

Approximating the universal thermal climate index using sparse regression with orthogonal polynomials

Sabin Roman¹, Gregor Skok², Ljupčo Todorovski^{2,1}, and Sašo Džeroski¹

¹ Department of Knowledge Technologies, Jožef Stefan Institute
`sabin.roman@ijs.si`

² Faculty of Mathematics and Physics, University of Ljubljana

1 **Abstract.** The Universal Thermal Climate Index (UTCI) is a measure
2 of thermal comfort that quantifies how humans experience environmental
3 conditions. Due to its robustness and versatility as a bioclimatic indica-
4 tor, it has been extensively employed across a wide range of studies in
5 bioclimatology and is increasingly used as an operational measure of
6 outdoor thermal comfort. At the same time, calculating the UTCI value
7 from the relevant environmental parameters is nominally not straight-
8 forward, which is why using a 6th-degree polynomial approximation has
9 become the standard way to calculate UTCI values. At the same time,
10 although it is computationally efficient, the error of this polynomial ap-
11 proximation can be substantial. The goal of this study was to develop
12 an improved version of the polynomial approximation – one that re-
13 tains comparable computational efficiency but is more robust in terms of
14 numerical stability and substantially more accurate, particularly in re-
15 ducing the frequency of larger errors. This goal was successfully achieved
16 using sparse orthogonal regression, namely sparse regression with an or-
17 thogonal polynomial basis, which not only substantially reduces the av-
18 erage errors (i.e., the mean error, the mean absolute error, and the root
19 mean square error) but also drastically reduces the frequency of large
20 errors. By leveraging Legendre polynomial bases, approximation models
21 could be constructed that efficiently populate a Pareto front of accuracy
22 versus complexity and exhibit stable, hierarchical coefficient structures
23 across varying model capacities. Training the new approximation models
24 over only 20% of the data, with the testing performed over the remaining
25 80%, highlights successful generalization, with the results also being ro-
26 bust under bootstrapping. The decomposition effectively approximates
27 the UTCI as a Fourier-like expansion in an orthogonal basis, yielding
28 results near the theoretical optimum in the L_2 (least squares) sense.

29 **Keywords:** Universal Thermal Climate Index · Sparse regression · Or-
30 thogonal polynomials

31 1 Introduction

32 The Universal Thermal Climate Index (UTCI) is a measure of thermal comfort
33 that quantifies how humans experience environmental conditions. It is derived

34 from an advanced thermo-physiological model (Pappenberger et al., 2015) and
35 expressed in units of temperature. The index accounts for multiple factors, in-
36 cluding air temperature, humidity, wind speed, radiation, and clothing insulation
37 (Bröde et al., 2012). A notable advantage of the UTCI compared to many other
38 bioclimatic indices is its ability to represent thermal conditions in terms that
39 are applicable to human strain under a wide range of climatic conditions (e.g.,
40 for both hot and cold conditions, Błażejczyk et al. (2012)). Based on the UTCI
41 value, the environmental conditions can be classified into one of the ten ther-
42 mal stress categories (Bröde et al., 2012), ranging from Extreme heat stress
43 ($UTCI > 43^{\circ}\text{C}$) to Extreme cold stress ($UTCI < -40^{\circ}\text{C}$).

44 Owing to its robustness and versatility as a bioclimatic indicator, the UTCI
45 has been extensively employed across a wide range of studies in bioclimatology
46 and related scientific disciplines. Its applications encompass diverse research ar-
47 eas, including the assessment of regional and local bioclimate characteristics, the
48 study of urban bioclimate, recreation, tourism, and sports, epidemiological and
49 health-related research, as well as the assessment and forecasting of bioclimatic
50 changes (Błażejczyk and Kuchcik, 2021). The UTCI has also seen growing adop-
51 tion across numerous countries as a standardized measure of outdoor thermal
52 comfort and is increasingly integrated into routine operational meteorological
53 forecasts. For example, within Europe, UTCI is used operationally in the Czech
54 Republic, Italy, Poland, Portugal, and Slovenia (Di Napoli et al., 2021a; Kuz-
55 manović et al., 2024).

56 At the same time, calculating the UTCI value from the relevant environ-
57 mental parameters is nominally not straightforward. Namely, the UTCI is based
58 on the Fiala multi-node model of human thermoregulation (Fiala et al., 2012).
59 However, running the complete Fiala model is computationally expensive and
60 requires expert knowledge to operate the complex simulation software (Bröde
61 et al., 2012). This is the reason the authors of Bröde et al. (2012) provided two
62 simplified approximate procedures for calculating the UTCI values that could be
63 used in operational settings. The first approximation is based on a 4-dimensional
64 look-up table of 104 643 accurate pre-calculated UTCI values that cover a wide
65 range of relevant combinations of the meteorological parameters. Using this look-
66 up table, interpolation from nearby data points can be used to determine ap-
67 proximate UTCI values for intermediate values of meteorological parameters.
68 The second approximation is based on a 6th-degree regression polynomial with
69 210 coefficients.

70 Each approximation has its benefits and weaknesses. The look-up table ap-
71 proach is more accurate, but storing the tabulated values and searching for
72 neighboring datapoints poses challenges to the implementation of this algorithm,
73 while also resulting in a longer execution time compared to the other approach
74 (Bröde, 2021a). In contrast, the polynomial approximation is less accurate, but
75 computationally faster and substantially easier to implement in various pro-
76 gramming languages and computational environments, as it relies on only the
77 most common, primitive mathematical operators and does not require storing
78 the tabulated values. At the same time, the motivation for improving the poly-

79 nomial approximation is not simply a matter of storage, since the size of the
 80 look-up table is modest in modern computational settings. Rather, an improved
 81 polynomial approximation remains attractive for several practical reasons:

- 82 (i) It is fully self-contained and does not depend on external tabulated data,
 83 which facilitates reproducibility and makes redistribution and integration
 84 into open-source software and operational tools more straightforward;
- 85 (ii) It is computationally more efficient than look-up-table-based interpolation,
 86 which has been reported to be slower by roughly three orders of magnitude
 87 (Bröde et al., 2012), an important consideration in large-scale applications
 88 such as numerical weather prediction and climate reanalysis;
- 89 (iii) It is simpler to implement and port across programming languages and com-
 90 putational environments, including constrained, embedded, or legacy sys-
 91 tems, because it requires only basic arithmetic operations and avoids the
 92 additional logic needed for multidimensional interpolation, data handling,
 93 and neighborhood search;
- 94 (iv) It provides a direct, continuous, and analytically defined mapping over the
 95 domain of validity, whereas the look-up table still requires interpolation,
 96 and in some cases extrapolation, for environmental states not explicitly rep-
 97 resented in the tabulated values;
- 98 (v) Its predictive behavior on unseen data can be assessed directly through a
 99 train–test evaluation framework; in the present case, training on 20% of the
 100 dataset and testing on the remaining 80% still yields very good predictive
 101 performance, indicating strong generalization.

102 For these reasons, the polynomial approximation is best viewed not as a universal
 103 replacement for the look-up-table approach, but as a complementary alternative
 104 that is particularly useful in applications where speed, portability, reproducibil-
 105 ity, and ease of deployment are important.

