



1 A comparison of gap-filling algorithms for eddy covariance 2 fluxes and their drivers

3

4 **Atbin Mahabbati¹, Jason Beringer¹, Matthias Leopold¹, Ian McHugh², James Cleverly³, Peter Isaac⁴,**
5 **Azizallah Izady⁵**

6 ¹School of Agriculture and Environment, The University of Western Australia, 35 Stirling Hwy,
7 Crawley, Perth WA, 6009, Australia

8 ²School of Ecosystem and Forest Sciences, The University of Melbourne, Richmond, VIC, 3121,
9 Australia

10 ³School of Life Sciences University of Technology Sydney Broadway NSW 2007

11 ⁴OzFlux Central Node, TERN Ecosystem Processes, Melbourne, VIC 3159, Australia

12 ⁵Water Research Center, Sultan Qaboos University, Muscat, Oman

13

14 *Correspondence to:* Atbin Mahabbati (atbin.m@hotmail.com)

15

16 **Abstract**

17 The errors and uncertainties associated with gap-filling algorithms of water, carbon and
18 energy fluxes data, have always been one of the prominent challenges of the global network of
19 microclimatological tower sites that use eddy covariance (EC) technique. To address this concern, and
20 find more efficient gap-filling algorithms, we reviewed eight algorithms to estimate missing values of
21 environmental drivers, and separately three major fluxes in EC time series. We then examined the
22 performance of mentioned algorithms for different gap-filling scenarios utilising data from five
23 OzFlux Network towers during 2013. The objectives of this research were a) to evaluate the impact
24 of training and testing window lengths on the performance of each algorithm; b) to compare the
25 performance of traditional and new gap-filling techniques for the EC data, for fluxes and their
26 corresponding meteorological drivers. The performance of algorithms was evaluated by generating
27 nine different training-testing window lengths, ranging from a day to 365 days. In each scenario, the
28 gaps covered the data for the entirety of 2013 by consecutively repeating them, where, in each step,
29 values were modelled by using earlier window data. After running each scenario, a variety of
30 statistical metrics was used to evaluate the performance of the algorithms. The algorithms showed
31 different levels of sensitivity to training-testing windows; The Prophet Forecast Model (FBP) revealed
32 the most sensitivity, whilst the performance of artificial neural networks (ANNs), for instance, did not
33 vary considerably by changing the window length. The performance of the algorithms generally
34 decreased with increasing training-testing window length, yet the differences were not considerable
35 for the windows smaller than 60 days. Gap-filling of the environmental drivers showed there was not
36 a significant difference amongst the algorithms, the linear algorithms showed slight superiority over



37 those of machine learning (ML), except the random forest algorithm estimating the ground heat flux
38 (RMSEs of 30.17 and 34.93 for RF and CLR respectively). For the major fluxes, though, ML algorithms
39 showed superiority (9 % less RMSE on average), except the Support Vector Regression (SVR), which
40 provided significant bias in its estimations. Even though ANNs, random forest (RF) and extreme
41 gradient boost (XGB) showed close performance in gap-filling of the major fluxes, RF provided more
42 consistent results with less bias, relatively. The results indicated that there is no single algorithm
43 which outperforms in all situations and therefore, but RF is a potential alternative for the ANNs as
44 regards flux gap-filling.

45

46 1. Introduction

47 To address the global challenges of climatological and ecological changes, environmental
48 scientists and policymakers are demanding data that are continuous in time and space. Besides, there
49 is a need for quantifying and reducing uncertainties in such data, including observations of carbon,
50 water and energy exchanges that are crucial components in national/international flux networks and
51 global earth observing systems. Satellites partially fill this gap as they provide excellent spatial
52 coverage but at a limited temporal resolution, and not measured at the point. As such, high-quality
53 long-term site observations of ecosystem process and fluxes are needed that are continuous in time
54 and space. The global eddy covariance (EC) flux tower networks (FLUXNET), consisted of its regional
55 counterparts (i.e. AmeriFlux, EUROFLUX, OzFlux, etc.), was established in the late 1990s to address
56 the global demand for such information (Beringer et al., 2016a; Hollinger et al., 1999; Tenhunen et al.,
57 1998). Despite the capability of EC to frequently validate process modelling analyses, field surveys
58 and remote sensing assessments (Hagen et al., 2006), there are some serious concerns regarding the
59 challenges associated with the technique, e.g. data gaps and uncertainties. Hence, filling data gaps
60 and reducing uncertainties through better gap-filling techniques are highly needed.

