AMERICAN METEOROLOGICAL SOCIETY

*Journal of Applied Meteorology and Climatology*

# EARLY ONLINE RELEASE

This is a preliminary PDF of the author-produced manuscript that has been peer-reviewed and accepted for publication. Since it is being posted so soon after acceptance, it has not yet been copyedited, formatted, or processed by AMS Publications. This preliminary version of the manuscript may be downloaded, distributed, and cited, but please be aware that there will be visual differences and possibly some content differences between this version and the final published version.

1

1 **Classification of Australian Thunderstorms using Multivariate Analyses of Large-Scale**

2 **Atmospheric Variables**

3

4 BRYSON C. BATES

5 *CSIRO Oceans and Atmosphere, Wembley, Western Australia, Australia*

6 *School of Earth and Environment, The University of Western Australia, Crawley, Western*

7 *Australia, Australia*

8

9 ANDREW J. DOWDY*

10 *Bureau of Meteorology, Melbourne, Victoria, Australia*

11

12 RICHARD E. CHANDLER

13 *Department of Statistical Science, University College London, London, UK*

14

15 ABSTRACT

16 Lightning accompanied by inconsequential rainfall (i.e. 'dry' lightning) is the primary

17 natural ignition source for wildfires globally. This paper presents a machine-learning and

18 statistical-classification analysis of 'dry' and 'wet' thunderstorm days in relation to

19 associated atmospheric conditions. The study is based on daily lightning flash count and

20 precipitation data from ground-based sensors and gauges, and a comprehensive set of

21 atmospheric variables based on the ERA-Interim reanalysis for the period from 2004 to 2013

22 at six locations in Australia. These locations represent a wide range of climatic zones

23 (temperate, subtropical to tropical). Quadratic surface representations and low-dimensional

24 summary statistics were used to characterize the main features of the atmospheric fields. Four

25 prediction skill scores were considered and ten-fold cross validation used to evaluate the

*Corresponding author (Email: andrew.dowdy@bom.gov.au; Tel: + 61 3 9669 4722)

26  performance of each classifier. The results were compared with those obtained by adopting

27  the approach used in an earlier study for the Pacific Northwest, United States. It was found

28  that: both approaches have prediction skill when tested against independent data, mean

29  atmospheric field quantities proved to be the most influential variables in determining dry

30  lightning activity and no single classifier or set of atmospheric variables proved to be

31  consistently superior to their counterparts for the six sites examined here.

32

33  **1.  Introduction**

34      Although human-caused wildfire ignitions are common in many regions of the world,

35  particularly in densely populated areas, fires ignited by lightning typically burn a larger area

36  than fires ignited by other sources. This is attributable to lightning occurrence in remote

37  locations and in large spatial and temporal clusters which hamper the response efforts of fire

38  management authorities (USDA Forest Service 1992; McRae 1992; Vazquez and Moreno

39  1998; Wotton et al. 2005; Wotton and Martell 2005; Kasischke et al. 2006; Dowdy and Mills

40  2012a). Lightning that occurs with relatively little precipitation (i.e., 'dry' lightning) has a

41  higher chance of igniting a fire than lightning accompanied by heavier precipitation ('wet'

42  lightning) (Rothermel, 1972; Wotton and Martell 2005; Dowdy and Mills 2012a). Therefore,

43  an improved understanding of dry lightning activity and the atmospheric conditions that

44  influence its occurrence is of importance for better preparedness and enhancing the ability to

45  respond to the impacts associated with wildfires ignited by lightning.

46      There are many physical factors that can influence lightning occurrence as demonstrated

47  in numerous previous climatological, dynamical modeling and seasonal prediction studies

48  including Weisman and Klemp (1982), Goodman et al. (2000), Burrows et al. (2005),

49  Williams et al. (2005), Deierling et al. (2008), Romero et al. (2007), Chronis et al. (2008),

50  Dai et al. (2009), Barthe et al. (2010), Romps et al. (2014), Magi (2015), Dowdy (2016),

51    Muñoz et al. (2016) and references therein. Although aspects of the microphysical processes

52    associated with lightning generation are not well understood in some cases, the role of ice in

53    facilitating charge separation within the cloud appears to be a critical factor in determining

54    whether or not lightning is produced (e.g., as indicated by laboratory experiments (Takahashi

55    and Miyawaki 2002) as well as observations (Lang et al. 2014)). Microphysical processes

56    such as ice formation are not well represented at the spatial and temporal scales of currently

57    available climate models and reanalyses, leading to the use of parameterization schemes for a

58    range of variables associated with convection. For example, the ERA-Interim reanalysis (Dee

59    et al. 2011) uses a convective parameterization based on a bulk mass flux scheme (as

60    originally described by Tiedtke 1989), with parameterizations also used to represent the

61    fallout of precipitation (e.g., Kuo and Raymond 1980) and factors such as virga (streaks of

62    water or ice particles that vaporize before reaching the ground) considered.

63        In addition to convective parameterization schemes, several studies have demonstrated

64    that statistical indicators of lightning activity can be found at relatively coarse spatial and

65    temporal scales (e.g., similar to the resolution of general circulation models (GCMs) and

66    reanalyses). For example, Romps et al. (2014) combined precipitation and Convective

67    Available Potential Energy (CAPE) based on GCM output for use as an indicator of

68    environments conducive to lightning activity, applying this indicator to examine the influence

69    of global warming on lightning strikes in the United States. A recent study based on

70    reanalyses demonstrated that even at spatial resolutions of 7.5° in latitude and longitude,

71    atmospheric conditions such as lower-tropospheric moisture content, temperature lapse rate

72    and CAPE can be strongly related to lightning activity (Dowdy 2016).

73        In contrast to the number of studies that have examined atmospheric conditions associated

74    with lightning activity in general, relatively few studies have focused specifically on dry

75    lightning. Notable early studies include Rorig and Ferguson (1999, hereafter designated as

76  RF99) and Rorig et al. (2007), demonstrating that a linear discriminant rule could separate

77  dry and wet lightning classes. The rule was composed of dewpoint depression at 850 hPa

78  (DD850) and temperature lapse from 850 to 500 hPa (TL850500), with dry lightning defined

79  as lightning accompanied by precipitation of less than one tenth of an inch (about 2.5 mm).

80  Dowdy and Mills (2012b) demonstrated that these two variables were also applicable in

81  southeast Australia, and that the average chance of a sustained fire ignition resulting from the

82  occurrence of lightning in that region is higher than average if the precipitation

83  accompanying the lightning is less than about 2 to 3 mm. Recent studies have examined a

84  somewhat wider range of variables in relation to the occurrence of dry lightning, including

85  studies in North America (Wallmann et al. 2010; Nauslar et al. 2013; Abatzoglou et al. 2016)

86  and Australia (Dowdy 2015), finding that some useful skill can be obtained for predicting the

87  occurrence of dry lightning based on several different methods. However, as dry lightning

88  activity remains relatively unstudied when compared with other aspects of thunderstorm

89  activity and associated convective processes, to date there have been no climatological

90  studies of the spatial and temporal variability of dry lightning activity, or the influence of

91  large-scale atmospheric drivers of dry lightning variability.

92      The approach presented in this paper (hereafter designated as BDC) represents a more

93  general approach to the two-category classification problem of dry and wet lightning days

94  than that of RF99. The paper has four objectives with a view to building on previous studies

95  of dry lightning occurrence. The first is to consider a wider range of atmospheric conditions

96  associated with dry lightning activity and precipitation occurrence than has been the case to

97  date. The second is to build on the suggestion put forward by Blouin et al. (2016) that a

98  comparison of classification methods (classifiers) may provide useful guidance for future

99  research. The third is to consider lightning, precipitation and atmospheric data from a wide

100 range of climatic zones. The fourth objective is to identify a subset of influential atmospheric

101 variables across climatic zones and different classifiers. The method of RF99 is used as a

102 benchmark for assessments of prediction accuracy and the applicability of the new approach

103 proposed in this paper. In this way, the paper provides a useful addition to the toolkit for

104 addressing questions related to lightning activity. The paper is divided into five sections.

105 Section 2 provides a description of the study sites, data and classifiers used. Results are

106 presented in Section 3. A summary and conclusions are given in Section 4. The quadratic

107 surface representations and low-dimensional summary statistics (LDSS) used to characterize

108 the main features of the atmospheric fields considered in this study are described in the

109 Appendix.

110

111 **2. Data and methods**

112 *a. Description of study sites and data*

113 The description of the daily lightning flash count datasets used herein parallels that of

114 Bates et al. (2015), and the text in the next two paragraphs is derived from there with minor

115 modifications. The data were collected from ground-based CIGRE 500 (Comité

116 Internationale des Grands Réseaux Electriques, 500 Hz peak transmission filter circuit)

117 sensors located at six weather stations operated by the Australian Bureau of Meteorology

118 (Figure 1 and Table 1). The sensors were selected because of their record length and quality,

119 and their locations in a variety of climatic settings including temperate, subtropical and

120 tropical sites. The records cover the period from January 2004 to at least December 2010

121 (Townsville) and at most February 2013 (Melbourne).

122

123                    **< Insert Figure 1 and Table 1 about here >**

124

125    Although the CIGRE 500 sensor was designed specifically to detect cloud-to-ground

126    flashes, it also responded to cloud-to-cloud flashes, with about 68% of the lightning flash

127    counts recorded being due to cloud-to-ground flashes. We considered the total number of

128    lightning flash counts since the CIGRE 500 sensor did not distinguish between intracloud and

129    cloud-to ground flash counts; and the ratio of intracloud to cloud-to-ground flashes can vary

130    significantly depending on thunderstorm type and intensity, region of occurrence and season

131    (Rakov and Uman, 2003). Estimates of the effective horizontal ranges of the sensor are 30

132    km for cloud-to-ground flashes and 15 km for cloud-to-cloud flashes (Kuleshov and

133    Jayaratne, 2004). As with other studies of this nature, these effective ranges should be taken

134    into consideration when interpreting results for specific purposes such as fire ignition from

135    cloud-to-ground lightning flashes. The electromechanical counters attached to the CIGRE

136    500 sensors were read manually each day between 0800 and 0900 h local time. Further

137    details can be found in Jayaratne and Kuleshov (2006), Kuleshov et al. (2009) and Bates et al.

138    (2015).

139    For a given weather station, thunderstorms were deemed to have occurred during a 24-h

140    period if at least one lightning flash count was registered by the CIGRE 500 sensor. They

141    were categorized as either 'dry' or 'wet' according to the concurrent daily precipitation

142    reading recorded by the storage gauge at the station. A thunderstorm was classified as dry if

143    the precipitation reading was less than 2.5 mm or wet otherwise. In Australia, daily

144    precipitation is nominally measured each day at 0900 h local time. Station data were obtained

145    from the SILO patch-point data set (Australian Bureau of Meteorology). There is a large

146    disparity in spatial scales between the detection range of the sensor and the diameter of a

147    precipitation gauge (15-30 km versus 203 mm). Thus it is possible for precipitation amounts

148    greater than 2.5 mm to occur within the sensor's detection limit but away from the station

149    gauge. However, the use of gridded station data has its own set of limitations in that the

150     interpolation involved is a form of smoothing that reduces precipitation variability. Thus, the

151     process of gridding can considerably increase (decrease) the frequency of low (high)

152     precipitation amounts (Ensor and Robeson, 2008), and this might have implications for the

153     classification of dry thunderstorms. The reduction in variability is dependent on the distance

154     from a grid point to the nearest gauge. A further concern is the relatively low density of

155     precipitation gauge networks in Australia. For example, the numbers of gauges within a 30

156     km radius of the Darwin, Townsville, Coffs Harbour and Port Hedland sites are 12, 8, 18 and

157     3, respectively. This low network density is likely to lead to excessive smoothing in some

158     instances and affect the distribution of daily precipitation amounts. Given the above, days

159     with station precipitation values flagged as interpolated were discarded.