106 Due to its simplicity and computational efficiency, the polynomial approx-
 107 imation has become the standard way of calculating the UTCI values. It has
 108 been incorporated into various biclimatic software packages and libraries (e.g.,
 109 the Bioklima software, Błażejczyk (2025), the Thermofeel Python library, Brim-
 110 icombe et al. (2022), and the pyThermalComfort Python library Tartarini and
 111 Schiavon (2020)), as well as numerical weather prediction and reanalysis sys-
 112 tems (e.g., the ALADIN model, Termonia et al. (2018), and the ERA5 reanal-
 113 ysis, Di Napoli et al. (2021b)). At the same time, the error of the polynomial
 114 approximation can be substantial. For example, when evaluated on the afore-
 115 mentioned look-up table of accurate UTCI values, the root-mean-square-error is
 116 about 1.1°C while the frequency of absolute errors larger than 2°C is about 8%,
 117 and the frequency of errors larger than 3°C is about 2%. This is problematic since
 118 an error of a few degrees Celsius can increase the likelihood of misclassification
 119 of the thermal stress category, some of which span only a 6°C interval.

120 The goal of this study is to develop an improved version of the polynomial approx-
 121 imation – one that has comparable computational complexity to the existing
 122 approximation but is more robust in terms of numerical stability and substan-
 123 tially more accurate, particularly in reducing the frequency of larger errors. To

Variable name	Description	Valid Range	Normalized range
Ta	Air temperature	-50 to +50 °C	[-1, 1]
va	Wind speed at 10 m	0.5 to 30.3 m/s	[-1, 1]
$Tr - Ta$	Mean Radiant-air temperature difference	-30 to +70 °C	[-1, 1]
rH	Relative humidity	5 to 100 %	[-1, 1]
pa	Water vapour pressure	0 to 5 kPa	Not used

Table 1: Description of variables used in this study, following Bröde et al. (2012). The normalized ranges map each variable to $[-1, 1]$, with respect to the interval of validity, suitable for use with Legendre polynomial bases. Although water vapor pressure (pa) is not used directly as an input for the new approximation, it can be computed from air temperature (Ta) and relative humidity (rH), and its effect is therefore accounted for through the inclusion of rH .

124 achieve this goal, symbolic and sparse regression techniques are used as tools for
 125 interpretable and efficient function approximation. We fit the UTCI offset using
 126 sparse regression on an orthogonal Legendre polynomial basis. To emphasize this
 127 key feature and distinguish it from standard sparse regression on monomials, we
 128 refer to this approach as sparse orthogonal regression.

129 We also note that the aim was not to derive an approximation that was as
 130 accurate as possible. For example, a sufficiently complex neural-network-based
 131 model would likely provide more accurate estimates of the UTCI values. How-
 132 ever, such a model would also require the use of machine-learning libraries, as
 133 well as suitable Graphics Processing Units, to function efficiently. This means
 134 that its implementation in various programming languages and computational
 135 environments would be substantially more difficult. On the other hand, replac-
 136 ing an existing polynomial approximation with a new one is fairly straightfor-
 137 ward, meaning that implementing the new approximation into existing bi climatic
 138 software packages/libraries and numerical weather prediction systems would be
 139 relatively easy.

140 2 Methods

141 Formally, the UTCI is defined as (Bröde et al., 2012)

$$142 \quad UTCI = Ta + \text{Offset}(Ta, va, Tr, rH \text{ or } pa), \quad (1)$$

143 where Ta is the air temperature and the Offset is the physiologically equivalent
 144 temperature difference, representing how other environmental factors modify the
 145 effect of the thermal stress on the human body. The Offset function represents
 146 the deviation of the UTCI from the actual air temperature and depends on Ta ,
 147 wind speed at 10 m (va), mean radiant temperature (Tr), which accounts for
 148 the effect of all incoming radiation, and humidity, which can be represented by
 149 either relative humidity (rH) or water vapour pressure (pa).

150 The dataset provided by Bröde et al. (2012) contains accurate values of the
 151 Offset function covering a wide range of environmental states. The variables

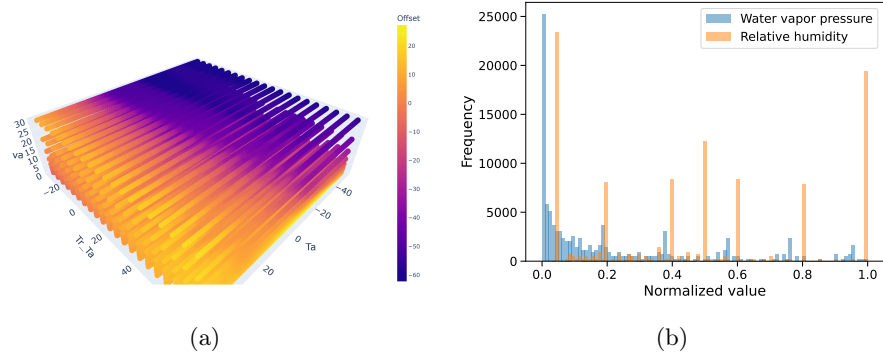


Fig. 1: (a) 3D plot of UTCI Offset Bröde et al. (2012) at 5% relative humidity, showing how wind speed (va), air temperature (Ta), and mean radiant temperature difference ($Tr - Ta$) combine to influence thermal stress. Color indicates the UTCI Offset magnitude across these environmental dimensions. (b) The different distributions of the water vapor pressure and relative humidity in the computed Offset dataset Bröde et al. (2012). The water vapor pressure is strongly peaked at zero, while the relative humidity is uniform across its range.

152 and their ranges are included in Table 1. The intervals of the environmental
 153 variables also represent the domain where the sixth-degree polynomial regression
 154 approximation is considered valid (Bröde et al., 2012). Using the approximation
 155 for conditions outside of these intervals can lead to large errors and unrealistic
 156 values of the Offset function and should be avoided (Bröde, 2021a).

157 In Fig. 1(a) we see how the UTCI Offset varies along the different environ-
 158 mental variables. Instead of the humidity (rH), the water vapor pressure (pa)
 159 can be used which is a nonlinear function of rH and the air temperature (Ta).
 160 However, the variables have different distribution, see Fig. 1(b), which impacts
 161 the extent that approximations of UTCI can generalize, discussed below.

162 Equation discovery aims to learn interpretable mathematical expressions, ei-
 163 ther differential or algebraic equations, from measurements of the variables of a
 164 given observed system (Todorovski and Džeroski, 1997). Positioned at the inter-
 165 section of symbolic machine learning and system identification, it is becoming
 166 increasingly relevant in environmental and climate science, where data-driven
 167 yet transparent models are essential (Steinmann et al., 2025; Roman, 2025b).
 168 Traditional modeling approaches rely on expert-derived formulations (Roman,
 169 2021, 2023; Roman and Bertolotti, 2022, 2023; Roman and Palmer, 2019), but
 170 the growing complexity and volume of climate data call for automated alterna-
 171 tives. Symbolic regression, which iteratively combines mathematical operators
 172 and variables to fit data, forms the core of equation discovery (Bridewell et al.,
 173 2005; Todorovski and Džeroski, 2006; Džeroski et al., 2007). Most methods em-
 174 ploy evolutionary or other (e.g., enumerative) search strategies to explore the
 175 space of candidate equations (Tanevski et al., 2016a, 2020; Mežnar et al., 2023).

