

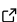
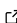
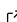
1 scores: A Python package for verifying and evaluating
2 models and predictions with xarray and pandas

3 Tennessee Leeuwenburg ¹¶, Nicholas Loveday ¹, Elizabeth E. Ebert¹,
4 Harrison Cook ¹, Mohammadreza Khanarmuei ¹, Robert J. Taggart ¹,
5 Nikeeth Ramanathan ¹, Maree Carroll ¹, Stephanie Chong ², Aidan
6 Griffiths³, and John Sharples¹

7 1 Bureau of Meteorology, Australia 2 Independent Contributor, Australia 3 Work undertaken while at the
8 Bureau of Meteorology, Australia ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#)).

9 Summary

10 scores is a Python package containing mathematical functions for the verification, evaluation
11 and optimisation of forecasts, predictions or models. It primarily supports the geoscience
12 communities; in particular, the meteorological, climatological and oceanographic communities.
13 In addition to supporting the Earth system science communities, it also has wide potential
14 application in machine learning and other domains such as economics.

15 scores not only includes common scores (e.g. Mean Absolute Error), it also includes novel
16 scores not commonly found elsewhere (e.g. Fixed Risk Multicategorical (FIRM) score, Flip-Flop
17 Index), complex scores (e.g. threshold-weighted continuous ranked probability score), and
18 statistical tests (such as the Diebold Mariano test). It also contains isotonic regression which
19 is becoming an increasingly important tool in forecast verification and can be used to generate
20 stable reliability diagrams. Additionally, it provides pre-processing tools for preparing data for
21 scores in a variety of formats including cumulative distribution functions (CDF). At the time
22 of writing, scores includes over 50 metrics, statistical techniques and data processing tools.

23 All of the scores and statistical techniques in this package have undergone a thorough scientific
24 and software review. Every score has a companion Jupyter Notebook tutorial that demonstrates
25 its use in practice.

26 scores primarily supports xarray datatypes for Earth system data, allowing it to work with
27 NetCDF4, HDF5, Zarr and GRIB data sources among others. scores uses Dask for scaling
28 and performance. It has expanding support for pandas.

29 The software repository can be found at <https://github.com/nci/scores/>.

30 Statement of Need

31 The purpose of this software is (a) to mathematically verify and validate models and predictions
32 and (b) to foster research into new scores and metrics.

33 Key Benefits of scores

34 In order to meet the needs of researchers and other users, scores provides the following key
35 benefits.

36 Data Handling

- 37 ■ Works with n-dimensional data (e.g., geospatial, vertical and temporal dimensions) for
38 both point-based and gridded data. scores can effectively handle the dimensionality,
39 data size and data structures commonly used for:
 - 40 – gridded Earth system data (e.g. numerical weather prediction models)
 - 41 – tabular, point, latitude/longitude or site-based data (e.g. forecasts for specific
42 locations).
- 43 ■ Handles missing data, masking of data and weighting of results.
- 44 ■ Supports xarray (Hoyer & Hamman, 2017) datatypes, and works with NetCDF4 (Unidata,
45 2024), HDF5 (The HDF Group & Koziol, 2020), Zarr (Miles et al., 2020) and GRIB
46 (World Meteorological Organization, 2024) data sources among others.

47 Usability

- 48 ■ A companion Jupyter Notebook (Jupyter Team, 2024) tutorial for each metric and
49 statistical test that demonstrates its use in practice.
- 50 ■ Novel scores not commonly found elsewhere (e.g. FIRM (Taggart et al., 2022), Flip-Flop
51 Index (Griffiths et al., 2019, 2021)).
- 52 ■ All scores and statistical techniques have undergone a thorough scientific and software
53 review.
- 54 ■ An area specifically to hold emerging scores which are still undergoing research and
55 development. This provides a clear mechanism for people to share, access and collaborate
56 on new scores, and be able to easily re-use versioned implementations of those scores.

