1   **Machine Learning for Wildfire Classification: Exploring Blackbox, eXplainable, Symbolic,**
2   **and SMOTE Methods**

3   M. Al-Bashiti[1], M.Z. Naser[1,2]
4   [1]School of Civil & Environmental Engineering and Earth Sciences (SCEEES), Clemson University, USA
5   [2]AI Research Institute for Science and Engineering (AIRISE), Clemson University, Clemson, SC 29634, USA
6   E-mail: malbash@clemson.edu, E-mail: mznaser@clemson.edu, Website: www.mznaser.com

7
8   **Abstract**

9   Whether triggered by natural or human-made events, wildfires are considered one of the most
10  traumatic events to our community and environment. Thus, properly predicting wildfires continues
11  to be an active area of research. This work showcases a statistical overview of the problem of
12  wildfires and then presents a dense data-driven ($D^3$) approach that leverages a variety of machine
13  learning (ML) techniques, namely, *blackbox* and *eXplainable* ML (i.e., deep learning (DL),
14  decision tree (DT), Stochastic Gradient Descent (SGD), Extreme Gradient Boosted Trees
15  (ExGBT), Logistic regression (LR)), and *symbolic* ML via genetic algorithms (GA) to classify and
16  predict wildfire breakouts. This approach was developed and validated using two databases
17  comprising more than 1.04 million geo-referenced wildfires that burned over 359,000 km$^2$ (88.7
18  million acres) between 1992 and 2015 in North America and Europe. Despite the complex nature
19  of wildfire formation and the interdependency of its governing factors, the findings of this $D^3$
20  analysis show the feasibility of utilizing ML in preciously classifying the expected size of wildfires
21  and predicting the possibility of the breakout of wildfires.

22
23  *Keywords:* Wildfires; Forests; Machine learning; Big data; explainable ML, Symbolic ML.
24
25  **Introduction**
26  The start of the twenty-first century marks a clear transition in which the number and intensity of
27  wildfires have exponentially risen [1]. While they can start naturally, wildfires are often caused by
28  humans with devastating consequences. On average, wildfires burn up to 1.11 billion acres of land
29  each year [2,3]. The united states wildfires have been significantly increasing from (140 to 250
30  wildfires) from (1980 to 2012)[3]. Although wildfires occur worldwide, they are most common
31  in regions with intense droughts and frequent lightning/thunderstorms.
32
33  This rise in wildfire occurrences mirrors the recent changes to our environment in which the
34  combination of dry conditions, extended high temperatures, and trapped emissions contribute to
35  some of these changes [4,5]. More specifically, climate change effects (and increased global
36  warming) generate heated conditions that draw moisture from the soil and dry out plants. Global
37  warming has not only led to the rise in wildfire occurrences, along with their intensity but has also
38  led to an increase in human and animal casualties, property losses, and environmental damage
39  [6,7]. This has been duly noted in the recent wildfires that broke in North America and Europe
40  over the past few years.
41
42  Wildfires require three components to breakout, known as the fire triangle. These include; a heat
43  source, fuel, and oxygen, heat sources, such as lightning, can supply enough heat to ignite a fire
44  that turns into flames when fuel or any flammable material is present [8]. Such ignition is bound
45  to spread and transport given the presence of favorable wind conditions [9,10].

46
47  Recent years have noted a surge in the amount of works that developed different approaches with
48  the power of data analytics to forecast the breakout of wildfires. Collectively, a number of
49  researchers [1,11–13] noted that there are three super high-tech approaches often used to predict
50  wildfire occurrences and stop them from surging. These approaches are grouped under physics-
51  based methods, statistical methods, and machine learning methods.
52
53  The first class of approaches, those lumped under physics-based methods, predicts fire breakout
54  by using a mathematical formulation that relies on fluid and heat transfer principles [14]. As such,
55  these approaches use novel software such as FireStation [15] and LANDIS-II [16] to model and
56  trace wildfire through geographical space and time. Due to the extensive use of software and the
57  need for detailed parameters on various inputs (i.e., fuel mass, characteristics of trees, air humidity,
58  soil moisture, etc.), predictions from such approaches heavily rely on assumptions used in the
59  analysis, are complex to set-up and computationally expensive [17].
60
61  The second approach, statistical methods, complements physics-based methods as they can also
62  be applied to model large/spatial areas while overcoming the simulation complexity. Further,
63  statistical methods can benefit from modern technologies (i.e., geographic information system
64  (GIS), etc.) and can be applied at different scales and resolution/roughness [18,19]. Some of the
65  notable statistical approaches include Poisson regression [20], Monte Carlo simulations [21],
66  weights of evidence [22], etc. Unfortunately, statistical methods could be sensitive to the type of
67  analyzed data and may require numerical manipulation to satisfy convergence criteria – especially
68  for those methods associated with nonlinear nature [23].
69
70  The third and most recent approach is one that leverages advancements in computer sciences. More
71  effectively, machine learning rises as an attractive approach given its good handling of complex
72  and high dimensional data, scalability, and affordability. Machine learning algorithms are applied,
73  tweaked, or created to understand the complex interaction of multi-variables associated with
74  wildfires [24–28]. While the open literature seems to favor the use of such algorithms (i.e., neural
75  networks [29], gradient boosting [30], k-nearest neighbors [31], etc.) and despite the convenience
76  of user-friendly and easy-to-use software that streamlines the development of machine learning by
77  employing pre-defined algorithms and training/validation procedures [32,33], we continue to lack
78  sufficient works on this front.
79
80  A recent look into some of the works in this area clearly shows the merits of applying machine
81  learning to enable modern and accurate prediction of wildfires [34,35]. In fact, Fig. 1 reinforces
82  this notion by presenting the publication trend in article publications between 2000-2020 as
83  obtained from a scientometrics analysis from the open-source scholarly database, *Dimensions*
84  [36,37]. As one can see, this search returns 8,716 papers. This trend of publication is expected to
85  continue to rise in the coming years as it capitalizes on the continued advancements in computer
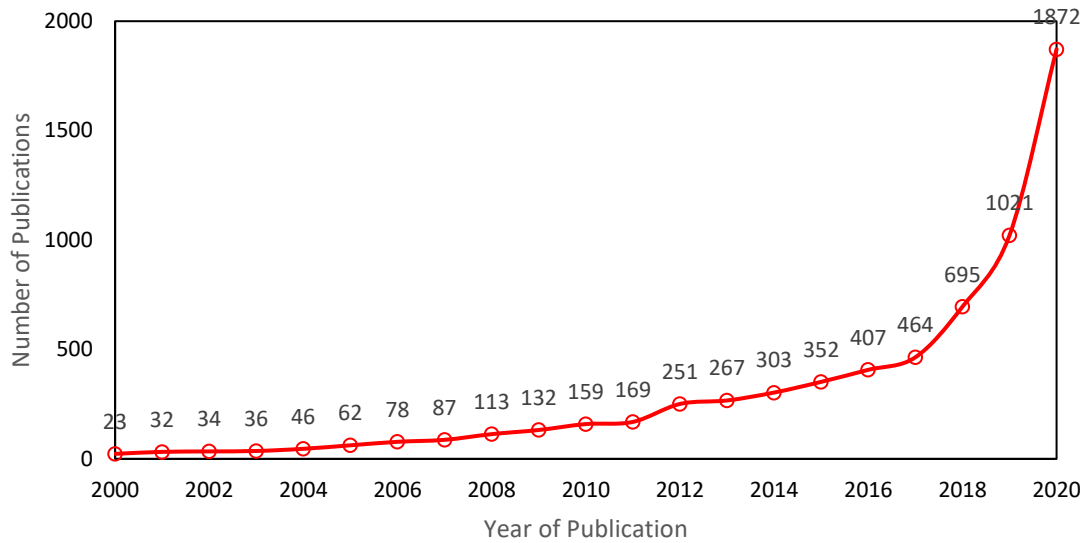86  science.

Fig. 1 Publication trend as obtained from the *Dimensions* academic database [36,37] [Note: Keywords "*Machine learning*" AND "*wildfire*"]

A deeper look into the majority of the noted publications displays that most works showcased the incorporation of one algorithm, often selected from researchers' familiarity with such algorithm (in a similar manner to opting to use a particular simulation software/package). However, our perspective is that reliance on a particular model, while it may produce favorable performance, can still generate biased models that could be overturned. In contrast, we would like to explore the use of multi-algorithmic search to identify suitable machine learning model candidates that can be used in parallel, thereby expanding a researcher's arsenal of predictive tools while adding an additional layer of redundancy.

In addition, the majority of the reviewed works adopt blackbox models that require the user to program and code the machine learning model. This may lead to dependence on computing stations and, most admittedly, a heavy reliance on the user's coding experience. To overcome these hurdles, we present the use of genetic algorithms as a means to augment the blackbox models and derive expressions that can substitute the need for algorithmic simulation. Simply put, the machine learning model will be run once to obtain the predictions, and then these predictions are fitted into an expression (or a form of a mapping function [38]) that can be substituted by hand or via a simple spreadsheet. The user would not need to code a new machine learning model to predict wildfire occurrence since the user can now use the derived expression directly.

In hopes of narrowing this knowledge gap and in pursuit of accelerating research efforts in this area, this work presents a statistical overview of the problem of wildfires and then deep dives to present a dense data-driven ($D^3$) approach that integrates different machine learning algorithms to realize modern wildfire assessment tools that have the capability to predict occurrence and size of wildfires. This approach was developed and validated using measured data points obtained from 1.04 million geo-referenced wildfires between 1992 and 2015 in North America and Europe. The

116  consequence of this $D^3$ analysis demonstrates the suitability and feasibility of exploiting intelligent
117  analysis tools to modernize wildfire disaster planning and optimal resource allocation.
118  **Statistical Overview**
119  Recent statistics have shown that the annual count of worldwide wildfires reaches 200,000 and
120  that these fires burn 3.5–4.5 million km$^2$ (equivalent to 0.86–1.11 billion acres) [2,3]. It is also
121  interesting to note that the average number of large wildfires occurring in the United States
122  increased from 140 to 160 to 250 in the periods of 1980-1989, 1990-1990, and 2000-2012,
123  respectively [3]. The United Nations Office for Disaster Risk Reduction (UNODRR) also supports
124  these statistics and reports that insured losses arising from wildfires around the world in 2017
125  totaled $14 billion, the highest ever in a single year [4]. Trends were most substantial for southern
126  and elevated regions, overlapping with tendencies toward amplified drought severity. The number
127  of large fires increased as well as the total fire area increased per year [80].
128
129  In the United States, wildfires in the Western region are greater and burn more land than their
130  counterparts in other regions. For example, nearly 26,000 wildfires burned approximately 9.5
131  million acres in the Western US, as compared with the over 33,000 fires that burned about 0.7
132  million acres in Eastern regions in 2020 [81]. This horrifying statistical information led to
133  extending the average length of fire season from 5 months in the early 1970s to slightly over seven
134  months nowadays (in the US) [3]. These prime conditions in forests for frequent and intense
135  wildfires as opposed to those experienced in the past decades [7,8]. On the European front, the
136  UNODRR also noted a similar observation and reported how the most damaging fires that occurred
137  in June and October of 2017 fell outside of the traditional fire season (July to September), thus,
138  indicating a shift towards a longer wildfire season [4].
139
140  The surge in the number of wildfires, along with their intensity, is expected to increase associated
141  casualties, property losses, and environmental damage [9,10]. This has been noted in the recent
142  wildfires that broke in North America and Europe over the past few years. For example, the last
143  year was one of the most destructive fire season in California, in which over 7,600 km$^2$ burned,
144  causing over $3.5 billion in damages. It was also in the same season that Mendocino Complex Fire
145  (which burned over 1,860 km2) became the largest single fire in California's history [11]. Within
146  the same timeframe, the Canadian province of British Columbia underwent its largest wildfire,
147  which caused the burning of an area equivalent to 1.3% of the total territory. This fire also led to
148  evacuating 40,000 people [11]. The past few years have also witnessed similar occurrences most
149  notably in Greece [12], Portugal [13], and most recently in Australia. In a nutshell, forest wildfires
150  pose a serious threat to our communities and need to be properly understood, predicted, and
151  mitigated [14].
152
153  To date, over 46 million homes in 70,000 urban, suburban and native communities are at risk of
154  wildfires in the US alone [39]. One should also be cognizant that in a severe fire season, wildfires
155  can burn thousands of structures (e.g., 10,488 buildings in the 2020 California wildfires [40] and
156  5,900 buildings in the 2020 Australian bushfires [41]), and such numbers are expected to rise given
157  the recent inertia for urban development and construction. While statistics on human losses are
158  often accessible [42,43], statistics on animal losses may not be as easily obtained. According to
159  the World Wide Fund for Nature [44] at least 1.25 billion animals were killed (and about 2.75

