Holm-Bonferroni or the Hochberg correction.

The main contribution in Niiler (2020) is the proposal of a new, less  
 conservative  $\alpha$  level correction which tries to take into account the corre-  
 lations between the test statistics  $\{t(g_1), \dots, t(g_M)\}$ . The proposed cor-  
 rection is given in equation (2) in Niiler (2020), but also presented here  
 for convenience:

$$\alpha' = \frac{\alpha}{M(1 - \hat{\rho}^2) + \hat{\rho}^2}, \quad (3)$$

where Niiler (2020) sets  $\alpha = 0.05$  and where  $\hat{\rho}$  denotes the sample auto-  
 correlation coefficient between adjacent test statistics  $t(g_j)$  and  $t(g_{j+1})$

$$\hat{\rho} = \frac{(M - 1)^{-1} \sum_{j=1}^{M-1} (t(g_j) - \bar{t})(t(g_{j+1}) - \bar{t})}{M^{-1} \sum_{j=1}^M (t(g_j) - \bar{t})^2}$$

45 with  $\bar{t} = M^{-1} \sum_{j=1}^M t(g_j)$ . Niiler (2020) motivates his proposal as follows:  
 46 While the case of perfect auto-correlation  $\hat{\rho} = 1$  leads to no  $\alpha$  level correc-  
 47 tion ( $\alpha' = \alpha$ ), the case of no auto-correlation  $\hat{\rho} = 0$  leads to the Bonferroni  
 48 correction ( $\alpha' = \alpha/M$ ).

49 **3. Main critique: Invalid  $\alpha$  level correction**

50 The auto-correlation coefficient  $\hat{\rho}$  is not a meaningful statistic in case  
 51 of non-stationary time series.<sup>1</sup> But even if  $t(g_1), \dots, t(g_M)$  were a station-  
 52 ary series of test statistics,  $\hat{\rho}$  would only measure the correlation between  
 53 adjacent test statistics  $t(g_j)$  and  $t(g_{j+1})$ . All the other pair-wise correla-  
 54 tions are not considered. Therefore, one cannot expect that this method  
 55 is able to control the family-wise type I error rate – except for trivial and  
 56 practically irrelevant special cases. Indeed, Niiler (2020) does not provide  
 57 any theoretical justification for his  $\alpha$  correction method in (3), and, as far  
 58 as I know, there is no similar correction in the statistical literature.

59 Niiler (2020) uses a Monte Carlo simulation study to demonstrate that  
 60 his  $\alpha$  level correction is able to control the family-wise type I error rate;  
 61 see Appendix C of the supplementary material of Niiler (2020). Simula-  
 62 tions, however, cannot replace theoretical considerations since the consid-  
 63 ered simulation scenarios may only reflect non-generalizing special cases  
 64 – which is exactly what happened in Niiler (2020).

The only reason, why Niiler (2020) was able to demonstrate that his  
 $\alpha$  level correction is able to control the family-wise type I error rate, is  
 the very specific choice of his simulation study. Niiler (2020) considers the  
 following overly simple type of random functions

$$\sin(x) + Z, \quad \text{with } x \in [0, 2\pi] \quad \text{and} \quad Z \sim \mathcal{N}(0, 1) \quad (4)$$

65 which are just random vertical shifts of a deterministic sinus function; see  
 Figure 1 (A) and (B).

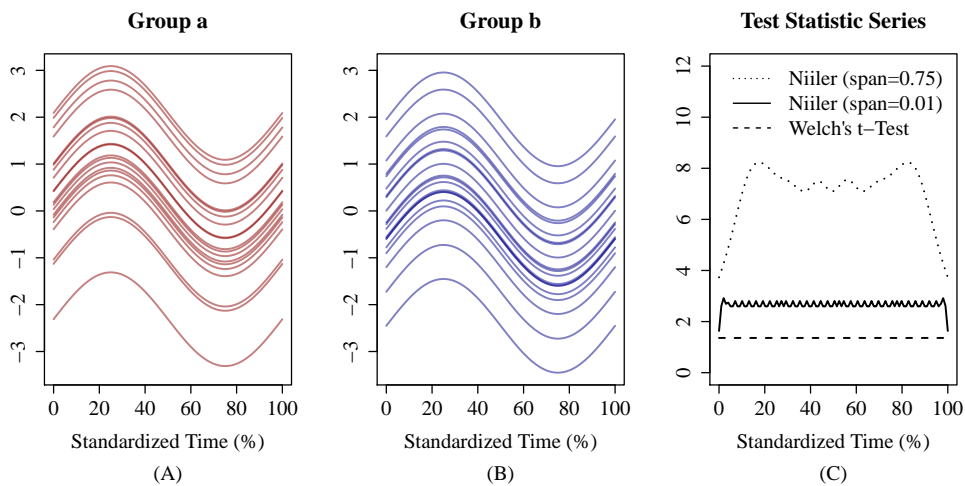


Figure 1: Plots (A) and (C): Exemplary functions from the sinus-shift random process in (4) as used in the Monte Carlo simulation study in Niiler (2020) for checking the family-wise type I error rate. Plot (C): test statistic series  $t(g_1), \dots, t(g_M)$  computed using the LOESS smoothing method suggested by Niiler (2020) with smoothing parameters  $\text{spar} = 0.01$  and  $\text{spar} = 0.75$ . Additionally, the classic Welch's  $t$ -test statistics series is shown.