57 Compatability

- 58 ■ Highly modular - provides its own implementations, avoids extensive dependencies and
59 offers a consistent API.
- 60 ■ Easy to integrate and use in a wide variety of environments. It has been used on
61 workstations, servers and in high performance computing (supercomputing) environments.
- 62 ■ Maintains 100% automated test coverage.
- 63 ■ Uses Dask (Dask Development Team, 2016) for scaling and performance.
- 64 ■ Expanding support for pandas (McKinney, 2010; The pandas development team, 2024).

65 Metrics, Statistical Techniques and Data Processing Tools Included in scores

66 At the time of writing, scores includes over 50 metrics, statistical techniques and data
67 processing tools. For an up to date list, please see the scores documentation.

68 The ongoing development roadmap includes the addition of more metrics, tools, and statistical
69 tests.

Table 1: A curated selection of the metrics, tools and statistical tests currently included in scores

	Description	A Selection of the Functions Included in scores
Continuous	Scores for evaluating single-valued continuous forecasts.	Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Additive Bias, Multiplicative Bias, Pearson's Correlation Coefficient, Flip-Flop Index (Griffiths et al., 2019, 2021), Quantile Loss, Murphy Score (Ehm et al., 2016).
Probability	Scores for evaluating forecasts that are expressed as predictive distributions, ensembles, and probabilities of binary events.	Brier Score (Brier, 1950), Continuous Ranked Probability Score (CRPS) for Cumulative Distribution Functions (CDFs) (including threshold-weighting, see Gneiting & Ranjan (2011)), CRPS for ensembles (Ferro, 2013; Gneiting & Raftery, 2007), Receiver Operating Characteristic (ROC), Isotonic Regression (reliability diagrams) (Dimitriadis et al., 2021).
Categorical	Scores for evaluating forecasts of categories.	Probability of Detection (POD), Probability of False Detection (POFD), False Alarm Ratio (FAR), Success Ratio, Accuracy, Peirce's Skill Score (Peirce, 1884), Critical Success Index (CSI), Gilbert Skill Score (Gilbert, 1884), Heidke Skill Score, Odds Ratio, Odds Ratio Skill Score, F1 Score, Symmetric Extremal Dependence Index (Ferro & Stephenson, 2011), Fixed Risk Multicategorical (FIRM) Score (Taggart et al., 2022).
Spatial	Scores that take into account spatial structure.	Fractions Skill Score (Roberts & Lean, 2008).
Statistical Tests	Tools to conduct statistical tests and generate confidence intervals.	Diebold-Mariano (Diebold & Mariano, 1995) with both the Harvey et al. (1997) and Hering & Genton (2011) modifications.
Processing Tools	Tools to pre-process data.	Data matching, discretisation, cumulative density function manipulation.

70 **Use in Academic Work**

71 In 2015, the Australian Bureau of Meteorology began developing a new verification system
 72 called Jive, which became operational in 2022. For a description of Jive see Loveday, Griffiths,
 73 et al. (2024). The Jive verification metrics have been used to support several publications
 74 (Foley & Loveday, 2020; Griffiths et al., 2017; Taggart, 2022a, 2022b, 2022c). scores has
 75 arisen from the Jive verification system and was created to modularise the Jive verification
 76 functions and make them available as an open source package.

77 scores has been used to explore user-focused approaches to evaluating probabilistic and
 78 categorical forecasts (Loveday, Taggart, et al., 2024).

79 Related Software Packages

80 There are multiple open source verification packages in a range of languages. Below is a
81 comparison of scores to other open source Python verification packages. None of these
82 include all of the metrics implemented in scores (and vice versa).

83 xskillscore (Bell et al., 2021) provides many but not all of the same functions as scores
84 and does not have direct support for pandas. The Jupyter Notebook tutorials in scores cover
85 a wider array of metrics.

86 climpred (Brady & Spring, 2021) uses xskillscore combined with data handling functionality,
87 and is focused on ensemble forecasts for climate and weather. climpred makes some design
88 choices related to data structure (specifically associated with climate modelling) which may
89 not generalise effectively to broader use cases. Releasing scores separately allows the differing
90 design philosophies to be considered by the community.