160 billion were harmed) during the 2020 Australian bushfires <u>alone</u>. At the time of this proposal, we
161 were not able to find a reliable source to report the number of animal losses to US wildfires.
162
163 **Methods**
164 *Development of Databases*
165 In order to effectively apply the $D^3$ approach, there is a need to compile observations on wildfires
166 in order to develop a proper wildfire database. As such, a literature survey was carried out and
167 resulted in identifying two publicly available databases comprising more than 1.04 million geo-
168 referenced wildfires that burned over 359,000 km$^2$ (88.7 million acres) between 1992-2015 in the
169 United States [45] and Portugal [46]. These databases cover well documented wildfires with
170 varying aspects and characteristics. These databases will be used for separate machine learning
171 analyses.
172
173 The first analysis aims to use the first database (to be referred to as the US database here) to create
174 machine learning classifiers that can predict the occurrence and expected size of a given wildfire
175 as a function of a set of variables (outlined below). In the second analysis, the second database
176 (aka. Portugal database) is used to create mathematical expressions that can identify the expected
177 size of a wildfire pending environmental features. Both databases, along with their variables, are
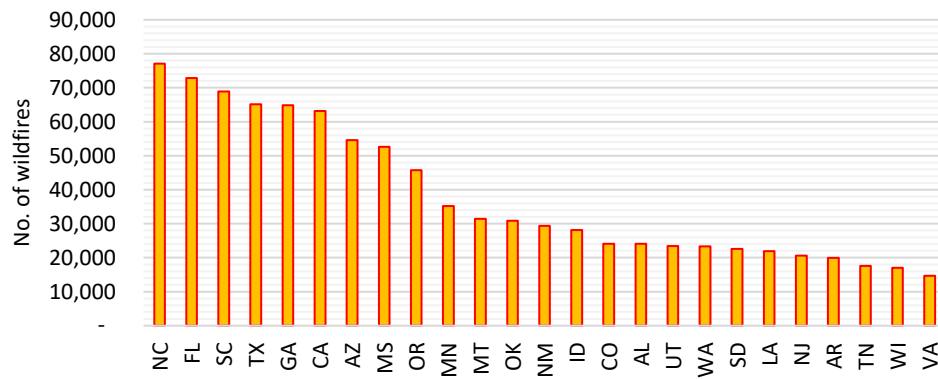178 described below.
179
180 <u>Database on wildfires occurring in the US</u>
181 The first database covers a spatial description of major wildfires that broke out within the United
182 States (US) from 1992 to 2015. The US area covers approximately 9,830,000 km$^2$. These fires
183 were obtained from the reports published by federal, state, and local fire organizations. The
184 observations were transformed to comply with the standards of the National Wildfire Coordinating
185 Group (NWCG) [47]. It is worth noting that this database was initially pre-processed to remove
186 redundant and incomplete observations. After this cleansing procedure, a total of 1.04 million (out
187 of 1.88 million) geo-referenced wildfire records that burned through 88.7 million acres during the
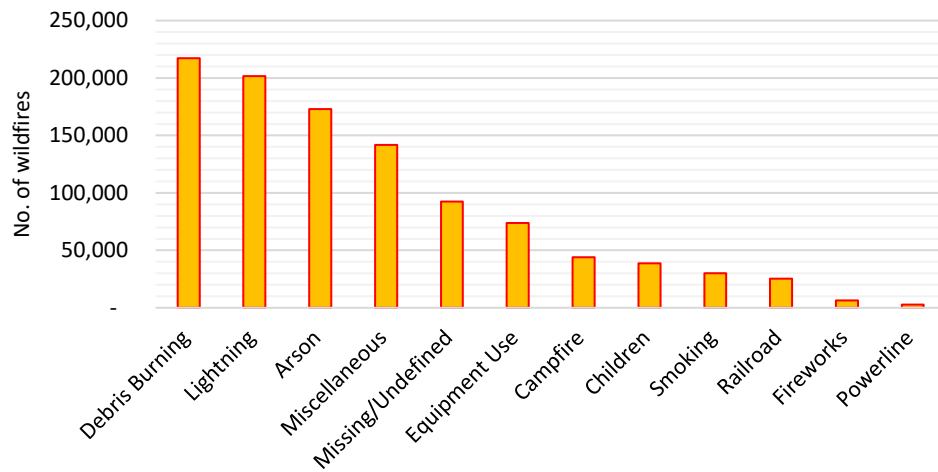188 aforementioned 24-year period were arrived at.
189
190 The same database contains 50 parameters (ranging from geographical location to fire breakout
191 cause and size etc.) and can be freely accessed at [45] or [48]. The database also contains six
192 variables: discovery day of wildfire (a numerical value ranging between 1-365), year of wildfire
193 (a numerical value ranging between 1992-2015), latitude and longitude of wildfire occurrence,
194 wildfire cause (in thirteen categories[*]), and state at which wildfire took place. Further, this database
195 has one predictor as "wildfire size," and this was divided into seven classes that are arranged
196 alphabetically; (A[†]=greater than 0 but less than or equal to 0.25 acres, B=0.26-9.9 acres, C=10.0-
197 99.9 acres, D=100-299 acres, E=300 to 999 acres, F=1,000 to 4,999 acres, and G=5,000+ acres).
198 Figure 2 shows further statistics and the geographic location of wildfires from this database.
199

---

[*] Categories include: arson, camp fire, debris burning, equipment use, fireworks, lightning, miscellaneous, powerline, railroad, smoking, structures, caused by children, and undefined.
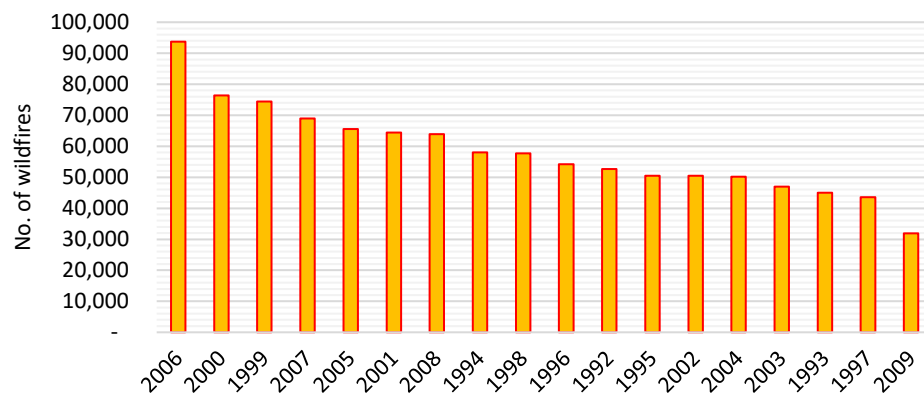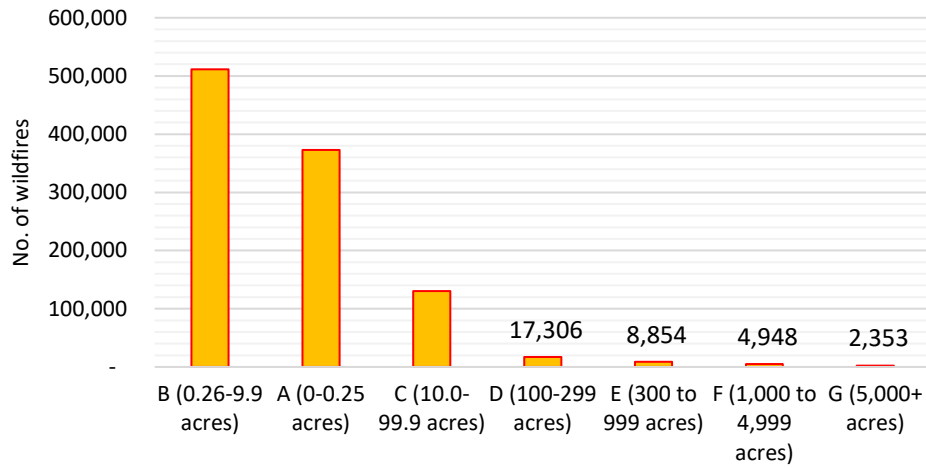[†] Due to its small area, this class was not examined further herein.

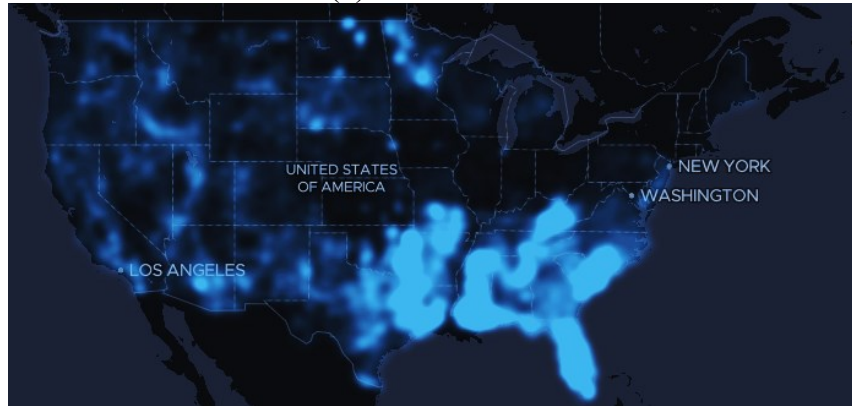(a) No. of wildfires in top 25 states in the US



(b) Cause of wildfires



(c) No. of wildfires per year

(d) Size of wildfires



(e) Spatial distribution of wildfires

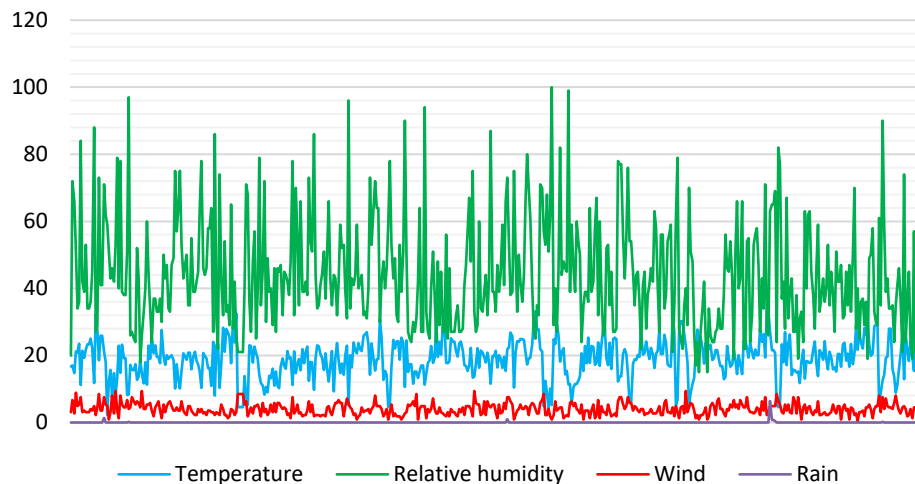Fig. 2 Statistics obtained from the first database (US)

Database on wildfires occurring in Portugal

The second database was prepared by Cortez and Morais [49,50], and this database was collected from the burned areas of Montesinho natural park, located in the northeast region of Portugal. The database contains 517 wildfires that occurred between January 2000 to December 2003.
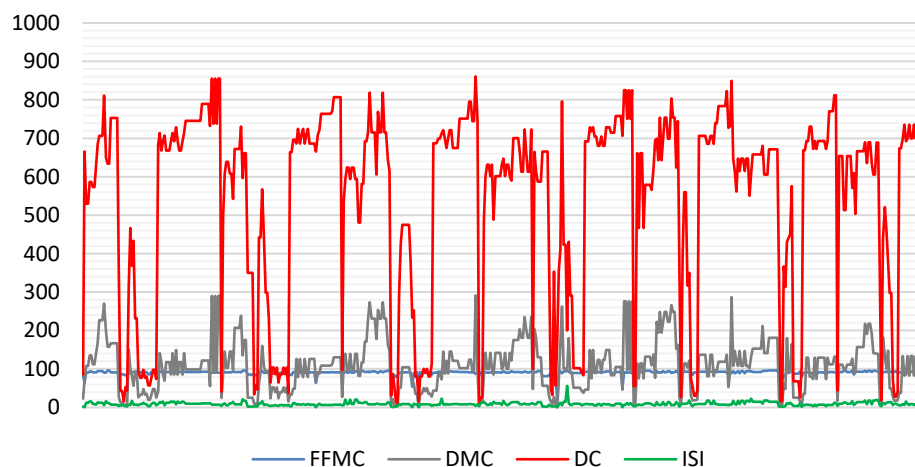
The database comprises the following attributes: geographic features, temporal variables, average monthly weather settings (e.g., temperature, relative humidity, wind speed, rain), as well as distinct weather-based indices. These indices include Fine Fuel Moisture Code (FFMC) which influences ignition and fire spread, Duff Moisture Code (DMC), Drought Code (DC), Initial Spread Index (ISI), which correlates with fire velocity spread, Buildup Index (BUI), and Fire Weather Index (FWI) following the Canadian system for rating fire danger[‡] [51]. As per suggestions laid out by Cortez and Morais (2008, 2007), the following four weather index were used as attributes from this database: FFMC, DMC, DC, and ISI (as the rest of the attributes make up these indices). The

---
[‡] FFMC: moisture content surface litter. DMC and DC represent the moisture content of shallow and deep organic layers (and hence affect fire intensity). BUI reflects upon the availability of fuel. FWI combines ISI with BUI and indicates the magnitude of fire intensity.
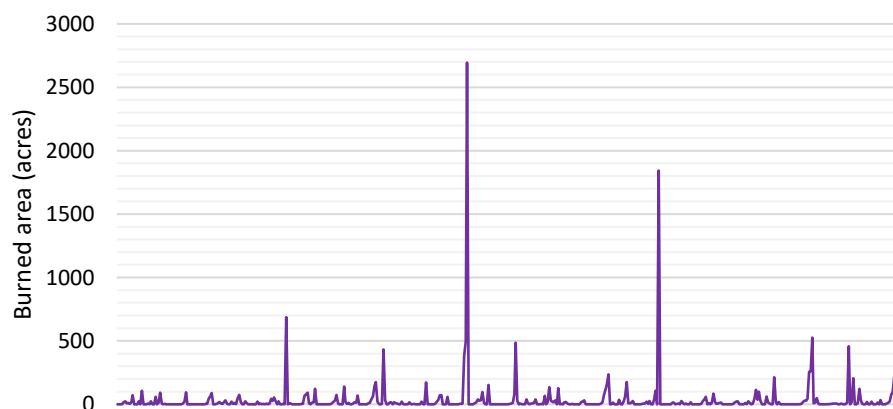
215    predictor in this database was the size of burned areas and ranged from 0-11 km$^2$ (0-2695.5 acres)
216    in a similar breakdown to classes that used in the first database is also followed herein. Figure 3
217    shows further statistics on this database.