61 Even though the EC is a common technique to measure fluxes of carbon, water and energy,
62 there are some challenges in providing robust, high-quality continuous observations. One of the
63 challenges regarding the technique, and therefore, the flux networks, is addressing data gaps and the
64 uncertainties associated with the gap-filling process, mainly when the gap windows are long (longer
65 than 12 consecutive days, as described by (Moffat et al., 2007)). These gaps happen very often due to
66 a variety of reasons, such as values out of range, spike detection or manual exclusion of date and time
67 ranges, instrument or power failure, herbivores, fire, eagles nests, cows, lightning, researchers on
68 leave, etc. (Beringer et al., 2016b). Since EC flux towers are often located in harsh climates, their data
69 are more susceptible to adverse weather (i.e. rain conditions), and they sometimes prevent quick
70 access to sites for repair and maintenance. As a result, this issue can, in turn, produce gaps which
71 might be relatively long (Isaac et al., 2017), and thus, problematic as follows. Firstly, loss of data is
72 considered a threat to scientific studies depending on the missing data quantity, pattern, mechanism
73 and nature (Altman and Bland, 2007; Molenberghs et al., 2014; Tannenbaum, 2010). That is because
74 using an incomplete dataset might lead to biased, invalid and unreliable results (Allison, 2000; Kang,



75 2013; Little, 2002). Second, continuous gap-filled data are required to calculate the annual or monthly
76 budgets of carbon or water balance components (Hutley et al., 2005).

77 Other than the challenges caused by missing data, there are several sources of errors and
78 uncertainties in the EC technique. Firstly, random error is associated with the stochastic nature of
79 turbulence, associated sampling errors (incomplete sampling of large eddies, uncertainty in the
80 calculated covariance between the vertical wind velocity and the scalar of interest), instrument errors,
81 and footprint variability (Aubinet et al., 2012). For instance, Dragoni et al. (2007) analysed an EC-based
82 data of Morgan-Monroe State Forest for eight years (1999-2006) and assessed that instrument
83 uncertainty was equal to 3 % of the total annual NEE. Another primary source of uncertainty in EC
84 measurements is systematic errors that are usually caused by methodological challenges and
85 instrument calibration problems (e.g. sonic anemometer errors, spikes, gas analyser errors, etc.).
86 Finally, one of the sources of uncertainties is data processing, especially data gap-filling (Isaac et al.,
87 2017; Moffat et al., 2007; Richardson et al., 2012; Richardson and Hollinger, 2007).

88

89 There are several uncertainties pertaining to gap-filling of missing values, including
90 measurement uncertainty (Richardson and Hollinger, 2007), lengths and timing the gaps (Falge et al.,
91 2001; Richardson and Hollinger, 2007) and the particular gap-filling algorithm that is used (Falge et
92 al., 2001; Moffat et al., 2007). However, there are two dominant issues of long data gaps and the choice
93 of a particular gap-filling algorithm (Aubinet et al., 2012). Firstly, long gaps can significantly increase
94 the total amount of uncertainty as the ecosystem behaviour might change because of different
95 agricultural periods or phenological phases (e.g. growing season, harvest period, bushfire, etc.). And
96 thereby show different responses under similar meteorological conditions (Aubinet et al., 2012; Isaac
97 et al., 2017; Richardson and Hollinger, 2007). Consequently, the period in which a long gap happens
98 is essential. For example, research undertaken by Richardson & Hollinger (2007) on data from a range
99 of FLUXNET sites revealed that a week data gap during spring green-up in a forest led to a higher
100 uncertainty over a three-week gap period during winter. Second, each gap-filling algorithm has its
101 strengths and weaknesses; for instance, Moffat et al. (2007) compared a couple of different commonly-
102 used gap-filling algorithms. They found that there was not a significant difference between the
103 performances of the algorithms with “good” reliability based on analysis of variance of RMSE.
104 Besides, the overall gap-filling uncertainty was within $\pm 25 \text{ g C m}^{-2} \text{ yr}^{-1}$ for most of the proper
105 algorithms, whereas, the other algorithms generated higher uncertainties of up to $\pm 75 \text{ g C m}^{-2} \text{ yr}^{-1}$. This
106 result is similar to the findings of Richardson & Hollinger (2007) who found uncertainties of up to
107 $\pm 30 \text{ g C m}^{-2} \text{ yr}^{-1}$ for long gaps by appropriate algorithms. Considering that the data provided by EC
108 tower networks are of use for research, government and policymakers, robust gap-filling is a need to
109 quantify and reduce uncertainties in flux estimations.