91 METplus (Brown et al., 2021) is a substantial verification system used by weather and climate
92 model developers. METplus includes a database and a visualisation system, with Python and
93 shell script wrappers to use the MET package for the calculation of scores. MET is implemented
94 in C++ rather than Python. METplus is used as a system rather than providing a modular
95 Python API.

96 Verif (Nipen et al., 2023) is a command line tool for generating verification plots whereas
97 scores provides a Python API for generating numerical scores.

98 Pysteps (Imhoff et al., 2023; Pulkkinen et al., 2019) is a package for short-term ensemble pre-
99 diction systems, and includes a significant verification submodule with many useful verification
100 scores. PySteps does not provide a standalone verification API.

101 PyForecastTools (Morley & Burrell, 2020) is a Python package for model and forecast
102 verification which supports darray rather than xarray data structures and does not include
103 Jupyter Notebook tutorials.

104 Acknowledgements

105 We would like to thank Jason West and Robert Johnson from the Bureau of Meteorology for
106 their feedback on an earlier version of this manuscript.

107 We would like to thank and acknowledge the National Computational Infrastructure (nci.org.au)
108 for hosting the scores repository within their GitHub organisation.

109 References

110 Bell, R., Spring, A., Brady, R., Huang, A., Squire, D., Blackwood, Z., Sitter, M. C., &
111 Chegini, T. (2021). *xarray-contrib/xskillscore: Metrics for verifying forecasts*. Zenodo.
112 <https://doi.org/10.5281/zenodo.5173153>

113 Brady, R. X., & Spring, A. (2021). climpred: Verification of weather and climate forecasts.
114 *Journal of Open Source Software*, 6(59), 2781. <https://doi.org/10.21105/joss.02781>

115 Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly*
116 *Weather Review*, 78(1), 1–3. [https://doi.org/10.1175/1520-0493\(1950\)078%3C0001:
117 vofeit%3E2.0.co;2](https://doi.org/10.1175/1520-0493(1950)078%3C0001:vofeit%3E2.0.co;2)

118 Brown, B., Jensen, T., Gotway, J. H., Bullock, R., Gilleland, E., Fowler, T., Newman, K.,
119 Adriaansen, D., Blank, L., Burek, T., & others. (2021). The Model Evaluation Tools (MET):
120 More than a decade of community-supported forecast verification. *Bulletin of the American*
121 *Meteorological Society*, 102(4), E782–E807. <https://doi.org/10.1175/bams-d-19-0093.1>

