

A system for airport weather forecasting based on circular regression trees

Pablo Rozas Larraondo^{a,*}, Iñaki Inza^b, Jose A. Lozano^{b,c}

^a*National Computational Infrastructure, Building 143, Australian National University, Ward Road, ACT, 2601, Australia*

^b*Intelligent Systems Group, Computer Science Faculty, University of the Basque Country, Paseo de Manuel Lardizabal, Donostia, 20018, Spain*

^c*Basque Center for Applied Mathematics (BCAM), Mazarredo 14, Bilbao, 48009, Spain*

Abstract

This paper describes a suite of tools and a model for improving the accuracy of airport weather forecasts produced by numerical weather prediction (NWP) products, by learning from the relationships between previously modelled and observed data. This is based on a new machine learning methodology that allows circular variables to be naturally incorporated into regression trees, producing more accurate results than linear and previous circular regression tree methodologies.

The software has been made publicly available as a Python package, which contains all the necessary tools to extract historical NWP and observed weather data and to generate forecasts for different weather variables for any airport in the world. Several examples are presented where the results of the proposed model significantly improve those produced by NWP and also by previous regression tree models.

Keywords: Circular Variables, Weather Forecasting, Meteorology, Regression Trees, Machine Learning

1 Software Availability

- 2 Name of software: AeroCirTree
- 3 Developer: Pablo Rozas Larraondo

*Corresponding author: Tel.: +61 (02) 6125 3211;

Email address: `pablo.larraondo@anu.edu.au` (Pablo Rozas Larraondo)

4 Contact Address: National Computational Infrastructure, Building 143, Aus-
5 tralian National University, Ward Road, ACT, 2601, Australia (pablo.larraondo@anu.edu.au)
6 Source: <http://github.com/prl900/AeroCirTree>
7 Programming Language: Python 3
8 Dependencies: Numpy, Pandas
9 Licence: GNU GPL v3

10 **1. Introduction**

11 Modern weather forecasting relies mostly on numerical models that sim-
12 ulate the evolution of the atmosphere, based on fluid dynamics and thermo-
13 dynamics equations. These equations are solved for the discrete points of a
14 regular grid covering the region of interest. Higher resolution models gener-
15 ate more detailed forecasts, but also require large computational resources
16 and longer running times. Operational models trade off resolution quality for
17 shorter processing times. The need for higher resolution forecasts has driven
18 numerous methodologies to generate more detailed outputs, which is known
19 as downscaling. Dynamic downscaling uses the output of a coarser model as
20 the initial condition of a higher resolution local model, which better resolves
21 sub-grid processes and topography [1]. Another approach is statistical down-
22 scaling, where historical observed data are used to enhance the output of a
23 numerical model. There are numerous methodologies for statistical down-
24 scaling based on different principles, such as analogues [2], interpolation [3]
25 or machine learning models [4, 5].

26 Aviation operations are highly affected by the weather and require the
27 best quality meteorological information to maximise efficiency and safety.
28 The International Civil Aviation Organization (ICAO) and the World Me-
29 teorological Organization (WMO) have established international standards
30 to ensure high quality meteorological reports [6]. To generate these reports,
31 national weather services across the world employ highly qualified personnel
32 who continuously observe and forecast conditions around the airport, such
33 as visibility, direction and speed of the wind or proximity of storm cells. Avi-
34 ation weather forecasters rely mainly on their knowledge of the airport and
35 the quality of the NWP used.

36 There are a number of tools that facilitate the process of generating air-
37 port weather forecasts [7, 8], being an area of active research at the moment.
38 Airports usually have long and regular series of high quality historical obser-
39 vation data that can be used to create statistical downscaling models to help

40 forecasters in their work. The effect of non-resolved surrounding mountains,
41 water bodies or local climate conditions can be incorporated by these models,
42 by studying the local effects produced by weather patterns in the past.

43 Circular variables are present in any directional measurement or variable
44 with an inherent periodicity. Weather data contain many parameters that
45 are represented as circular variables, such as wind direction, geographical
46 coordinates or timestamps. Most of the current regression machine learn-
47 ing algorithms focus on modelling the relationships between linear variables.
48 Circular variables have a different nature to linear variables, so traditional
49 methodologies are not able to represent their content thoroughly, leading to
50 suboptimal results in most cases. The model presented in this article builds
51 upon the concept of circular regression trees introduced by Lund [9]. Our
52 model is computationally more efficient and generates contiguous splits for
53 circular variables, which results in improved accuracy when compared to its
54 precursor.

55 Circular regression trees can better represent circular variables, as they
56 consider more possibilities for splitting the space than linear regression trees
57 do. Circular regression trees can define subsets of data around the origin
58 $0, 2\pi$ radians point. For example, when predicting an event that shows a high
59 correlation with the winter months in the northern hemisphere, a circular tree
60 would be able to isolate the months from December to March in one group.
61 On the other hand, a linear tree would most likely consider splits starting
62 or ending at the beginning of the year, failing to create a group containing
63 these months.

64 This paper introduces **AeroCirTree**, a system based on the described
65 circular regression tree model, which is able to generate improved airport
66 weather forecasts for any airport in the world. This software presents a
67 general solution where all the necessary tools required to extract historical
68 weather data, train models and generate new forecasts are made available.
69 This system is intended to help aviation weather forecasters to produce better
70 quality reports and for machine learning researchers to build upon more
71 sophisticated models.

