

**1 A New Heuristic Method for the Optimal Selection of Tolerance r for Entropy
2 Indices**

3 Dominique Makowski¹

4 ¹ School of Social Sciences, Nanyang Technological University, Singapore

5 Author Note

**6 Correspondence concerning this article should be addressed to Dominique Makowski,
7 HSS 04-18, 48 Nanyang Avenue, Singapore (dom.makowski@gmail.com).**

8

Abstract

9 The tolerance threshold r is a key parameter of several entropy algorithms (e.g., *SampEn*).
10 Unfortunately, the gold standard method to estimate its optimal value - the one that
11 maximizes *ApEn* - is computationally costly, prompting users to rely to cargo-cult
12 rules-of-thumb such as $0.2 * \text{SD}$. In this simulation study, we attempted at validating a
13 new heuristic, based on the embedding dimension m and the signal's length n (optimal $r =$
14 $\text{SD} * 0.281(m-1) + 0.005(\log(n)) - 0.02(m-1 * \log(n))$), which was superior to other
15 existing heuristics. All the methods of optimal tolerance r estimation used in this study
16 are available in the *NeuroKit2* Python software (Makowski et al., 2021).

17 *Keywords:* chaos, complexity, entropy, tolerance, r, physiology

18 Word count: 925

19 A New Heuristic Method for the Optimal Selection of Tolerance r for Entropy
20 Indices

21 Introduction

22 Complexity analysis is a growing approach to physiological signals, including cardiac
23 (e.g., Heart Rate Variability, [1]) and brain activity [2]. It is an umbrella term for the usage
24 of various complexity indices that quantify concepts such as chaos, entropy, fractal
25 dimension, randomness, predictability, and information. Importantly, some of the most
26 popular indices of entropy (e.g., *ApEn*, *SampEn*, their fuzzy and multiscale variations) and
27 recurrence quantification analysis (RQA), rely on a similar set of parameters. Namely,
28 these are the delay τ , the embedding dimension m , and the tolerance r , which are critical
29 to accurately capture the space in which complexity becomes quantifiable. Unfortunately,
30 despite the existence of methods to estimate optimal values for these parameters depending
31 on the signal at hand, their choice often relies on simple heuristics and cargo-cult
32 conventions.

33 Such is the case of the tolerance threshold r , which typically corresponds to the
34 minimal distance required to assigning two points in a state-space as belonging to the same
35 state. It directly impacts the amount of “recurrences” of a system and the measure of its
36 tendency to revisit past states, which is the base metric for the calculation of the
37 aforementioned entropy indices. Despite its importance, it is often selected as a function of
38 the standard deviation (SD) of the signal, with (in)famous magic values including 0.1 or
39 $0.2 * SD$ [3]. One of the reason for the longevity of such approach is 1) the past literature
40 (as many past studies used it, it becomes easier to justify the choice of the same values)
41 and 2) the fact that other approaches to estimate the optimal r are computationally costly.

42 The aim of the present study is to investigate the relationship between different
43 methods for optimal tolerance r estimation. The ground-truth method used is the
44 tolerance value corresponding to a maximal value of Approximate Entropy - *ApEn* [4–6].

⁴⁵ As this method is computationally expensive, the objective is to see whether any heuristic
⁴⁶ proxies can be used to satisfactorily approximate $r_{maxApEn}$.

⁴⁷ **Methods**

⁴⁸ For $n = 5760$ combinations of different signal types and lengths, as well as noise
⁴⁹ types and intensities (the procedure used was the same as in [7], and the data generation
⁵⁰ code is available at
⁵¹ <https://github.com/DominiqueMakowski/ComplexityTolerance>), we computed
⁵² the Approximate Entropy ($ApEn$), which peak is used to estimate the optimal tolerance
⁵³ level for time-delay embedding spaces ranging from 1 to 9 embedding dimensions m .

⁵⁴ The aim of the analysis is to establish a new heuristic based only on the signal's SD
⁵⁵ and the embedding dimension m ; and compare all of these approximations with other
⁵⁶ existing heuristics such as 0.2 SD , *Chon* [6], and the *Schötzl* method
⁵⁷ $(1.3334 + 0.5627 * \log(\text{dimension}))$ implemented in the package *nolds* [8].