176 Recent advances integrate probabilistic grammars to incorporate prior knowl-
 177 edge and constrain the search to physically meaningful expressions (Brence et al.,

206 In this work, we employ sparse regression to identify compact, interpretable
 207 models of the UTCI, emphasizing its suitability for high-dimensional input spaces
 208 with redundant or weakly relevant features. While sparse modeling is well-
 209 established in statistical learning, its application to orthogonal polynomial bases-
 210 particularly in the context of bioclimatic indices—remains unexplored. By lever-
 211 aging the structure of orthogonal polynomials, we obtain improved numerical
 212 stability and additive expansions that facilitate coefficient interpretability. To
 213 our knowledge, this is the first application of sparse regression using orthogo-
 214 nal bases to approximate the UTCI, addressing both predictive accuracy and
 215 model parsimony. Our results show that this approach surpasses the standard
 216 sixth-degree polynomial approximation in both accuracy and efficiency.

217 3 Results and discussion

218 Table 2 presents a detailed comparison of model performance across a range of
 219 polynomial degrees for both standard (non-sparse) linear regression and sparse
 220 regression techniques, evaluated in the context of approximating the UTCI. The
 221 standard approximation (Bröde et al., 2012) is a sixth-degree regression poly-
 222 nomial model with four variables, consisting of 210 terms and achieving a root
 223 mean squared loss of 1.12°C . This serves as the benchmark to be matched or
 224 improved upon. It is important to note that the standard approximation does
 225 not directly employ the relative humidity (rH), but the water vapor pressure
 226 (pa), which can be derived from the relative humidity (rH) and air temperature
 227 (Ta). As we noted above, in the dataset, the relative humidity is well represented
 228 across its entire range, see Fig. 1(b), while the water vapor pressure is strongly
 229 peaked close to zero. Optimization employing the water vapor pressure (pa) as
 230 an independent variable (instead of rH) is thus poorly conditioned and leads
 231 to instability in the regression coefficients, both in simple and sparse regression.
 232 While using the pa (instead of rH) can achieve better accuracy (lower loss),
 233 it comes at the price of losing parameter consistency across optimizations with
 234 different polynomial degrees. For this reason, we report our results employing
 235 the relative humidity (rH) instead of the water vapor pressure (pa), see Table 1.

236 The regression methods are applied to polynomial basis expansions of increas-
 237 ing degree, evaluated on the basis of root mean squared test loss and number
 238 of active parameters. Unlike many studies in the literature where models are
 239 trained on the majority of the data and evaluated on a relatively small test
 240 set, our approach inverts this paradigm: training is conducted on only 20% of
 241 the available data, while performance is assessed on the remaining 80%. Despite
 242 this stringent evaluation setting, the models achieve comparable performance on
 243 both training and test sets, underscoring their strong generalization capabilities.
 244 This performance stability is further validated through bootstrapping, which
 245 reveals minimal variance in both loss metrics and selected features across resam-
 246 pled datasets. The reported performance metrics—such as train/test loss and
 247 number of parameters—remain stable when the model training and evaluation
 248 process is repeated on multiple random re-samplings (bootstrapped subsets) of

Method	Polynomial degree						
	4th	6th	8th	10th	12th	14th	16th
Standard	1.12 (210)						
Linear regression	2.1 (70)	1.3 (210)	0.92 (495)	Train: 0.67 Test: 0.71 (1001)	0.54 0.62 (1820)	0.44 0.66 (3060)	0.36 1.74 (4845)
Sparse orthogonal regression	2.1 (65)	1.38 (124)	1.03 (176)	0.88 (209)	0.69 (355)	0.63 (400)	0.6 (424)

Table 2: Root mean squared train loss [°C], test loss [°C] and the number of parameters (shown in parenthesis) in approximating the UTCI Offset. The baseline reference, labeled as “Standard,” corresponds to the sixth-degree regression polynomial model with four variables (Bröde et al., 2012). Unless otherwise stated the test loss equals the train loss. Where two loss values are reported (train loss on the top and test loss below), they indicate a notable train-test discrepancy, typically suggesting overfitting. Training is done with 20% of the data and testing is performed with 80%. Results are robust under bootstrapping.

249 the data. This suggests that the results are not sensitive to specific data splits and
 250 that the models generalize well across different subsets of the dataset, indicating
 251 reliability and consistency in the reported findings. These findings demonstrate
 252 the robustness and reliability of the proposed framework.

253 To make the fitted model class explicit, let \tilde{T}_a , \tilde{v}_a , $\widetilde{\Delta T}_r$, and \widetilde{rH} denote the
 254 normalized versions of T_a , v_a , $T_r - T_a$, and rH , respectively, each mapped to the
 255 interval $[-1, 1]$ according to the ranges in Table 1. In this formulation, relative
 256 humidity is retained as an input variable in order to account for the effect of
 257 water vapor. The approximation of the UTCI offset can then be written in the
 258 general form

$$259 \quad \widehat{\text{Offset}}(T_a, v_a, T_r - T_a, rH) = \sum_{\alpha \in \mathcal{A}_p} c_{\alpha} \prod_{j=1}^4 P_{\alpha_j}(x_j), \quad (2)$$

260 where $P_n(\cdot)$ denotes the Legendre polynomial of degree n , $(x_1, x_2, x_3, x_4) =$
 261 $(\tilde{T}_a, \tilde{v}_a, \widetilde{\Delta T}_r, \widetilde{rH})$, $\alpha = (\alpha_1, \alpha_2, \alpha_3, \alpha_4)$ is a multi-index, and

$$262 \quad \mathcal{A}_p = \{\alpha \in \mathbb{N}_0^4 : \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 \leq p\}$$

263 is the set of all basis terms up to total polynomial degree p . Thus, the model is a
 264 linear combination of products of Legendre polynomials in the four normalized
 265 environmental variables. For a given maximum degree p , the full candidate basis
 266 contains $\binom{p+4}{4}$ terms, which yields the sequence 70, 210, 495, ... reported in Ta-
 267 ble 2 for degrees 4, 6, 8, ... Sparse orthogonal regression restricts this expansion

268 by retaining only a subset of the candidate terms,

$$269 \quad \widehat{\text{Offset}}(T_a, v_a, T_r - T_a, rH) = \sum_{\alpha \in S_p} c_\alpha \prod_{j=1}^4 P_{\alpha_j}(x_j), \quad (3)$$

270 where $S_p \subseteq \mathcal{A}_p$ is selected by the Lasso regularization. The number of active
 271 parameters therefore depends on two factors: the maximum polynomial degree,
 272 which determines the size of the candidate pool, and the regularization strength,
 273 which determines how many of those candidate terms are retained in the final
 274 model. This is the reason why the number of parameters changes across poly-
 275 nomial degrees and also along the Pareto fronts shown in Fig. 3. In this sense,
 276 the approximation can be viewed as a Fourier-like decomposition in an orthogo-
 277 nal polynomial basis, where lower-order terms capture the dominant structure of
 278 the UTCI offset and higher-order terms provide progressively finer corrections. A
 279 key advantage of the orthogonal basis is that it yields order-by-order consistency,
 280 see Fig. 4: when higher-degree terms are introduced, the coefficients associated
 281 with lower-order structure remain much more stable than in regressions based
 282 on ordinary monomials.