72 The paper is structured as follows: Section 2 contains the methodology
73 used to create the model. Section 3 contains an introduction to the observed
74 and numerical weather datasets used to develop and test the system. Sec-
75 tion 4 presents results where the proposed model is compared with other
76 regression tree methodologies. This section also contains a discussion of
77 the results, providing the reader with deeper insight into the novelty of the

78 proposed model. Section 5 provides a high level description of the model im-
 79 plementation, including its key components and their functionality as well as
 80 examples on how to use the software. Section 6 concludes this paper, revisit-
 81 ing the research highlights and proposing some ideas on future developments
 82 to carry this work forward.

83 2. Methodology

84 Because of their simplicity, training speed and performance, regression
 85 trees are a popular and effective technique for modelling linear variables.
 86 Classification and Regression Trees (CART) [10] is one of the most popular
 87 versions of regression trees.

88 Linear regression trees recursively partition the space, finding the best
 89 split at each non-terminal node. Each split divides the space in two sets
 90 using a cost function, which is usually based on a metric for minimising the
 91 combined variance of the resulting children nodes.

92 Figure 1 contains an example of a regression tree based on two linear
 93 variables x_1 and x_2 . On the right side, there is a graphical representation of
 94 how the space is divided by creating splits on these two variables.

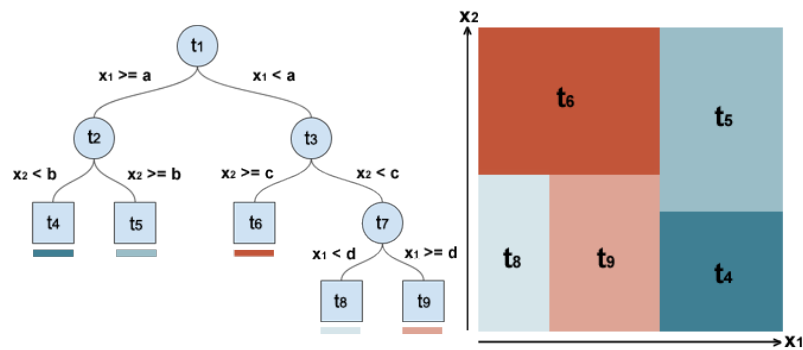


Figure 1: Example of a classic linear regression tree and a representation of how the space is divided.

95 Circular variables are numerical variables whose values are constrained
 96 into a cyclical space - for example, a variable measuring angles in radians,
 97 spans between 0 and 2π , where both values represent the same point in space.
 98 Although these variables can be included in a linear regression tree, they have
 99 to be treated as linear variables, which is an oversimplification and normally
 100 leads to suboptimal results [9].

101 A circular variable defines a circular space. A circular space is cyclic
 102 in the sense that it is not bounded; for instance, the notion of a minimum
 103 and maximum value does not apply. The distance between two values in the
 104 space becomes an ambiguous concept, as it can be measured in clockwise and
 105 anticlockwise directions, yielding different results. Also, this space cannot be
 106 split in two halves by selecting a value, as the ' $<$ ' and ' $>$ ' operators are not
 107 applicable.

108 In order to split a circular variable, at least two different values need
 109 to be defined. These two values describe two complementary sectors, each
 110 containing a portion of the data. Circular regression trees use this splitting
 111 approach for incorporating circular variables into regression trees.

112 There are many examples of circular variables. Any variable representing
 113 directional data or a periodic event is circular. More specifically, in the field
 114 of airport weather forecasting, wind direction, the time of the day or the date
 115 are examples of circular variables.

116 Lund [9] proposes a methodology that allows circular variables to be
 117 incorporated into regression trees. Figure 2 contains a similar representation
 118 to the previous example, but considering one circular variable α and a linear
 119 one x_1 . On the right side, there is a chart representing how the space is
 120 partitioned using polar coordinates.

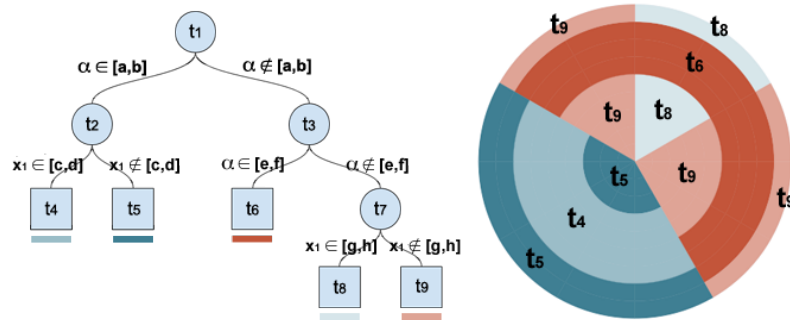


Figure 2: Example of Lund's original proposal of circular regression tree and a representation of how the space is divided.

121 The methodology presented in this work builds upon the concept of cir-
 122 cular regression trees, presenting an alternative that improves computational
 123 performance and the accuracy of its results. Figure 3 shows how the space
 124 is partitioned using the proposed methodology.