⁵⁸ **Results**

⁵⁹ **Maximum Approximate Entropy.** **Figure 1** shows the normalized value of
⁶⁰ Approximate Entropy $ApEn$ as a function of tolerance r and embedding dimension m . As
⁶¹ expected, the value of $ApEn$ peaks at certain values of r (hence its usage as an indicator of
⁶² the optimal tolerance). The location of this peak seems strongly impacted by the
⁶³ embedding dimension m , getting more variable - and on average larger - as m increases.

⁶⁴ **New Heuristic.** Selecting the tolerance based on the signal's SD alone makes
⁶⁵ hardly sense, as the embedding dimension has a strong impact on it. In order to validate a
⁶⁶ new heuristic to approximate the optimal r value based on the embedding dimension m
⁶⁷ and the signal's length, we compared different regression specifications using the BIC-based
⁶⁸ Bayes Factor test. The model including the log-transformed signal length (in samples) and
⁶⁹ the embedding dimension m minus 1 (with no intercept) performed significantly better

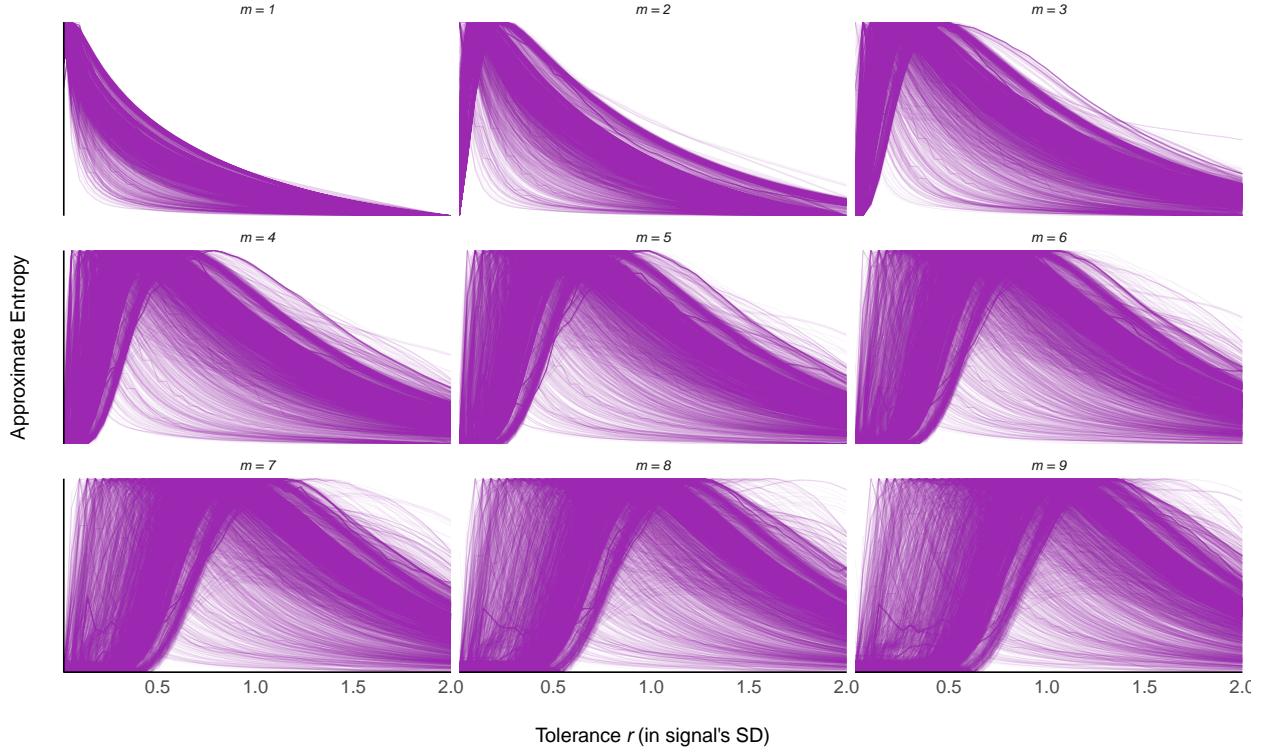


Figure 1. Approximate Entropy $ApEn$ as a function of tolerance r (expressed in signal SD) and embedding dimension m .