110

111 To manage the missing data problem, several methods have been typically used to fill data
112 gaps in both fluxes and their meteorological drivers. Due to computational constraints of complex



113 algorithms, early works to impute EC data gaps used interpolation methods based mostly on linear
114 regression or temporal autocorrelation (Falge et al., 2001; Lee et al., 1999). These approaches were
115 replaced quickly by more sophisticated methods such as non-linear regressions (Barr et al., 2004; Falge
116 et al., 2001; Moffat et al., 2007; Richardson et al., 2006); lookup tables (Falge et al., 2001; Law et al.,
117 2002; Zhao and Huang, 2015); artificial neural networks (ANNs) (Aubinet et al., 1999; Beringer et al.,
118 2016a; Cleverly et al., 2013; Hagen et al., 2006; Isaac et al., 2017; Kunwor et al., 2017; Moffat et al., 2007;
119 Papale and Valentini, 2003; Pilegaard et al., 2001; Staebler, 1999); mean diurnal variation (Falge et al.,
120 2001; Moffat et al., 2007; Zhao and Huang, 2015), multiple imputations (Hui et al., 2004; Moffat et al.,
121 2007), etc. Each of these methods has its pros and cons as follows: a) Interpolation methods such as
122 the Mean Diurnal Variation (MDV), do not need any drivers, yet, their accuracy is lower than other
123 approaches (Aubinet et al., 2012). Moreover, this method may provide biased results on extremely
124 clear or cloudy days (Falge et al., 2001). MDV is not recommended when a gap is longer than two
125 weeks, for it cannot consider the non-linear relations between the drivers and the flux, and thus leads
126 to a high level of uncertainty (Falge et al., 2001). And b) The Lookup table, especially its modified
127 version, Marginal Distribution Sampling (MDS), has provided performance close to ANNs, and are
128 more reliable and consistent than the other algorithms so far. Hence, MDS was chosen as one of the
129 standard gap-filling methods in EUROFLUX (Aubinet et al., 2012). Nevertheless, one of the concerns
130 regarding this algorithm is that the independent variables, here meteorological drivers, might be auto-
131 correlated. c) ANNs have commonly been used to gap-fill EC fluxes since 2000 and because of their
132 robust and consistent results are considered as a standard gap-filling algorithm in several networks,
133 e.g. ICOS, FLUXNET, OzFlux, etc. (Aubinet et al., 2012; Beringer et al., 2017; Isaac et al., 2017). Despite
134 their reliable performance, ANNs –and generally all other ML algorithms– face some challenges. Over-
135 fitting, for instance, is a big concern and can happen when the number of degrees of freedom is high,
136 while the training window is not long enough respectively, or the quality of the training dataset is
137 low. This challenge becomes acute when the gaps happen within a period when the ecosystem
138 behaviour is changing and thereby showing different response under similar meteorological
139 conditions. Furthermore, there is a desire to have the training windows short so that the algorithm
140 can track the ecosystem behaviour shift. Yet, this increases the risk of over-fitting depending on the
141 algorithm. In other words, the training window length should be neither too short to cause over-
142 fitting, and nor too long to lead algorithms to ignore ecological condition changes. Besides, long gaps
143 are considered as one of the primary uncertainty sources of CO₂ flux in the FLUXNET (Aubinet et al.,
144 2012). As a result, studying the effects of the gap lengths, as well as the window length whereby an
145 algorithm is trained are both critical challenges associated with the environmental data gap-filling.