(a) Weather conditions



(b) Weather indices

(c) Size of burned areas

218    Fig. 3 Statistics obtained from the second database (Montesinho natural park, Portugal) [Note:
219                        The horizontal represents each individual wildfire]

220    *Description of Machine Learning Algorithms*
221    A dense data-driven ($D^3$) approach that leverages machine learning to uncover hidden patterns
222    within the two datasets described above is presented herein.
223
224    The overarching goal of the $D^3$ approach is to draw conclusions that could be mapped into a
225    solution (or set of solutions) to the wildfire occurrence phenomenon being investigated as part of
226    this study. To attain at such a solution, *key features* governing a wildfire breakout and spread in
227    addition to the *governing relation* that connects these features, are to be identified first. As
228    researchers, our domain knowledge alludes to the fact that a wildfire can breakout once/if several
229    conditions converge. Such conditions may include weather and climate factors (i.e., temperature,
230    humidity, etc.), spatial factors (topology, ignition agents, etc.), and fuel conditions (i.e., vegetation
231    type, heterogeneity of landscape, etc.), among others. The interaction of these features determines
232    how a wildfire can break out and how it will spread, intensify, and potentially be controlled.
233
234    Hence, the rationale behind adopting $D^3$ is that since wildfires behavior can be observed (say in
235    the databases from actual fires as collected in Sec. 3.0), then a governing relation connecting the
236    cycle of a wildfire to its key features can be obtained through $D^3$. Such a relation can be arrived at
237    via machine learning models, as well as could be converted into a mathematical expression via
238    symbolic ML. A systematic analysis of such a magnitude will require special computational
239    treatment, and this is where $D^3$ shines. The following algorithms are used herein; deep learning
240    (DL), decision tree (DT), Stochastic Gradient Descent (SGD), Extreme Gradient Boosted Trees
241    (ExGBT), Logistic regression (LR), and genetic algorithms (GA), and these are further described
242    below.
243
244    <u>Deep learning (DL)</u>
245    The architecture of a DL algorithm follows that of the brain and consists of a similar topology or
246    layout (see Fig. 4). Such topology is characterized by layers. The outermost layer receives the data
247    (representing attributes, say wildfire cause, metrological conditions, etc.) to be analyzed. For this
248    reason, this layer is denoted as the *input layer*. The inputs are then fed into the next set of layers,
249    the *hidden layers*. These layers, or in some cases one layer, house processing units called neurons.
250    The neurons analyze input data via a series of generated weightages (connections). It is through
251    this analysis that the algorithm learns and recognizes any relevant patterns impeded by input data
252    points. This recognition is then mapped into patterns using transformative operations and
253    functions. This aforenoted process sums up the training process of a typical DL algorithm. Once
254    this process passes fitness requirements (whether a pre-defined number of iterations and/or until
255    satisfying a set of fitness metrics), the training is completed, and the algorithm is set into the testing
256    stage.
257
258    The most frequently adopted optimization technique in DL is called Leveberg-Marquard. This
259    technique assesses the error by evaluating the mean squared error (MSE) [52]. In this optimization
260    method, if $z$ is the experimental dataset, then MSE is evaluated using Eq. 1:

$$MSE = \frac{1}{z}\sum_{i=0}^{z}(e_i)^2 = \frac{1}{z}\sum_{i=0}^{z}(m_i - p_i)^2 \tag{1}$$

261 where, $z$ = the total number of datasets, $e_i$ = the error for each input set, $m_i$ = the measured output,
262 and $p_i$ = the estimated output.
263
264 In this development, a pre-sensitivity analysis inferred that adopting a *ReLu* activation function for
265 DL with an initial learning rate of 0.001, 3 hidden layers (with 256, 128, and 64 units) led to
266 achieving an optimal DL architecture. The final outcome within the hidden layers is then
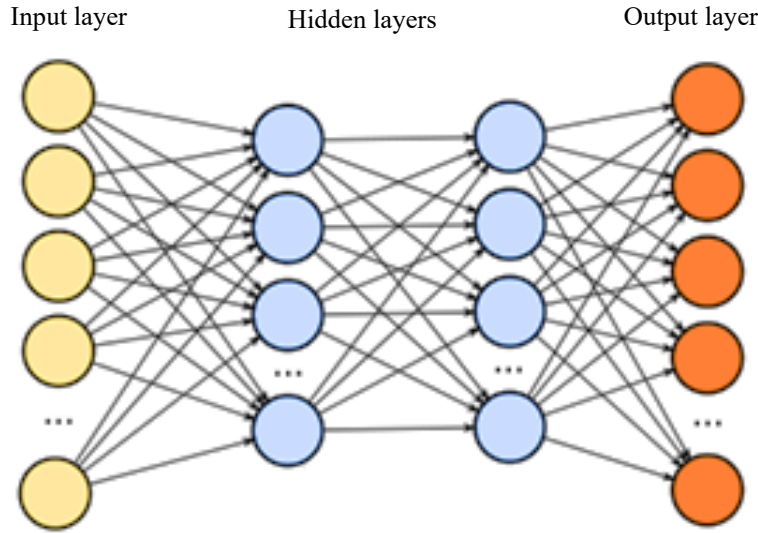267 forwarded to the output layer for visualization.
268



Fig. 4 A typical DL topology

272 Decision tree (DT)
273 A DT algorithm is attractive in a classification-like problem similar to that tackled herein to
274 classify the expected size of a wildfire. A key advantage to DT is its ability to create a diagram-
275 like depiction of all likely decisions [53]. The DT analysis starts by separating the database into
276 branch-like shapes. Then, a random decision tree is created at a root node and then grows into
277 other tree-like components (i.e., leaves, etc.). The created DT was optimally designed to have a
278 maximum depth of 45, with a confidence level = 0.1, minimum leaf size, and maximum size for
279 split equals 2 and 4, respectively [54,55]. A DT analysis may utilize additional measures such as
280 Gini impurity to facilitate the analysis and processing of data points. For example, for a node $t$,
281 Gini index $g(t)$ is defined as [56]:
282
283 $g(t) = \sum_{j\neq i}p(j|t)p(i|t)$ (2)
284
285 where $i$ and $j$ are target field categories.
286
287 $p(j,t) = \frac{p(j,t)}{p(t)}; p(j,t) = \frac{\pi(j)N_j(t)}{N_j};$ and $p(t) = \sum_j p(j,t)$ (3)
288

289    Stochastic Gradient Descent (SGD)
290    SGD regularizes linear models such as support vector machines and logistic regression with
291    stochastic gradient descent (SGD) learning in classification problems [57]. SGD adopts a plain
292    stochastic gradient descent learning process with a loss penalty function as shown in Eq (4).
293

294          $E(w, b) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, f(x_i)) + \alpha R(w)$        (4)
295

296          Where, $L$ is a loss function that measures model, $R$ is a regularization term; $\alpha > 0$ is a non-
297          negative hyperparameter that controls the regularization strength. The developed algorithm
298          was incorporated with *LogLoss* as a loss function, *ElasticNet* as a regularization function,
299          and $\alpha = 2.2 \times 10^{-5}$.
300

301    Extreme Gradient Boosted Trees (ExGBT)
302    This algorithm [58] re-samples data points into a series of tree, where each tree boostraps a sample
303    some data points in each iteration. ExGBT fits each successive tree to the residual errors from all
304    the previous trees and focuses on the most difficult cases to predict to increase its prediction
305    accuracy (see Eq. 5). The developed algorithm incorporated a learning rate of 0.05, a maximum
306    tree depth of 5, a subsample feature of 0.8, and a minimum interval for early stopping of 200.
307

308          $Y = \sum_{k=1}^{M} f_k(x_i), f_k \in F = \{f_x = w_{q(x)}, q : R^p \rightarrow T, w \in R^T\}$    (5)
309

310          where, $M$ is additive functions, $T$ is the number of leaves in the tree, $w$ is a leaf weights
311          vector, $w_i$ is a score on $i$-th leaf, and $q(x)$ represents the structure of each tree that maps an
312          observation to the corresponding leaf index [59].
313

314    Logistic regression (LR)
315    The regularized LR algorithm aims to maximize the likelihood of observing a phenomenon
316    through its capability to estimate coefficients for identified features to measure the comparative
317    influence of each feature on the phenomenon [60]. Therefore, LR is noted to be a successful
318    algorithm for classification problems [61]. LR, and just like other algorithms, can suffer from
319    overfitting. To avoid this, LR's loss function can be modified with a penalty term to
320    shrink/penalize the estimates of the coefficients. L2 penalty is used herein as it is proven effective
321    during a pre-sensitivity study [62]. The used algorithm has a true fit intercept and approximates
322    the multi-linear regression function:
323

324          $logit(p) = \beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n$        (6)
325

326          where, $p$ is the probability of the presence of a phenomenon. The logit transformation is
327          defined as the logged odds:
328

329          $odds = \frac{p}{1-p}$        (7)
330

331          and,

332

333
$$logit(p) = \ln\left(\frac{p}{1-p}\right) + L2_{(penalty)} \tag{8}$$

334 The developed algorithm incorporated a learning rate of 0.05, a maximum tree depth of 5,
335 subsample feature of 0.8, and a minimum interval for early stopping of 200.

336

337 Genetic Algorithms (GA)
338 This algorithm is an evolutionary method that was initially presented by Holland [63] and Koza
339 [64]. GA leverages the concept of the natural selection process to arrive at hidden relations between
340 attributes and expected outcomes in a symbolic format. In GA, a set of expressions are numerically
341 derived from mapping to mathematical expressions that can be used to represent the size of
342 wildfires [65,66].

343

344 The GA analysis starts by creating a populace of arbitrary expressions. These expressions consist
345 of a tree-like formation that houses mathematical operations (addition, multiplication, etc.) and/or
346 mathematical functions (power, log, etc.). In some scenarios, a GA-based expression may also
347 contain conditional and logic functions. The GA-based expression is configured into a tree with
348 hierarchical form, which can then be transformed into a *Karva-expression* as shown in Fig. 5. Once
349 a set of the suitable formula is generated, the algorithm then assesses the fitness (i.e., accuracy) of
350 each expression. Only the fittest expression is then selected for the next stage of analysis. In this
351 stage, the expression is then manipulated by bio-inspired transformative operations, i.e.,
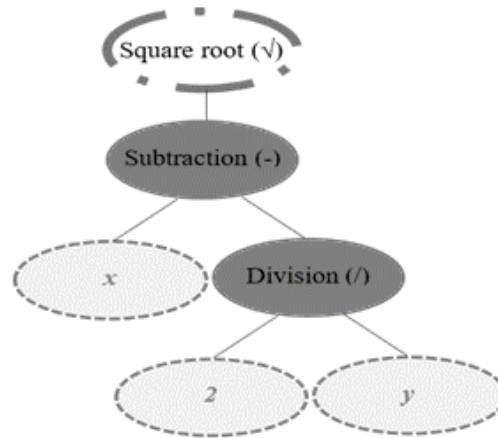352 *reproduction*, *crossover*, and *mutation* [64,67].



353
354 Fig. 5 Representation of a typical GA

355
356 The first, reproduction, the operation ensures that fittest expressions have higher primality of
357 selection to the following stages of analysis. The second, crossover, operation allows the exchange
358 of genetic code (i.e., mathematical functions) between evolved expressions. The third mutation,
359 an operation, can randomly select a function from an expression to mutate into another function
360 [68]. Similar to other algorithms, the GA analysis also terminates once the fitness of a fit
361 expression is achieved or by satisfying a convergence condition. Figure 6 demonstrates a typical
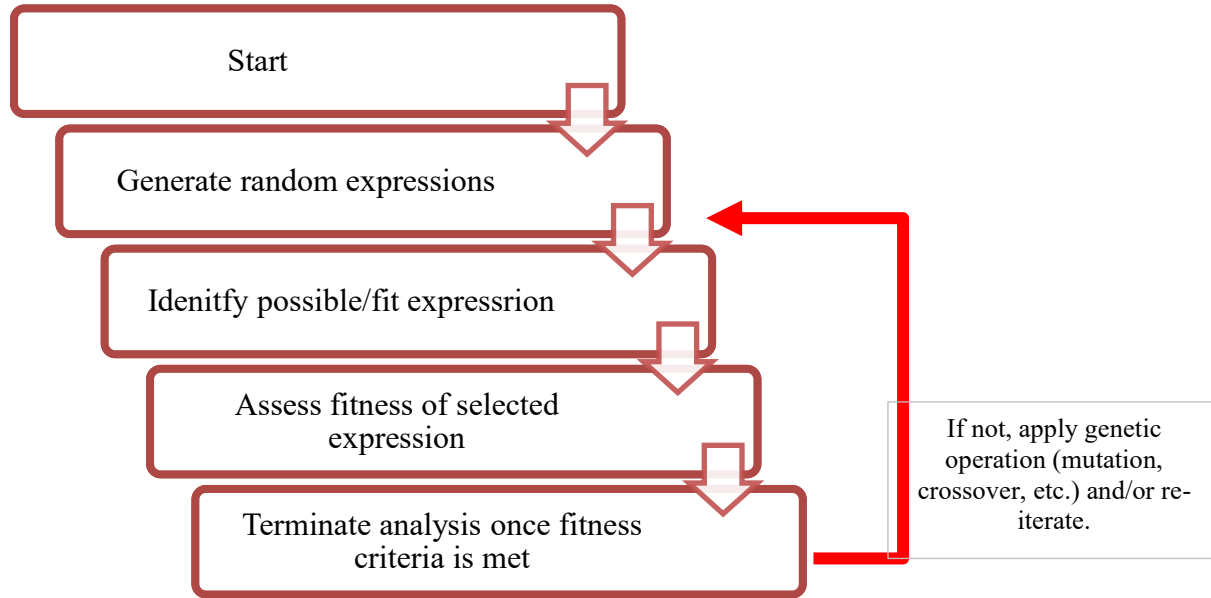362 flow of GA analysis.

363

364
365    Fig. 6 A flowchart of GA analysis
366

367    **Results and Discussion**
368    Now that the databases are compiled, these databases can be analyzed using the selected
369    algorithms. This analysis resembles a classification problem where each machine learning model
370    is expected to correctly label the examined fires (given each fire's set of variables). To start this
371    analysis, first, each dataset was first randomly shuffled to minimize biases arising from a specific
372    wildfire attribute. Then, the ML algorithms are trained using 10 k-fold cross validation [69,70].
373    The analysis was conducted by using the aforenoted algorithms in Matlab [71], Python [72], and
374    GMDH environments [73,74][§].
375
376    The outcome of each machine learning model is then cross checked against that of the ground
377    truth. In this pursuit, specific classification metrics are used. The first is a composite metric known
378    as the confusion matrix, and the second is the *LogLoss* error [75].
379
380    The outcome of the $D^3$ analysis is listed in Table 1 by means of the confusion matrix. This matrix
381    lists the fitness of the applied algorithms in classifying the wildfires as well as two fitness metrics
382    (Accuracy and *LogLoss* error). The Accuracy (ACC) metric evaluates the ratio of a number of
383    correct predictions to the total number of samples.
384    $$ACC = \frac{TP+TN}{P+N} = \frac{TP+TN}{TP+TN+FP+FN}$$    (9)
385
386    where, *P* (denotes the number of real positives), *N* (denotes the number of real negatives), *TP*
387    (denotes the true positives), TN (denotes t the rue negatives), *FP* (denotes the false positives), and

---

[§] In this study, the databases were kept in their original datapoints without any preprocessing to minimize their imbalanced-nature to examine the raw effectiveness of the selected algorithms when applied "as is". A future study will explore different treatment techniques for imbalanced data for the same algorithms. Incorporating such techniques and results can significantly push the size of this paper beyond the limitation of a standard article.