283 Linear regression without any sparsity constraints shows improved perfor-
 284 mance at higher degrees, with test loss reducing as model capacity increases.
 285 However, this comes with a dramatic increase in the number of parameters; it
 286 reaches over 1800 coefficients by degree 12. Furthermore, the discrepancy be-
 287 tween train and test losses at higher degrees (e.g., 0.62°C vs. 0.54°C at degree
 288 12) indicates overfitting, despite the improved predictive accuracy. The resulting
 289 models are also substantially more complex, raising concerns regarding inter-
 290 pretation and generalization. Sparse regression with standard polynomial bases
 291 show similar performance at low degrees but fails to converge beyond the 6th de-
 292 gree. This indicates that enforcing sparsity in a poorly conditioned basis becomes
 293 increasingly difficult as model complexity grows.

294 In contrast, sparse regression using an orthogonal Legendre basis (or sparse
 295 orthogonal regression) exhibits superior stability and accuracy across all degrees.
 296 It outperforms the baseline 6th-degree polynomial fit from degree 8th onward,
 297 achieving a test loss of 0.88°C at degree 10 with only 209 parameters—almost the
 298 same count as the original benchmark model, but with improved generalization.
 299 As the degree increases to 16, the loss reduces further to 0.60°C using 424 param-
 300 eters—a fraction of those used by the corresponding standard regression model.
 301 The orthogonality of the Legendre basis likely contributes to better numerical
 302 conditioning, facilitating sparse model discovery even at high degrees. These re-
 303 sults emphasize the importance of basis selection and regularization strategy in
 304 symbolic regression tasks. Sparse methods, when combined with well-structured
 305 bases like Legendre polynomials, offer a promising path toward accurate, com-
 306 pact, and interpretable models in high-dimensional settings.

307 Furthermore, optimization of nonlinear objective functions using gradient-
 308 based algorithms can be computationally intensive, especially in high-dimensional
 309 spaces where convergence is slow and local minima may hinder performance.
 310 In contrast, the regression-based approach proposed in this article—particularly

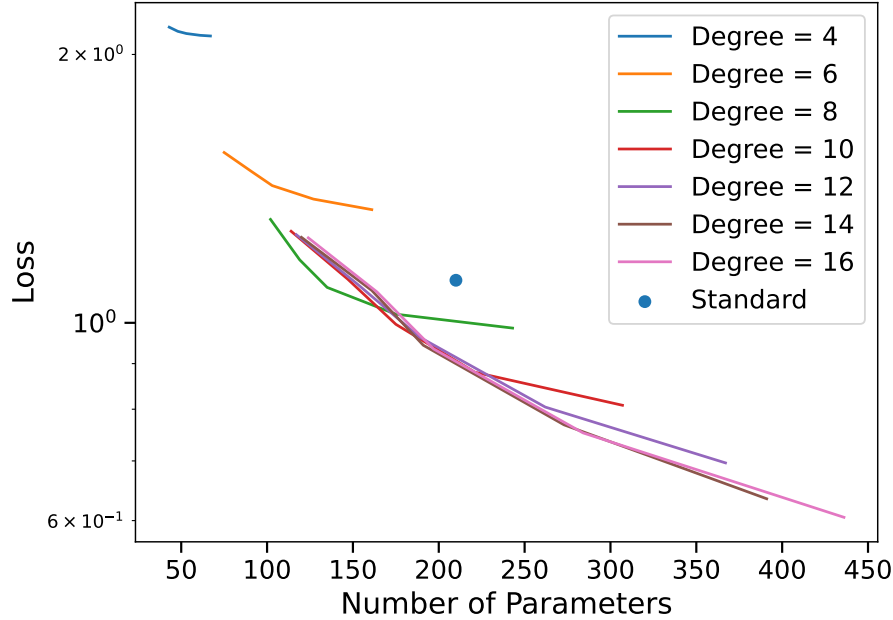


Fig. 3: Loss versus number of parameters for different polynomial degrees. The regularization parameter was varied in the lasso regression to yield a Pareto front in model accuracy and complexity for each degree.

311 through sparse regression with orthogonal polynomials—offers significantly faster
 312 computation. By framing the problem as a structured regression task rather
 313 than a nonlinear optimization, the method avoids costly iterative procedures
 314 and scales efficiently with dimensionality, making it highly suitable for rapid
 315 modeling of complex environmental indices like the UTCI.

316 Fig. 3 illustrates the relationship between model complexity (measured by the
 317 number of parameters) and prediction accuracy (log-scaled loss) for sparse re-
 318 gression models using Legendre polynomial bases of varying degrees. Each curve
 319 corresponds to a fixed polynomial degree, ranging from 4 to 16, with points re-
 320 flecting models of increasing complexity obtained through regularization. A clear
 321 trend is observed: for a given polynomial degree, increasing the number of pa-
 322 rameters generally results in improved model accuracy (i.e., lower loss). However,
 323 diminishing returns set in, and the rate of improvement flattens. More notably,
 324 the envelope formed by the lowest loss at each level of complexity across all de-
 325 grees traces an emergent Pareto front (Smits and Kotanchek, 2005). This front
 326 captures the trade-off between model simplicity and predictive performance.

327 Higher-degree models (e.g., degrees 12–16) dominate this frontier at higher
 328 parameter counts, offering better loss with only marginal increases in complexity.
 329 In contrast, lower-degree models saturate quickly, highlighting their limited ex-
 330 pressivity. The Pareto front thus reflects the optimal set of models that balance

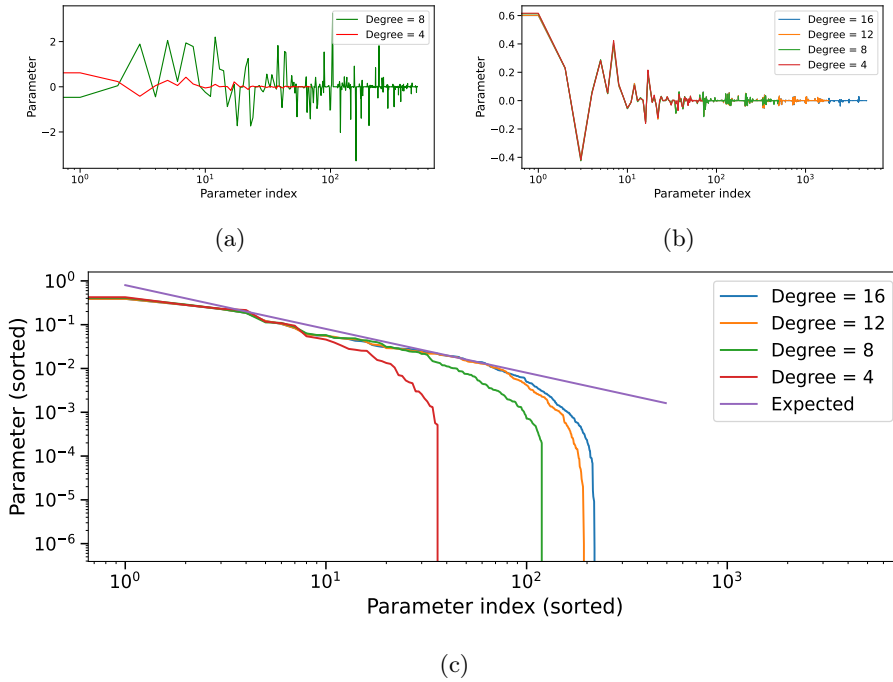


Fig. 4: Parameters (or polynomial coefficients) and how they change for different polynomial degrees for (a) simple regression and (b) sparse regression (using Legendre basis). (c) Sorted sparse-regression coefficients (Legendre basis) versus parameter index on a logarithmic x-axis show a clear, Fourier-like decay with order—approximately $1/n$ —that is stable across model capacities (degrees 4, 8, 12, 16), indicating a hierarchical structure where lower-order terms dominate and higher-order terms provide incremental refinement.