146

147 Apart from the limitations and disadvantages of the mentioned algorithms, gap-filling of fluxes
148 (i.e. NEE) experiences some other challenges that make it necessary to find or develop new gap-filling
149 algorithms. That is because the current methods are not flexible enough to perform well in special
150 occasions or extreme values (Kunwor et al., 2017), and there is almost no room to optimise them to
151 improve their outcome (Moffat et al., 2007). Moreover, even using the best available algorithm, such



152 as ANNs, the model (gap-filling) uncertainty still accounts for a sizable proportion of the total
153 uncertainties, especially when the gaps are relatively long. Since the 2000s when MDS and ANNs were
154 chosen as the most reliable gap-filling methods for EC flux observations, many new ML and
155 optimisation algorithms have been developed and used in varieties of scientific fields. Some of which
156 have shown superiority over ANNs, either individually or as a part of a hybrid or ensemble model,
157 e.g. (Gani et al., 2016). As a result, comparing the cutting-edge algorithms with the current standard
158 ones can show whether there is any room to improve the gap-filling process within the field.
159 According to the concerns mentioned above, this paper had two objectives. a) To find out the impact
160 of different window lengths on the performance of each algorithm. And b) evaluate the performance
161 of traditional and new gap-filling techniques for the OzFlux Network, separately for fluxes and their
162 meteorological drivers, particularly soil moisture, for this has always been a challenging variable to
163 gap-fill for a couple of reasons, such as of the biology and heterogeneity of soil parameters. To address
164 these objectives, we utilised eight different algorithms (Extreme Gradient Boost (XGB), Random Forest
165 Algorithm (RF), Artificial Neural Networks (ANNs), Classic Linear Regression (CLR), Support Vector
166 Regression (SVR), Elastic net regularisation (ELN), Panel Data (PD) and Prophet Forecast Model
167 (FBP)) to fill the gaps of environmental drivers and the major fluxes. We then assessed their relative
168 performance to evaluate potentially better ways to fill EC flux data. To test the approaches, we used
169 five flux towers from the OzFlux network. To evaluate the performance of these algorithms, nine
170 scenarios for gaps were planned – from a day to a whole year - and applied to the datasets, and
171 different common performance metrics (e.g. RMSE, MBE, etc.), as well as visual graphs were used.

172

173 2. Materials and methods

174

175 To address the first objective of this research, data of nine different window lengths were
176 considered to train and test the algorithms, i.e. 1, 5, 10, 20, 30, 60, 90, 180 and 365 days. To address the
177 second objective, we chose eight different algorithms to fill the gaps, including a wide variety of
178 different approaches, e.g. from a simple algorithm like CLR to cutting-edge ML algorithms, such as
179 XGB. The data used in this paper came from five EC towers of the OzFlux Network, i.e. Alice Springs
180 Mulga, Calperum, Gingin, Howard Springs and Tumbarumba from 2011 to 2013, with a time
181 resolution of 30 minutes. Additionally, data coming from three additional sources outside of the
182 network were also used as ancillary data to help the algorithms fill the gaps of environmental drivers.

183 2.1. Data

184 The data used for this research came from OzFlux, which is the regional Australian and New
185 Zealand flux tower network that aims to provide a continental-scale national research facility to
186 monitor and assess Australia's terrestrial biosphere and climate (Beringer et al., 2016a). As described
187 in (Isaac et al., 2017), all OzFlux towers continuously measure and record 28 environmental features
188 at resolutions up to 10 Hz, and use a 30 min averaging period, with a few exceptions (data are available
189 from (<http://data.ozflux.org.au/portal>)). Besides, the network acquires additional data from the
190 Australian Bureau of Meteorology (BoM), the European Centre for Medium-Range Weather



191 Forecasting (ECMWF), and the Moderate Resolution Imaging Spectroradiometer (MODIS) on the
 192 TERRA and AQUA satellites (Isaac et al., 2017). These additional data, also known as ancillary data,
 193 provide alternative data for gap-filling flux tower datasets (Isaac et al., 2017). As explained in (Isaac
 194 et al., 2017), OzFlux uses the BoM automated weather station (AWS) datasets to gap-fill the
 195 meteorological data, the BoM weather forecasting model (ACCESS-R) for radiation and soil data from
 196 2011 onward, and MODIS MOD13Q1 for Normalised Difference Vegetation Index (NDVI) and
 197 Enhanced Vegetation Index (EVI). Moreover, the data provided by BIOS2, a physically-based model-
 198 data integration environment for tracking Australian carbon and water (Haverd et al., 2015), were also
 199 used as another ancillary source for varieties of environmental features. Current ACCESS-R and
 200 MODIS data are available from the BoM OPeNDAP (<http://www.opendap.org/>) server and TERN-
 201 AusCover data (<http://www.auscover.org.au/>), respectively.