388    *FN* (denotes false negatives) – and hence the composite nature of the matrix. And, LogLoss error
389    metric measures where the prediction input is a probability value.
390    $LLE = -\sum_{c=1}^{M} A_i \log P$                                      (10)
391
392    where, M: number of classes, c: class label, y: binary indicator (0 or 1) if c is the correct
393    classification for a given observation. It is worth noting that an accuracy closer to unity and a
394    LogLoss error close to zero imply favorable predictive performance.
395
396    *Blackbox ML*
397    The first analysis adopts six *blackbox* machine learning algorithms and six variables: discovery
398    day of wildfire, year of wildfire, latitude, and longitude of wildfire occurrence, wildfire cause, and
399    state at which wildfire took place to predict the expected size of the wildfire. A closer look at Table
400    1 shows that all models achieved a comparable accuracy that centers around 80% and LogLoss
401    error ranging between (0.42-0.61).
402
403    These results show a couple of interesting observations. For a start, regardless of the machine
404    learning model type, or search mechanism, it is clear that the adopted models have a good grasp
405    on predicting wildfire occurrences (with minimal tuning, as noted in Sec. 4.). Secondly, the DL,
406    DT, SGD, ExGBT, and LR algorithms achieved comparable performance in accuracy, with DL,
407    DT, and ExGBT ranking top three. Recent works on the front of wildfires have also noted the
408    predictive capacity of such algorithms [76–78]. Thus, we can comfortably say that adopting these
409    three models as independent and redundant models to identify wildfire breakouts can be of merit.
410
411    In all cases, whenever a wildfire class is mistakenly classified with an error larger than 20%, this
412    error is highlighted in red. In addition, it is clear that SGD and LR suffered in predicting individual
413    wildfire sizes. It is clear that the imbalanced nature of the used database on US wildfires adversely
414    affected these algorithms.
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432

433    Table 1 Performance of ML algorithms

| DL | True C | True D | True E | True F | True G | Accuracy | LogLoss |
|---|---|---|---|---|---|---|---|
| Pred. C | 84.3 | 23.8 | 51.7 | 30 | 54.3 | | |
| Pred. D | 7.5 | 76.0 | 0.0 | 0.0 | 0.0 | | |
| Pred. E | 5.4 | 0.0 | 47.8 | 3.9 | 0.0 | 0.825 | 0.418 |
| Pred. F | 1.7 | 0.0 | 0.4 | 49.4 | 0.2 | | |
| Pred. G | 0 | 0.0 | 0.0 | 16.5 | 45.4 | | |

| DT | True C | True D | True E | True F | True G | | |
|---|---|---|---|---|---|---|---|
| Pred. C | 84.4 | 23.6 | 52.6 | 28.7 | 31.8 | | |
| Pred. D | 8 | 76.3 | 0.0 | 0.0 | 0.0 | | |
| Pred. E | 5.4 | 0.0 | 47.3 | 3.8 | 4.5 | 0.822 | 0.433 |
| Pred. F | 1.9 | 0.0 | 0.0 | 50.6 | 40.9 | | |
| Pred. G | 1.1 | 0.0 | 0.0 | 16.7 | 22.7 | | |

| SGD | True C | True D | True E | True F | True G | | |
|---|---|---|---|---|---|---|---|
| Pred. C | 79.5 | 99.4 | 100 | 44.4 | 0.0 | | |
| Pred. D | 10.5 | 0.5 | 0.0 | 0.0 | 0.0 | | |
| Pred. E | 5.4 | 0.0 | 0.0 | 0.0 | 0.0 | 0.794 | 0.544 |
| Pred. F | 3 | 0.0 | 0.0 | 44.4 | 0.0 | | |
| Pred. G | 1.4 | 0.0 | 0.0 | 11.1 | 0.0 | | |

| ExGBT | True C | True D | True E | True F | True G | | |
|---|---|---|---|---|---|---|---|
| Pred. C | 82.8 | 25 | 37.9 | 29.8 | 54.5 | | |
| Pred. D | 8.7 | 74.9 | 0.0 | 0.0 | 0.0 | | |
| Pred. E | 5.6 | 0.0 | 62.1 | 0.0 | 0.0 | 0.818 | 0.451 |
| Pred. F | 1.7 | 0.0 | 0.0 | 49.7 | 1.8 | | |
| Pred. G | 1 | 0.0 | 0.0 | 16.3 | 43.6 | | |

| LR | True C | True D | True E | True F | True G | | |
|---|---|---|---|---|---|---|---|
| Pred. C | 79.5 | 0.0 | 0.0 | 83.3 | 0.0 | | |
| Pred. D | 10.5 | 0.0 | 0.0 | 0.0 | 0.0 | | |
| Pred. E | 5.4 | 0.0 | 0.0 | 0.0 | 0.0 | 0.797 | 0.611 |
| Pred. F | 3 | 0.0 | 0.0 | 0.0 | 0.0 | | |
| Pred. G | 1.4 | 0.0 | 0.0 | 16.6 | 0.0 | | |

434

435    *Symbolic ML*
436    While the above algorithms are used as more of a standard assessment "tool" for the first database,
437    the second analysis uses GA to arrive at *symbolic* expressions that can be substituted into to
438    estimate wildfire class in the second database, given the availability of information regarding four
439    weather-based indices (FFMC, DMC, DC, and ISI). This decision can be rationalized by the notion
440    that GA, unlike the other blackbox models, can output an expression that a user can
441    apply/substitute directly instead of running a coded model.

15

442 For practicality, and since the second database included a number of wildfires that were of low
443 size and intensity, GA-expressions were derived for wildfire classes C, D and E (and greater).
444 Table 2 lists these expressions along with their performance. Table 2 shows that GA managed to
445 properly derive simple expressions that can be used to predict the size of a given wildfire. The
446 predictivity of these expressions was established through the correlation coefficient[**], $R$ – see Eq.
447 11. As one can see, these equations are highly nonlinear and represent the complex nature of the
448 phenomenon on hand.

449

450 $$R = \frac{\sum_{i=1}^{n}(A_i-\overline{A}_i)(P_i-\overline{P}_i)}{\sqrt{\sum_{i=1}^{n}(A_i-\overline{A}_i)^2 \sum_{i=1}^{n}(P_i-\overline{P}_i)^2}}$$ (11)

451 where, $A$ is actual data points, and $P$ is for predicted data points.

452

453 The derived expressions can come in handy in assessing the projected size of a wildfire knowing
454 the magnitude of the previously identified weather indices. Having such tools can come in handy
455 in a variety of scenarios, especially those associated with abrupt wildfire breakout and those that
456 may require a quick judgment call to allocate proper resources to fight the wildfire. For
457 transparency and completion, we expect future works to be able to devise improved expressions
458 with higher accuracy – especially once the dataset is massaged for imbalanced data.

459

460 Table 2 GA-derived expressions to predict wildfire class via weather indices.

| Class | Expression | $R$ |
|-------|-----------|-----|
| C | $Class\ C = Step\ (0.0815FFMC + 0.0208DC + 7.047\tan(1.522DMC) + 3.852\tan(193.6DMC) + \tan(136.5DMC) + \tan(\tan(0.2517DMC)) - \tan(5.203FFMC \times DMC) - 0.02121ISI - 0.0647DMC)$ | 0.82 |
| D | $Class\ D = Logistic\ (0.0326DMC + 0.00656FFMC) + (-1.352 - \frac{\tan(8.933DMC)}{0.02ISI+3.96\times10^{-10}DC\times DMC^3} - \tan(1.524DMC) - \tan(1.522DMC) - \tan(5.446 \times 10^{-5} \times maximum(0.005DMC, -0.7548\tan(1.522DMC))) - 0.01207DC - 0.04806ISI)$ | 0.86 |
| E | $Class\ E = Logistic\ (0.04039DMC + \tan(\tan(3.442DC)) - \cosh(sin(tan(0.005DMC)) \times \cos(0.161DC^2))) - 0.00117DC - 0.1261ISI - 0.1525FFMC)$ | 0.87 |
| In each case, a value of 1.0 indicates that the outcome of a given expression agrees with the identified class. | | |

461

462 *Explainable ML*
463 To supplement the GA analysis and to combat the blackbox nature of the traditional algorithm
464 wherein, for example, the models listed in Table 1 do not articulate how the correct or poor
465 predictions were arrived at, we apply the *explainability* method SHapley Additive exPlanations
466 (SHAP) [79] to the ExGBT to better analyze its performance when applied to the second database.
467 In parallel, the initial phase of analysis noted a need to improve the model given the imbalances

---

[**] We also recommend the adoption of other companion metrics.

468    of the data in the categories (E, F). Thus, the Synthetic Minority Oversampling TEchnique
469    (SMOTE) was applied [80]. SMOTE copies data from the small classes with the lower data point
470    and adds it to the dataset to create a balanced dataset and better resemble or match the number of
471    examples in most classes. Note that such a technique does not affect the model accuracy and only
472    provides the model with different copies of the samples from the same category.
473
474    We start by re-validating this model against the second database. The confusion matrix (see Fig.
475    7) is also used to validate the model. This matrix shows exactly how many errors have been made
476    by the model by comparing the testing set class with the predicted results and the training set class
477    with the predicted results. It is clear that the model was able to classify over 90% of the samples
478    correctly on the training set only; however, for the testing set, the model achieved 84%. Please
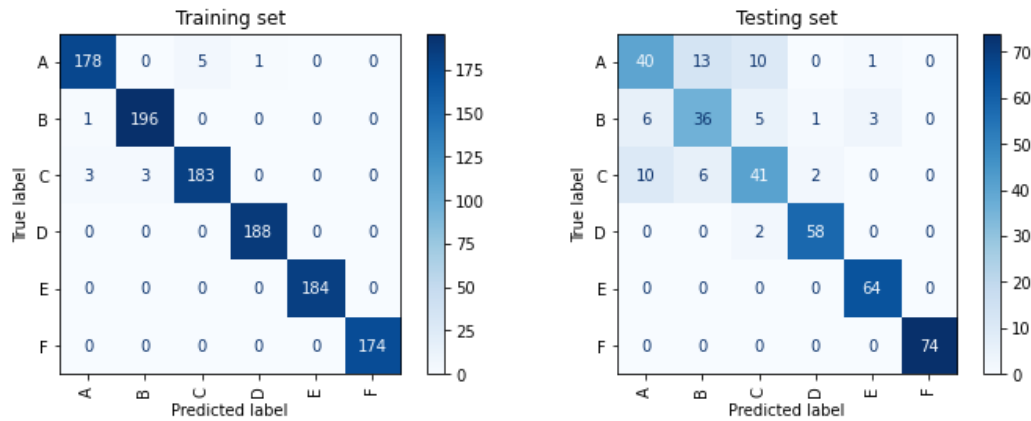479    note that our Python code is provided in the Appendix.



480    Fig. 7 Confusion matrix on the second database
481
482    Now that the model is validated, let us examine the results of our analysis use SHAP *feature*
483    *importance plot,* which shows all the features stacked in horizontal lines representing the effect of
484    each feature on the predicted class of the occurred wildfire (see Fig. 8). For instance, the
485    temperature was found to be the most influential overall. However, for a specific class, that is not
486    true. Taking class E as an example, the most important factor that affects class E wildfire is ISI.
487    Similarly, the temperature was not seen to be o high importance for classes E and B. Also,
488    temperature and wind is the most important factor in predicting the occurrence of class A wildfires
489    but not for class C. Another example is that by looking at the DMC, we can conclude that it highly
490    influences classes E and D. An inherent issue with such a plot is that it does not explain how each
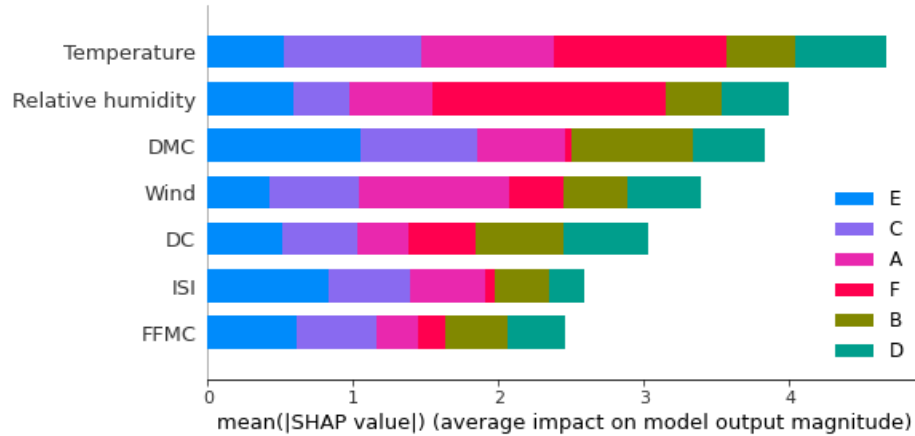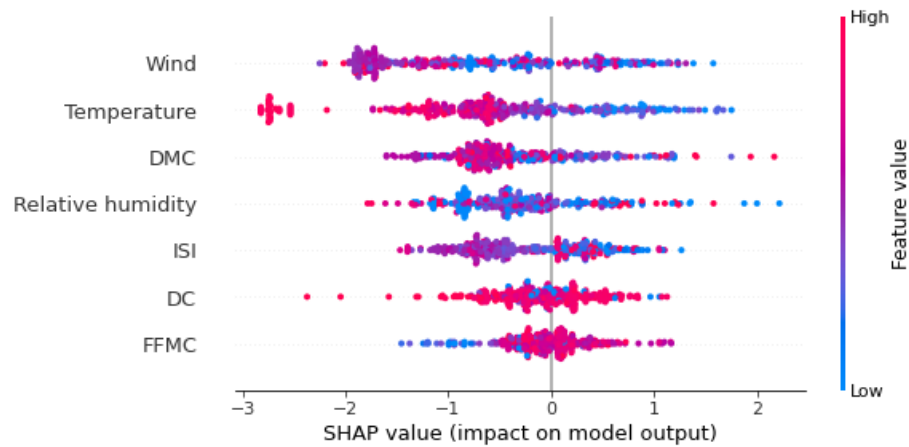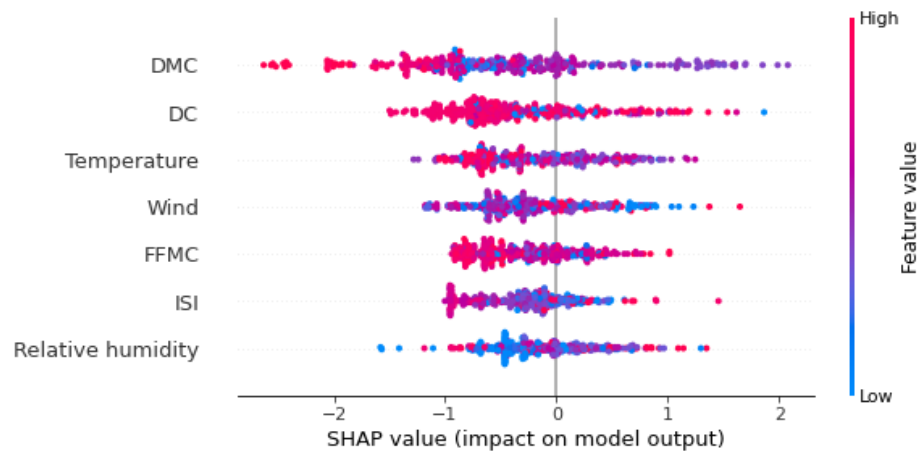491    feature affects the model, positively or negatively.
492