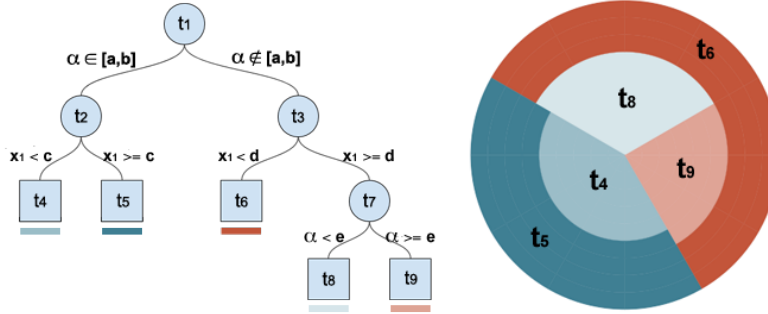


Figure 3: Example of the proposed circular regression tree and a representation of how the space is divided.

125 Visually comparing Figure 2 and Figure 3, it is evident that regions are
 126 split differently. The novelty of this methodology, when compared to the origi-
 127 nal version proposed by Lund, is that it always generates contiguous splits.
 128 In doing so, we avoid an excessive fragmentation of the space, and the splits
 129 provide a better generalisation for its child nodes. The original methodology
 130 uses the ' \in ' and ' \notin ' operators to generate all the splits for circular variables.
 131 This usually generates partitions in which the subsets defined by the \in clause
 132 are surrounded by the complementary \notin subset. Our methodology uses these
 133 operators to create just the first split of a circular variable and, after that,
 134 uses the '<' and '>'' operators to create the subsequent splits. This change
 135 also results in a reduction of the search space for possible splits. The pro-
 136 posed algorithm for generating circular trees has, as a consequence, $\mathcal{O}(n)$
 137 cost instead of $\mathcal{O}(n^2)$, when compared to Lund's original proposal. The only
 138 exception is when computing the first split of a circular variable, which has
 139 a computational cost of $\mathcal{O}(n^2)$, as it has to consider all the different splits
 140 around the circle.

141 3. Software and datasets

142 **AeroCirTree** is a collection of Python scripts which provides the tools to
 143 train and test the three previously described regression tree methodologies
 144 using airport weather data. It uses NWP variables as the input and generates
 145 a more accurate value for the selected output variable by learning from the
 146 observed values for a certain location. Once the model has been trained, it

147 can be used to improve the accuracy of the forecasted output value provided
148 by new incoming NWP data.

149 It is worth noting that regression tree models are presented in this work as
150 a method to statistically downscale the output of NWP for specific locations.
151 They are not used to predict future values of a time-series but to improve the
152 values produced by NWP. Analysis data from the NWP model and observed
153 data are used to train the regression trees. These trees can account for biases
154 and systematic errors of the NWP model. Trained models can be applied
155 to any forecasting horizon produced by the NWP to correct systematic and
156 random errors.

157 The `AeroCirTree` software presented in this work offers a general im-
158 plementation of a regression tree. `AeroCirTree` allows its users to train
159 linear regression trees as well as circular versions using non-contiguous or
160 contiguous splits, as we propose. To determine which methodology is used,
161 each variable in the input or output can be tagged as being either [`linear`,
162 `circular`] using a configuration file. An extra tag, `contiguous`, which can
163 be set to [`true`, `false`], indicates the split methodology applied to circular
164 variables. Different values of these tags indicate different versions of regres-
165 sion trees. For example, classic linear regression trees can be generated by
166 tagging all their input variables as `linear` and `contiguous=true`. Lund’s
167 proposal of circular tree would require the circular input variables to be
168 tagged as `circular` and `contiguous=false`. Lastly, our proposed methodol-
169 ogy would require the same circular input variables to be `contiguous=true`.

170 `AeroCirTree` makes use of two weather datasets. The first is the output of
171 a global NWP, called the Global Forecast System model (GFS) [11], which is
172 run operationally by the National Oceanic and Atmospheric Administration
173 (NOAA). The second uses Meteorological Aerodrome Reports (METARs)
174 [6], which contain periodic meteorological observations from airports around
175 the world.

176 Each of these datasets contains several variables describing different weather
177 parameters, such as the temperature, humidity, wind speed or cloud cover
178 at the different locations they represent. The GFS model represents data
179 using a regular grid which covers the whole world with a spatial resolution
180 of approximately 50 km and a temporal resolution of 3 hours. NOAA main-
181 tains an Operational Model Archive and Distribution System (NOMADS) to
182 publish the GFS data. This archive contains the GFS outputs for the last
183 10 years.

184 METARs are weather text reports that encode observed meteorological

185 parameters at airport runways using a well defined code. METARS are pro-
186 duced with an hourly or half-hourly frequency and are also made publicly
187 available through the WMO Global Telecommunication System (GTS). The
188 National Centers for Environmental Prediction (NCEP) maintains a sys-
189 tem called Meteorological Assimilation Data Ingest System (MADIS), which
190 archives all the METAR reports that have been produced in the world for
191 the last 10 years. Each report is uniquely identified by its header, which
192 contains the International Civil Aviation Organization (ICAO) airport code
193 and a UTC time stamp.

194 The provided `AeroCirTree` software contains a command line utility that
195 extracts the information from these two datasets for any given airport and
196 date range. The output is presented as a convenient `csv` file containing
197 the values of the different variables as a time series. All operations, such
198 as locating the airport coordinate in the GFS grid, parsing and extracting
199 METARs or homogenising variable units, are handled by the software, so the
200 user can easily get a clean dataset for the desired airport. This `csv` file is the
201 input used to train new models.

202 4. Experiments and results

203 The hypothesis of this study is that our proposed methodology for gener-
204 ating regression trees provides better generalisation and accuracy than pre-
205 vious non-contiguous circular regression trees when using circular variables
206 and the equivalent classic linear methodologies.