66 For this special case, the Welch's  $t$ -test statistics  $t(g_1), \dots, t(g_M)$  are  
 67 all exactly equal to each other (see Fig. 1 (C)) which demonstrates that  
 68

<sup>1</sup>Biomechanical curve data are usually not stationary.

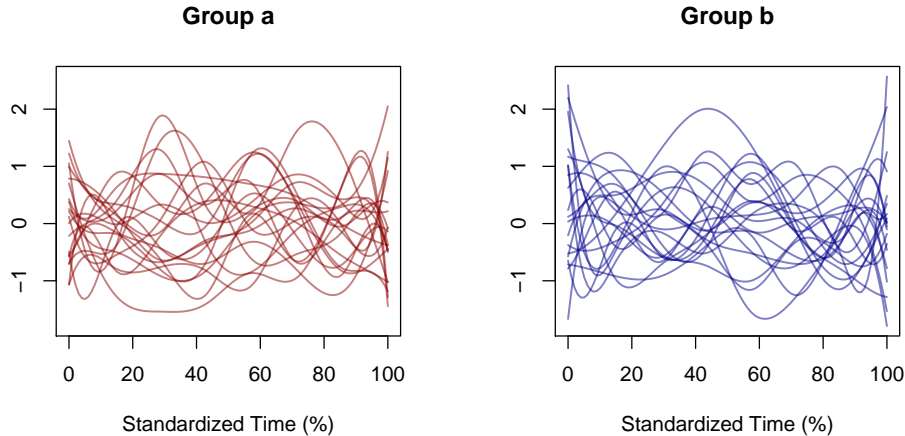


Figure 2: Exemplary functions from the B-splines random process in (5) as used in my Monte Carlo simulation for checking the family-wise type I error rate.

69 there is no multiple testing problem in this special case. All  $t$ -tests are  
70 either *simultaneously* significant or not, and, therefore, the family-wise  
71 type I error rate coincides with the type I error rate of a single  $t$ -test  
72 which makes the stochastic process in (4) unsuitable for checking  $\alpha$  level  
73 correction methods.

In the following, I consider a practically more relevant situation by drawing random functions from

$$f(t) = \sum_{k=1}^{10} Z_k B_{k,10}(t), \quad \text{with } t \in [0\%, 100\%], \quad (5)$$

74 where  $Z_k \sim \mathcal{N}(0, 1)$  and where  $B_{k,10}(t)$  denotes the  $k$ th cubic B-spline  
75 function based on equidistant knots in  $[0, 100]$ ; see Figure 2. B-spline func-  
76 tions have compact supports which guarantees that the random functions  
77 in (5) consist of independent as well as dependent stochastic components  
78 making this case suitable for checking  $\alpha$  level correction methods.

79 To check the family-wise type I error rate of the  $\alpha$  level correction  
80 in Niiler (2020), I simulate  $N_a$  and  $N_b$  many functions from (5). Both  
81 groups  $a$  and  $b$  have the same population mean (zero) such that all  $M$   
82 null hypotheses  $H_0(g_j): \mu_a(g_j) = \mu_b(g_j)$ ,  $j = 1, \dots, M$ , are fulfilled. Under  
83 this scenario, the family-wise type I error rate must be smaller or equal  
84 to the pre-selected significance level  $\alpha = 0.05$ ; otherwise, the  $\alpha$  correction  
85 method is invalid and cannot be used in practice.

86 Table 1 shows my simulation results based on 10,000 Monte Carlo repli-  
87 cations. To check the effect of different choices for the number of sampling  
88 grid points  $M$ , I consider the values  $M \in \{50, 75, 100\}$ . To check the ef-  
89 fect of different sample sizes, I consider the values  $N_a = N_b \in \{10, 20, 50\}$ .  
90 The LOESS smoother used by Niiler (2020) involves setting a smoothing  
91 parameter, where I use a small smoothing parameter  $\text{span} = 0.01$  and a  
92 relatively large smoothing parameter  $\text{span} = 0.75$ . I compare the infer-  
93 ence method of Niiler (2020) with a Bonferroni adjusted series of Welch's  
94  $t$ -tests and with the random field theory based method SPM1d (Pataky,  
95 2016). For the latter method, I use the R package `ffscb` which contains

96 SPM1d based bands, but also more general simultaneous confidence bands  
 as proposed by Liebl & Reimherr (2020).

Table 1: False positive (type I error) rates under the null-hypothesis for a pre-selected significance level  $\alpha = 0.05$