160     With future applications in mind, the study was designed to be conducted at spatial and

161     temporal resolutions similar to that of current general circulation models and reanalyses.

162     Atmospheric information was obtained from the ERA-Interim reanalysis archive (Dee et al.

163     2011). The spatial and temporal resolution of the dataset used is 0.75 degrees (in both latitude

164     and longitude) and 6 hours, respectively. For each CIGRE 500 site, atmospheric data were

165     extracted for the 49 reanalysis grid points closest to the sensor's location. The aim of the grid

166     was to capture the presence of a thunderstorm over or in the proximity of a sensor. The

167     lightning and precipitation series were synchronized with the ERA-Interim series for 0600

168     UTC (1600 h Eastern Australia Time) within the 24-hour period represented by the lightning

169     and precipitation data. This is because the diurnal variation in temperature lapse rate over

170     land, due to solar radiation, produces conditions that are more favorable for lightning activity

171     to occur during the late-afternoon period in general than at other times of the day or night

172     (Christian et al. 2003; Dowdy and Mills 2009; Allen et al. 2011). Thus, the synchronization

173     ensures that the atmospheric variables for each daily lightning flash count correspond to the

174     time at which the lightning is most likely to have occurred. Since the use of a single time

175    point can be viewed as reductive, the possibility that atmospheric variables at other times of

176    day may also be relevant was considered. However, the additional information was found to

177    be largely redundant because correlations within a 24-hour period are invariably high (e.g.

178    correlations between individual variables at 0600 UTC and 1200 UTC, spanning the time

179    period during which most deep convective processes occur in Australia, are typically greater

180    than 0.95 and greater than 0.85 in every case examined). The atmospheric variables

181    considered herein are listed in Table 2.

182

183                            **< Insert Table 2 about here >**

184

185    The set of atmospheric variables examined here represents a wider range than has typically

186    been examined in previous studies, particularly those studies focused on climate-scale

187    analyses rather than finer-resolution numerical weather prediction or radar observational

188    studies. This is because there have been very few studies that have specifically examined dry

189    lightning activity and the atmospheric conditions that influence its occurrence. Consequently,

190    the literature on dry and general lightning activity was combed and physical understanding

191    used to reduce the number of variables as far as possible. The variables listed in Table 2

192    represent a broad variety of physical processes that can be associated with deep convection,

193    including both dynamical and thermodynamical processes. The variables comprise various

194    measures of temperature lapse, moisture content, vertical motion and water phase state,

195    including at a range of different pressure levels (to allow potential variations in height

196    between dry and wet thunderstorm characteristics to be distinguished).

197    *b.  Variable selection*

198    To identify the dominant large-scale controls on lightning activity from among the

199    variables listed in Table 2 is a challenging statistical problem: there are dependencies among

200    the variables leading to collinearity and, moreover, the processes controlling lightning

201    activity are complex so that the variables must be considered concurrently rather than in

202    isolation. Regression-based approaches, notably those based on generalized linear models,

203    are ideally suited to this kind of problem (e.g. Yan et al. 2002, Chandler 2005). However, an

204    additional complication in the present application is that the explanatory variables are spatial

205    fields over a 7×7 grid, rather than individual values. In principle, this can be handled using

206    modern statistical techniques such as functional regression (e.g. Morris, 2015). However, in

207    their current state of development such methods are most effective when the number of

208    candidate variables is relatively small. The current state of knowledge is insufficient to

209    identify a small number of candidate variables from the list in Table 2 with high confidence.

210    The strategy adopted here is therefore to use a combination of approaches that are designed to

211    isolate the most influential variables from many candidates.

212       To handle the spatial nature of the atmospheric variables listed in Table 2, the daily fields

213    for each variable were reduced to a set of five LDSS designed to capture the main synoptic

214    features. This was done by fitting quadratic surfaces to each daily field (see Appendix) and

215    using the fitted surfaces to derive physically-interpretable daily summaries (overall means,

216    vertical and horizontal gradients, and curvature). Note that the intention is not to provide

217    highly accurate descriptions of the fields, but rather to provide indices that broadly describe

218    the synoptic structure. The use of LDSS reduces the dimensionality of the problem from 49

219    grid point values per atmospheric variable per day to 5. Other dimension reduction techniques

220    are available, notably principal component (empirical orthogonal) analysis which was

221    explored as an alternative to the LDSS considered here. It was found that five or more

222    components were necessary to explain 70 to 80% of the variance for each data set. Only the

223    first component had any predictive power in terms of discriminating between dry and wet

224    lightning. Although the loadings for this component would often indicate a contrast between

225 two sets of variables, a defensible interpretation of the contrast proved elusive. Moreover, its

226 predictive skill was lower than that obtained with LDSS.

227     At this point, a LDSS of an atmospheric variable will be referred to as a potential

228 candidate variable. As an initial screening procedure, for each potential candidate variable,

229 comparative boxplots of each LDSS were used to contrast its values for dry and wet lightning

230 cases. Two variable selection criteria were considered. First, potential candidate variables

231 where the $75^{th}$ ($25^{th}$) percentile for one lightning type was below (above) the $25^{th}$ ($75^{th}$)

232 percentile for the other were deemed informative in terms of discriminating between dry and

233 wet lightning days. These variables were reserved for further analysis. Second, depending on

234 the number of such candidate variables found, they were supplemented by including

235 additional candidate variables where the median in one lightning type was above the $75^{th}$

236 percentile or below the $25^{th}$ percentile of the other (see, e.g., Figure 2). The resulting

237 candidate variables formed the columns of an atmospheric data matrix. This approach could

238 be criticized as ad hoc: it is natural to ask whether alternative techniques, such as automatic

239 variable selection procedures, would be preferable. The main reason for the approach taken

240 here is that manual inspection of boxplots can provide checks on the data, as well as

241 preliminary insights that may aid subsequent interpretation and that cannot be obtained from

242 an automated analysis. In any case, the aim is merely to carry out a very preliminary

243 screening of the data so as to focus subsequently on quantities that may have some predictive

244 power in discriminating between dry and wet lightning.

245

246                     **< Insert Figure 2 about here >**

247

248     Many of the candidate variables are measured on very different scales and thus are not

249 commensurable in terms of magnitude or variability. This means that some variables could

250    dominate or influence the results of the classification analysis because of their measurement

251    units alone (Everitt and Hothorn, 2011). Thus, the columns of the data matrix were

252    standardized to zero mean and unit variance prior to further analysis. This process places

253    candidate variables on the same relative scale without disturbing the shape of the distribution

254    of the data. It facilitates interpretation of the results of a discriminant or regression analysis,

255    and helps to concentrate precisely on the conditions that are present during thunderstorms

256    because it focuses on the relative variations of each variable within its own physical limits.

257    The colldiag function from the perturb package in the R computing environment (Hendrickx

258    2012; R Core Team 2015) was used to detect the presence of collinearity in the data matrix.

259    Colldiag is an implementation of the regression collinearity diagnostic procedures found in

260    Belsley et al. (1980). It computes the condition indices of the data matrix and provides the

261    variance decomposition proportions associated with each condition index. As a rule of thumb,

262    variables with proportions greater than 0.99 were considered sources of severe collinearity.

263    Thus the corresponding columns were removed to form a reduced data matrix. A second

264    proportion threshold of 0.8 was used to assess the degree of the sensitivity to threshold

265    selection. It was found that the results obtained from the procedures described below showed

266    only a slight sensitivity. Therefore, the results obtained using the proportion threshold of 0.8

267    will not be reported here.

268    *c.   Multivariate statistical analysis*

269      Two machine-learning and three statistical methods were used for classification:

270    classification and regression trees (CART); random forests (RF); linear discriminant analysis

271    (LDA), quadratic discriminant analysis (QDA) and logistic regression (LR). Detailed

272    descriptions of CART, RF and LR can be found in Faraway (2016), and LDA and QDA in

273    Everitt and Dunn (2001). The R packages used in this work were: DiscriMiner (Sanchez

274    2013); MASS (Venables and Ripley, 2002); randomForest (Liaw and Wiener 2002); and tree

275    (Ripley, 2014). CART uses binary recursive partitioning to divide the data space, splitting it

276    along the coordinate axes of the candidate variables to give increasingly homogenous subsets

277    and hence the maximal separation of the classes until it is infeasible to continue. The measure

278    of node heterogeneity is the deviance (a quality-of-fit statistic)**.** The partitioning leads to a set

279    of decision rules in the form of a binary tree. The tree is 'pruned' to identify a parsimonious

280    tree with acceptable misclassification rates. Cross validation can be used to determine an

281    appropriate tree size. RF is an ensemble learning algorithm which generates a large number

282    of CART from bootstrap samples of the original data. An estimate of the misclassification

283    rate can be obtained by using each tree to predict the data not in the bootstrap sample and

284    averaging the predictions over all trees. The randomForest package can be used to produce

285    variable importance plots which reveal how important each variable is in classifying the data

286    and contributing to the homogeneity of the nodes. LDA is derived from an underlying model

287    in which the distributions of the variables on dry and wet lightning days are both multivariate

288    normal, with possibly different means and a common covariance matrix. LDA is somewhat

289    robust with respect to minor violations of these assumptions. Although serious violations will

290    often result in unreliable estimates of the coefficients, the procedure can still be a good

291    heuristic. The discriminant function is a linear combination of the candidate variables, the

292    coefficients of which are estimated by ordinary least squares so that the ratio of the between-

293    classes variance and the within-classes variance is maximized. This function takes the value

294    zero at the decision boundary. If the value of the discriminant function is negative the

295    variable vector is assigned to one class, if positive it is assigned to the other class. Given that

296    the variables are standardized, the coefficients indicate the relative importance of each

297    variable in predicting class assignment. QDA is a generalization of LDA in which the two

298    classes need not have the same covariance matrix, but the assumption of multivariate

299    normality still applies. The interpretation of the coefficients in terms of the relative

300 importance of each variable is more difficult to assess than for LDA as the discriminant

301 function contains quadratic as well as linear and constant terms. The LR model can be written

302 as

303

304
$$\text{logit}(\pi_i) = \ln[\pi_i/(1-\pi_i)] = \beta_0 + \sum_{j=1}^{p} \beta_j X_j \qquad (1)$$

305

306 where $\pi_i$ is the probability of occurrence of class $i$ $(i = 1, 2)$, $\pi_i/(1-\pi_i)$ is the odds ratio for

307 class $i$, $p$ is the number of columns in the data matrix and $\beta_0, \ldots, \beta_p$ are the regression

308 coefficients which are determined via maximum likelihood estimation. (Obviously, with only

309 two categories it is only necessary to estimate the coefficients for one of the categories since

310 $\pi_2 = 1 - \pi_1$.) Classification on the basis of the variables is then done by setting a threshold $\tau$

311 say, and allocating a day to category 1 if $\pi_1 > \tau$. For each site, a receiver operating

312 characteristic (ROC) curve was used to select the threshold $\tau$ by minimizing the distance

313 from the curve to the point representing perfect classification accuracy: this was done to

314 account for the fact that the sample sizes for dry and wet lightning days were noticeably

315 unequal for several sites (Table 1). Experiments using Youden's (1950) Index indicated that

316 threshold estimates were not sensitive to the selection technique used. With LR, by contrast

317 with LDA and QDA, there is no formal requirement for multivariate normality of the

318 explanatory variables within each category of the response variable, and the use of binary or

319 categorical variables is acceptable. A combination of stepwise selection and analysis of

320 deviance was used to determine the significance of variables in the LR models. Further

321 details on the above classifiers can be found in Breiman (2001), Venables and Ripley (2002)

322 and Hilbe (2009).