- 122 Dask Development Team. (2016). *Dask: Library for dynamic task scheduling*. <http://dask.pydata.org>
123
- 124 Diebold, F. X., & Mariano, R. S. (1995). Comparing predictive accuracy. *Journal of Business & Economic Statistics*, 13(3), 253–263. <https://doi.org/10.3386/t0169>
125
- 126 Dimitriadis, T., Gneiting, T., & Jordan, A. I. (2021). Stable reliability diagrams for probabilistic classifiers. *Proceedings of the National Academy of Sciences*, 118(8), e2016191118.
127 <https://doi.org/10.1073/pnas.2016191118>
128
- 129 Ehm, W., Gneiting, T., Jordan, A., & Krüger, F. (2016). Of quantiles and expectiles: Consistent scoring functions, Choquet representations and forecast rankings. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 78(3), 505–562.
130 <https://doi.org/10.1111/rssb.12154>
131
- 133 Ferro, C. A. T. (2013). Fair scores for ensemble forecasts: Fair scores for ensemble forecasts. *Quarterly Journal of the Royal Meteorological Society*, 140(683), 1917–1923. <https://doi.org/10.1002/qj.2270>
134
- 136 Ferro, C. A. T., & Stephenson, D. B. (2011). Extremal dependence indices: Improved verification measures for deterministic forecasts of rare binary events. *Weather and Forecasting*, 26(5), 699–713. <https://doi.org/10.1175/WAF-D-10-05030.1>
137
- 139 Foley, M., & Loveday, N. (2020). Comparison of single-valued forecasts in a user-oriented framework. *Weather and Forecasting*, 35(3), 1067–1080. <https://doi.org/10.1175/waf-d-19-0248.1>
140
- 142 Gilbert, G. K. (1884). Finley's tornado predictions. *American Meteorological Journal*, 1(5), 166–172.
143
- 144 Gneiting, T., & Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477), 359–378. <https://doi.org/10.1198/016214506000001437>
145
- 147 Gneiting, T., & Ranjan, R. (2011). Comparing density forecasts using threshold-and quantile-weighted scoring rules. *Journal of Business & Economic Statistics*, 29(3), 411–422. <https://doi.org/10.1198/jbes.2010.08110>
148
- 150 Griffiths, D., Foley, M., Ioannou, I., & Leeuwenburg, T. (2019). Flip-Flop Index: Quantifying revision stability for fixed-event forecasts. *Meteorological Applications*, 26(1), 30–35. <https://doi.org/10.1002/met.1732>
151
- 153 Griffiths, D., Jack, H., Foley, M., Ioannou, I., & Liu, M. (2017). *Advice for automation of forecasts: A framework*. Bureau of Meteorology. <https://doi.org/10.22499/4.0021>
154
- 155 Griffiths, D., Loveday, N., Price, B., Foley, M., & McKelvie, A. (2021). Circular Flip-Flop Index: Quantifying revision stability of forecasts of direction. *Journal of Southern Hemisphere Earth Systems Science*, 71(3), 266–271. <https://doi.org/10.1071/es21010>
156
- 158 Harvey, D., Leybourne, S., & Newbold, P. (1997). Testing the equality of prediction mean squared errors. *International Journal of Forecasting*, 13(2), 281–291. [https://doi.org/10.1016/S0169-2070\(96\)00719-4](https://doi.org/10.1016/S0169-2070(96)00719-4)
159
- 161 Hering, A. S., & Genton, M. G. (2011). Comparing spatial predictions. *Technometrics*, 53(4), 414–425. <https://doi.org/10.1198/tech.2011.10136>
162
- 163 Hoyer, S., & Hamman, J. (2017). xarray: N-D labeled arrays and datasets in Python. *Journal of Open Research Software*, 5(1). <https://doi.org/10.5334/jors.148>
164
- 165 Imhoff, R. O., De Cruz, L., Dewettinck, W., Brauer, C. C., Uijlenhoet, R., Heeringen, K.-J. van, Velasco-Forero, C., Nerini, D., Van Genderachter, M., & Weerts, A. H. (2023). Scale-dependent blending of ensemble rainfall nowcasts and numerical weather prediction in the open-source pysteps library. *Quarterly Journal of the Royal Meteorological Society*,
166
167
168