202

203 The datasets were used in this research came from five towers amongst the OzFlux Network
 204 between 2011 and 2013, each representative of a different climate and land cover of Australian
 205 ecological conditions; i.e. Alice Springs Mulga: Tropical and Subtropical Desert, Calperum: steppe,
 206 Gingin: Mediterranean, Howard Springs: Tropical Savanna, Tumbarumba: Oceanic (Table 1)
 207 (Beringer et al. 2016). The datasets included 15 meteorological drivers as well as three major fluxes
 208 recorded (Table 2) based upon EC technique at a 30-minute temporal resolution, except for
 209 Tumbarumba, which was hourly. Additionally, relevant ancillary datasets for the mentioned towers
 210 were used to follow the OzFlux Network gap-filling protocol. Each dataset was quality checked at
 211 three levels based on the OzFlux Network protocol described in (Isaac et al., 2017) and applied using
 212 PyFluxPro ver. 0.9.2. To address the underestimation of canopy respiration by EC measurements at
 213 night, we used the CPD method of (Barr et al., 2013) to reject nightly records when the friction velocity
 214 fell below the threshold value of each site. After dismissing the inappropriate measurements, overall
 215 coverage of 72-88 % and 21-48 % were achieved for diurnal and nocturnal records, respectively.

216

217 *Table 1. The information of the five towers that their data were used, including their name, location, dominant species and*
 218 *climate.*

Site	Location	Species	Climate	Latitude, Longitude (degree)
Alice Springs Mulga [AU-ASM]	Pine Hill cattle station, near Alice Springs, Northern Territory	Semi-arid mulga (Acacia aneura) ecosystem	Tropical and Subtropical Desert Climate (Bwh)	-22.2828° N, 133.2493° E
Calperum [AU-Cpr]	Calperum Station, 25 km NW of Renmark, South Australia	Recovering Mallee woodland	Steppe Climate (Bsk)	-34.0027° N, 140.5877° E
Gingin [AU-Gin]	Swan Coastal Plain 70 km north of Perth, Western Australia	Coastal heath Banksia woodland	Mediterranean Climate (Csa)	-31.3764° N, 115.7139° E
Howard Springs [AU-How]	E of Darwin, NT	Tropical savanna (wet)	Tropical Savanna Climate (Aw)	-12.4943° N, 131.1523° E
Tumbarumba [AU- Tum]	Near Tumbarumba, NSW	Wet temperate sclerophyll eucalypt	Oceanic climate (Cfb)	-35.6566° N, 148.1517° E



219

220 *Table 2. List of variables and their units used in this research, including the three main fluxes and their environmental drivers.*

List of variables	Units
Drivers:	
Ah	Absolute Humidity (g m^{-3})
Fa	Available energy (W m^{-2})
Fg	Ground heat flux (W m^{-2})
Fld	Downwelling long-wave radiation (W m^{-2})
Flu	Upwelling long-wave radiation (W m^{-2})
Fn	Net radiation (W m^{-2})
Fsd	Downwelling short-wave radiation (W m^{-2})
Fsu	Upwelling short-wave radiation (W m^{-2})
ps	Surface pressure (kPa)
Sws	Soil water content (m m^{-1})
Ta	Air temperature (C)
Ts	Soil temperature (C)
Ws	Wind speed (m s^{-1})
Wd	Wind direction (deg)
Precip	Precipitation (mm)
Fluxes:	
Fc	CO_2 flux ($\mu\text{mol m}^{-2} \text{s}^{-1}$)
Fh	Sensible heat flux (W m^{-2})
Fe	Latent heat flux (W m^{-2})