70 ($BF_{10} > 1000$) than any other version, with an explained variance of 92.64%. Based on this
 71 simple regression model, we can derive the following approximation (assuming a
 72 standardized signal with an SD of 1):

$$\hat{r} = 0.2811(m - 1) + 0.0049(\log(n)) - 0.02(m - 1 \times \log(n)) \quad (1)$$

73 It is to note that shorter signals require larger tolerance values, and the impact of
 74 length lowers as the signal gets longer. Also, for a dimension m of 2 (and short signal
 75 lengths), this equation returns close results to the *0.2 SD* heuristic, which is not entirely
 76 surprising as the latter was initially derived under such conditions (that are typical of some
 77 applications, such as heart-rate intervals).

78 **Heuristics Comparison.**

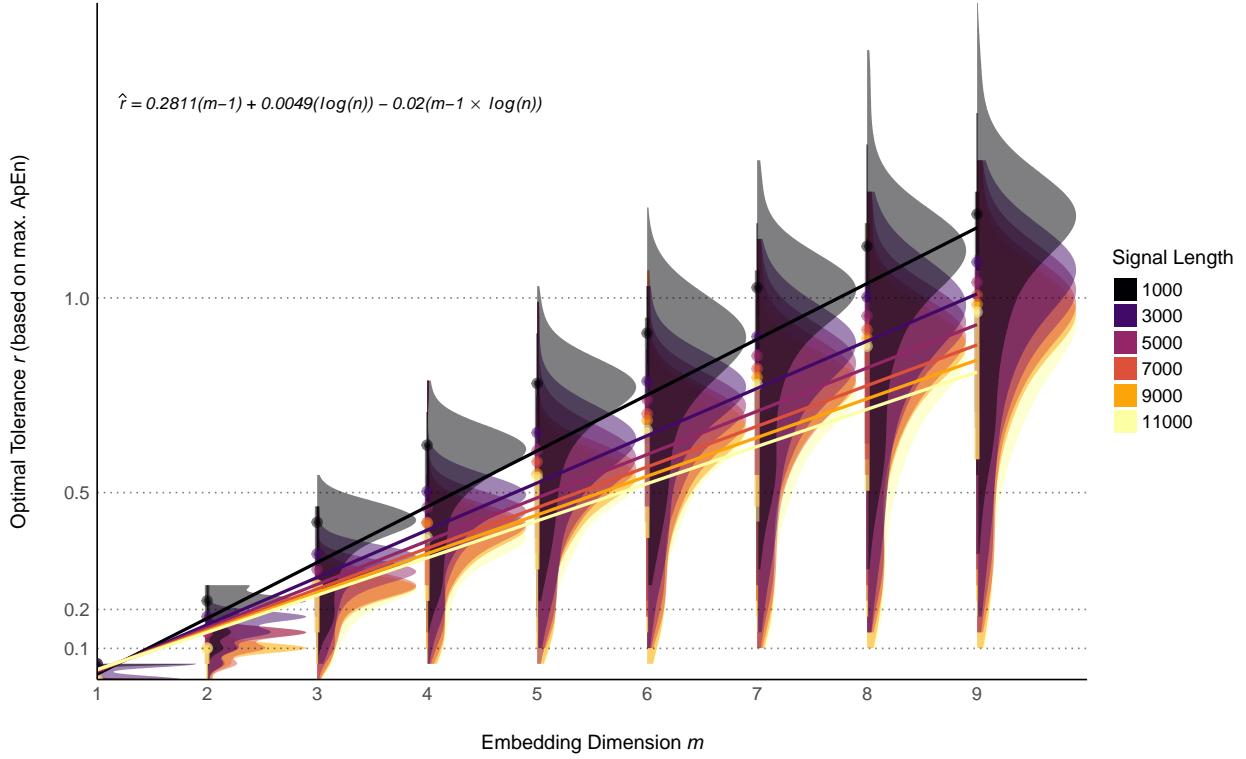


Figure 2. Optimal tolerance values approximated by a new heuristic model based on the embedding dimension m and the signal length n (in samples). The density plots show the true optimal tolerance values as based on max. ApEn.