207 The next sections go through the required steps to extract the necessary
208 data, train the models and generate the forecasts. The last section contains
209 an analysis of the proposed model accuracy and a comparison with the results
210 provided by the GFS raw output, Lund’s methodology and classic linear
211 regression trees.

212 4.1. Data extraction and model training

213 To compare the differences in performance between methodologies, we use
214 weather data coming from simulated NWP and observed data from different
215 airports. Regression trees are trained using NWP as input and the observed
216 speed of the wind as the output variable. It is worth noting that regression
217 tree models are not used to forecast wind speeds into the future. These
218 models are used to statistically downscale NWP data, correcting biases and
219 systematic errors.

220 We choose to forecast the observed speed of the wind at 5 different lo-
221 cations in Europe. Data from the airports of Berlin Tegel (EDDT), Lon-
222 don Heathrow (EGLL), Barcelona El Prat (LEBL), Paris Charles de Gaulle
223 (LFPG) and Milano Malpensa (LIMC) are used to train the different models
224 and to analyse the results. The models are trained using three-hourly data
225 for the years 2011, 2012 and 2013, providing approximately 8760 samples per
226 airport.

227 Each model generates the required partitions to predict the observed wind
228 speed using the following GFS parameters as input variables: relative humid-
229 ity, speed and direction of the 10-meter wind as well as the time of the day
230 associated with the values. Wind speed is one of the most important weather
231 variables affecting airport operations. This variable is also highly dependent
232 on another variable, wind direction, which is circular. The reason for in-
233 cluding these two variables in our experiments is that, in conjunction, they
234 can represent local topography effects non resolved by weather models. Sur-
235 face relative humidity is used as an indicator for phenomena such as rain
236 or fog conditions. Lastly, time of the day, also a circular variable, is highly
237 correlated with the daily patterns of the wind.

238 The stop criterium for all the considered trees is based on the number
239 of elements in a node. Splits are recursively performed until the number of
240 data entries in a node falls below a certain value. Then, the splitting process
241 is stopped and the node is denoted as a leaf. This value receives the name
242 “maximum leaf size”. Large values of “maximum leaf size” generate shallow
243 trees, whereas small values generate deep trees with a larger number of nodes.
244 For each airport, different versions of the model are generated using different
245 maximum tree leaf sizes. The maximum leaf size values considered in this
246 experiment are: 1000, 500, 250, 100 and 50. This is the content of the config
247 file used to train our proposed model for the comparison defining a maximum
248 leaf size of 100 (*please refer to Section 5.2 for more details on how these files*
249 *are used and defined.*):

```
250  
251 {"output":{"name":"metar_wind_spd","type":"linear"},  
252   "input":[{"name":"gfs_wind_spd","type":"linear"},  
253           {"name":"gfs_wind_dir","type":"circular"},  
254           {"name":"gfs_rh","type":"linear"},  
255           {"name":"time","type":"circular"}]},  
256   "contiguous":true
```

257 "max_leaf_size":100}

258 4.2. Experimental analysis

259 Following the process described in the previous sections, data from 2011
260 to 2013 is extracted for the 5 selected airports. For each airport and value
261 of maximum leaf size, three different models are generated: classic linear
262 regression tree (using the \mathbf{u} , \mathbf{v} components of the wind speed and time of the
263 day), Lund's and our proposed circular regression tree.

264 To evaluate the differences in accuracy between these three methodolo-
265 gies, a 5-fold cross validation procedure is used. This validation process
266 ensures that models are tested using data that has not been used at training
267 time. In order to avoid differences in the results caused by different partitions
268 in the validation process, the same 5-fold partition is used to validate all the
269 methodologies for the different values of the "maximum leaf size" parame-
270 ter. The error in forecasting is defined as the difference between the speed
271 of the wind predicted by the tree, which is the mean of the target values
272 contained in the corresponding leaf, and the observed METAR wind speed
273 value. The Refined Index of Agreement (RIA) [12] is used to measure the
274 differences in accuracy between methodologies. This index provides greater
275 separation when comparing models that perform relatively well and is less
276 sensitive to errors concentrated in outliers when compared to other methods
277 such as absolute or root mean squared error. The RIA can be expressed as

$$RIA = 1 - \frac{\sum_{i=1}^n |P_i - O_i|}{2 \sum_{i=1}^n |O_i - \bar{O}|}$$

278 Where O_i represents the observations and P_i the predictions produced by
279 the model.

280 Table 1 contains the resulting RIA values for each tree methodology as
281 well as the reference value of the 10-meter wind speed produced by GFS in
282 the airports previously referenced. Higher values of RIA indicate better accu-
283 racy in the results. Similar results using different combinations of input and
284 output variables combining linear and circular variables are made available,
285 as a text file, at the main code repository.

286 Looking at the RIA values contained in Table 1, it can be noted that
287 the use of regression tree models significantly improves the level of accuracy
288 from the output of the GFS model. The level of improvement is highly de-
289 pendent on the selected airport. This may be due to the fact that each grid

290 point of the GFS model contains a representation of the weather in an area
291 of approximately 50 square kilometres, and some locations and variables are
292 better represented by this simplification than others. For example, airports
293 surrounded by mountains will benefit more from statistical models than air-
294 ports located on large plains.