493

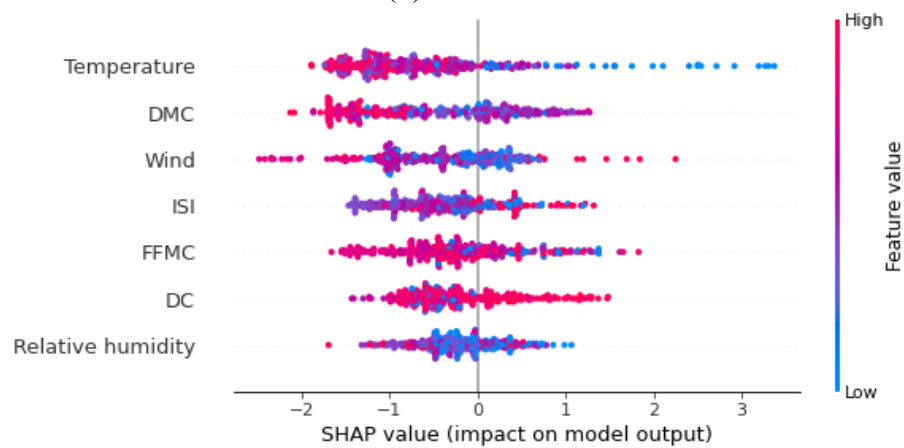Fig. 8 Feature importance plot for all the classes (A-F).

495    To explain multi-class models with two or more classes, one needs to generate new features that
496    are uniquely dedicated to these classes. This can be seen in terms of the *SHAP Summary Plot*,
497    which represents both the feature importance and the direction each feature affects the model's
498    class of wildfire. For instance, Fig. 9a represents class A's feature importance and the direction of
499    each feature's effect in a specific class. One can see those high values of wind temperature and
500    DMC negatively affects the occurrence of a wildfire in class A. the sub-figures shown in Fig. 9
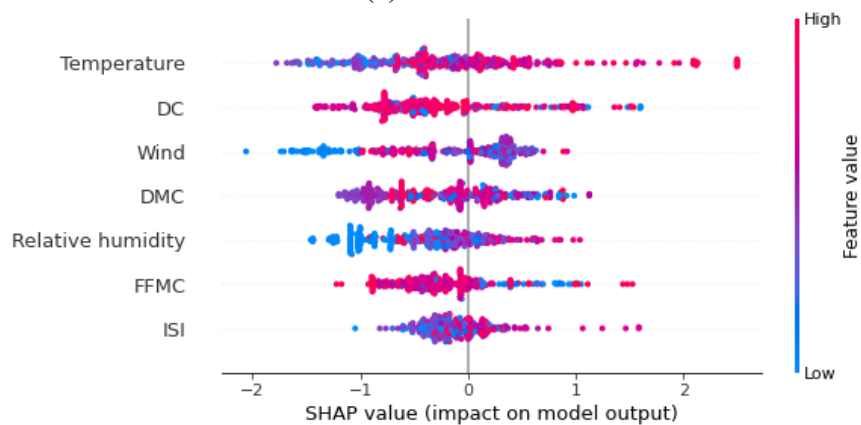501    represent the other classes B, C, D, E, and F, respectively.
502



503
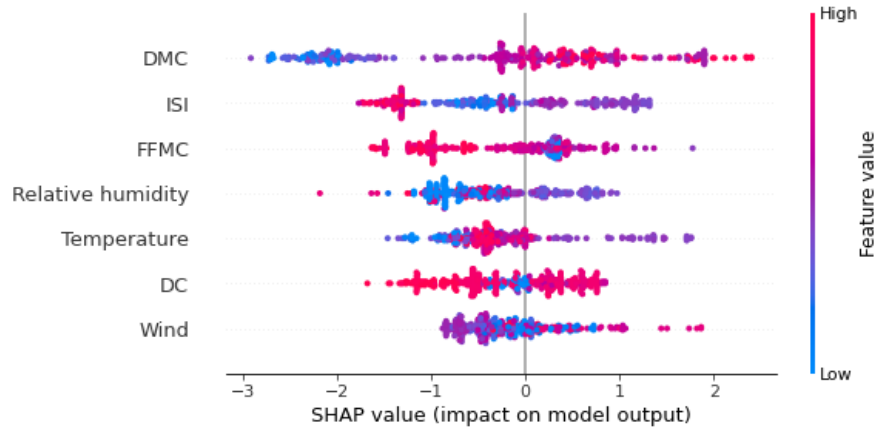504                                 (a)  Class A
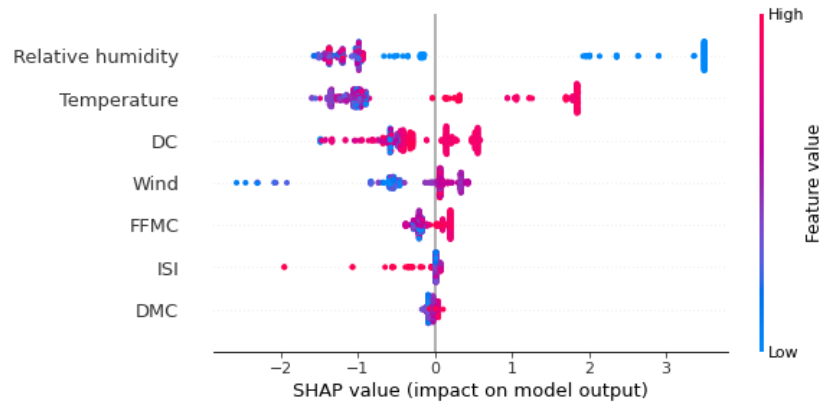
(b) Class B



(c) Class C



(d) Class D

19

(e) Class E



(f) Class F

Fig. 9 SHAP summary plot for all the classes (A-F).

**Closing Remarks on Wildfires Predictive and Classification Tools**

It goes without saying that the accuracy in predictions obtained applying the dense data-driven ($D^3$) approach relies on the presence of information on correctly identified wildfires as well as properly documented parameters such as weather indices. While this study presents results obtained on two databases, one in the US and another one from Portugal, the reader should keep in mind that the presented approach can also be extended to other parts of the world as well as to encompasses a variety of input parameters [81]. This work infers that $D^3$ approach can lead to developing *support tools* that can aid the human-heavy decision making process, and we hope to explore such aspects in future work.

For example, if authorities are preparing for a wildfire season in the state of California, then they could possibly use the DL or DT tool to gauge the expected size of a wildfire, given that they input attributes comprising of: the discovery day of wildfire, latitude, and longitude of expected incident occurrence, and wildfire causes. Based on the outcome of the developed tools, the authorities will be able to estimate how many resources are expected to be allocated and deployed for such wildfires. One instance is given here as an illustration. In this scenario, a wildfire is expected to breakout on the 201st day of a given year in California in a location with longitude and latitude of -123.0 and 40.0, respectively. Based on the analysis from DL and DT tools, these tools show how

20

535   that such wildfire is expected to be primarily of "B" size fire (based on observations from 1992-
536   1994 and 2001-2009) with the potential to grow into a size "C" and beyond (based on observations
537   collected between 1994-2001). While this estimation heavily relies on previous wildfire incidents
538   still, it can be helpful to gauge the size and intensity of future wildfires with ease and in
539   combination with currently used methods that utilize qualitative metrics and methods such as that
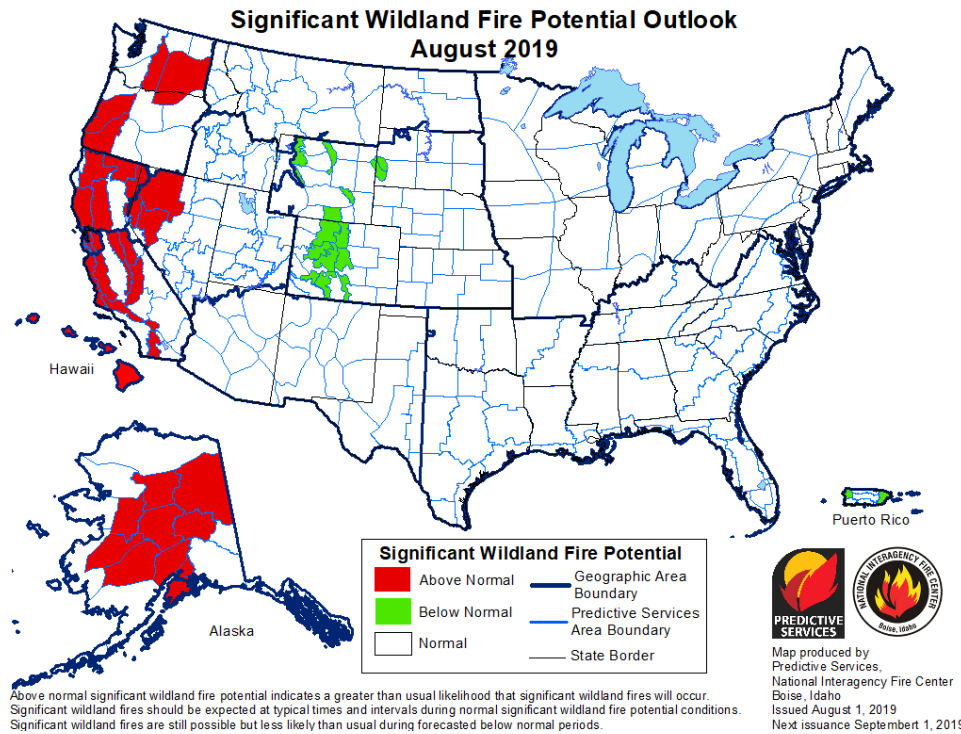540   shown in Fig. 10 [82,83].



541
542   Fig. 10 Predictions for wildland fire potential as obtained by the National Interagency Fire
543   Center for the Month of August 2019

544
545   Similarly, the GA-derived expressions and explainable model for the case of wildfires occurring
546   in Portugal can also be used to estimate the size of wildfires, given insights into weather indices.
547   In this scenario, these expressions can be used to alarm authorities, occupants, and commuters in
548   areas with high vulnerability to wildfire breakouts. This can also turn handy for preparedness and
549   ensuing awareness in particular regions prone to wildfires. In all cases, the developed tools can be
550   used as either predictive methods (i.e., to evaluate if a wildfire is expected to break out) or as
551   classification methods (i.e., to estimate the size of an ongoing wildfire).

552
553   Finally, one should note that machine learning algorithms are adaptable and can improve by
554   collecting new observations for analysis [84–86]. The proposed expressions/tools can also be
555   designed to account for other attributes than those applied here. For example, future works are
556   invited to explore adding attributes covering weather conditions, the magnitude of resource
557   allocations (i.e., number of first responders, evacuation crews, etc.), expected damage to the
558   environment (i.e., air quality, smoke/toxicity levels, fire spread, etc.) as well as to infrastructure
559   (number of collapsed structures or bridges, etc.).
560

561 **Conclusions**
562 This work shows the merit of leveraging computational intelligence in order to develop predictive
563 tools that are able to accurately predict the breakout and size of wildfires. More specifically, this
564 paper explores the integration of deep learning (DL), decision tree (DT), Stochastic Gradient
565 Descent (SGD), Extreme Gradient Boosted Trees (ExGBT), Logistic regression (LR), and genetic
566 algorithms (GA) to gauge expected size of a wildfire given knowledge on existing geographical
567 and environmental condition as well as human-based factors. In lieu of the above, the following
568 conclusions can also be drawn from the findings of this study:

569 • Recent incidents have noted the increasing frequency and intensity of modern wildfires.
570 As such, there is a need to properly predict the occurrence and size of such wildfires.
571 • Deep learning and decision tree algorithms seem to properly capture the wildfire
572 phenomenon with accuracy exceeding 80%. On the other hand, genetic algorithms can
573 also derive appropriate expressions that can be easily implemented into spreadsheets to
574 predict the expected size of wildfires with good accuracy (*R* exceeding 80%). All these
575 tools can potentially be implemented in practice to predict and classify wildfire sizes
576 • The use of explainable and symbolic ML can lead to realizing different types of
577 transparent and equation-like tools to predict wildfires.
578 • The performance of the utilized algorithms herein (together with those to be developed
579 in the near future) can be further enhanced with further training against properly
580 documented wildfire observations, as well as historical information, etc.
581

582 *Conflict of Interest:* The authors declare that he has no conflict of interest.
583

588

589 **References**
590 [1]  A. Jaafari, E.K. Zenner, M. Panahi, H. Shahabi, Hybrid artificial intelligence models
591      based on a neuro-fuzzy system and metaheuristic optimization algorithms for spatial
592      prediction of wildfire probability, Agric. For. Meteorol. (2019).
593      https://doi.org/10.1016/j.agrformet.2018.12.015.
594 [2]  Insurance Information Institute, Facts + Statistics: Wildfires | III, (2020).
595      https://www.iii.org/fact-statistic/facts-statistics-wildfires (accessed August 13, 2019).
596 [3]  Union of Concerned Scientists, Infographic: Western Wildfires and Climate Change |
597      Union of Concerned Scientists, (2013). https://www.ucsusa.org/global-warming/science-
598      and-impacts/impacts/infographic-wildfires-climate-change.html (accessed August 2,
599      2020).
600 [4]  J. Sherry, T. Neale, T.K. McGee, M. Sharpe, Rethinking the maps: A case study of
601      knowledge incorporation in Canadian wildfire risk management and planning, J. Environ.
602      Manage. (2019). https://doi.org/10.1016/j.jenvman.2018.12.116.
603 [5]  T.B. Paveglio, C.M. Edgeley, A.M. Stasiewicz, Assessing influences on social
604      vulnerability to wildfire using surveys, spatial data and wildfire simulations, J. Environ.