331 accuracy and sparsity, guiding model selection under complexity constraints. The
 332 use of Legendre polynomials ensures numerical stability and encourages efficient
 333 basis representations, which supports the recovery of compact yet accurate mod-
 334 els in this sparse regression setting.

335 The Fig. 4(a) and (b) we visualize the behavior of regression coefficients obtained
 336 from simple regression and sparse regression with orthogonal Legendre
 337 polynomials. Both plots use a logarithmic x-axis to indicate the parameter index
 338 and reveal how coefficients evolve as higher-degree polynomial terms are intro-
 339 duced. In Fig. 4(a), each line corresponds to simple regression solutions using
 340 polynomial bases of increasing degree. The x-axis denotes the index of poly-
 341 nomial terms (sorted or sequential), while the y-axis shows the corresponding
 342 coefficient values. A key observation is that the coefficients of lower-degree terms
 343 (left side of the plot) are not stable across model orders. As higher-degree terms
 344 are added, previously estimated lower-order coefficients shift significantly, often
 345 changing sign and magnitude.

346 Figure 4(b) presents coefficient values for sparse regression using Legendre
 347 polynomials, with colors indicating contributions from different polynomial de-
 348 grees. Here, a contrasting pattern emerges: coefficients associated with lower-
 349 degree terms remain stable as higher-degree terms are added. New coefficients
 350 primarily emerge in the higher-order region of the x-axis, without disturbing
 351 the existing ones. This stability results from the orthogonality of the Legendre
 352 basis, which decorrelates the polynomial terms and enables additive refinement
 353 without re-tuning existing coefficients.

354 The contrast between the Figs. 4(a) and (b) underscores the advantage of
 355 orthogonal polynomial bases in sparse regression. Simple regression results in
 356 unstable, entangled coefficient estimates that shift with basis expansion, com-
 357 plicating interpretability and reuse. Sparse regression with ordinary polynomial
 358 bases fails to converge for higher degrees. In contrast, sparsity and orthogo-
 359 nal polynomials yield stable, hierarchical models where lower-order structure is
 360 preserved and higher-order terms incrementally enrich the representation. This
 361 behavior is particularly valuable for symbolic regression and interpretable mod-
 362 eling, where each term ideally reflects a distinct, meaningful contribution to the
 363 model output.

364 In Fourier analysis, the magnitude of coefficients typically decays as $1/n$
 365 (where n is the order of the term) for functions of bounded variation (Stein and
 366 Shakarchi, 2011) – a class that includes many naturally occurring signals and
 367 is a reasonable assumption for observational data. This decay reflects the fact
 368 that higher-order (or higher-frequency) components contribute less to the overall
 369 structure of such functions. A similar trend is observed in sparse regression
 370 using orthogonal polynomial bases, see Fig. 4(c). When coefficients are sorted
 371 by magnitude, they exhibit a clear decreasing pattern, analogous to the Fourier
 372 case, with lower-order terms capturing the dominant structure and higher-order
 373 terms refining the approximation in a controlled manner.

374 This suggests that through the use of sparse regression with an orthogonal
 375 polynomial basis, we have achieved a Fourier-like decomposition of the UTCI Off-
 376 set in the Legendre basis (instead of the trigonometric one). This has a number
 377 of theoretical advantages: due to the orthogonality of the basis functions, the de-
 378 composition minimizes the L_2 distance (least squares) between approximation
 379 and function, guaranteeing the best possible polynomial fit for a given model
 380 complexity (Stein and Shakarchi, 2011). Additionally, the coefficients are un-
 381 correlated and hierarchically structured, ensuring that lower-order components
 382 remain stable as higher-order terms are added—enhancing both interpretability
 383 and numerical robustness.

384 Based on the analysis results and one of the initial goals (that the new ap-
 385 proximation should have comparable computational complexity to the existing
 386 one), we selected the sparse regression model based on tenth-degree Legendre
 387 polynomials as the most suitable approximation. The final version of the new
 388 polynomial, which has 209 coefficients, was calculated using the whole dataset
 389 of tabulated values.

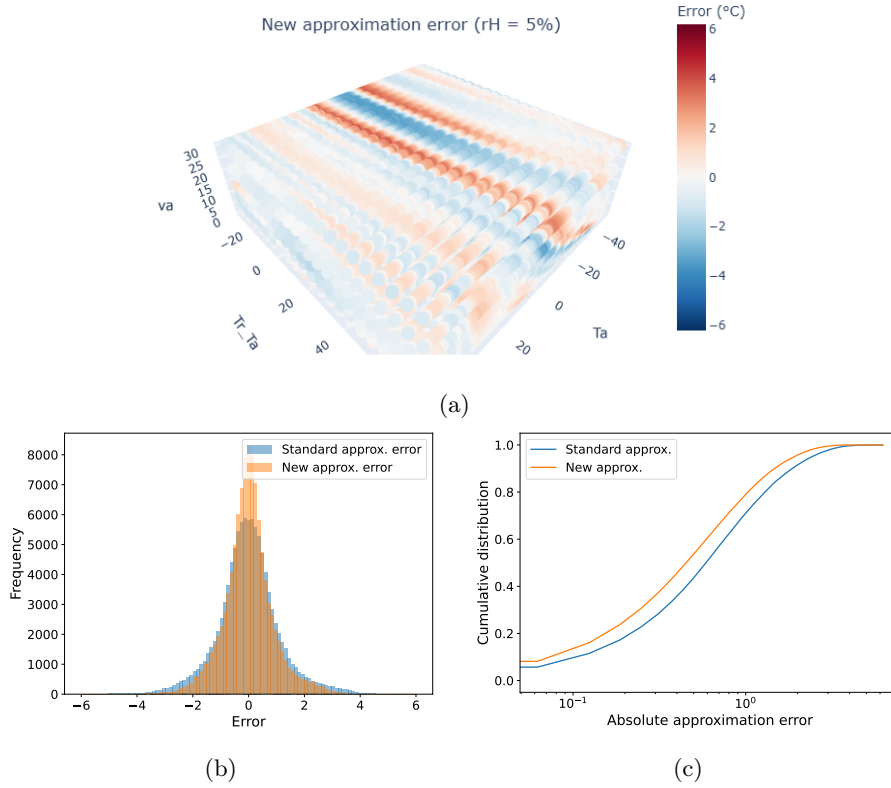


Fig. 5: (a) Spatial distribution of the UTCI Offset error (approximation minus reference) for the new sparse-model-based approximation at a fixed relative humidity of 5%, showing small, smoothly varying discrepancies. (b) Comparison of error histograms for the standard UTCI approximation and the new approximation based on the tenth-degree Legendre polynomials. (c) Cumulative distributions of the absolute errors of the two approximations.