295 Comparing the differences in accuracy between the three regression tree
296 models shown in Table 1, the use of the proposed model provides better re-
297 sults in most of the cases. The level of improvement also varies significantly
298 between different airport locations. Results are analysed considering the case
299 of shallow and deep trees. For shallow trees, the two circular models show
300 very similar behaviour outperforming the linear approach. As the maximum
301 leaf size parameter gets smaller, we see an improvement in accuracy for all
302 three models. Deeper trees still show better results for the circular mod-
303 els, but Lund’s proposal starts showing signs of premature over-fitting when
304 compared to the other two models. In the case of the deepest tree (maximum
305 leaf size equal to 50), all three models show a deterioration of performance,
306 with Lund’s being the most noticeable case.

307 In the case of Paris Charles de Gaulle (LFPG), shallow circular trees
308 show an improvement of around 4 to 5% when compared to the classic linear
309 tree version. This improvement is maintained by our proposed model when
310 considering deeper trees. However, Lund’s model does not improve at the
311 same rate. A more systematic analysis of the results of this test is offered at
312 the end of the section, providing the statistical significance of the differences
313 between methodologies.

314 Figure 4 and Figure 5 show a graphical representation of the evolution
315 of the RIA when predicting wind speed for the airports of London Heathrow
316 (EGLL) and Barcelona El Prat (LEBL) respectively. All the regression tree
317 methodologies improve their accuracy as the maximum leaf size decreases,
318 showing signs of overfitting for the smallest leaf size case. The value of the
319 GFS wind speed value at the closest grid point is shown as a reference to
320 represent the relative improvement achieved by each model.

321 As introduced in Section 2, the circular methodologies have the benefit of
322 considering extra partitions for circular variables, those that cross the origin,
323 when compared to linear methods. The benefits of using circular trees are
324 more noticeable for the case of shallow trees, the ones with larger values of
325 maximum leaf size. The first split of a circular variable normally happens at
326 one of the first nodes of the tree, near the root node. Splits that happen at
327 the top part of a tree have a major impact on its performance, because they

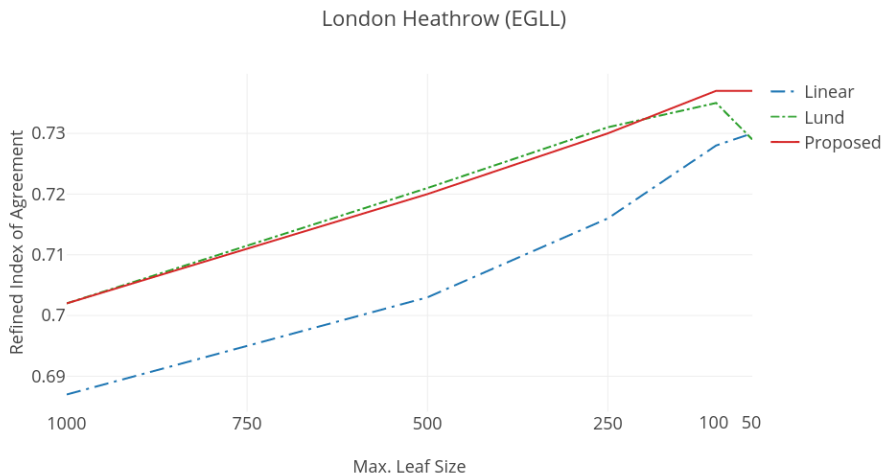


Figure 4: RIA values for the airport of London Heathrow (EGLL), comparing the accuracy of the output for different maximum leaf sizes.

328 divide a bigger proportion of the dataset. For shallow trees, finding a good
 329 partition at these levels is critical, whereas deeper trees can improve poor
 330 partitions by creating new ones.

331 Non-contiguous circular regression trees generate partitions that seem
 332 to provide a poorer generalisation for subsequent splits than the other two
 333 methodologies. The good results shown by Lund’s method for shallow trees
 334 quickly deteriorate for deeper trees. The proposed methodology, based on
 335 contiguous circular trees, achieves a similar performance to Lund’s method
 336 for shallow trees and also better results than the other two methodologies for
 337 deeper ones. Moreover, as mentioned in Section 2, the proposed methodology
 338 is more efficient computationally than the non-contiguous version.

339 In order to evaluate the results, the methodology proposed by Demsar
 340 [13] is used to assess the statistical significance of the differences between
 341 methods. The null hypothesis of similarity is rejected for linear and both
 342 circular regression trees. This justifies the use of post-hoc bivariate tests,
 343 Nemenyi in our case, which assess the statistical difference between pairs
 344 of algorithms. The results of these tests can be graphically expressed using
 345 Critical Difference (CD) diagrams. The Nemenyi test pairwise compares
 346 every methodology. The accuracy of any two methodologies is considered
 347 significantly different if the corresponding average rank differs by at least the