605    Manage. (2018). https://doi.org/10.1016/j.jenvman.2018.02.068.

606 [6]  J. Handmer, R. Betts, Estimating the economic, social and environmental impacts of
607    wildfires in Australia Catherine Stephenson1, Environ. Hazards. (2013).
608    https://doi.org/10.1080/17477891.2012.703490.

609 [7]  D. Paton, P.T. Buergelt, S. McCaffrey, F. Tedim, Wildfire Hazards, Risks, and Disasters,
610    2014. https://doi.org/10.1016/C2012-0-03331-5.

611 [8]  K.-M. Hung, L.-M. Chen, T.-W. Chen, A Novel Hierarchical Wildfire Alarm System
612    Based on Vegetation Features, J. Comput. 32 (2021) 137–151.
613    https://doi.org/10.53106/199115992021083204011.

614 [9]  A. Tohidi, N.B. Kaye, Stochastic modeling of firebrand shower scenarios, Fire Saf. J.
615    (2017). https://doi.org/10.1016/j.firesaf.2017.04.039.

616 [10] R.A. Anthenien, S.D. Tse, A. Carlos Fernandez-Pello, On the trajectories of embers
617    initially elevated or lofted by small scale ground fire plumes in high winds, Fire Saf. J.
618    (2006). https://doi.org/10.1016/j.firesaf.2006.01.005.

619 [11] F. Hejazi, I. Toloue, M.S. Jaafar, J. Noorzaei, Optimization of earthquake energy
620    dissipation system by genetic algorithm, Comput. Civ. Infrastruct. Eng. 28 (2013) 796–
621    810. https://doi.org/10.1111/mice.12047.

622 [12] M.A. Haq, Planetscope Nanosatellites Image Classification Using Machine Learning,
623    Comput. Syst. Sci. Eng. (2022). https://doi.org/10.32604/csse.2022.023221.

624 [13] M.A. Haq, P. Baral, S. Yaragal, B. Pradhan, Bulk processing of multi-temporal modis
625    data, statistical analyses and machine learning algorithms to understand climate variables
626    in the indian himalayan region, Sensors. (2021). https://doi.org/10.3390/s21217416.

627 [14] W. Mell, M.A. Jenkins, J. Gould, P. Cheney, A physics-based approach to modelling
628    grassland fires, Int. J. Wildl. Fire. (2007). https://doi.org/10.1071/WF06002.

629 [15] A.M.G. Lopes, M.G. Cruz, D.X. Viegas, Firestation - An integrated software system for
630    the numerical simulation of fire spread on complex topography, Environ. Model. Softw.
631    (2002). https://doi.org/10.1016/S1364-8152(01)00072-X.

632 [16] B.R. Sturtevant, R.M. Scheller, B.R. Miranda, D. Shinneman, A. Syphard, Simulating
633    dynamic and mixed-severity fire regimes: A process-based fire extension for LANDIS-II,
634    Ecol. Modell. (2009). https://doi.org/10.1016/j.ecolmodel.2009.07.030.

635 [17] A. Bar Massada, A.D. Syphard, T.J. Hawbaker, S.I. Stewart, V.C. Radeloff, Effects of
636    ignition location models on the burn patterns of simulated wildfires, Environ. Model.
637    Softw. (2011). https://doi.org/10.1016/j.envsoft.2010.11.016.

638 [18] M. Van Kreveld, Geographic information systems, in: Handb. Discret. Comput. Geom.
639    Third Ed., 2017. https://doi.org/10.1201/9781315119601.

640 [19] M. Finney, I.C. Grenfell, C.W. McHugh, Modeling containment of large wildfires using
641    generalized linear mixed-model analysis, For. Sci. (2009).
642    https://doi.org/10.1093/forestscience/55.3.249.

643 [20] I. Kochi, P.A. Champ, J.B. Loomis, G.H. Donovan, Valuing mortality impacts of smoke
644    exposure from major southern California wildfires, J. For. Econ. (2012).
645    https://doi.org/10.1016/j.jfe.2011.10.002.

646 [21] H. Xue, F. Gu, X. Hu, Data assimilation using sequential monte carlo methods in wildfire
647    spread simulation, ACM Trans. Model. Comput. Simul. (2012).
648    https://doi.org/10.1145/2379810.2379816.

649 [22] T.E. Dilts, J.S. Sibold, F. Biondi, A weights-of-evidence model for mapping the

650    probability of fire occurrence in lincoln county, Nevada, Ann. Assoc. Am. Geogr. (2009).
651    https://doi.org/10.1080/00045600903066540.

652  [23]  J.T. Abatzoglou, T.J. Brown, A comparison of statistical downscaling methods suited for
653    wildfire applications, Int. J. Climatol. (2012). https://doi.org/10.1002/joc.2312.

654  [24]  M. Castelli, L. Vanneschi, A. Popovič, Predicting burned areas of forest fires: An artificial
655    intelligence approach, Fire Ecol. (2015). https://doi.org/10.4996/fireecology.1101106.

656  [25]  M. Rodrigues, J. De la Riva, An insight into machine-learning algorithms to model
657    human-caused wildfire occurrence, Environ. Model. Softw. (2014).
658    https://doi.org/10.1016/j.envsoft.2014.03.003.

659  [26]  D. Tien Bui, Q.T. Bui, Q.P. Nguyen, B. Pradhan, H. Nampak, P.T. Trinh, A hybrid
660    artificial intelligence approach using GIS-based neural-fuzzy inference system and
661    particle swarm optimization for forest fire susceptibility modeling at a tropical area, Agric.
662    For. Meteorol. (2017). https://doi.org/10.1016/j.agrformet.2016.11.002.

663  [27]  F. Zhang, P. Zhao, J. Thiyagalingam, T. Kirubarajan, Terrain-influenced incremental
664    watchtower expansion for wildfire detection, Sci. Total Environ. (2019).
665    https://doi.org/10.1016/j.scitotenv.2018.11.038.

666  [28]  Q. Zhao, S. Yu, F. Zhao, L. Tian, Z. Zhao, Comparison of machine learning algorithms
667    for forest parameter estimations and application for forest quality assessments, For. Ecol.
668    Manage. (2019). https://doi.org/10.1016/j.foreco.2018.12.019.

669  [29]  J.A. Blackard, D.J. Dean, Comparative accuracies of artificial neural networks and
670    discriminant analysis in predicting forest cover types from cartographic variables,
671    Comput. Electron. Agric. (1999). https://doi.org/10.1016/S0168-1699(99)00046-0.

672  [30]  S. Sachdeva, T. Bhatia, A.K. Verma, GIS-based evolutionary optimized Gradient Boosted
673    Decision Trees for forest fire susceptibility mapping, Nat. Hazards. (2018).
674    https://doi.org/10.1007/s11069-018-3256-5.

675  [31]  D.T. Bui, K.T.T. Le, V.C. Nguyen, H.D. Le, I. Revhaug, Tropical forest fire susceptibility
676    mapping at the Cat Ba National Park area, Hai Phong City, Vietnam, using GIS-based
677    Kernel logistic regression, Remote Sens. (2016). https://doi.org/10.3390/rs8040347.

678  [32]  D.E. King, Dlibml: A Machine Learning Toolkit, J. Mach. Learn. Res. (2009).

679  [33]  R. Collobert, K. Kavukcuoglu, C. Farabet, Torch7: A Matlab-like Environment for
680    Machine Learning, 2011.

681  [34]  J. Yao, M. Brauer, S. Raffuse, S.B. Henderson, Machine Learning Approach to Estimate
682    Hourly Exposure to Fine Particulate Matter for Urban, Rural, and Remote Populations
683    during Wildfire Seasons, Environ. Sci. Technol. (2018).
684    https://doi.org/10.1021/acs.est.8b01921.

685  [35]  J. Rogan, J. Franklin, D. Stow, J. Miller, C. Woodcock, D. Roberts, Mapping land-cover
686    modifications over large areas: A comparison of machine learning algorithms, Remote
687    Sens. Environ. (2008). https://doi.org/10.1016/j.rse.2007.10.004.

688  [36]  Dimensions, Dimensions.ai, (2021). https://www.dimensions.ai/.

689  [37]  M. Thelwall, Dimensions: A competitor to Scopus and the Web of Science?, J. Informetr.
690    (2018). https://doi.org/10.1016/j.joi.2018.03.006.

691  [38]  M.Z. Naser, Mapping functions: A physics-guided, data-driven and algorithm-agnostic
692    machine learning approach to discover causal and descriptive expressions of engineering
693    phenomena, Measurement. 185 (2021) 110098.
694    https://doi.org/10.1016/J.MEASUREMENT.2021.110098.

695   [39]   D. Bailey, WUI Fact Sheet, 2013.
696          http://www.iawfonline.org/pdf/WUI_Fact_Sheet_08012013.pdf.
697   [40]   CALFIRE, 2020 Fire Season - CALFIRE, (2020). https://www.fire.ca.gov/incidents/2020/
698          (accessed November 16, 2020).
699   [41]   D. Tin, A.J. Hertelendy, G.R. Ciottone, What we learned from the 2019–2020 Australian
700          Bushfire disaster: Making counter-terrorism medicine a strategic preparedness priority,
701          Am. J. Emerg. Med. (2020). https://doi.org/10.1016/j.ajem.2020.09.069.
702   [42]   National Interagency Fire Center, Wildland Fire Fatalities by Year, 2017.
703          https://www.nifc.gov/safety/safety_documents/Fatalities-by-Year.pdf (accessed
704          November 16, 2020).
705   [43]   BBC, Australia fires: A visual guide to the bushfire crisis - BBC News, BBC News.
706          (2020). https://www.bbc.com/news/world-australia-50951043 (accessed November 16,
707          2020).
708   [44]   WWF, Australian Bushfires - WWF-Australia - WWF-Australia, (2020).
709          https://www.wwf.org.au/what-we-do/bushfire-recovery/bushfires#gs.lm0sfk (accessed
710          November 16, 2020).
711   [45]   K. Short, Spatial wildfire occurrence data for the United States, 1992-2015, 2017.
712          https://doi.org/https://doi.org/10.2737/RDS-2013-0009.4.
713   [46]   P. Cortez, A. Morais, A Data Mining Approach to Predict Forest Fires using
714          Meteorological Data, in: Proc. 13th Port. Conf. Artif. Intell., 2007.
715   [47]   The National Wildfire Coordinating Group, NWCG | NWCG is an operational group
716          designed to coordinate programs of the participating wildfire management agencies.,
717          (n.d.). https://www.nwcg.gov/ (accessed August 13, 2019).
718   [48]   R. Tatman, 1.88 Million US Wildfires, Kaggle.Com. (2017).
719          https://www.kaggle.com/rtatman/188-million-us-wildfires (accessed August 6, 2019).
720   [49]   P. Cortez, A. Morais, Forest Fires Data Set, UCI Mach. Learn. Repos. (2008).
721          https://archive.ics.uci.edu/ml/datasets/forest+fires (accessed August 6, 2019).
722   [50]   P. Cortez, A. Morais, Forest Fires Data Set Portugal | Kaggle, (2007).
723          https://www.kaggle.com/datasets/ishandutta/forest-fires-data-set-portugal (accessed July
724          11, 2022).
725   [51]   B.J. Stocks, T.J. Lynham, B.D. Lawson, M.E. Alexander, C.E. Van Wagner, R.S.
726          McAlpine, D.E. Dubé, The Canadian Forest Fire Danger Rating System: An Overview,
727          For. Chron. (1989). https://doi.org/10.5558/tfc65450-6.
728   [52]   M. Naser, G. Abu-Lebdeh, R. Hawileh, Analysis of RC T-beams strengthened with CFRP
729          plates under fire loading using ANN, Constr. Build. Mater. 37 (2012) 301–309.
730          https://doi.org/10.1016/j.conbuildmat.2012.07.001.
731   [53]   S.R. Safavian, D. Landgrebe, A survey of decision tree classifier methodology, IEEE
732          Trans. Syst. Man. Cybern. 21 (1991) 660–674. https://doi.org/10.1109/21.97458.
733   [54]   SciKit, Decision Tree, (2020). https://scikit-
734          learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesClassifier.html?highlight
735          =decision tree#sklearn.ensemble.ExtraTreesClassifier.decision_path (accessed January 22,
736          2021).
737   [55]   D. Che, Q. Liu, K. Rasheed, X. Tao, Decision Tree and Ensemble Learning Algorithms
738          with Their Applications in Bioinformatics, in: Springer, New York, NY, 2011: pp. 191–
739          199. https://doi.org/10.1007/978-1-4419-7046-6_19.