323    Although the approach of RF99 used only LDA to discriminate between dry and wet

324    lightning, the four other classifiers considered herein (CART, RF, QDA and LR) were

325    applied to the means of the DD850 and TL850500 fields (designated mu.DD850 and

326    mu.TL850500, see Appendix) to ascertain whether a higher classification performance could

327    be achieved. This extended approach is hereafter designated as E-RF99. The analysis was

328    conducted in parallel with an identical study of a much larger set of variables to determine the

329    extent to which it is possible to improve on the RF99 variable pair. Four measures of

330    prediction skill were considered: the hit rate for dry lightning (HR), the false alarm ratio for

331    dry lightning (FAR), the Brier (1950) score (BS) and for LR the area under the ROC (AUC).

332    HRs, FARs, BSs, ROC curves and AUCs were determined using the verification package in

333    R (NCAR 2015). For a perfect classification, HR=1, FAR=0, BS=0 and AUC=1. HR values

334    near 0, FAR and BS values near 1, and AUC values near 0.5 indicate poor classification

335    performance. For the convenience of the reader, in what follows, a list of the acronyms and

336    abbreviations used in this paper and their meaning is given in Table 3.

337

338                          **< Insert Table 3 about here >**

339

340    *d.   Cross validation experiments*

341    Initial assessments of the prediction skill of the five classifiers (CART, RF, LDA, QDA

342    and LR) were based on the data matrices for the six CIGRE 500 sites. As this can lead to

343    optimistic bias in the estimated skill scores, ten-fold cross validation experiments were used

344    to assess how well the results generalized to an independent dataset. Here the lightning,

345    precipitation and candidate variable data were partitioned into ten subsamples of equal size.

346    From these subsamples, a single subsample was retained for testing the model, and the

347    remaining nine subsamples used for training (model fitting). The process is then repeated ten

348    times with each of the subsamples used exactly once for validation. The R packages used in

349    the cross validation experiments were cvTools (Alfons 2012) and verification (NCAR 2015).

350

351    **3.   Results**

352    *a.   Preliminary analyses*

353        Lightning activity at the six sites occurs primarily during the warmer months of the year

354    (November to April). However, the most severe fire weather conditions in Australia occur at

355    different times of the year, generally ranging from summer in the south to winter (i.e., the dry

356    season) in the north. There are some regional variations to this, particularly along the eastern

357    seaboard (including Coffs Harbour) where the peak fire weather conditions occur somewhat

358    earlier (around October) than at other similar latitudes in Australia (Luke and McArthur,

359    1978). Further details on the lightning climatology of Australia may be found in Kuleshov et

360    al. (2009), Dowdy and Kuleshov (2014), Bates et al. (2015) and references therein, and will

361    not be repeated here.

362        The proportions of dry and wet lightning days for the six CIGRE 500 sites are reported in

363    Table 1 and illustrated in Figure 1. Darwin is one of the most lightning prone areas in

364    Australia. The number of lightning days for Darwin is markedly higher than those for the

365    remaining sites, even for the case of Townsville which is in the same climatic zone. Perth has

366    the lowest number of lightning days by a wide margin. Port Hedland has the highest

367    proportion of dry lightning days, reflecting its desert environment. For the tropical and

368    subtropical sites, the proportion of dry lightning days exceeds that of wet lightning days. This

369    is somewhat surprising, and may in part be explained by the use of a single precipitation

370    gauge to characterize rainfall over the detection range of the CIGRE 500 sensor (Section 2).

371    To a lesser extent, it might reflect the effects of the precipitation threshold of 2.5 mm on

372    lightning day classification. For example, the proportions of dry and wet lightning days at

373    Darwin are essentially equal if the precipitation threshold is reset to 2.0 mm.

374        Median adjusted $R^2$ values for the fitted quadratic surfaces varied from variable to variable

375    with 5.2 to 16% below 0.5 across the six CIGRE 500 sites and 47 to 68% above 0.75. The

376    highest values were obtained for GPH500 and GPH700 ($> 0.97$), and the lowest for W and

377    MING, (0.09 to 0.43). This pattern was consistent across all sites. Thus, overall, the quadratic

378    surfaces described in the Appendix gave a reasonable representation of the main features of

379    the atmospheric fields considered herein.

380    *b.    Classification analyses*

381        Scatterplots of the skill scores obtained from the five classifiers are shown in Figure 3.

382    The HR and FAR are for dry lightning and the radii of the circles represent the magnitude of

383    the BS. For each CIGRE 500 site, the convex hull of the five data points obtained using only

384    mu.DD850 and mu.TL850500 as candidate variables is displayed to facilitate their

385    delineation. (The convex hull of a set of points is the smallest convex set enclosing the

386    points.) The plots reveal six key features. First, for any site, approach (E-RF99 or BDC) and

387    classifier, the HR exceeds the FAR (note the differences in the axis scales). Thus, both

388    approaches and all five classifiers have some skill in discriminating dry lightning from wet

389    lightning. Apart from Port Hedland, it is also evident that the RF99 approach (denoted by

390    filled squares enclosed by green circles) provides lower HRs. Second, for Darwin and

391    Townsville the approach of BDC often provides higher HRs than those for E-RF99 but at

392    expense of higher FARs for some classifiers. Third, for Coffs Harbour, Melbourne and Perth,

393    the approach of BDC produced simultaneously higher hit rates and lower FARs than those for

394    E-RF99. Fourth, in the case of Port Hedland, the HRs and FARs obtained for a given

395    classifier and the two approaches considered herein (E-RF99 and BDC) are similar despite

396    the differences in candidate variable sets: the BDC candidate variable set included terms such

397    as mu.TOTP, mu.CONVP and mu.TCW (see Appendix for details regarding their derivation).

398    Fifth, except for Darwin, the application of LR to the BDC candidate variable set produced

399    low FARs. Sixth, the approach of BDC produces similar or lower BSs than those for E-RF99.

400

401    **< Insert Figure 3 about here >**

402

403    *c.  Influential variables*

404    The relative frequency histogram of influential variables across the six CIGRE 500 sites

405    and four classifiers with easily interpreted decision rules or boundaries (CART, RF, LDA and

406    LR) is shown in Figure 4. Overall, 16 out the 28 variables are means, nine are magnitudes of

407    gradient vectors, two are vertical gradients and one is SEASON (Table 2). The seven most

408    frequent variables are associated with atmospheric water content (mu.TOTP, mu.CONVP,

409    gd.TOTP and mu.TCWV) and instability and lifting potential (mu.CBH, mu.DD700 and

410    vg.T). Thus, five of the seven most frequent variables are means. In terms of the raw

411    atmospheric variables listed in Table 2, not one of the variables used by RF99 (DD850 and

412    TL850500) is present in this subset. Additionally, DD850 does not appear to be as influential

413    as DD700, with DD700 and DD850 having relative frequencies of 0.0879 and 0.0220. As

414    shown in Fig. 2c, high values of DD700 are typically associated with dry lightning rather

415    than wet lightning, with a physical interpretation of this being that relatively dry air results in

416    an increased likelihood of precipitation evaporating before reaching the ground (i.e., virga).

417    The absence of CAPE and the near absence of W in Figure 4 suggest that these variables are

418    not informative in terms of discriminating between dry and wet lightning conditions. This is

419    unlikely to be the case for lightning activity studies involving discrimination between

420    lightning and non-lightning days.

421    The dominance of the mean terms in the set of influential variables could be related to

422    temporal variations in the timing of thunderstorms with respect to a given location, noting

423    that although our analyses is based on afternoon values of the atmospheric variables (as this

424    is when lightning most frequently occurs in these regions), lightning can also occur at other

425    times of the day and night. The apparent influence of mu.TOTP and mu.CONVP must give

426    rise to concern that information about precipitation has been used twice: once as daily

427    precipitation readings at ground-based storage gauges were used to classify lightning days as

428    either dry or wet; and twice as mu.TOTP and mu.CONVP values at 0600 UTC were derived

429    from modeled precipitation and used as explanatory variables. However, scatter plots and

430    quantile-quantile plots of mu.TOTP and mu.CONVP against the precipitation readings (not

431    shown) revealed little evidence of relationships for all six CIGRE 500 sites. Except for

432    Melbourne, robust estimates of the correlation coefficients ranged from 0.1 to 0.3. For

433    Melbourne, the estimates were about 0.4. This lack of a simple relationship, and the positions

434    of mu.TOTP and mu.CONVP in the histogram depicted in Figure 4, suggest that the

435    construction of these variables captures useful additional information about atmospheric

436    conditions that cannot be obtained from the other potential candidate variables considered.

437    Some evidence for this conjecture is provided in Davies et al. (2013). For one of the tropical

438    sites considered herein (Darwin), they created two concurrent long-term data sets that

439    described the large-scale atmosphere and the characteristics of small-scale convection. They

440    found that estimates of convective precipitation have a strong relationship with dynamical

441    variables such as moisture convergence and vertical velocity at mid-levels. Wind rather than

442    moisture convergence was used in the current study (Table 2), and vertical velocities in

443    reanalyses can suffer from large inaccuracies (Abalos et al., 2015). The latter may have also

444    contributed to the position of mu.W and gd.W in Figure 4.

445

446      **< Insert Figure 4 about here >**

447

448      Figure 5 displays the relative frequency histograms of the most-frequent atmospheric

449      variables on a site-by-site basis. Here the maximum frequency for any variable is limited to

450      four (the number of classifiers with interpretable decision rules or boundaries). Furthermore,

451      the minimum count for any variable can be zero as not every one of the variables was found

452      to be influential for every site. Colored bars indicate the seven most-frequent variables

453      depicted in Figure 4. For the sake of clarity, white bars indicate additional variables that have

454      a frequency of at least two. The variables depicted in Figure 5 are primarily associated with

455      atmospheric water content and instability and lifting potential. Comparison of Figures 5a-d

456      and 5e-f indicates a marked difference in the shapes of the histograms for sites located in

457      western Australia (Perth and Port Hedland) and those in central and eastern Australia

458      (Darwin, Townsville, Coffs Harbour and Melbourne). In the case of Perth (Figure 5e), five of

459      the seven most-influential variables across all sites and classifiers depicted in Figure 4 have

460      zero frequencies and the frequencies of the remaining two (mu.CBH and gd.TOTP) are low.