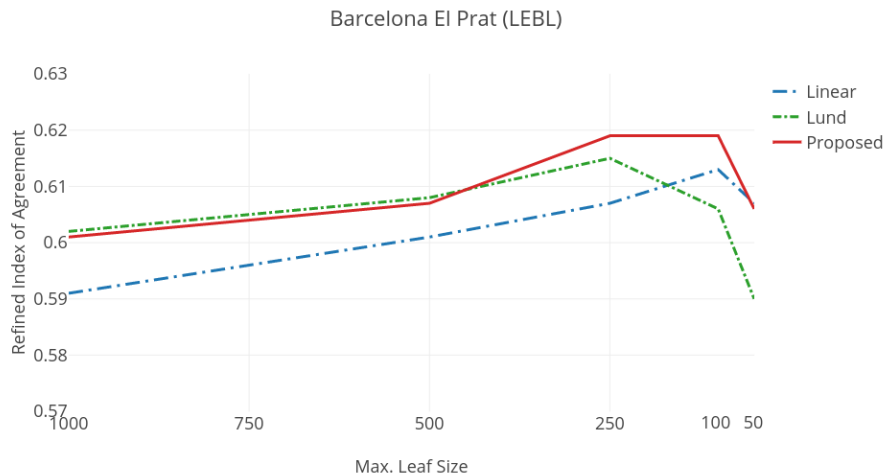


Figure 5: RIA values for the airport of Barcelona El Prat (LEBL) comparing the accuracy of the output for different maximum leaf sizes.

critical difference.

Figure 6 represents the RIA results of the Nemenyi test ($\alpha = 0.05$) making use of CD diagrams for the maximum leaf sizes of 1000, 100 and 50, as they represent both extremes of the proposed range.

CD diagrams connect the groups of algorithms for which no significant differences were found, or in other words, those whose distance is less than the fixed critical difference, shown above the graph. Note that algorithms ranked with lower values in CD diagrams imply higher RIA scores. These tests have been performed using the `scmamp` R package, which is publicly available at the Comprehensive R Archive Network (CRAN) [14].

As can be seen in the CD diagrams in Figure 6, for shallow trees, both circular methodologies outperform the linear approach (maximum leaf size 1000). As the experiment progresses into deeper trees (maximum leaf size 100), the proposed methodology statistically outperforms the other two in the considered datasets. Even for the case of maximum leaf size 50, when all the methods show a deterioration in accuracy, the proposed methodology shows the best results. Lund’s methodology, on the other hand, reveals a major degradation in accuracy for the smallest maximum leaf size. These results corroborate our experimental hypothesis: the proposed circular regression tree is able to generate models that provide better generalizations for circular

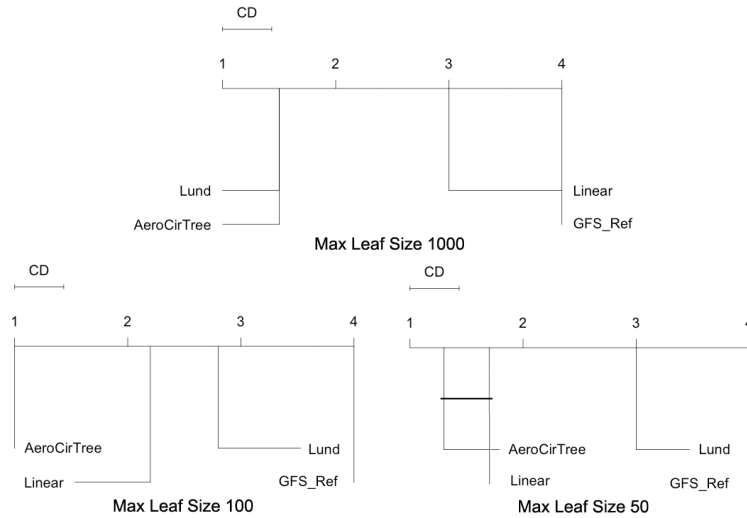


Figure 6: Critical Differences comparing the three methodologies for shallow and deep trees. $\alpha = 0.05$

368 variables.

369 5. Design and use of the software

370 `AeroCirTree` is a Python 3 package implementing regression trees and
 371 a set of command line tools to extract weather data and train tree models
 372 for any airport in the world. Users will normally use the provided package
 373 by using three scripts, named `aerocirtree_extract`, `aerocirtree_train`
 374 and `aerocirtree_test`, which fetch historical time-series weather data, train
 375 models and test results respectively, for any airport in the world.

376 5.1. Implementation design

377 The proposed circular regression tree has been implemented as a Python
 378 package. Most of its functionality is contained in two classes, called `Data`
 379 and `Node`. A tree is modelled as a nested structure of `Node` instances. Each
 380 `Node` in the tree contains an instance of the `Data` class, which represents the
 381 subset of the dataset contained in that node. The `Data` object is built around
 382 the Python Pandas `DataFrame` class.

383 `Node` contains two class attributes of type `Node`, named `left_child` and
 384 `right_child`, defining a recursive structure. Each non-terminal node in a

385 tree contains two `Node` instances which constitute its left and right children.
386 On the other hand, terminal nodes or leaves are characterised by having the
387 contents of its children set to the `None` value.

388 `Node` defines also the `.Split()` method which creates a split generating
389 two new instances of the `Node` class. Each of these two new `Node` instances
390 contains one part of the original `Data` and is assigned to the `left_child` and
391 `right_child` attributes. A tree is built by recursively calling the `.Split()`
392 method on each of the children `Node` until the stop criteria is satisfied. The
393 stop criteria can be configured to be a minimum number of elements or
394 variance value for the `Data` contents of a node.

395 Each column of a node's `Data` has to be tagged as linear or circular to
396 designate the nature of the data it represents. By tagging columns, we can
397 dynamically train different tree versions and compare their results. Classic
398 regression trees consider all the variables as linear, whereas our proposed
399 methodology allows some of the variables to be treated as circular. For
400 example, by tagging all variables as linear, we will get a classic regression
401 tree.