740  [56]  J.-S.S. Chou, C.-F.F. Tsai, A.-D.D. Pham, Y.-H.H. Lu, Machine learning in concrete
741        strength simulations: Multi-nation data analytics, Constr. Build. Mater. 73 (2014) 771–
742        780. https://doi.org/10.1016/j.conbuildmat.2014.09.054.
743  [57]  SciKit, SGD Classifier , (2020). https://scikit-
744        learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html (accessed
745        January 22, 2021).
746  [58]  Y. Freund, R.E. Schapire, A Decision-Theoretic Generalization of On-Line Learning and
747        an Application to Boosting, J. Comput. Syst. Sci. (1997).
748        https://doi.org/10.1006/jcss.1997.1504.
749  [59]  Gradient boosted tree (GBT), (2019). https://software.intel.com/en-us/daal-programming-
750        guide-details-24 (accessed April 9, 2019).
751  [60]  SciKit, Logistic Regression , (2020). https://scikit-
752        learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html?highlig
753        ht=logistic regression#sklearn.linear_model.LogisticRegression (accessed January 22,
754        2021).
755  [61]  Y. Jafari Goldarag, A. Mohammadzadeh, A.S. Ardakani, Fire Risk Assessment Using
756        Neural Network and Logistic Regression, J. Indian Soc. Remote Sens. (2016).
757        https://doi.org/10.1007/s12524-016-0557-6.
758  [62]  F. Bunea, Honest variable selection in linear and logistic regression models via l1 and
759        l1+l2 penalization, Electron. J. Stat. (2008). https://doi.org/10.1214/08-EJS287.
760  [63]  D.E. Goldberg, J.H. Holland, Genetic Algorithms and Machine Learning, Mach. Learn.
761        (1988). https://doi.org/10.1023/A:1022602019183.
762  [64]  J.R. Koza, A genetic approach to finding a controller to back up a tractor-trailer truck, in:
763        Proc. 1992 Am. Control Conf., 1992.
764  [65]  M.Z. Naser, Heuristic machine cognition to predict fire-induced spalling and fire
765        resistance of concrete structures, Autom. Constr. 106 (2019) 102916.
766        https://doi.org/10.1016/J.AUTCON.2019.102916.
767  [66]  M.Z. Naser, AI-based cognitive framework for evaluating response of concrete structures
768        in extreme conditions, Eng. Appl. Artif. Intell. 81 (2019) 437–449.
769        https://www.sciencedirect.com/science/article/pii/S0952197619300466 (accessed April 1,
770        2019).
771  [67]  A.H. Alavi, A.H. Gandomi, M.G. Sahab, M. Gandomi, Multi expression programming: A
772        new approach to formulation of soil classification, Eng. Comput. 26 (2010) 111–118.
773        https://doi.org/10.1007/s00366-009-0140-7.
774  [68]  C. Ferreira, Gene Expression Programming: a New Adaptive Algorithm for Solving
775        Problems, Ferreira, C. (2001). Gene Expr. Program. a New Adapt. Algorithm Solving
776        Probl. Complex Syst. 13. (2001). https://www.semanticscholar.org/paper/Gene-
777        Expression-Programming%3A-a-New-Adaptive-for-
778        Ferreira/3232b2a24c2584ca8e81cb5bf6f55aef34f0aefe (accessed March 16, 2019).
779  [69]  S.L. Oh, V. Jahmunah, C.P. Ooi, R.S. Tan, E.J. Ciaccio, T. Yamakawa, M. Tanabe, M.
780        Kobayashi, U. Rajendra Acharya, Classification of heart sound signals using a novel deep
781        WaveNet model, Comput. Methods Programs Biomed. (2020).
782        https://doi.org/10.1016/j.cmpb.2020.105604.
783  [70]  J. Abawajy, A. Kelarev, X. Yi, H.F. Jelinek, Minimal ensemble based on subset selection
784        using ECG to diagnose categories of CAN, Comput. Methods Programs Biomed. (2018).

785        https://doi.org/10.1016/j.cmpb.2018.01.019.

786  [71]  D. Searson, D. Searson, GPTIPS Genetic Programming &amp; Symbolic Regression for
787        MATLAB User Guide, (2009).
788        http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.177.494 (accessed January 22,
789        2019).

790  [72]  F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
791        P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M.
792        Brucher, M. Perrot, É. Duchesnay, E. Duchesnay, Scikit-learn: Machine learning in
793        Python, J. Mach. Learn. Res. 12 (2011) 2825–2830.

794  [73]  GMDH, GMDH Shell DS, (2019). https://gmdhsoftware.com/.

795  [74]  M. Sheikholeslami, F. Bani Sheykholeslami, S. Khoshhal, H. Mola-Abasia, D.D. Ganji,
796        H.B. Rokni, Effect of magnetic field on Cu-water nanofluid heat transfer using GMDH-
797        type neural network, Neural Comput. Appl. (2014). https://doi.org/10.1007/s00521-013-
798        1459-y.

799  [75]  M.Z. Naser, A.H. Alavi, · Amir, H. Alavi, Error Metrics and Performance Fitness
800        Indicators for Artificial Intelligence and Machine Learning in Engineering and Sciences,
801        Archit. Struct. Constr. 1 (2021). https://doi.org/https://doi.org/10.1007/s44150-021-00015-
802        8.

803  [76]  D. Radke, A. Hessler, D. Ellsworth, Firecast: Leveraging deep learning to predict wildfire
804        spread, in: IJCAI Int. Jt. Conf. Artif. Intell., 2019. https://doi.org/10.24963/ijcai.2019/636.

805  [77]  Z. Langford, J. Kumar, F. Hoffman, Wildfire mapping in interior alaska using deep neural
806        networks on imbalanced datasets, in: IEEE Int. Conf. Data Min. Work. ICDMW, 2019.
807        https://doi.org/10.1109/ICDMW.2018.00116.

808  [78]  P. Jain, S.C.P. Coogan, S.G. Subramanian, M. Crowley, S.W. Taylor, M.D. Flannigan, A
809        review of machine learning applications in wildfire science and management, Environ.
810        Rev. (2020). https://doi.org/10.1139/er-2020-0019.

811  [79]  S.M. Lundberg, S.I. Lee, A unified approach to interpreting model predictions, in: Adv.
812        Neural Inf. Process. Syst., 2017.

813  [80]  N. V. Chawla, K.W. Bowyer, L.O. Hall, W.P. Kegelmeyer, SMOTE: Synthetic minority
814        over-sampling technique, J. Artif. Intell. Res. (2002). https://doi.org/10.1613/jair.953.

815  [81]  Y.O. Sayad, H. Mousannif, H. Al Moatassime, Predictive modeling of wildfires: A new
816        dataset and machine learning approach, Fire Saf. J. (2019).
817        https://doi.org/10.1016/J.FIRESAF.2019.01.006.

818  [82]  G.A. Trunfio, Predicting Wildfire Spreading Through a Hexagonal Cellular Automata
819        Model, in: 2004. https://doi.org/10.1007/978-3-540-30479-1_40.

820  [83]  C. Vega-García, E. Chuvieco, Applying local measures of spatial heterogeneity to
821        Landsat-TM images for predicting wildfire occurrence in Mediterranean landscapes,
822        Landsc. Ecol. (2006). https://doi.org/10.1007/s10980-005-4119-5.

823  [84]  Y. Xie, M. Peng, Forest fire forecasting using ensemble learning approaches, Neural
824        Comput. Appl. (2019). https://doi.org/10.1007/s00521-018-3515-0.

825  [85]  A. Erdil, E. Arcaklioglu, The prediction of meteorological variables using artificial neural
826        network, Neural Comput. Appl. (2013). https://doi.org/10.1007/s00521-012-1210-0.

827  [86]  M. Liu, S.M. Lo, B.Q. Hu, C.M. Zhao, On the use of fuzzy synthetic evaluation and
828        optimal classification for computing fire risk ranking of buildings, Neural Comput. Appl.
829        (2009). https://doi.org/10.1007/s00521-009-0244-4.

830   **Appendix**
831   Here is our code. The database can be found at [46] and [50].
832
833   ```
      import sklearn
834   from sklearn.model_selection import train_test_split
835   import pandas as pd
836   import numpy as np
837   import shap
838   import xgboost as xgb
839   from matplotlib import pyplot
840   from sklearn.metrics import accuracy_score
841   from sklearn.metrics import plot_confusion_matrix
842   from sklearn.metrics import accuracy_score
843   from imblearn.over_sampling import SMOTE
844   from imblearn.over_sampling import BorderlineSMOTEIn
845   wildfire=pd.read_excel('Portugal ABCDEF.xlsx')
846   wildfire
      ```
847

| | FFMC | DMC | DC | ISI | Temperature | Relative humidity | Wind | Class |
|---|---|---|---|---|---|---|---|---|
| **0** | 83.0 | 23.3 | 85.3 | 2.3 | 16.7 | 20 | 3.1 | A |
| **1** | 63.5 | 70.8 | 665.3 | 0.8 | 17.0 | 72 | 6.7 | A |
| **2** | 90.1 | 108.0 | 529.8 | 12.5 | 14.7 | 66 | 2.7 | A |
| **3** | 94.8 | 227.0 | 706.7 | 12.0 | 23.3 | 34 | 3.1 | A |
| **4** | 94.8 | 227.0 | 706.7 | 12.0 | 25.0 | 36 | 4.0 | A |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| **512** | 92.5 | 121.1 | 674.4 | 8.6 | 18.2 | 46 | 1.8 | E |
| **513** | 91.0 | 129.5 | 692.6 | 7.0 | 18.8 | 40 | 2.2 | E |
| **514** | 89.2 | 103.9 | 431.6 | 6.4 | 22.6 | 57 | 4.9 | E |
| **515** | 94.8 | 222.4 | 698.6 | 13.9 | 27.5 | 27 | 4.9 | F |
| **516** | 92.5 | 121.1 | 674.4 | 8.6 | 25.1 | 27 | 4.0 | F |

848                           517 rows × 8 columns

849   ```
      x=wildfire.drop(['Class'],axis=1)
850   y=wildfire['Class']
851   oversampled = SMOTE(sampling_strategy='auto',
852                  random_state=5,k_neighbors = 1
853                  )
854   x, y = oversampled.fit_resample(x, y)
855   x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.250,random_state=10)
856
857   y_train.value_counts()
858   wildfire.info()
859   wildfire.isnull().any()
      ```
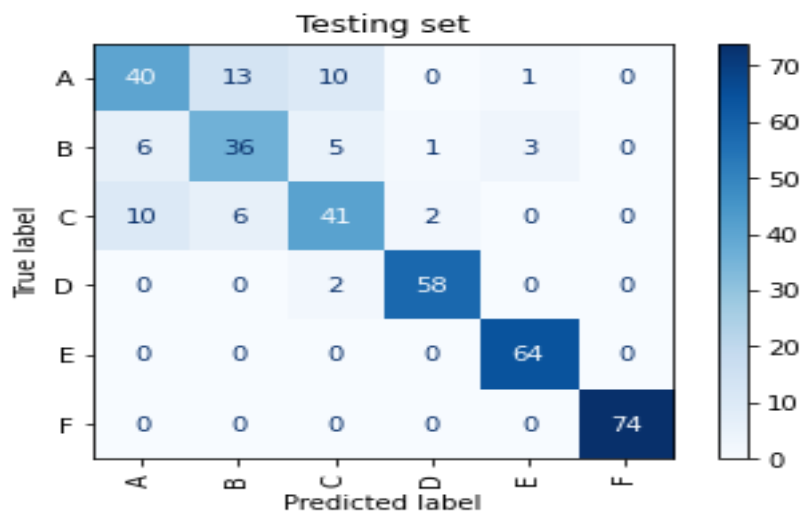
```
860
861    FFMC            False
862    DMC             False
863    DC          False
864    ISI         False
865    Temperature        False
866    Relative humidity    False
867    Wind            False
868    Class           False
869    dtype: bool
870
871    xgbc=xgb.XGBClassifier(objective='multi:softprob',
872                    learning_rate =0.6,
873                     n_estimators=800,
874                     max_depth=6,
875                     min_child_weight=0,
876                     gamma=0.2,
877                     subsample=0.9,
878                     colsample_bytree=0.7,
879                     nthread=40,
880                      seed=230)
881    xgbc.fit(x_train,y_train)
882    predictions = xgbc.predict(x_test)
883    accuracy = accuracy_score(y_test, predictions)
884    print("Accuracy: %.2f%%" % (accuracy * 100.0))
885
886    Accuracy: 84.14%
887
888    class_names = ['A', 'B', 'C','D','E','F']
889    disp = plot_confusion_matrix(xgbc, x_test, y_test, display_labels=class_names, cmap=pyplot.cm.Blues,
890    xticks_rotation='vertical')
891    pyplot.title('Testing set')
```
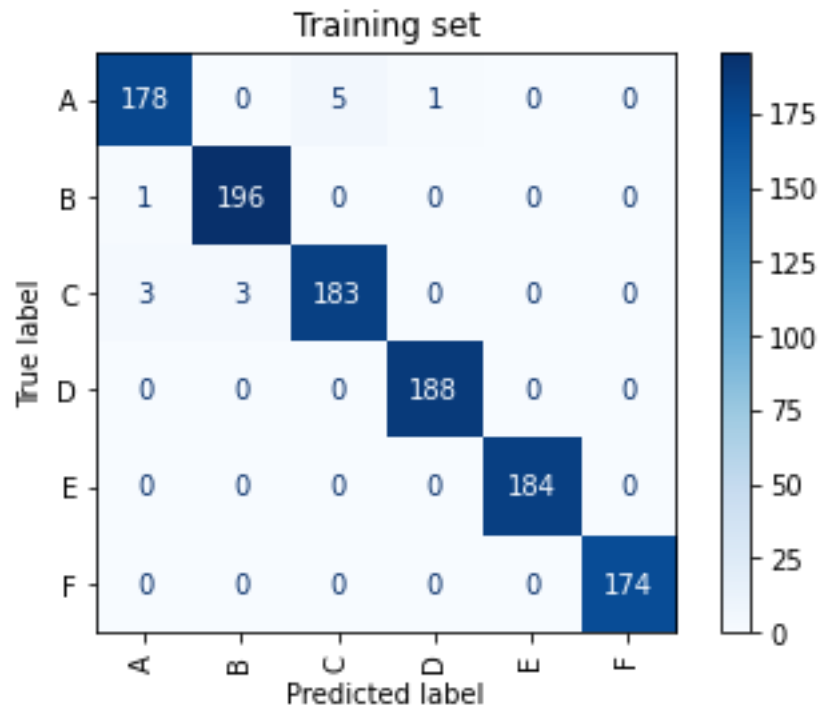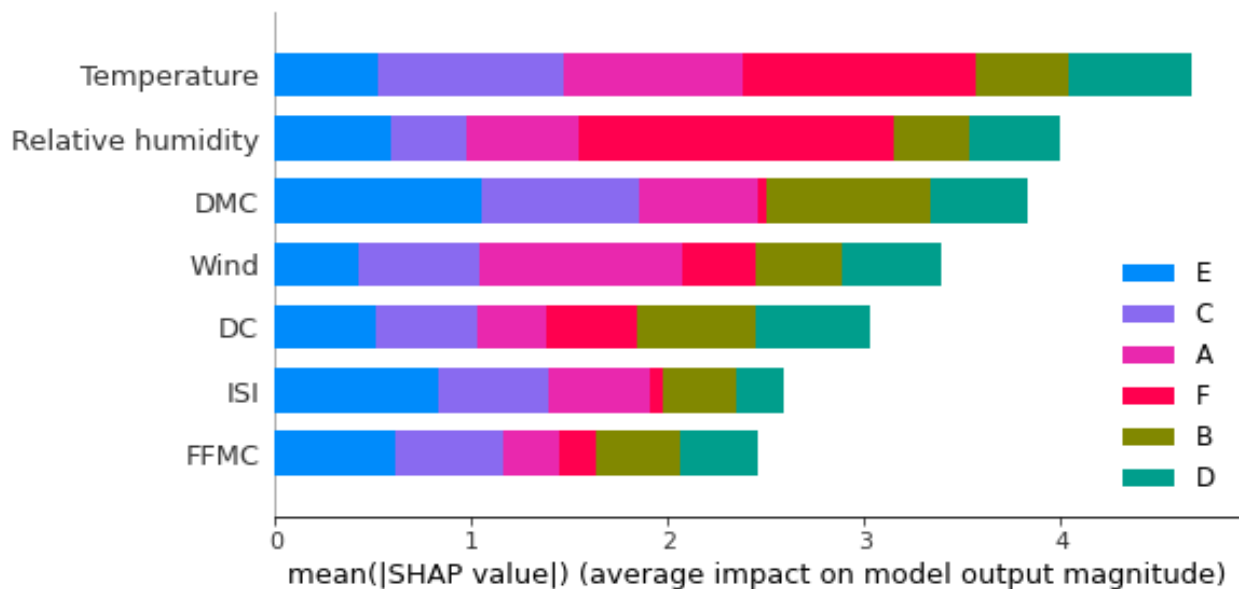