461      It is the only site not to include both mu.TOTP and mu.CONVP amongst its set of influential

462      variables. The most common variables across the four classifiers for Perth are indicated by

463      white bars. Three of these four variables (mu.TGM7001000, gd.GPH500 and gd.GPH700)

464      are not included in the variable sets for the other sites (cf. Figure 4). These variables are

465      potential indicators of convective systems associated with fronts. The fourth variable

466      (mu.TL850500) is selected for Coffs Harbour by LR only. For Port Hedland (Figure 5f),

467      three of the seven most frequent variables in Figure 4 have zero frequencies. It is the only site

468      to not include mu.CBM amongst its set of influential variables. The remaining four variables

469      (mu.TOTP, mu.CONVP, gd.TOTP and m.TCWV) characterize atmospheric water content.

470      There are four additional variables (gd.TOTP, gd.CONVP, mu.MING and mu.T2) that are

471  not depicted in Figure 5f since they have a frequency of one. Port Hedland is also different to

472  the other sites in that it has a notably higher HRs and lower FARs (Figure 3). This is because

473  the ratio of dry to wet lightning proportions for Port Hedland is 4.3 which is much higher

474  than that for the other sites where it is between 1.1 and 1.8 (Table 1). With the exceptions of

475  Townsville and Coffs Harbour, the frequencies of vg.T are zero for the four remaining sites.

476  Sharp temperature gradients are a potential indicator of troughs, and convergence along

477  troughs can lead to showers and thunderstorms. The so-called inland (or easterly) trough is

478  located on the inland side of the Great Dividing Range in Australia, forming a boundary

479  between the moist air near the coast and dry air inland. It typically extends through central

480  Queensland and into central New South Wales and is active during the months from

481  September to May. Furthermore, the frequency of mu.ICE is greater than zero for Melbourne

482  alone. Ice water content and lightning activity are highly correlated (Xu et al., 2010 and

483  references therein), and this variable may provide information about the low (high) lightning

484  flash rates associated with dry (wet) lightning. These results, and those illustrated in Figure 4,

485  suggest that the optimal variable sets for lightning classification problems may vary between

486  different climatic zones.

487

488  **< Insert Figure 5 about here >**

489

490  *d.  Cross validation experiments*

491  Scatterplots of the mean skill scores obtained from the cross validation experiments are

492  shown in Figure 6. Again, the HR and FAR are for dry lightning and the radii of the circles

493  represent the magnitude of the BS. The radii of the circles have been placed on the same scale

494  as those shown in Figure 3. The plots in Figure 6 reveal five key features. First, in all cases

495  the mean HR exceeds the mean FAR. This indicates that both approaches (E-RF99 and BDC)

496    and the classifiers considered herein have prediction skill when tested with independent data.

497    While this is also true for the approach of RF99, the mean HRs are relatively low compared

498    to those of either the E-RF99 or BDC approach. Second, the plots confirm the earlier finding

499    that the approach of BDC generally provides either higher hit rates, or simultaneously higher

500    hit rates and lower FARs, than that of E-RF99. Third, the mean FARs obtained using QDA

501    are not always robust. This is particularly evident for Townsville, Melbourne, Perth and Port

502    Hedland. This reflects the method's sensitivity to outliers. Fourth, when tested with

503    independent data, applying LR to the BDC variable set often produced the lowest or

504    competitive mean FARs. Fifth, the approach of BDC often produces competitive or lower

505    BSs when tested with independent data than that of E-RF99.

506

507                          **< Insert Figure 6 about here >**

508

509        A scatterplot of AUC values obtained from cross-validation of the LR models is shown in

510    Figure 7. For all sites and both approaches (E-RF99 and BDC), the AUC values are greater

511    than 0.5 indicating that prediction skill is better than climatology. However, the AUCs for the

512    BDC approach are greater than those for E-RF99. The lowest AUC values were obtained for

513    Darwin and the highest for Port Hedland (E-RF99 approach) and Perth (BDC approach).

514

515                          **< Insert Figure 7 about here >**

516

517    **4.  Summary and conclusions**

518        Daily lightning flash count and precipitation data from ground-based sensors and gauges,

519    atmospheric information from the ERA-Interim reanalysis and five classification techniques

520    (classifiers) were used to distinguish between 'dry' and 'wet' thunderstorm days for the

521     period from 2004 to 2013 at six locations in Australia. The locations of the lightning flash

522     (CIGRE 500) sensors represent a range of climatic settings (including temperate, subtropical

523     and tropical regions). The earlier approach of Rorig and Ferguson (1999, RF99), which used

524     two atmospheric variables and one classifier (linear discriminant analysis) for one region in

525     the United States (the Pacific Northwest), was used as a benchmark to test whether the

526     inclusion of additional atmospheric information and a wider range of classifiers resulted in a

527     notable improvement in prediction accuracy for the climatic settings considered herein.

528         With future applications in mind, the study was designed to be conducted at the spatial

529     resolution of current GCMs and reanalyses. Quadratic surfaces and determination of low-

530     dimensional summary statistics (LDSS) were used to capture the main features of the

531     atmospheric fields. Five classifiers were considered: classification and regression trees

532     (CART); random forests (RF); linear discriminant analysis (LDA), quadratic discriminant

533     analysis (QDA) and logistic regression (LR). Four prediction skill scores were considered,

534     with a focus on dry lightning since it is the primary cause of wildfire ignition. Ten-fold cross

535     validation was used to estimate the prediction accuracy of the classifiers. The study findings

536     can be summarized as follows:

537     1)  The use of LDSS captured useful and interpretable information in terms of the large-scale

538         spatial structure of thunderstorms. While it can be argued that the LDSS are somewhat

539         crude, our results suggest that there is value in their application to the problem of

540         thunderstorm classification.

541     2)  The approach outlined in this paper (BDC) and an extended version of that of Rorig and

542         Ferguson (1999, herein designated as E-RF99) have prediction skill when tested against

543         independent data for a wide range of climatic zones.

544     3)  Overall, while five LDSS were used to better capture the main features of the atmospheric

545         fields used, the mean field proved to be the most useful. The seven most-frequent

546    variables across the six sites and five classifiers considered are associated with

547    atmospheric water content (mu.TOTP, mu.CONVP, gd.TOTP and mu.TCWV) and

548    instability and lifting potential (mu.CBH, mu.DD700, and vg.T). The preceding lists of

549    variables contain five spatial means, two gradient terms and variables derived from

550    convective parameterizations. The results presented herein suggest that the latter may

551    provide unique information that is not contained in ground-based precipitation data.

552    4) Despite the finding above, the set of influential atmospheric variables varied from site-to-

553    site and between classifiers. This result needs to be tested using data from dense

554    monitoring networks in different countries and a wide variety of climatic zones. The

555    question of whether it is legitimate to use the same atmospheric variables and statistical

556    classification techniques at different locations within the same climatic zone will be the

557    subject of future research.

558    5) No single classifier proved to be consistently superior to its counterparts across the six

559    sites considered. However, LR often produced lower FARs while the predictive accuracy

560    of QDA was compromised by the presence of outliers in the variables.

561    6) Although the BDC variable selection approach requires more effort than that of E-RF99,

562    with the exception of the Port Hedland site it produced either higher hit rates or

563    simultaneously higher hit rates and lower false alarm ratios for dry lightning than that of

564    E-RF99. It also tended to produce lower Brier (1950) scores and higher AUCs for LR

565    models.

566    Although a number of previous studies have examined lightning and thunderstorm activity

567    at the spatial and temporal scales of current reanalyses and GCMs, very few of these studies

568    have considered 'dry' and 'wet' systems separately. The results presented here are intended

569    to lead to an improved ability to classify deep convective systems in terms of their likelihood

570    of being 'dry' or 'wet', as well as enhanced capability to understand the observed

571  climatological characteristics of these systems. It is envisaged that the approach of this study

572  will find application in future studies involving finer-scale reanalyses and GCM runs as they

573  become available. Such work might lead to classification decision rules and boundaries that

574  are less dependent on model parameterizations.  Given the importance of dry thunderstorms

575  for the ignition of wildfires by lightning, as well as wet thunderstorms in relation to a range

576  of associated hazards (including extreme rainfall), a greater understanding of dry and wet

577  thunderstorms could have significant benefits for improved resilience to the impacts of

578  lightning and thunderstorms throughout the world.

579

588

589  APPENDIX

590  Representation of atmospheric variables

591      Most of the daily atmospheric variable information is available at single pressure level or

592  is defined as a mean or difference for fixed pressure levels and hence can be considered as a

593  function of two spatial dimensions: $z = f(x, y)$. An exception is convective mass flux (CMF)

594  which, by definition, has a constant value across all 49 grid points for a given day and UTC

595  time. Other variables such as air temperature, minimum geostrophic vorticity, vertical

596    velocity, specific humidity, and zonal and meridional wind are defined for specific

597    atmospheric pressure levels ($p$) at each grid point (Table 2). These variables can be

598    considered as a function of three spatial dimensions: $z = f(x,\, y,\, p)$. For each day, quadratic

599    surfaces were fitted to the atmospheric fields for 0600 UTC using ordinary least squares. A

600    quadratic surface in two spatial dimensions is defined by

601

602    $$z = f(x, y) = c_1 + c_2 x + c_3 x^2 + c_4 y + c_5 xy + c_6 y^2 \qquad\qquad \text{(A.1)}$$

603

604    and the corresponding surface in three spatial dimensions by

605

606    $$z = f(x, y, p)\ = c_1 + c_2 x + c_3 x^2 + c_4 y + c_5 xy + c_6 y^2 + c_7 p +$$

607    $$c_8 xp + c_9 yp + c_{10} p^2 \qquad\qquad \text{(A.2)}$$

608

609    Instead of fitting (A.1) and (A.2) directly, the linear and quadratic terms were replaced by

610    orthogonal polynomials in order to ensure that: the intercept and linear and quadratic

611    regression coefficients are independent of each other (i.e. they do not change when higher-

612    order terms are added); the estimates of the intercept and regression coefficients are placed on

613    the same scale; and it allows the decomposition of relationships into general components of

614    magnitude as well as into linear and nonlinear rates of change. The estimates were calculated

615    in a coordinate system centered on the CIGRE 500 sensor (i.e., a $7 \times 7$ grid described in

616    Section 2a). The adjusted $R^2$ was used as a goodness-of-fit measure for the quadratic surfaces.

617    Let $\theta_1, \ldots, \theta_{10}$ denote the orthogonal polynomial regression coefficients. Five low-

618    dimensional summary statistics (LDSS) for the above surfaces were used to facilitate physical

619    interpretation: the intercept which is equivalent to the mean across the domain ($mu = \theta_1$); the

620    magnitude of the gradient vector (gd) and its direction (dr) in the $x - y$ plane; Gaussian

621    curvature (gc); and vertical gradient $(vg = c_7)$. The magnitude of the gradient vector and its

622    direction in terms of linear rate of change are defined by $gd = \sqrt{\theta_2^2 + \theta_4^2}$ and $dr = \tan^{-1}(\theta_4/\theta_2)$

623    . Given the use of orthogonal polynomial regression, the values of gd and dr are the same as

624    those that would have been obtained had a linear surface been fitted to the data. Gaussian

625    curvature is an intrinsic geometric property of a surface which is independent of the

626    coordinate system used to describe it. It is defined by

627

628    $$gc = \det(\mathbf{H}) = \lambda_1 \lambda_2 \tag{A.3}$$

629

630    where $\det(\bullet)$ denotes the determinant, $\mathbf{H}$ is the Hessian matrix given by

631

632    $$\mathbf{H} = \begin{bmatrix} \dfrac{\partial^2 z}{\partial x^2} & \dfrac{\partial^2 z}{\partial x \partial y} \\ \dfrac{\partial^2 z}{\partial y \partial x} & \dfrac{\partial^2 z}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2\theta_3 & \theta_5 \\ \theta_5 & 2\theta_6 \end{bmatrix} \tag{A.4}$$

633

634    and $\lambda_1$ and $\lambda_2$ are the eigenvalues of $\mathbf{H}$ (also the maximum and minimum principal

635    curvatures).