390 Fig. 5(a) shows the spatial distribution of the Offset errors for the new ap-
 391 proximation at a fixed relative humidity of 5%. The errors are small and smoothly
 392 varying, indicating good agreement across the input space. Fig. 5(b) presents
 393 a comparison of error histograms for both the standard and new approxima-
 394 tions. The sparse-model-based approximation produces a narrower, more sharply
 395 peaked distribution centered at zero, highlighting a reduction in error variance
 396 and suggesting better generalization. Fig. 5(c) shows the cumulative distribution
 397 of absolute errors for the two approximations. The curve for the new approxi-
 398 mation rises more steeply and reaches higher cumulative values at lower error
 399 thresholds, indicating that a larger proportion of predictions fall within smaller
 400 error margins.

401 Table 3 summarizes the most relevant properties of the two approximations.
 402 The results show a clear improvement in accuracy: the new approximation not

	Standard approximation	New approximation
Polynomial degree	6th	10th
Basis functions	monomials	Legendre
Number of coefficients	210	209
Mean Error	$1.7 \cdot 10^{-3} \text{ }^\circ\text{C}$ (0.35 $^\circ\text{C}$)	$-2.7 \cdot 10^{-15} \text{ }^\circ\text{C}$ (0.22 $^\circ\text{C}$)
Mean Absolute Error	0.81 $^\circ\text{C}$ (1.33 $^\circ\text{C}$)	0.64 $^\circ\text{C}$ (0.71 $^\circ\text{C}$)
Root Mean Square Error	1.17 $^\circ\text{C}$ (2.77 $^\circ\text{C}$)	0.88 $^\circ\text{C}$ (0.96 $^\circ\text{C}$)
Freq. of abs. errors larger than 2 $^\circ\text{C}$	8.4 % (15.5 %)	4.2 % (5.0 %)
Freq. of abs. errors larger than 3 $^\circ\text{C}$	2.2 % (6.3 %)	0.50 % (0.60 %)
Freq. of abs. errors larger than 4 $^\circ\text{C}$	0.34 % (3.8 %)	0.011 % (0.10 %)
Freq. of abs. errors larger than 5 $^\circ\text{C}$	0.038 % (3.3 %)	0.00096 % (0 %)

Table 3: Comparison of properties of the standard (Bröde et al., 2012) and new polynomial approximations of UTCI Offset function. The values outside of the parentheses reflect the evaluation of the approximations on the full dataset of 104 643 accurate Offset values provided by Bröde et al. (2012). The values shown in the parentheses reflect the evaluation using the independent dataset of 1000 accurate UTCI values (Bröde, 2021b), which were not used during the development of the new approximation. Both approximations are only valid for the intervals of environmental variables available in the full dataset (Table 1).

403 only substantially reduces the average errors (i.e, the mean error, the mean
404 absolute error, and the root mean square error) but also drastically reduces
405 the frequency of large deviations compared to the standard approximation. For
406 example, the frequency of absolute errors larger than 2 $^\circ\text{C}$ is halved from 8% to
407 4%, the frequency of errors larger than 3 $^\circ\text{C}$ reduces from 2% to 0.5%, while the
408 frequency of errors larger than 4 $^\circ\text{C}$ reduces from 0.3% to 0.01%. These results
409 clearly show the added benefits of the new approximation and confirm that
410 the sparse regression approach can achieve comparable or improved predictive
411 accuracy while maintaining interpretability and model parsimony.

412 We also evaluated the new approximation on the independent dataset of 1000
413 accurate UTCI values, which were not used during the development of the ap-
414 proximation. This dataset was prepared by the authors of the Bröde et al. (2012)
415 paper, and is freely available on a Zenodo repository (Bröde, 2021b). Similarly
416 to the evaluation of the new approximation on the full dataset, evaluation on
417 the independent dataset shows a substantial reduction of the mean errors and
418 a drastic reduction in the frequency of large errors compared to the standard
419 approximation (Table 3).

420 Since the new approximation was determined using the full dataset of accu-
421 rate Offset values (Bröde et al., 2012), it is, same as the standard approximation,
422 only valid for the intervals of environmental variables available in this dataset
423 (Table 1). Using the approximation for conditions outside of these intervals can
424 potentially lead to large errors or unrealistic results and should be avoided.

425 4 Conclusions

426 The goal of this study was to develop an improved version of the polynomial
 427 approximation – one that would have comparable computational complexity to
 428 the existing approximation but would be more robust in terms of numerical
 429 stability and substantially more accurate, particularly in reducing the frequency
 430 of larger errors. This goal was successfully achieved using sparse regression with
 431 an orthogonal polynomial basis.

432 Sparse regression methods, such as LASSO, helped reduce overfitting and im-
 433 prove interpretability. As we have shown, the choice of basis functions is crucial:
 434 orthogonal polynomials like Legendre polynomials offer better numerical stabil-
 435 ity and conditioning than monomials. They enable hierarchical models where
 436 higher-order terms don't affect lower-order estimates, making them especially
 437 useful in sparse, interpretable models. Empirical results support these theoretic-
 438 al advantages.

439 Using sparse regression with an orthogonal polynomial basis (or sparse or-
 440 thogonal regression), we have:

- 441 (a) Achieved substantially better accuracy – compared to the standard approx-
 442 imation, the new approximation not only substantially reduces the average
 443 errors (i.e, the mean error, the mean absolute error, and the root mean square
 444 error) but also drastically reduces the frequency of large errors.
- 445 (b) Retained a comparable computational complexity – the number of coeffi-
 446 cients is almost the same for both approximations, meaning the computa-
 447 tional complexity is comparable.
- 448 (c) Found a Pareto front for different model complexities – loss curves reveal
 449 that sparse models with orthogonal bases efficiently populate a Pareto front,
 450 balancing complexity and accuracy.
- 451 (d) Determined coefficients consistent over models with different capacities - co-
 452 efficient plots for models built on orthogonal bases show the progressive in-
 453 clusion of higher-order components without disrupting lower-order structure,
 454 in contrast to models using simple regression and ordinary polynomials.
- 455 (e) Achieved successful generalization – training the model over only 20% of the
 456 data, while testing was performed over the other 80%, highlights successful
 457 generalization. The results are also robust under bootstrapping.
- 458 (f) Essentially decomposed the UTCI in a Fourier expansion with a Legendre-
 459 polynomial basis, with parameters scaling as expected. Thus, we are arguably
 460 close to the theoretical optimum results for a robust approximation in the
 461 L_2 metric (or least squares).

462 Sparse orthogonal regression provides an effective framework for constructing
 463 accurate and numerically stable polynomial approximations of the UTCI. Our
 464 main contribution is therefore not methodological novelty in sparse regression it-
 465 self, but the use of an orthogonal polynomial basis as a practical approximation
 466 strategy with favorable numerical properties, including order-by-order consis-
 467 tency and stable low-order truncations. In addition, the results obtained from
 468 random train–test splits, together with their robustness under bootstrapping,

469 show that using only 20% of the data for training is not a requirement of the
470 method, but a deliberately stringent test of generalization. The comparable per-
471 formance on the remaining 80% of the data indicates that the approach remains
472 accurate, robust, and efficient even under a severe limitation in the number of
473 training data points, while remaining well suited for practical applications that
474 require portability and ease of implementation.