892

```
893    disp = plot_confusion_matrix(xgbc, x_train, y_train, display_labels=class_names, cmap=pyplot.cm.Blues,
894    xticks_rotation='vertical')
895    pyplot.title('Training set')
896
```
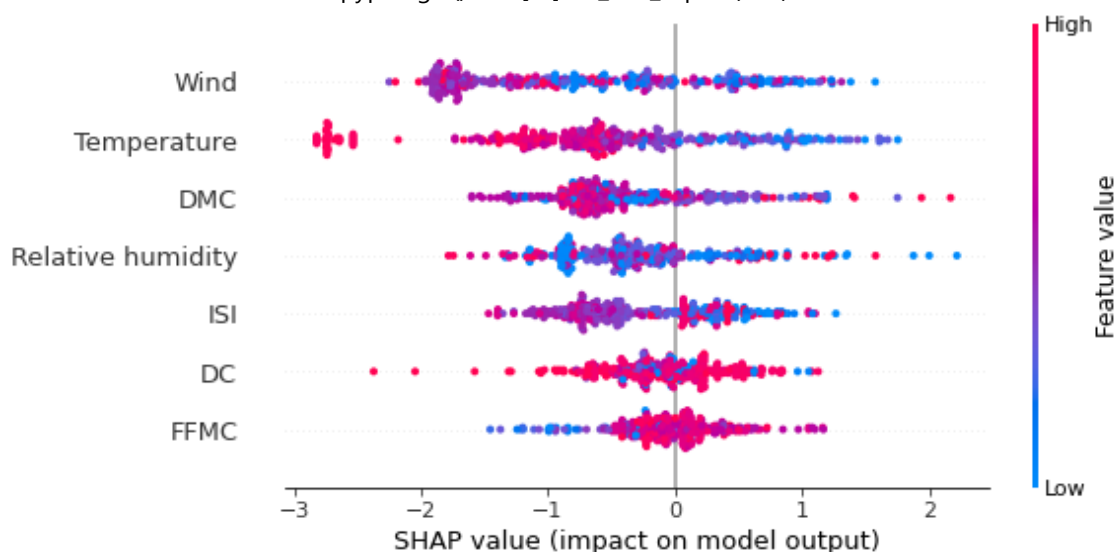


```
897
898    shap_values = shap.TreeExplainer(xgbc).shap_values(x_test)
899            shap.summary_plot(shap_values, x_test,class_names = class_names, plot_type='bar')
```
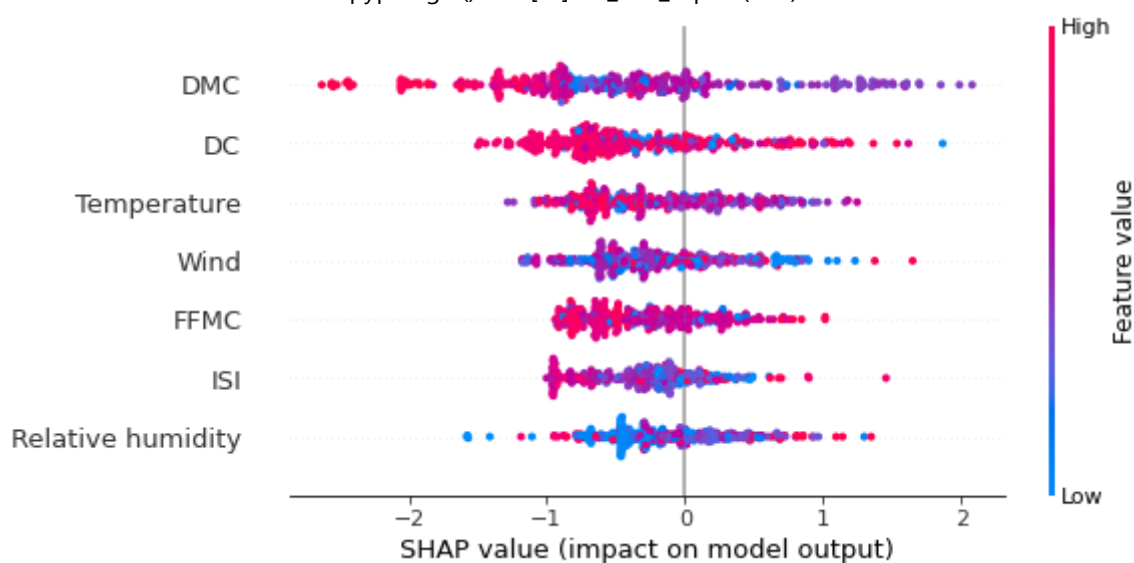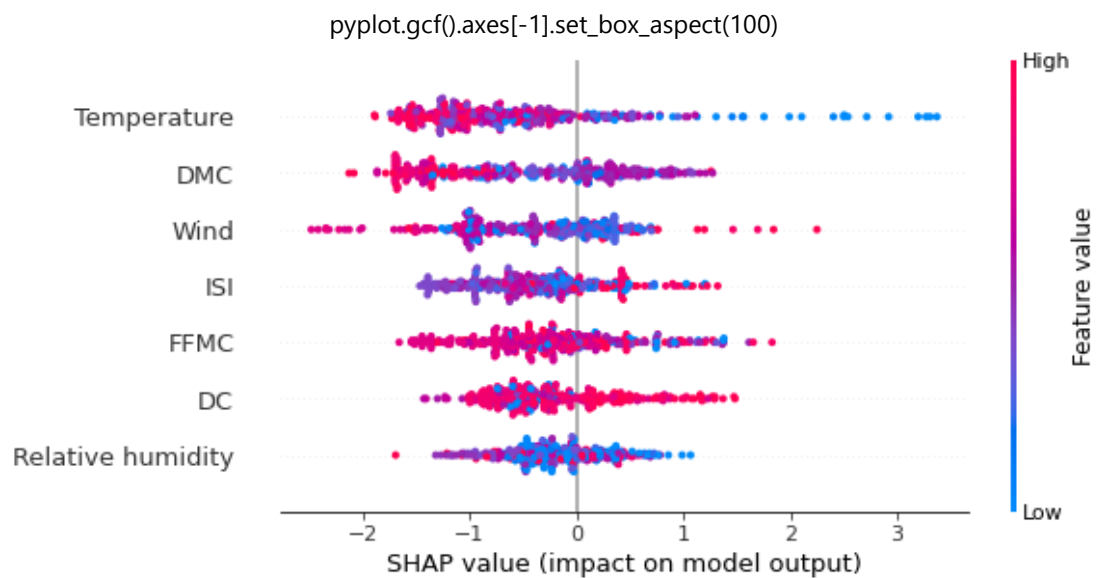


```
900
901
902    shap.summary_plot(shap_values[0], x_test, class_names=class_names,show=False)
903    pyplot.gcf().axes[-1].set_box_aspect(50)
```

904 pyplot.gcf().axes[-1].set_aspect(100)
905        pyplot.gcf().axes[-1].set_box_aspect(100)



906
907 shap.summary_plot(shap_values[1], x_test, class_names=class_names,show=False)
908 pyplot.gcf().axes[-1].set_box_aspect(50)
909 pyplot.gcf().axes[-1].set_aspect(100)
910        pyplot.gcf().axes[-1].set_box_aspect(100)



911
912
913 shap.summary_plot(shap_values[2], x_test, class_names=class_names,show=False)
914 pyplot.gcf().axes[-1].set_box_aspect(50)
915 pyplot.gcf().axes[-1].set_aspect(100)

916

pyplot.gcf().axes[-1].set_box_aspect(100)



917
918    shap.summary_plot(shap_values[3], x_test, class_names=class_names,show=False)
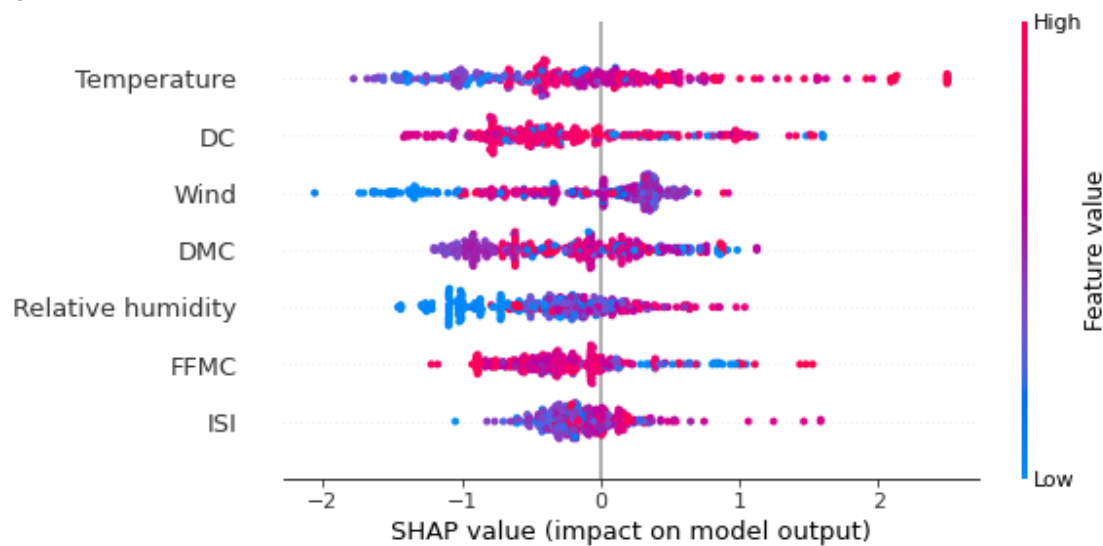919    pyplot.gcf().axes[-1].set_box_aspect(50)
920    pyplot.gcf().axes[-1].set_aspect(100)
921    pyplot.gcf().axes[-1].set_box_aspect(100)



922
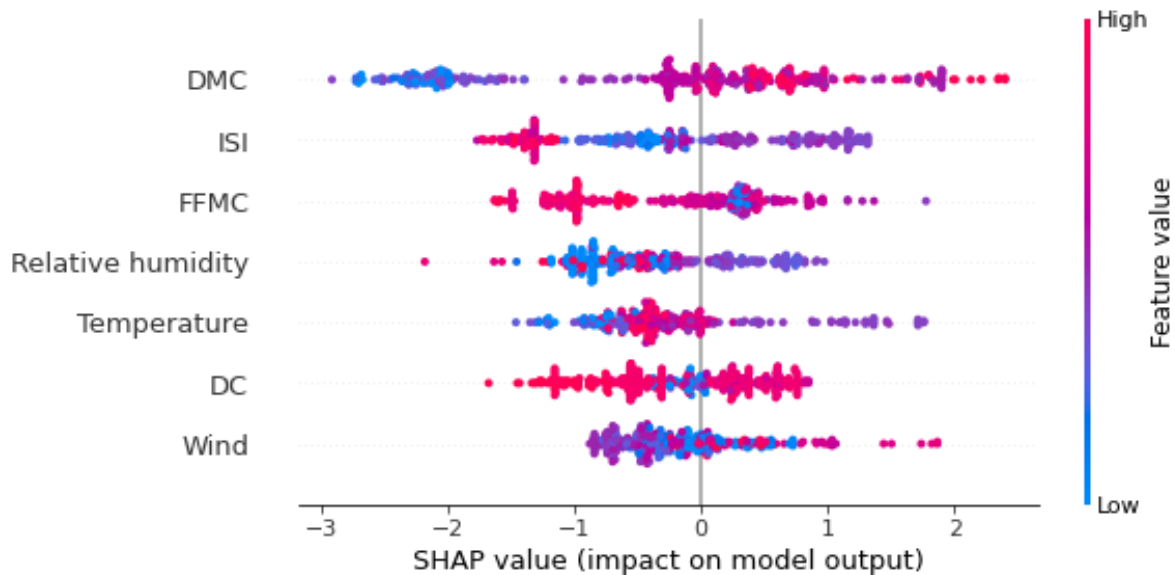923    shap.summary_plot(shap_values[4], x_test, class_names=class_names,show=False)
924    pyplot.gcf().axes[-1].set_box_aspect(50)
925    pyplot.gcf().axes[-1].set_aspect(100)
926    pyplot.gcf().axes[-1].set_box_aspect(100)

927
928     shap.summary_plot(shap_values[5], x_test, class_names=class_names,show=False)
929     pyplot.gcf().axes[-1].set_box_aspect(50)
930     pyplot.gcf().axes[-1].set_aspect(100)
931                              pyplot.gcf().axes[-1].set_box_aspect(100)



932
933
934