636

637                                    REFERENCES

638

639    Abalos, M., B. Legras, F. Ploeger, and W. J. Randel, 2015: Evaluating the advective Brewer-

640    Dobson circulation in three reanalyses for the period 1979-2012. *J. Geophys. Res. Atmos.*,

641    **120**, 7534-7554, doi:10.1002/2015JD023182.

642

643  Abatzoglou, J. T., C. A. Kolden, J. K. Balch, and B. A. Bradley, 2016: Controls on

644  interannual variability in lightning-caused fire activity in the western US. *Environ. Res. Lett*.,

645  **11**, 1-11, doi:10.1088/1748-9326/11/4/045005.

646

647  Alfons, A., 2012: cvTools: Cross-validation tools for regression models. R package version

648  0.3.2. http://CRAN.R-project.org/package=cvTools.

649

650  Allen, J. T., D. J. Karoly, and G. A. Mills, 2011: A severe thunderstorm climatology for

651  Australia and associated thunderstorm environments. . *Aust. Meteorol. Oceanogr. J.*, **61**, 143-

652  158.

653

654  Barthe, C., W. Deierling, and M. C. Barth, 2010: Estimation of total lightning from various

655  storm parameters: A cloud-resolving study. *J. Geophys. Res*., **115**, D24202,

656  doi:10.1029/2010JD014405.

657

658  Bates, B. C., R. E. Chandler, and A. J. Dowdy, 2015: Estimating trends and seasonality in

659  Australian monthly lightning flash counts. *J. Geophys. Res. Atmos*., **120**, 3973–3983,

660  doi:10.1002/2014JD023011.

661

662  Belsley, D., E. Kuh, and R. Welsch, 1980: *Regression Diagnostics Identifying Influential*

663  *Data and Sources of Collinearity*. Wiley, 292 pp.

664

665     Blouin, K. D., M. D. Flannigan, X. Wang, and B. Kochtubajda, 2016: Ensemble lightning

666     prediction models for the province of Alberta, Canada. *Intl. J. Wildland Fire*, **25**, 421-432,

667     http://dx.doi.org/10.1071/WF15111.

668

669     Breiman, L., 2001: Random forests. *Mach. Learn.*, **45**, 5-32, doi:10.1023/A:1010933404324.

670

671     Brier, G. W., 1950: Verification of forecasts expressed in terms of probability. *Mon. Wea.*

672     *Rev.,* **78,** 1–3, doi: http://dx.doi.org/10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2.

673

674     Burrows, W. R., C. Price, and L. J. Wilson, 2005: Warm season lightning probability

675     prediction for Canada and the northern United States. *Wea. Forecasting*, **20**, 971-988.

676

677     Chandler, R. E., 2005: On the use of generalized linear models for interpreting climate

678     variability. *Environmetrics,* **16**, 699-715, doi:10.1002/env.731.

679

680     Christian, H. J., R. J. Blakeslee, D. J. Boccippio, W. L. Boeck, D. E. Buechler, K. T. Driscoll,

681     S. J. Goodman, J. M. Hall, W. J. Koshak, D. M. Mach, and M. F. Stewart, 2003: Global

682     frequency and distribution of lightning as observed from space by the Optical Transient

683     Detector. *Journal of Geophysical Research: Atmospheres*, 108(D1).

684

685     Chronis T. G., S. J. Goodman, D. Cecil, D. Buechler, F. J. Robertson, J. Pittman, and R. J.

686     Blakeslee, 2008: Global lightning activity from the ENSO perspective. *Geophys. Res. Lett.*,

687     **35**, L19804, doi: (2008).10.1029/2008GL034321.

688

689  Dai, J., Y. Wang, L. Chen, L. Tao, J. Gu, J. Wang, X. Xu, H. Lin, and Y. Gu, 2009: A

690  comparison of lightning activity and convective indices over some monsoon-prone areas of

691  China. *Atmos. Res.*, **91**, 438-452, doi:10.1016/j.atmosres.2008.08.002.

692

693  Davies, L., C. Jakob, P. May, V. V. Kumar, and S. Xie, 2013: Relationships between the

694  large-scale atmosphere and the small-scale convective state for Darwin, Australia. *J.*

695  *Geophys. Res. Atmos.*, **118**, 11,534-11,545, doi:10.1002/jgrd.50645.

696

697  Dee, D. P., and Coauthors, 2011: The ERA-Interim reanalysis: configuration and

698  performance of the data assimilation system. *Q. J. R. Meteorol. Soc.*, **137**, 553-597,

699  doi:10.1002/qj.828.

700

701  Deierling, W., W. A. Petersen, J. Latham, S. Ellis, and H. J. Christian, 2008: The relationship

702  between lightning activity and ice fluxes in thunderstorms. *J. Geophys. Res.*, **113**, D15210,

703  doi:10.1029/2007JD009700.

704

705  Dowdy, A. J., 2015: Large-scale modelling of environments favourable for dry lightning

706  occurrence. *MODSIM2015, 21st International Congress on Modelling and Simulation.*

707  *Modelling and Simulation Society of Australia and New Zealand*, Gold Coast, Queensland,

708  Australia, T. Weber, M. J. McPhee, and R. S. Anderssen, Eds, 1524-1530, ISBN: 978-0-

709  9872143-5-5. [Available online at www.mssanz.org.au/modsim2015/G4/dowdy.pdf.]

710

711  Dowdy, A. J., 2016: Seasonal forecasting of lightning and thunderstorm activity in tropical

712  and temperate regions of the world. *Sci. Rep.*, **6**, doi:10.1038/srep20874.

713

714    Dowdy, A. J., and Y. Kuleshov, 2014: Lightning climatology of Australia: temporal and

715    spatial variability. *Aust. Meteorol. Oceanogr. J.*, **64**, 103-108.

716

717    Dowdy, A. J., and G. A. Mills, 2009: Atmospheric States Associated with the Igntion of

718    Lightning-Attributed Fires. CAWCR Tech. Rep. No. 019, 34 pp.

719

720    Dowdy, A. J., and G. A. Mills, 2012a: Characteristics of lightning-attributed fires in south-

721    east Australia. *Int. J. Wildland Fire*, **21**, 521–524, doi.org/10.1071/WF10145.

722

723    Dowdy, A. J., and G. A. Mills, 2012b: Atmospheric and fuel moisture characteristics

724    associated with lightning-attributed fires. *J. Appl. Meteor. Climatol.,* **51**, 2025-2037,

725    doi:10.1175/JAMC-D-11-0219.1.

726

727    Ensor, L. A., and Robeson, S. M., 2008: Statistical characteristics of daily precipitation:

728    Comparisons of gridded and point datasets. *J. Appl. Meteor. Climatol.*, **47**, 2468-2476,

729    doi:10.1175/2008JAMC1757.1.

730

731    Everitt, B .S., and G. Dunn, 2001: *Applied Multivariate Data Analysis.* Second edition.

732    Arnold, 342 pp.

733

734    Everitt, B., and T. Hothorn, 2011: *An Introduction to Applied Multivariate Analysis with R.*

735    Springer, 273 pp.

736

737    Faraway, J. J., 2016: *Extending the linear model with R: generalized linear, mixed effects and*

738    *nonparametric regression models.* Second edition. Chapman and Hall / CRC Press, 399 pp.

739

740 Goodman S. J., D. E. Buechler, K. Knupp, D. Driscoll, and E. E. McCaul Jr., 2000: The

741 1997–98 El Niño event and related wintertime lightning variations in the southeastern United

742 States. *Geophys. Res. Lett*. **27**, 541–544, doi:10.1029/1999GL010808.

743

744 Hendrickx, J., 2012: perturb: Tools for evaluating collinearity. R package version 2.05.

745 http://CRAN.R-project.org/package=perturb.

746

747 Hilbe, J. M., 2009: *Logistic regression models*. Chapman and Hall/CRC, 656 pp.

748

749 Jayaratne, E. R., and Y. Kuleshov, 2006: Geographical and seasonal characteristics of the

750 relationship between lightning ground flash density and rainfall within the continent of

751 Australia. *Atmos. Res*., **79**, 1–14, doi:10.1016/j.atmosres.2005.03.004.

752

753 Kasischke, E. R., T. S. Rupp, and D. L. Verbyla, 2006: Fire trends in the Alaskan Boreal

754 Forest. *Alaska's Changing Boreal Forest*, F. S. Chapin, M. Oswood, K. van Cleve, L.

755 Viereck, and D. Verbyla, Eds., Oxford University Press, 285-301.

756

757 Kuleshov, Y., and E. R. Jayaratne, 2004: Estimates of lightning ground flash density in

758 Australia and its relationship to thunder-days. *Aust. Meteorol. Mag*., **53**, 189–196.

759

760 Kuleshov, Y., D. Mackerras, and M. Darveniza, 2009: Spatial distribution and frequency of

761 thunderstorms and lightning in Australia. *Lightning: Principles, Instruments and*

762 *Applications*, H. D. Betz, U. Schumann, and P. Laroche, Eds., Springer, 187–207,

763 doi:10.1007/978-1-4020-9079-0_8.

764

765    Kuo, H. L., and W. H. Raymond, 1980: A quasi-one-dimensional cumulus cloud model and

766    parametrization of cumulus heating and mixing effects. *Mon. Wea. Rev*., **108**, 991–1009, doi:

767    http://dx.doi.org/10.1175/1520-0493(1980)108<0991:AQODCC>2.0.CO;2.

768

769    Lang, T. J., S. A. Rutledge, B. Dolan, P. Krehbiel, W. Rison, and D. T.  Lindsey, 2014:

770    Lightning in wildfire smoke plumes observed in Colorado during summer 2012. *Mon. Wea.*

771    *Rev*., **142**, 489-507, doi:http://dx.doi.org/10.1175/MWR-D-13-00184.1.

772

773    Liaw, A., and M. Wiener, 2002: Classification and Regression by randomForest. *R News*, **2**,

774    18-22.

775

776    Luke, R. H., and A. G. McArthur, 1978: *Bushfires in Australia*. Australian Government

777    Publishing Service, 368 pp.

778

779    Magi, B. I., 2015: Global lightning parameterization from CMIP5 climate model output. *J.*

780    *Atmos. Oceanic Technol.,* **32**, 434-452*,* doi:10.1175/JTECH-D-13-00261.1.

781

782    McRae, R. H. D., 1992: Prediction of areas prone to lightning ignition. *Intl. J. Wildland Fire*,

783    **2**, 123-130, doi:10.1071/WF9920123.

784

785    Morris, J. S., 2015: Functional regression. *Annu. Rev. Stat. Appl*., **2**, 321-359,

786    doi:10.1146/annurev-statistics-010814-020413.