475 We have also prepared an easy-to-use Python function for the new approxi-
476 mation (please refer to the Code and data availability section on how to obtain
477 the code). The code relies only on basic mathematical operations, which makes
478 it easy to adapt to other programming languages, such as Fortran or C++. We
479 also implemented a check to see if the environmental state falls within the do-
480 main of validity of the approximation. If this is not the case, the code produces a
481 warning that the resulting UTCI values could have large errors or be unrealistic.

482 **Funding**

483 This publication is supported by the European Union’s Horizon Europe re-
484 search and innovation programme under the Marie Skłodowska-Curie Postdoc-
485 toral Fellowship Programme, SMASH co-funded under the grant agreement
486 No. 101081355. The operation (SMASH project) is co-funded by the Republic
487 of Slovenia and the European Union from the European Regional Development
488 Fund. The authors acknowledge the financial support of the Slovenian Research
489 Agency via the Gravity project *AI for Science*, GC-0001 and of the Slovenian
490 Research And Innovation Agency (research core funding No. P1-0188).

491 **Author Contributions**

492 S.R. - Conceptualization, Data curation, Formal analysis, Investigation, Method-
493 ology, Software, Validation, Visualization, Writing (original draft preparation),
494 G.S. - Conceptualization, Resources, Validation, Software, Writing (review and
495 editing), L.T. - Conceptualization, Methodology, Project administration, Super-
496 vision, Validation, Writing (review and editing), S.D. - Conceptualization, Fund-
497 ing acquisition, Project administration, Resources, Supervision, Writing (review
498 and editing).

499 **Code and data availability**

500 The Python function code for the calculation of the new approximation as well
501 as all the code needed to reproduce the results of the analysis is published on
502 Zenodo (Roman, 2025a).

503 **Conflict of Interest Statement**

504 The authors declare no conflicts of interest.

Bibliography

- 506 Atanasova, N., Recknagel, F., Todorovski, L., Džeroski, S., and Kompare, B.:
507 Computational assemblage of ordinary differential equations for chlorophyll-a
508 using a lake process equation library and measured data of Lake Kasumi-
509 gaura, *Ecological Informatics: Scope, Techniques and Applications*, pp. 409–
510 427, 2006a.
- 511 Atanasova, N., Todorovski, L., Džeroski, S., Remec, Š. R., Recknagel, F., and
512 Kompare, B.: Automated modelling of a food web in lake Bled using measured
513 data and a library of domain knowledge, *Ecological Modelling*, 194, 37–48,
514 2006b.
- 515 Atanasova, N., Todorovski, L., Džeroski, S., and Kompare, B.: Application of
516 automated model discovery from data and expert knowledge to a real-world
517 domain: Lake Glumsø, *ecological Modelling*, 212, 92–98, 2008.
- 518 Atanasova, N., Džeroski, S., Kompare, B., Todorovski, L., and Gal, G.: Auto-
519 mated discovery of a model for dinoflagellate dynamics, *Environmental Mod-
520 elling & Software*, 26, 658–668, 2011.
- 521 Blazejczyk, K., Epstein, Y., Jendritzky, G., Staiger, H., and Tinz, B.: Compar-
522 ison of UTCI to selected thermal indices, *International Journal of Biometeo-
523 rology*, 56, 515–535, <https://doi.org/10.1007/s00484-011-0453-2>, 2012.
- 524 Brence, J., Todorovski, L., and Džeroski, S.: Probabilistic grammars for equation
525 discovery, arXiv preprint arXiv:2012.00428, 2020.
- 526 Brence, J., Džeroski, S., and Todorovski, L.: Dimensionally-consistent equation
527 discovery through probabilistic attribute grammars, *Information Sciences*, 632,
528 742–756, 2023.
- 529 Bridewell, W., Asadi, N. B., Langley, P., and Todorovski, L.: Reducing over-
530 fitting in process model induction, in: *Proceedings of the 22nd international
531 conference on Machine learning*, pp. 81–88, 2005.
- 532 Brimicombe, C., Napoli, C. D., Quintino, T., Pappenberger, F., Cornforth, R.,
533 and Cloke, H. L.: Thermofeel: A python thermal comfort indices library,
534 *SoftwareX*, 18, 101005, <https://doi.org/10.1016/j.softx.2022.101005>,
535 2022.
- 536 Bröde, P.: Issues in UTCI Calculation from a Decade’s Experience, pp. 13–
537 21, Springer International Publishing, Cham, [https://doi.org/10.1007/
538 978-3-030-76716-7_2](https://doi.org/10.1007/978-3-030-76716-7_2), 2021a.
- 539 Bröde, P.: UTCI-Test-Data, <https://doi.org/10.5281/zenodo.5503967>,
540 2021b.
- 541 Bröde, P., Fiala, D., Błażejczyk, K., Holmér, I., Jendritzky, G., Kampmann,
542 B., Tinz, B., and Havenith, G.: Deriving the operational procedure for the
543 Universal Thermal Climate Index (UTCI), *International Journal of Biometeo-
544 rology*, 56, 481–494, <https://doi.org/10.1007/s00484-011-0454-1>, URL
545 <http://link.springer.com/10.1007/s00484-011-0454-1>, 2012.

- 546 Brunton, S. L., Proctor, J. L., and Kutz, J. N.: Discovering governing equations
547 from data by sparse identification of nonlinear dynamical systems, *Proceedings*
548 *of the national academy of sciences*, 113, 3932–3937, 2016.
- 549 Błażejczyk, K.: BioKlima - Universal tool for bioclimatic and thermophysiological
550 studies, URL <https://www.igipz.pan.pl/bioklima-crd.html>, [Accessed
551 date: 10.10.2025.], 2025.
- 552 Błażejczyk, K. and Kuchcik, M.: UTCI applications in practice (methodological
553 questions), *Geographia Polonica*, 94, <https://doi.org/10.7163/GPol.0198>,
554 2021.
- 555 Čerepnalkoski, D., Taškova, K., Todorovski, L., Atanasova, N., and Džeroski, S.:
556 The influence of parameter fitting methods on model structure selection in
557 automated modeling of aquatic ecosystems, *Ecological Modelling*, 245, 136–
558 165, 2012.
- 559 Di Napoli, C., Messeri, A., Novák, M., Rio, J., Wieczorek, J., Morabito, M.,
560 Silva, P., Crisci, A., and Pappenberger, F.: The Universal Thermal Climate
561 Index as an Operational Forecasting Tool of Human Biometeorological Con-
562 ditions in Europe, in: *Applications of the Universal Thermal Climate In-*
563 *dex UTCI in Biometeorology*, pp. 193–208, Springer International Publishing,
564 Cham, https://doi.org/10.1007/978-3-030-76716-7_10, URL [https://](https://link.springer.com/10.1007/978-3-030-76716-7_10)
565 link.springer.com/10.1007/978-3-030-76716-7_10, 2021a.
- 566 Di Napoli, C., Barnard, C., Prudhomme, C., Cloke, H. L., and Pappenberger, F.:
567 ERA5-HEAT: A global gridded historical dataset of human thermal comfort
568 indices from climate reanalysis, *Geoscience Data Journal*, 8, [https://doi.](https://doi.org/10.1002/gdj3.102)
569 [org/10.1002/gdj3.102](https://doi.org/10.1002/gdj3.102), 2021b.
- 570 Džeroski, S., Langley, P., and Todorovski, L.: Computational discovery of sci-
571 entific knowledge, in: *Computational discovery of scientific knowledge: Intro-*
572 *duction, techniques, and applications in environmental and life sciences*, pp.
573 1–14, Springer, 2007.
- 574 Fiala, D., Havenith, G., Brode, P., Kampmann, B., et al.: UTCI- Fiala multi-
575 node model of human heat transfer and temperature regulation, *Int J Biomete-*
576 *orol.*, 56, <https://doi.org/10.1007/s00484-011-0424-7>, 2012.
- 577 Jeraž, M., Džeroski, S., Todorovski, L., and Debeljak, M.: Application of machine
578 learning methods to palaeoecological data, *Ecological modelling*, 191, 159–169,
579 2006.
- 580 Kuzmanović, D., Banko, J., and Skok, G.: Improving the operational forecasts
581 of outdoor Universal Thermal Climate Index with post-processing, *Internation-*
582 *al Journal of Biometeorology*, 68, 965–977, [https://doi.org/10.1007/](https://doi.org/10.1007/s00484-024-02640-6)
583 [s00484-024-02640-6](https://doi.org/10.1007/s00484-024-02640-6), 2024.
- 584 Mežnar, S., Džeroski, S., and Todorovski, L.: Efficient generator of mathematical
585 expressions for symbolic regression, *Machine Learning*, 112, 4563–4596, 2023.
- 586 Omejc, N., Gec, B., Brence, J., Todorovski, L., and Džeroski, S.: Probabilistic
587 grammars for modeling dynamical systems from coarse, noisy, and partial
588 data, *Machine Learning*, 113, 7689–7721, 2024.
- 589 Pappenberger, F., Jendritzky, G., Staiger, H., Dutra, E., Di Giuseppe, F.,
590 Richardson, D. S., and Cloke, H. L.: Global forecasting of thermal health haz-
591 ards: the skill of probabilistic predictions of the Universal Thermal Climate In-