$N_a = N_b$	$M$	Niiler (2020)		$t$ -test (Bonferroni)	SPM1d
		(span= 0.01)	(span= 0.75)		
10	50	0.10	0.82	0.00	0.05
10	75	0.60	0.92	0.00	0.05
10	100	0.81	0.97	0.00	0.06
20	50	0.19	0.88	0.01	0.05
20	75	0.73	0.95	0.00	0.05
20	100	0.88	0.98	0.00	0.05
50	50	0.25	0.90	0.01	0.05
50	75	0.77	0.96	0.01	0.05
50	100	0.91	0.98	0.00	0.05

97  
 98 The only two methods that are able to control the type I error rate in  
 99 this simulation study are the point-wise  $t$ -tests with Bonferroni correction  
 100 and the random fields theory based method SPM1d. While the Bonferroni  
 101 correction is overly conservative, SPM1d is able to exploit the significance  
 102 level  $\alpha = 0.05$  very well. The inference procedure proposed by Niiler  
 103 (2020) fails to control the type I error rate severely and, therefore, leads  
 104 to invalid inferences.

105 **Further issues.** Niiler (2020) misses several further issues. For instance,  
 106 the standard errors computed by the `loess` function in R are invalid for  
 107 smoothing functional data when  $M$  becomes large (see Liebl, 2019). More-  
 108 over, nonparametric smoothing methods like LOESS have biased estimates  
 109 and can suffer from boundary problems both leads to distorted estimation  
 110 results as indicated in Figure 2 (C). To get valid inference in finite samples,  
 111 one would need bias and boundary corrections – both are not considered  
 112 in Niiler (2020).

#### 113 4. Conclusion

114 The development of simultaneous inference methods for functional  
 115 data is an active research field in statistics (see Degras (2011), Cao et al.  
 116 (2012), Wang et al. (2020), Pini & Vantini (2017), Choi & Reimherr  
 117 (2018), Liebl & Reimherr (2020), and many others). New steps forward  
 118 in this literature are often published in the most prestigious statistical  
 119 science journals. In my humble opinion, the work of Niiler (2020) fails to  
 120 make a contribution to this literature.

121 Every statistical testing procedure can lead to invalid inferences when  
 122 misapplied. However, this is different here. The proposed  $\alpha$  level testing  
 123 procedure of Niiler (2020) will lead to false inferences even when “cor-  
 124 rectly” applied. This is a serious issue since the method is already applied  
 125 in the literature (see Shoja et al., 2020).

126 **References**

- 127 Cao, G., Yang, L., & Todem, D. (2012). Simultaneous inference for the  
128 mean function based on dense functional data. *Journal of Nonparametric*  
129 *Statistics*, *24*, 359–377.
- 130 Choi, H., & Reimherr, M. (2018). A geometric approach to confidence  
131 regions and bands for functional parameters. *Journal of the Royal Sta-*  
132 *tistical Society: Series B (Statistical Methodology)*, *80*, 239–260.
- 133 Cleveland, W. S., & Devlin, S. J. (1988). Locally weighted regression: An  
134 approach to regression analysis by local fitting. *Journal of the American*  
135 *Statistical Association*, *83*, 596–610.
- 136 Degras, D. A. (2011). Simultaneous confidence bands for nonparametric  
137 regression with functional data. *Statistica Sinica*, *21*, 1735–1765.
- 138 Liebl, D. (2019). Inference for sparse and dense functional data with  
139 covariate adjustments. *Journal of Multivariate Analysis*, *170*, 315–335.
- 140 Liebl, D., & Reimherr, M. (2020). Fast and fair simultaneous confidence  
141 bands for functional parameters. *Preprint arXiv:1910.00131*, (pp. 1–  
142 37).
- 143 Niiler, T. (2020). Comparing groups of time dependent data using locally  
144 weighted scatterplot smoothing alpha-adjusted serial t-tests. *Gait &*  
145 *Posture*, *76*, 58–63.
- 146 Pataky, T. C. (2016). Rft1d: Smooth one-dimensional random field up-  
147 crossing probabilities in python. *Journal of Statistical Software*, *71*,  
148 1–22.
- 149 Pini, A., & Vantini, S. (2017). Interval-wise testing for functional data.  
150 *Journal of Nonparametric Statistics*, *29*, 407–424.
- 151 R Core Team (2021). *R: A Language and Environment for Statistical*  
152 *Computing*. R Foundation for Statistical Computing Vienna, Austria.  
153 URL: <https://www.R-project.org/>.
- 154 Shoja, O., Farsi, A., Towhidkhah, F., Feldman, A. G., Abdoli, B., &  
155 Bahramian, A. (2020). Visual deprivation is met with active changes  
156 in ground reaction forces to minimize worsening balance and stability  
157 during walking. *Experimental Brain Research*, *238*, 369–379.
- 158 Wang, Y., Wang, G., Wang, L., & Ogden, R. T. (2020). Simultaneous  
159 confidence corridors for mean functions in functional data analysis of  
160 imaging data. *Biometrics*, *76*, 427–437.