787

788  Muñoz, Á. G., J. Díaz-Lobatón, X. Chourio, and M. J. Stock, 2016: Seasonal prediction of

789  lightning activity in North Western Venezuela: Large-scale versus local drivers. *Atmos. Res*.,

790  **172-173**, 147-162, doi:10.1016/j.atmosres.2015.12.018

791

792  Nauslar, N. J., M. L. Kaplan, J. Wallmann, and T. J. Brown, 2013: A forecast procedure for

793  dry thunderstorms. *J. Oper. Meteor*., **1**, 200-214,

794  doi:http://dx.doi.org/10.15191/nwajom.2013.0117.

795

796  NCAR - Research Applications Laboratory, 2015: verification: Weather Forecast Verification

797  Utilities. R package version 1.42. http://CRAN.R-project.org/package=verification.

798

799  R Core Team (2015). R: A language and environment for statistical computing. R Foundation

800  for Statistical Computing, Vienna, Austria. URL http://www.R-project.org/.

801

802  Rakov, V. A. and M. A. Uman, 2003: *Lightning: Physics and Effects*, Cambridge University

803  Press, 687 pp.

804

805  Ripley, B., 2014: tree: Classification and regression trees. R package version 1.0-35.

806  http://CRAN.R-project.org/package=tree.

807

808  Romero, R., M. Gayà, and C. A. Doswell, 2007: European climatology of severe convective

809  storm environmental parameters: A test for significant tornado events. *Atmos. Res*., **83**, 389-

810  404, doi:10.1016/j.atmosres.2005.06.011.

811

812    Romps, D. M., J. T. Seeley, D. Vollaro, and J. Molinari, 2014: Projected increase in lightning

813    strikes in the United States due to global warming. *Science*, **346**, 851-854,

814    doi:10.1126/science.1259100.

815

816    Rorig, M. L., and S. A. Ferguson, 1999: Characteristics of lightning and wildland fire ignition

817    in the Pacific Northwest. *J. Appl. Meteor.*, **38**, 1565-1575,

818    doi:http://dx.doi.org/10.1175/1520-0450(1999)038<1565:COLAWF>2.0.CO;2.

819

820    Rorig, M. L., S. J. McKay, S. A. Ferguson, and P. Werth, 2007: Model-generated predictions

821    of dry thunderstorm potential. *J. Appl. Meteor. Climatol.*, **46**, 605-614,

822    doi:http://dx.doi.org/10.1175/JAM2482.1.

823

824    Rothermel, R. C., 1972: A mathematical model for predicting fire spreads in wildland fuels.

825    U.S. Forest Service Res. Paper INT-115, 40 pp.

826

827    Sanchez, G., 2013: DiscriMiner: Tools of the Trade for Discriminant Analysis. R package

828    version 0.1-29. http://CRAN.R-project.org/package=DiscriMiner.

829

830    Takahashi, T., and K. Miyawaki, 2002: Reexamination of riming electrification in a wind

831    tunnel. *J. Atmos. Sci.*, **59**, 1018-1025, doi: http://dx.doi.org/10.1175/1520-

832    0469(2002)059<1018:ROREIA>2.0.CO;2.

833

834    Tiedtke, M., 1989: A comprehensive mass flux scheme for cumulus parameterization in

835    large-scale models. *Mon. Wea. Rev.*, **117**, 1779–1800, doi:http://dx.doi.org/10.1175/1520-

836    0493(1989)117<1779:ACMFSF>2.0.CO;2.

837

838 USDA Forest Service, 1992: 1984-1990 Wildfire Statistics. State and Private Forestry, Fire

839 and Aviation Management Staff, USDA Forest Service (Washington, DC), 240 pp.

840

841 Vazquez A., and J. M. Moreno, 1998: Patterns of lightning- and people-caused fires in

842 peninsula Spain. *Intl. J. Wildland Fire*, **8**, 103-115, doi:10.1071/WF9980103.

843

844 Venables, W. N. and B. D. Ripley, 2002: *Modern Applied Statistics with S*. Fourth Edition.

845 Springer, 495 pp.

846

847 Wallmann, J., R. Milne, C. Smallcomb, and M. Mehle, 2010: Using the 21 June 2008

848 California lightning outbreak to improve dry lightning forecast procedures. *Wea.*

849 *Forecasting*, **25**, 1447-1462, doi:http://dx.doi.org/10.1175/2010WAF2222393.1.

850

851 Weisman, M. L., and J. B. Klemp, 1982: The dependence of numerically simulated

852 convective storms on vertical wind shear and buoyancy. *Mon. Wea. Rev*., **110**, 504-520,

853 doi:http://dx.doi.org/10.1175/1520-0493(1982)110<0504:TDONSC>2.0.CO;2.

854

855 Williams, E., V. Mushtak, D. Rosenfeld, S. Goodman, and D. Boccippio, 2005:

856 Thermodynamic conditions favorable to superlative thunderstorm updraft, mixed phase

857 microphysics and lightning flash rate. *Atmos. Res*., **76**, 288-306,

858 doi:10.1016/j.atmosres.2004.11.009.

859

860 Wotton B. M., and D. L. Martell, 2005: A lightning fire occurrence model for Ontario. *Can.*

861 *J. For. Res*., **35**, 1389-1401, doi:10.1139/x05-071.

862

863    Wotton, B. M., B. J. Stocks, and D. L. Martell, 2005: An index for tracking sheltered forest

864    floor moisture within the Canadian Forest Fire Weather Index System. *Intl. J. Wildland Fire*,

865    **14**, 169-182, http://dx.doi.org/10.1071/WF04038.

866

867    Xu, W., E. J. Zipser, C. Liu, and H. Jiang, 2010: On the relationship between lightning

868    frequency and thundercloud parameters of regional precipitation systems. *J. Geophys. Res.,*

869    *Atmos*., **115**, D12203, doi:10.1029/2009JD013385.

870

871    Yan, Z., S. Bate, R. E. Chandler, V. Isham, and H. S. Wheater, 2002: An analysis of daily

872    maximum windspeed in northwestern Europe using generalized linear models. *J. Climate*, **15,**

873    2073-2088, doi:http://dx.doi.org/10.1175/1520-0442(2002)015<2073:AAODMW>2.0.CO;2.

874

875    Youden, W. J., 1950: Index for rating diagnostic tests. *Cancer* **3**, 32–35, doi:10.1002/1097-

876    0142(1950)3:1<32::AID-CNCR2820030106>3.0.CO;2-3.

877

878

879

880

881     TABLE 1. Site and data details for CIGRE 500 lightning flash counters. Daily lightning flash count records cover the period from January 2004

882     to at least December 2010 (Townsville) and at most February 2013 (Melbourne).

| Site No. | Location | Altitude (m) | Köppen classification | No. of lightning days | Proportion dry lightning days | Proportion wet lightning days |
|---|---|---|---|---|---|---|
| 1 | Darwin | 30 | Tropical savanna climate (Aw) | 1350 | 0.53 | 0.47 |
| 2 | Townsville | 4 | Tropical savanna climate (Aw) | 286 | 0.53 | 0.47 |
| 3 | Coffs Harbour | 5 | Humid subtropical climate (Cfa) | 501 | 0.58 | 0.42 |
| 4 | Melbourne | 113 | Marine west coast (Cfb) | 570 | 0.64 | 0.36 |
| 5 | Perth | 15 | Mediterranean (Csa) | 148 | 0.55 | 0.45 |
| 6 | Port Hedland | 6 | Subtropical desert (BWh) | 401 | 0.81 | 0.19 |

883

884

885

886

887 TABLE 2. Abbreviations, full names, units of measure and specifications for atmospheric variables.

| Abbreviation | Full name | Specification |
|---|---|---|
| *Instability and lifting potential* | | |
| CAPE | Convective available potential energy (J kg$^{-1}$) | As provided in ERA-Interim reanalysis (maximum CAPE based on lifting parcels within a near-surface layer) |
| CBH | Cloud base height (m) | Based on temperature and dewpoint at a height of 2 m with lifting to condensation level using an idealized constant lapse rate |
| CMF | Convective mass flux (Pa$^2$ s$^{-1}$ K$^{-1}$) | 500 hPa: calculated as the product of air density, fraction of grid points covered by updrafts within the 7x7 gridded region, and the vertical velocity averaged across all updrafts. |

| CONV1000850 | Mean low-level horizontal wind convergence ($s^{-1}$) | Mean value at 850 and 1000 hPa pressure levels |
|---|---|---|
| DD | Dewpoint depression (°C) | 500, 700 and 850 hPa |
| DDIV | Density-weighted mean upper-level divergence minus density-weighted mean low-level divergence ($s^{-1}$) | {300, 400} – {850, 1000} hPa |
| EPTL | Mean low-level equivalent potential temperature minus mean mid-level equivalent potential temperature (°C) | Mean value at 1000 and 850 hPa – mean value at 700 and 500 hPa |
| TD850T500 | Cross totals index (°C) | 850 and 500 hPa |
| TGD | Direction of thickness gradient (rad) | {500, 700}, {500, 1000} and {700, 1000} hPa |
| TGM | Magnitude of thickness gradient ($m^2\,s^{-2}$) | {500, 700}, {500, 1000} and {700, 1000} hPa |
| THETA_W1000 | Wet-bulb potential temperature (°C) | 1000 hPa |
| THETA_W850500 | Wet-bulb potential temperature difference (°C) | 850 – 500 hPa |
| THK7001000 | Geopotential thickness ($m^2\,s^{-2}$) | 700 – 1000 hPa geopotential heights |
| TL850500 | Temperature lapse (°C) | 850 – 500 hPa |

888

| TL850700 | Temperature lapse (°C) | 850 – 700 hPa |
|---|---|---|
| TTI | Total totals index (°C) | 850 and 500 hPa |
| W | Vertical velocity (Pa s$^{-1}$) | 200, 300, 500, 700, 850 and 1000 hPa |
| *Atmospheric water content* | | |
| CONVP | Convective precipitation (m) | As provided in ERA-Interim reanalysis |
| ICE | Total column ice water (kg m$^{-2}$) | As provided in ERA-Interim reanalysis |
| SH | Specific humidity (kg kg$^{-1}$) | 500, 700 and 850 hPa |
| TCWV | Total column water vapor (kg m$^{-2}$) | As provided in ERA-Interim reanalysis |
| TOTP | Total precipitation (m) | As provided in ERA-Interim reanalysis |
| *Wind speed* | | |
| MVWS | Maximum vertical wind shear (m s$^{-1}$) | 300 to 850 hPa |
| S06 | Vertical wind shear between 0 and 6 km (m s$^{-1}$) | 1000 and 500 hPa |
| U | Zonal wind velocity (m s$^{-1}$) | 300, 500, 700, 850 and 1000 hPa |

| V | Meridional wind velocity (m s$^{-1}$) | 300, 500, 700, 850 and 1000 hPa |
|---|---|---|
| *General atmospheric state and variability* | | |
| SEASON | Season-of-year | DJF, MAM, JJA and SON |
| T | Air temperature (°C) | 2 meters, 500, 700 and 850 hPa |
| MSLP | Mean sea level pressure (Pa) | As provided in ERA-Interim reanalysis |
| GPH | Geopotential height (m$^2$ s$^{-2}$) | 500 and 700 hPa |
| MING | Minimum geostrophic vorticity (s$^{-2}$) | Laplacian of geopotential at 500, 700 and 850 hPa |

889

890

891                     TABLE 3. Frequently used acronyms and abbreviations

| Acronym or abbreviation | Full name |
|---|---|
| AUC | Area under receiver operating characteristic curve |
| BDC | Approach of Bates, Dowdy and Chandler (this paper) |
| BS | Brier (1950) score |
| CART | Classification and regression trees |
| E-RF99 | Extended approach of Rorig and Ferguson (1999) |
| FAR | False alarm ratio for dry lightning |
| GCM | General circulation model |
| HR | Hit rate for dry lightning |
| LDA | Linear discriminant analysis |
| LDSS | Low-dimensional summary statistics |
| LR | Logistic regression |
| QDA | Quadratic discriminant analysis |
| RF | Random forests |
| RF99 | Approach of Rorig and Ferguson (1999) |
| ROC | Receiver operating characteristic |

892

893     LIST OF FIGURES

894         FIG. 1. Locations of CIGRE 500 lightning flash counters (filled circles) and relative

895     proportions of dry lightning and wet lightning days in daily lightning flash count series. Key

896     to numerals is given in Table 1. Widths of gray rectangles indicate proportions of dry

897     lightning days and heights proportions of wet lightning days.