402 This implementation is generic and can be applied to data from any field
403 if made available in `csv` format.

404 5.2. User guide

405 `AeroCirTree` also provides a series of scripts to extract weather data,
406 train and test regression tree models. These scripts make use of the previously
407 described package to train specific models for any airport in the world.

408 Here is an example that shows how to extract the data for the airport of
409 London Heathrow from the 1st of January 2016 to the 1st of June 2016:

```
410  
411 $ ./aerocirtree_extract --airport EGLL --start_date 20160101\  
412 --end_date 20160601  
413 metar_press , metar_rh , metar_temp , metar_wind_spd , gfs_press ,\  
414 gfs_rh , gfs_temp , gfs_wind_dir , gfs_wind_spd , time , date  
415 1025.0 , 75.5 , 6.0 , 2.57 , 1016 , 92 , 3 , 280 , 3 , 45.0 , 0  
416 1024.0 , 80.92 , 5.0 , 4.12 , 1016 , 96 , 3 , 290 , 3 , 90.0 , 0  
417 1024.0 , 80.92 , 5.0 , 2.57 , 1015 , 97 , 4 , 300 , 3 , 135.0 , 0  
418 1024.0 , 86.99 , 6.0 , 2.57 , 1016 , 93 , 6 , 340 , 3 , 180.0 , 0  
419 ...
```

420 Note that the values of time and date are transformed to their numerical
421 values as circular variables, where the origin [0-360] corresponds to 00:00
422 hours and the 1st of January respectively. The output of this command can
423 be redirected to a local file. These files are used as the input required to
424 train tree models.

425 Once a dataset is available for a given airport, a model can be trained by
426 defining its input and target variables. The output variable has to be one
427 of the observed variables coming from the METAR reports and the input
428 variables are the GFS forecasted variables or a subset of them.

429 Doing it this way, when new forecast data from the GFS is available, the
430 model can be used to generate an enhanced forecast of the target variable.
431 The different options to create a model are specified through a configura-
432 tion file. This configuration file contains a JSON object with three fields:
433 “output”, “input” and “max_leaf_size”. The name of the target variable pro-
434 duced by the tree is specified in “output”. Input variables are listed in the
435 “input” field along with a tag to treat them as either circular or linear. The
436 max_leaf_size parameter specifies the value to control the depth of the result-
437 ing tree. For example, to specify a model to forecast temperature using GFS
438 relative humidity, wind direction as a circular variable and a maximum leaf
439 size of 100, a file with the following content should be specified:

```
440  
441 { "output": { "name": "metar_temp", "type": "linear" },  
442   "input": [ { "name": "gfs_wind_dir", "type": "circular" },  
443             { "name": "gfs_rh", "type": "linear" } ],  
444   "contiguous": true  
445   "max_leaf_size": 100 }
```

446 To train a model we use `aerocirtree_train`, which receives as arguments
447 the paths of a file containing the data and a configuration file. Supposing
448 the output of the data extracted in the previous section has been saved
449 in a file named `EGLL.csv` and the presented configuration file is saved as
450 `Model_A.json`, a model can be trained by running:

```
451  
452 $ ./aerocirtree_train --data EGLL.csv --config Model_A.json
```

453 This command learns the specified model and saves it using a name that
454 combines both input file names and using the extension `.mod`. The previous
455 model would be saved on disk with the file name `EGLL_Model_A.mod`.

456 Finally, `aerocirtree_test` can be used to run the model on new data.
457 This script receives the path to a saved model file and input csv as arguments.
458 The script returns the resulting model outputs for each line of the input file.

459 For example, supposing we want to test our previously trained model
460 `EGLL_Model_A.mod` with new data contained in the file `EGLL.csv`, we could
461 run:

```
462  
463 $ ./aerocirtree_test --data EGLL_new.csv --model EGLL\_Model\_A.mod
```

464 This command computes the resulting temperature values for each of the
465 input values at the airport of London Heathrow.

466 6. Conclusions

467 This work presents a software application for forecasting the weather in
468 any airport of the world. It also proposes a new circular regression tree
469 methodology which offers better accuracy when compared to classic linear
470 methods, and also better accuracy and computational efficiency than Lund's
471 original proposal of circular regression trees.

472 This software contains a library that implements a general version of
473 regression trees as well as the command line tools to train, test and download
474 new airport datasets. These tools have been designed so users can create
475 their own forecasts and also so that they can experiment and explore the
476 differences between models, input variables and airports. Scripts and libraries
477 are written in a simple way so users can read the code to understand what
478 the program is doing and also modify parts of it. `AeroCirTree` comes with a
479 GNU GPLv3 licence so anyone can use, modify and share this program for
480 any purpose.

481 The model proposed in this work is based on a new methodology to
482 build a basic circular regression tree. Regression trees have evolved with the
483 introduction of many different techniques that improve both their accuracy
484 and efficiency. Well known techniques that modify standard regression trees
485 such as pruning, balancing, smoothing [10, 15] or random forests [16] and
486 ensembles [17] can be also applied to circular regression trees and can improve
487 the accuracy of results when compared to basic regression trees. Future work
488 could implement the ideas presented in the referred publications offering more
489 advanced models.

490 **Acknowledgements**

491 We would like to thank the National Computational Infrastructure (NCI)
492 at the Australian National University and the University of the Basque Coun-
493 try for their support and advice in carrying out this research work.