- 592 dex (UTCI), *International Journal of Biometeorology*, 59, 311–323, <https://doi.org/10.1007/s00484-014-0843-3>, URL <http://link.springer.com/10.1007/s00484-014-0843-3>, 2015.
- 593
594
- 595 Reid, S., Tibshirani, R., and Friedman, J.: A study of error variance estimation
596 in lasso regression, *Statistica Sinica*, pp. 35–67, 2016.
- 597 Roman, S.: Historical dynamics of the Chinese dynasties, *Heliyon*, 7, 2021.
- 598 Roman, S.: Theories and models: Understanding and Predicting Societal Col-
599 lapse, in: *The Era of Global Risk: An Introduction to Existential Risk Studies*,
600 pp. 27–54, Open Book Publishers, URL <https://doi.org/10.11647/OBP.0336.02>, 2023.
- 601
- 602 Roman, S.: Code for Approximating the universal thermal climate index (UTCI)
603 using sparse regression with orthogonal polynomials, <https://doi.org/10.5281/zenodo.16880382>, 2025a.
- 604
- 605 Roman, S.: Maximum Entropy Models for Unimodal Time Series: Case Studies of
606 Universe 25 and St. Matthew Island, in: *International Conference on Discovery*
607 *Science*, pp. 32–44, Springer, 2025b.
- 608 Roman, S. and Bertolotti, F.: A master equation for power laws, *Royal Society*
609 *open science*, 9, 220 531, 2022.
- 610 Roman, S. and Bertolotti, F.: Global history, the emergence of chaos and in-
611 ducing sustainability in networks of socio-ecological systems, *Plos one*, 18,
612 e0293 391, 2023.
- 613 Roman, S. and Palmer, E.: The Growth and Decline of the Western Roman
614 Empire: Quantifying the Dynamics of Army Size, Territory, and Coinage,
615 *Cliodynamics*, 10, 2019.
- 616 Simidjievski, N., Todorovski, L., and Džeroski, S.: Learning ensembles of popula-
617 tion dynamics models and their application to modelling aquatic ecosystems,
618 *Ecological Modelling*, 306, 305–317, 2015.
- 619 Simidjievski, N., Todorovski, L., and Džeroski, S.: Modeling dynamic systems
620 with efficient ensembles of process-based models, *PloS one*, 11, e0153 507, 2016.
- 621 Smits, G. F. and Kotanchek, M.: Pareto-front exploitation in symbolic regression,
622 in: *Genetic programming theory and practice II*, pp. 283–299, Springer, 2005.
- 623 Stein, E. M. and Shakarchi, R.: *Fourier analysis: an introduction*, vol. 1, Prince-
624 ton University Press, 2011.
- 625 Steinmann, P., Verstegen, J., Van Voorn, G., Roman, S., and Ligtenberg, A.:
626 Scenario search: finding diverse, plausible and comprehensive scenario sets for
627 complex systems, *Socio-Environmental Systems Modelling*, 7, 18 823–18 823,
628 2025.
- 629 Tanevski, J., Todorovski, L., and Džeroski, S.: Learning stochastic process-based
630 models of dynamical systems from knowledge and data, *BMC systems biology*,
631 10, 1–17, 2016a.
- 632 Tanevski, J., Todorovski, L., and Džeroski, S.: Process-based design of dynamical
633 biological systems, *Scientific reports*, 6, 34 107, 2016b.
- 634 Tanevski, J., Todorovski, L., and Džeroski, S.: Combinatorial search for select-
635 ing the structure of models of dynamical systems with equation discovery,
636 *Engineering Applications of Artificial Intelligence*, 89, 103 423, 2020.

- 637 Tartarini, F. and Schiavon, S.: pythermalcomfort: A Python package for ther-
638 mal comfort research, *SoftwareX*, 12, 100578, [https://doi.org/10.1016/j.
639 softx.2020.100578](https://doi.org/10.1016/j.softx.2020.100578), 2020.
- 640 Termonia, P., Fischer, C., Bazile, E., Bouyssel, F., Brozkova, R., Bénard, P.,
641 Bochenek, B., Degrauwe, D., Derková, M., Khatib, R., Hamdi, R., Mašek, J.,
642 Pottier, P., Pristov, N., Seity, Y., Smolikova, P., Španiel, O., Tudor, M., Wang,
643 Y., and Joly, A.: The ALADIN System and its canonical model configurations
644 AROME CY41T1 and ALARO CY40T1, *Geoscientific Model Development*,
645 11, <https://doi.org/10.5194/gmd-11-257-2018>, 2018.
- 646 Todorovski, L. and Džeroski, S.: Declarative bias in equation discovery, in: *Pro-
647 ceedings of the International Conference on Machine Learning*, pp. 376–384,
648 1997.
- 649 Todorovski, L. and Džeroski, S.: Theory revision in equation discovery, in: *In-
650 ternational Conference on Discovery Science*, pp. 389–400, Springer, 2001.
- 651 Todorovski, L. and Džeroski, S.: Integrating knowledge-driven and data-driven
652 approaches to modeling, *Ecological Modelling*, 194, 3–13, 2006.
- 653 Todorovski, L., Džeroski, S., and Kompare, B.: Modelling and prediction of phy-
654 toplankton growth with equation discovery, *Ecological Modelling*, 113, 71–81,
655 1998.