898         FIG. 2. Examples of comparative boxplots of potential candidate variables for the Coffs

899     Harbour CIGRE 500 site: (a) mu.CBH, (b) mu.CONVP, (c) mu.DD700, (d) mu.DD850, (e)

900     mu.ICE, and (f) mu.TOTP.

901         FIG. 3. Skill scores obtained using the methods of E-RF99 (filled squares) and BDC

902     (filled triangles) for five classifiers and six CIGRE 500 sites. Radii of the circles are

903     proportional to the Brier score. Dashed lines represent the convex hull of the false alarm ratio

904     (FAR) and hit rate (HR) values for dry lightning obtained using the methods of E-RF99.

905         FIG. 4. Relative frequency histogram of selected variables across six CIGRE 500 sites and

906     four classifiers: classification and regression trees (CART), random forests (RF), linear

907     discriminant analysis (LDA) and logistic (LR).

908         FIG. 5. Relative frequency histograms of influential variables for discriminating dry

909     lightning from wet lightning days for each CIGRE 500 site across four classifiers:

910     classification and regression trees (CART), random forests (RF), linear discriminant analysis

911     (LDA) and logistic (LR).

912         FIG. 6. Mean skill scores obtained from cross validation experiments using the methods of

913     E-RF99 (filled squares) and BDC (filled triangles) for five classifiers and six CIGRE 500

914     sites. Radii of the circles are proportional to the Brier score. Dashed lines represent the

915     convex hull of the mean false alarm ratio (FAR) and hit rate (HR) values for dry lightning

916     obtained using the methods of E-RF99.

917  FIG. 7. Scatter plot of mean area under receiver operating characteristic curve (AUC)

918 values obtained from cross-validation of logistic regression (LR) models. Key to numerals is
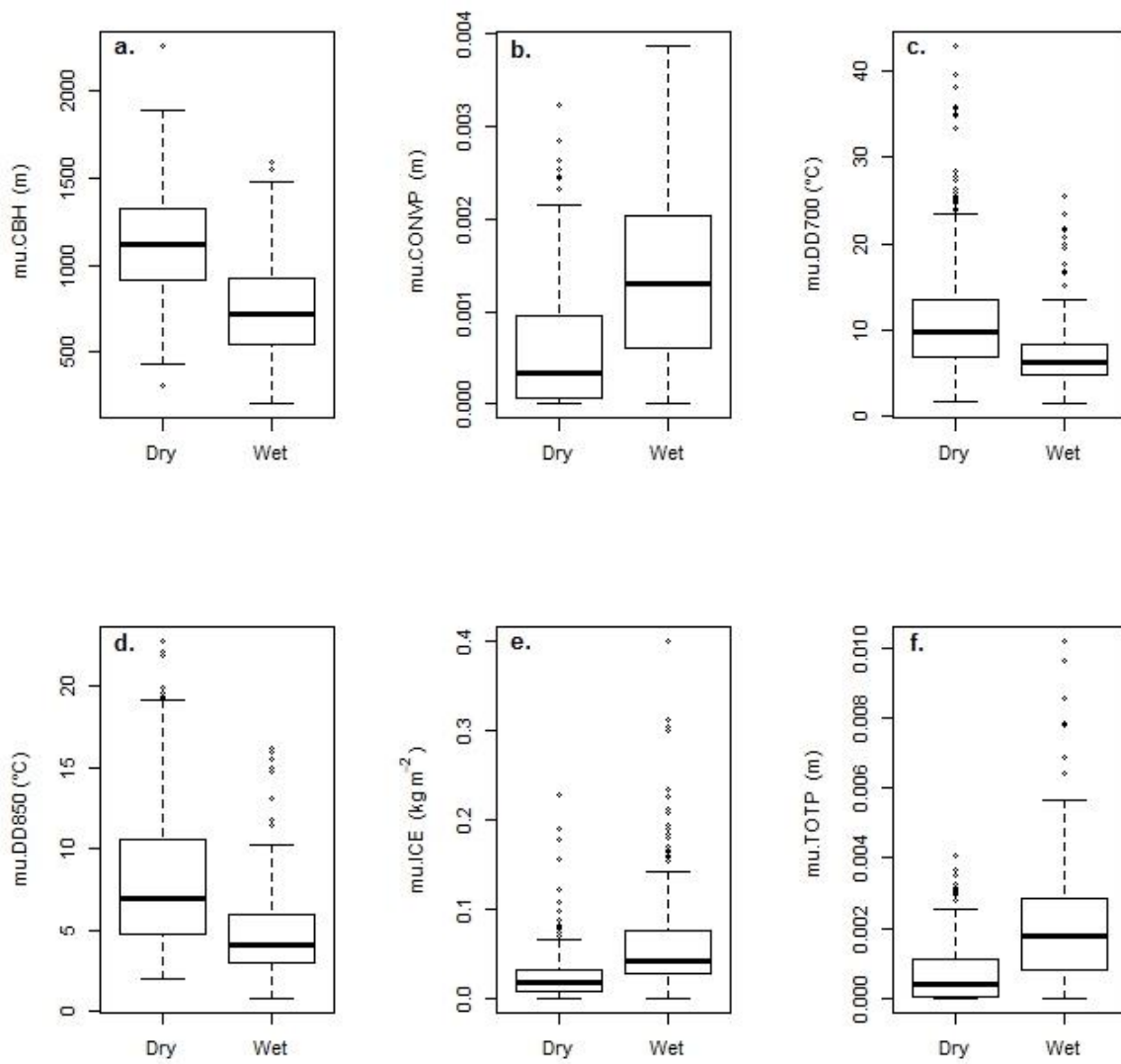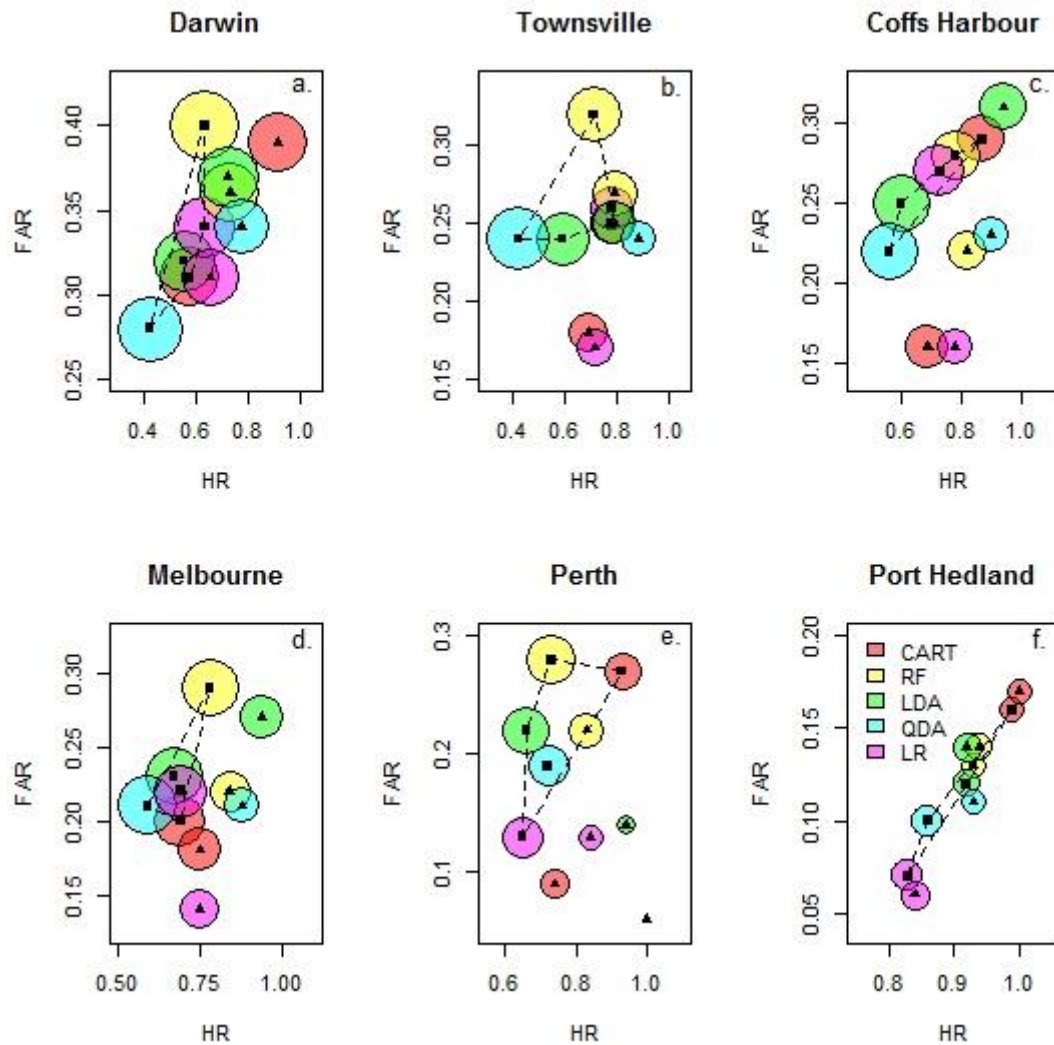
919 given in Table 1.

920

921

922

923     FIG. 1. Locations of CIGRE 500 sensors (filled circles) and relative proportions of dry

924    lightning and wet lightning days in daily lightning flash count series. Key to numerals is

925    given in Table 1. Widths of gray rectangles indicate proportions of dry lightning days and

926    heights proportions of wet lightning days.

927

928

929

930    FIG. 2. Examples of comparative boxplots of potential candidate variables for the Coffs

931    Harbour CIGRE 500 site: (a) mu.CBH, (b) mu.CONVP, (c) mu.DD700, (d) mu.DD850, (e)

932    mu.ICE, and (f) mu.TOTP.