494 We are grateful for the support of the Basque Government (IT609-13), the
495 Spanish Ministry of Economy and Competitiveness (TIN2016-78365-R) and
496 a University-Society Project (15/19 Basque Government and UPV/EHU).

497 Jose A. Lozano is also supported by BERC program 2014-2017 (Basque
498 Gov.) and Severo Ochoa Program SEV-2013-0323 (Spanish Ministry of Econ-
499 omy and Competitiveness).

500 **References**

501 [1] A. C. Carvalho, A. Carvalho, H. Martins, C. Marques, A. Rocha, C. Bor-
502 rego, D. X. Viegas, A. I. Miranda, Fire weather risk assessment under
503 climate change using a dynamical downscaling approach, *Environmental*
504 *Modelling & Software* 26 (2011) 1123-1133.

505 [2] M. Bannayan, G. Hoogenboom, Weather analogue: A tool for real-
506 time prediction of daily weather data realizations based on a modified
507 k-nearest neighbor approach, *Environmental Modelling & Software* 23
508 (2008) 703-713.

509 [3] C. C. F. Plouffe, C. Robertson, L. Chandrapala, Comparing interpola-
510 tion techniques for monthly rainfall mapping using multiple evaluation
511 criteria and auxiliary data sources: A case study of sri lanka, *Environ-*
512 *mental Modelling & Software* 67 (2015) 57-71.

513 [4] P. Rozas-Larraondo, I. Inza, J. A. Lozano, A method for wind speed
514 forecasting in airports based on nonparametric regression., *Weather and*
515 *Forecasting* 29 (2014) 1332-1342.

516 [5] T. Salameh, P. Drobinski, M. Vrac, P. Naveau, Statistical downscaling of
517 near-surface wind over complex terrain in southern france, *Meteorology*
518 *and Atmospheric Physics* 103 (2009) 253-265.

519 [6] WMO, *Manual on Codes International Codes VOLUME I.1.*, Geneva,
520 1995.

- 521 [7] J. E. Ghirardelli, B. Glahn, The meteorological development laboratorys
522 aviation weather prediction system, *Weather and Forecasting* 25 (2010)
523 1027-1051.
- 524 [8] A. J. M. Jacobs, N. Maat, Numerical guidance methods for decision
525 support in aviation meteorological forecasting, *Weather and Forecasting*
526 20 (2005) 82-100.
- 527 [9] U. J. Lund, Tree-based regression for a circular response., *Communica-*
528 *tions in Statistics - Theory and Methods* 31 (2002) 1549-1560.
- 529 [10] L. Breiman, J. Friedman, R. Olshen, C. Stone, *Classification and regres-*
530 *sion trees*, Wadsworth Books, 1984.
- 531 [11] K. Campana, P. Caplan, Technical procedure bulletin for t382 global
532 forecast system., 2005.
- 533 [12] C. J. Willmott, S. M. Robeson, K. Matsuura, A refined index of model
534 performance, *International Journal of Climatology* 32 (2012) 2088–2094.
- 535 [13] J. Demsar, Statistical comparisons of classifiers over multiple data sets.,
536 *Journal of Machine Learning Research* 7 (2006) 1–30.
- 537 [14] B. Calvo, G. Santafe, scmamp: Statistical comparison of multiple algo-
538 rithms in multiple problems., *The R Journal* (2016).
- 539 [15] J. R. Quinlan, *C4.5: Programs for Machine Learning*, Morgan Kauf-
540 mann Publishers, Inc., 1993.
- 541 [16] L. Breiman, Random forests., *Machine Learning* 45.1 (2001) 5–32.
- 542 [17] P. Bühlmann, *Handbook of Computational Statistics*, Springer Berlin
543 Heidelberg, 2012.

Table 1: Comparison of the RIA values when forecasting the observed METAR wind speed for the different airports using the direct output of GFS, a classic linear regression tree, Lund’s circular tree and the proposed model.

<i>Airport</i>	<i>Method</i>	RIA per Max Leaf Size				
		1000	500	250	100	50
EDDT	GFS (ref.)	0.669	0.669	0.669	0.669	0.669
	Linear	0.684	0.695	0.710	0.716	0.713
	Lund	0.700	0.713	0.720	0.715	0.702
	AeroCirTree	0.700	0.712	0.717	0.721	0.714
EGLL	GFS (ref.)	0.653	0.653	0.653	0.653	0.653
	Linear	0.687	0.703	0.716	0.728	0.730
	Lund	0.702	0.721	0.731	0.735	0.729
	AeroCirTree	0.702	0.720	0.730	0.737	0.737
LEBL	GFS (ref.)	0.362	0.362	0.362	0.362	0.362
	Linear	0.591	0.601	0.607	0.613	0.607
	Lund	0.602	0.608	0.615	0.606	0.590
	AeroCirTree	0.601	0.607	0.619	0.619	0.606
LFPG	GFS (ref.)	0.604	0.604	0.604	0.604	0.604
	Linear	0.674	0.691	0.702	0.711	0.707
	Lund	0.704	0.716	0.719	0.706	0.691
	AeroCirTree	0.704	0.712	0.715	0.714	0.707
LIMC	GFS (ref.)	0.401	0.401	0.401	0.401	0.401
	Linear	0.517	0.519	0.519	0.509	0.496
	Lund	0.521	0.520	0.518	0.500	0.482
	AeroCirTree	0.522	0.521	0.521	0.513	0.501