- 149(753), 1335–1364. <https://doi.org/10.1002/qj.4461>
- 170 Jupyter Team. (2024). *Jupyter interactive notebook*. GitHub. <https://github.com/jupyter/notebook>
171 <https://github.com/jupyter/notebook>
- 172 Loveday, N., Griffiths, D., Leeuwenburg, T., Taggart, R., Pagano, T. C., Cheng, G., Plastow,
173 K., Ebert, E., Templeton, C., Carroll, M., Khanarmuei, M., & Nagpal, I. (2024). *The Jive*
174 *verification system and its transformative impact on weather forecasting operations*. arXiv.
175 <https://doi.org/10.48550/arXiv.2404.18429>
- 176 Loveday, N., Taggart, R., & Khanarmuei, M. (2024). A user-focused approach to evaluating
177 probabilistic and categorical forecasts. *Weather and Forecasting*. <https://doi.org/10.1175/waf-d-23-0201.1>
178 <https://doi.org/10.1175/waf-d-23-0201.1>
- 179 McKinney, W. (2010). Data structures for statistical computing in Python. In S. van der Walt
180 & J. Millman (Eds.), *Proceedings of the 9th Python in Science Conference* (pp. 56–61).
181 <https://doi.org/10.25080/Majora-92bf1922-00a>
- 182 Miles, A., Kirkham, J., Durant, M., Bourbeau, J., Onalan, T., Hamman, J., Patel, Z., shikharsg,
183 Rocklin, M., dussin, raphael, Schut, V., Andrade, E. S. de, Abernathey, R., Noyes, C.,
184 sbalmer, bot, pyup.io, Tran, T., Saalfeld, S., Swaney, J., ... Banihirwe, A. (2020). *Zarr-*
185 *developers/zarr-python: v2.4.0* (Version v2.4.0). Zenodo. <https://doi.org/10.5281/zenodo.3773450>
186 <https://doi.org/10.5281/zenodo.3773450>
- 187 Morley, S., & Burrell, A. (2020). *Drsteve/PyForecastTools: Version 1.1.1* (Version v1.1.1).
188 Zenodo. <https://doi.org/10.5281/zenodo.3764117>
- 189 Nipen, T. N., Stull, R. B., Lussana, C., & Seierstad, I. A. (2023). Verif: A weather-prediction
190 verification tool for effective product development. *Bulletin of the American Meteorological*
191 *Society*, 104(9), E1610–E1618. <https://doi.org/10.1175/bams-d-22-0253.1>
- 192 Peirce, C. S. (1884). The numerical measure of the success of predictions. *Science*, ns-4(93),
193 453–454. <https://doi.org/10.1126/science.ns-4.93.453.b>
- 194 Pulkkinen, S., Nerini, D., Pérez Hortal, A. A., Velasco-Forero, C., Seed, A., Germann, U., &
195 Foresti, L. (2019). Pysteps: An open-source Python library for probabilistic precipitation
196 nowcasting (v1.0). *Geoscientific Model Development*, 12(10), 4185–4219. <https://doi.org/10.5194/gmd-12-4185-2019>
197 <https://doi.org/10.5194/gmd-12-4185-2019>
- 198 Roberts, N. M., & Lean, H. W. (2008). Scale-selective verification of rainfall accumulations
199 from high-resolution forecasts of convective events. *Monthly Weather Review*, 136(1),
200 78–97. <https://doi.org/10.1175/2007MWR2123.1>
- 201 Taggart, R. (2022a). *Assessing calibration when predictive distributions have discontinuities*.
202 <http://www.bom.gov.au/research/publications/researchreports/BRR-064.pdf>
- 203 Taggart, R. (2022b). Evaluation of point forecasts for extreme events using consistent scoring
204 functions. *Quarterly Journal of the Royal Meteorological Society*, 148(742), 306–320.
205 <https://doi.org/10.1002/qj.4206>
- 206 Taggart, R. (2022c). Point forecasting and forecast evaluation with generalized huber loss.
207 *Electronic Journal of Statistics*, 16(1), 201–231. <https://doi.org/10.1214/21-ejs1957>
- 208 Taggart, R., Loveday, N., & Griffiths, D. (2022). A scoring framework for tiered warnings and
209 multicategorical forecasts based on fixed risk measures. *Quarterly Journal of the Royal*
210 *Meteorological Society*, 148(744), 1389–1406. <https://doi.org/10.1002/qj.4266>
- 211 The HDF Group, & Koziol, Q. (2020). *HDF5-version 1.12.0*. <https://doi.org/10.11578/dc.20180330.1>
212 <https://doi.org/10.11578/dc.20180330.1>
- 213 The pandas development team. (2024). *Pandas-dev/pandas: pandas* (Version v2.2.2). Zenodo.
214 <https://doi.org/10.5281/zenodo.10957263>

- 215 Unidata. (2024). *Network common data form (NetCDF)*. UCAR/Unidata Program Center.
216 <https://doi.org/10.5065/D6H70CW6>
- 217 World Meteorological Organization. (2024). *WMO no. 306 FM 92 GRIB (edition 2)*. World
218 Meteorological Organization. <https://codes.wmo.int/grib2>

DRAFT