933

934

935

936     FIG. 3. Skill scores obtained using the methods of E-RF99 (filled squares) and BDC

937     (filled triangles) for five classifiers and six CIGRE 500 sites. Radii of the circles are

938     proportional to the Brier score. Dashed lines represent the convex hull of the false alarm ratio

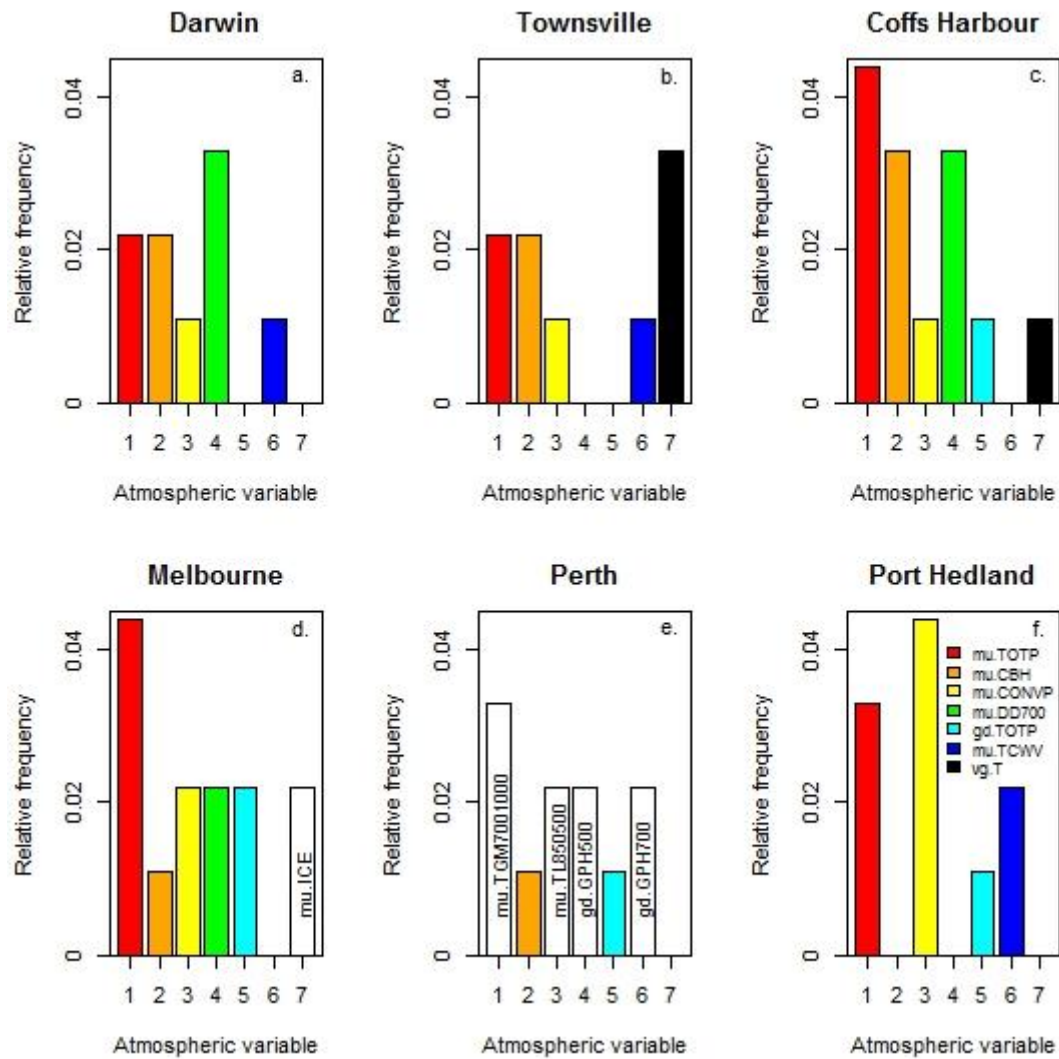939     (FAR) and hit rate (HR) values for dry lightning obtained using the methods of E-RF99.

940

48



FIG. 4. Relative frequency histogram of selected variables across six CIGRE 500 sites and

four classifiers: classification and regression trees (CART), random forests (RF), linear

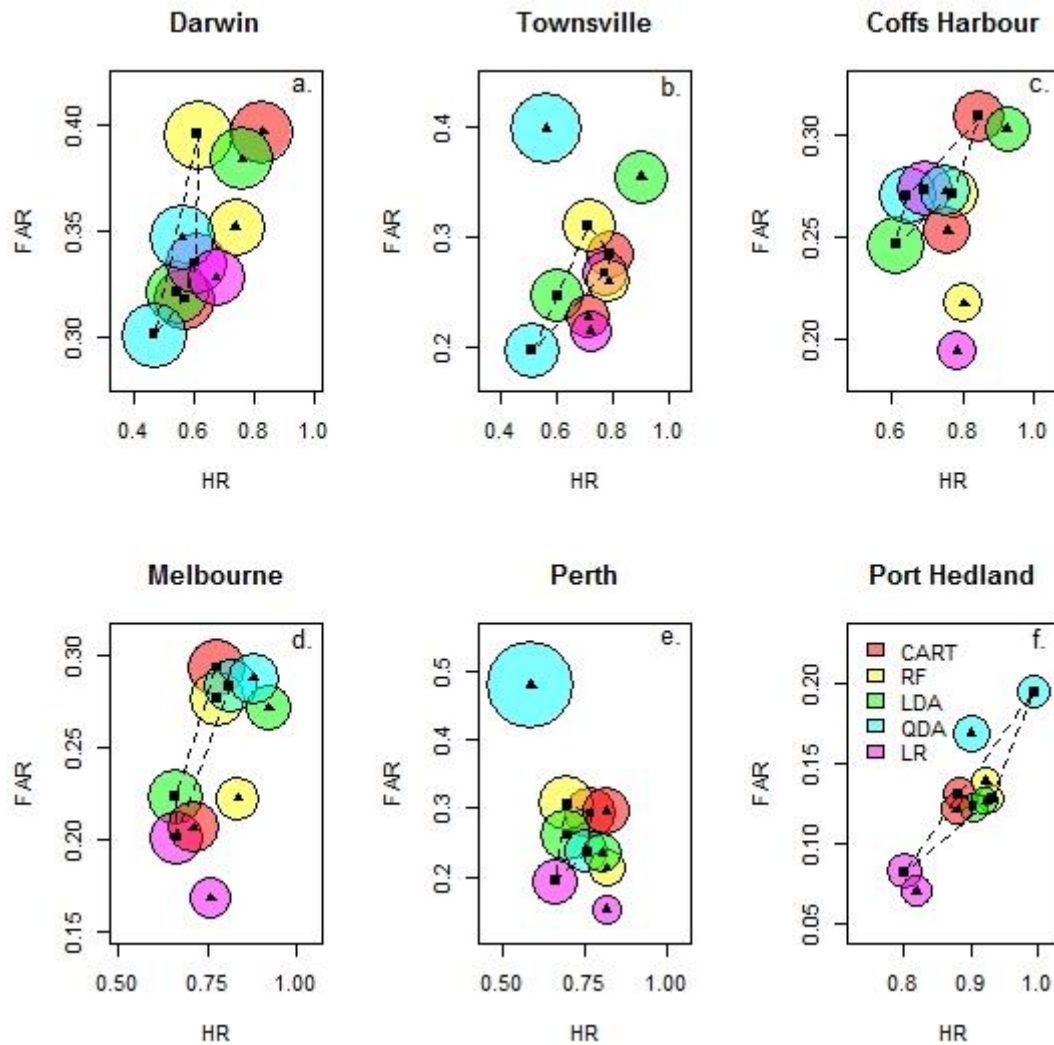discriminant analysis (LDA) and logistic (LR).

941

942

943

944

945

946

FIG. 5. Relative frequency histograms of influential variables for discriminating dry
lightning from wet lightning days for each CIGRE 500 site across four classifiers:
classification and regression trees (CART), random forests (RF), linear discriminant analysis
(LDA) and logistic (LR).

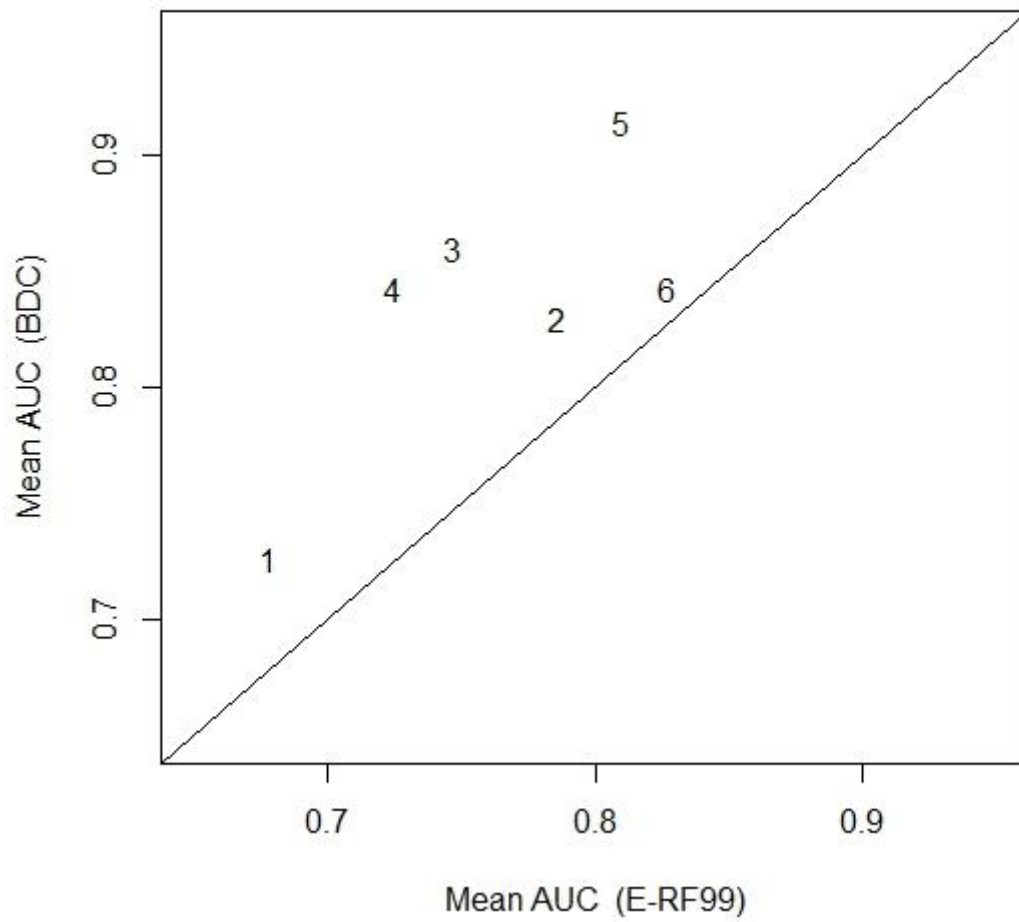FIG. 6. Mean skill scores obtained from cross validation experiments using the methods of E-RF99 (filled squares) and BDC (filled triangles) for five classifiers and six CIGRE 500 sites. Radii of the circles are proportional to the Brier score. Dashed lines represent the convex hull of the mean false alarm ratio (FAR) and hit rate (HR) values for dry lightning obtained using the methods of E-RF99.

FIG. 7. Scatter plot of mean area under receiver operating characteristic curve (AUC)

values obtained from cross-validation of logistic regression (LR) models. Key to numerals is

given in Table 1.