Table 1

Comparison of Model Performance Indices

Model	BIC	R2	BF
Makowski	-41375.73	0.77	1.00
Scholzel	-24637.03	0.68	0.00e+00
Chon	24614.17	0.18	0.00e+00
SD	35119.90	0.00	0.00e+00

79 We compared together different methods to approximate $r_{maxApEn}$ (see **Table 1**) by
 80 comparing $r_{maxApEn}$ (our ground-truth) to the value estimated by different methods. Our
 81 heuristic method presented in this study, based on the signal's SD, the embedding
 82 dimension and the log-transformed length, surpassed any other models ($BF_{10} > 1000$, $R^2 =$

83 0.77), and was followed by the *Schötzel* adjustment ($R^2 = 0.68$).

84 **Discussion**

85 The tolerance threshold r is a key parameter of several entropy algorithms, including
86 widely popular ones like *SampEn*. The current gold standard method to estimate the
87 optimal r is to compute Approximate Entropy (*ApEn*) over a range of different r values
88 and to select the one corresponding to the maximum *ApEn* value. Unfortunately, this
89 method is computationally costly.

90 In this study, we have shown that a simple heuristic approximation based on the
91 embedding dimension m and the log-transformed signal length was a valid approximation
92 of $r_{maxApEn}$, showing superior performance to other heuristic methods. We suggest the use
93 of this method as a default alternative to the *0.2 SD* rule of thumb.

94 While we believe that our data generation procedure was able to generate a wide
95 variety of signals, and that our results are to some extent generalizable, future studies
96 could attempt at refining the estimation procedures for specific signals (for instance, EEG,
97 or heart rate data). All the methods of optimal tolerance r estimation used in this study,
98 including our new proposal, are available in the *NeuroKit2* open-source Python software,
99 `complexity_tolerance()` function [9].

100 References

- 101 1. Pham, T.; Lau, Z.J.; Chen, S.; Makowski, D. Heart Rate Variability in Psychology:
102 A Review of HRV Indices and an Analysis Tutorial. *Sensors* **2021**, *21*, 3998.
- 103 2. Lau, Z.J.; Pham, T.; Annabel, S.; Makowski, D. Brain Entropy, Fractal Dimensions
104 and Predictability: A Review of Complexity Measures for EEG in Healthy and Neu-
ropsychiatric Populations. **2021**.
- 105 3. Pincus, S.M.; Viscarello, R.R. Approximate Entropy: A Regularity Measure for Fetal
106 Heart Rate Analysis. *Obstetrics and gynecology* **1992**, *79*, 249–255.
- 107 4. Chen, X.; Solomon, I.; Chon, K. Parameter Selection Criteria in Approximate Entropy
108 and Sample Entropy with Application to Neural Respiratory Signals. *Am. J. Physiol.
Regul. Integr. Comp. Physiol.*, *to be published* **2008**.
- 109 5. Lu, S.; Chen, X.; Kanters, J.K.; Solomon, I.C.; Chon, K.H. Automatic Selection of
110 the Threshold Value r for Approximate Entropy. *IEEE Transactions on Biomedical
Engineering* **2008**, *55*, 1966–1972.
- 111 6. Chon, K.H.; Scully, C.G.; Lu, S. Approximate Entropy for All Signals. *IEEE engi-
112 neering in medicine and biology magazine* **2009**, *28*, 18–23.
- 113 7. Makowski, D.; Te, A.S.; Pham, T.; Lau, Z.J.; Chen, S.H.A. The Structure of Chaos:
114 An Empirical Comparison of Fractal Physiology Complexity Indices Using NeuroKit2.
Entropy **2022**, *24*, 1036, doi:10.3390/e24081036.
- 115 8. Schölzel, C. *Nonlinear Measures for Dynamical Systems*; Zenodo, 2019;
116
- 117 9. Makowski, D.; Pham, T.; Lau, Z.J.; Brammer, J.C.; Lespinasse, F.; Pham, H.;
118 Schölzel, C.; Chen, S. NeuroKit2: A Python Toolbox for Neurophysiological Sig-
nal Processing. *Behavior research methods* **2021**, *53*, 1689–1696.