221

222 The datasets whereby each environmental variable was gap-filled are shown in Table 3. For each of
 223 these variables, the same variable of the ancillary source was used to fill the gaps. For instance, to gap-
 224 fill Ah, the Ah records of AWS, ACCESS-R and BIOS2 were used. To gap-fill the missing values of
 225 fluxes, i.e. Fc, Fh and Fe, eight drivers were used as follows: Ta, Ws, Sws, Fg, VPD, Fn, q and Ts based
 226 on trial and error. Different libraries of Python Programming Language (ver. 3.6.4) were utilised for
 227 training and testing the algorithms, i.e. xgboost for XGB, fbprophet for FBP, statsmodels for PD and
 228 sklearn for the rest of algorithms. Each algorithm was tuned up individually using grid search, and
 229 the number of nodes, layers, irritations, etc. were chosen therefor.

230

231

232 *Table 3. The ancillary sources whereby each environmental driver was gap-filled.*

List of variables (γ)	Ancillary Source
Drivers:	
Ah	AWS, ACCESS-R, BIOS2
Fa	ACCESS-R, BIOS2
Fg	ACCESS-R, BIOS2
Fld	ACCESS-R, BIOS2
Flu	ACCESS-R, BIOS2
Fn	ACCESS-R, BIOS2
Fsd	ACCESS-R, BIOS2
Fsu	ACCESS-R, BIOS2
ps	AWS, ACCESS-R
Sws	ACCESS-R, BIOS2



Ta	AWS, ACCESS-R, BIOS2
Ts	ACCESS-R, BIOS2
Ws	AWS, ACCESS-R
Wd	AWS, ACCESS-R
Precip	AWS, ACCESS-R, BIOS2

233

234

235 *2.2. Gap-filling algorithms*

236

237

238

239

240

241

242

243

244

Artificial Neural Networks (ANN)

245

246

247

248

249

250

251

252

253

254

255

Classical Linear Regression (CLR)

256

257

258

259

260

A classical linear regression is an equation developed to estimate the value of the dependent variable (y) based on independent values (x_i). In contrast, each x_i has its specific coefficient and an overall intercept value. In this method, these coefficients are determined by minimising the squared residuals (errors) of estimated vs observed values, called least squares. A CLR algorithm can be formulated as follows (Freedman, 2009):

$$y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_i X_i + \varepsilon \quad (1)$$

261

262

263

where y is the dependent variable, α is the interception, X_s are independent variables, and β_i is coefficient of X_i , and ε is the error term. We chose this algorithm as a baseline to find out how better more complicated algorithms can estimate dependent variables comparatively.



264 **Random Forests (RF)**

265 Random forest, a supervised ML algorithm, used for both classification and regression,
266 consists of multiple trees constructed systematically by pseudorandomly selecting subsets of
267 components of the feature vector, that is, trees constructed in randomly chosen subspaces (Ho, 1998).
268 RF algorithm has been developed to control the overcome over-fitting problem, a commonplace
269 limitation of its preceding decision tree-based methods (Ho, 1995, 1998).
270 Sklearn.ensemble.RandomForestRegressor was used to apply this method in Python, and the
271 hyperparameters used were 5 and 1000 for “max_depth” and “n_estimators”, respectively based on
272 grid search.

273

274 **Support Vector Regression (SVR)**

275 As a non-linear method, support vector regression was developed based on Vanpik’s concept
276 of support vectors theory (Drucker et al., 1997). An SVR algorithm is trained by trying to solve the
277 following problem:

278

279 minimise $\frac{1}{2} \|w\|^2$

280 subject to $\begin{cases} y_i - \langle w, x_i \rangle - b \leq \varepsilon, \\ \langle w, x_i \rangle + b - y_i \leq \varepsilon, \end{cases}$

281 where x_i and y_i are training sample and target value in a row. The inner product plus intercept
282 $\langle w, x_i \rangle + b$ is the prediction for that sample, and ε is a free parameter that serves as a threshold.
283 sklearn.svm.SVR was used to apply this method in Python, and the hyperparameters that used were
284 1 and 0.001 for “C” and “gamma”, respectively based on grid search.

285 **Elastic net regularisation (ELN)**

286 The elastic net is a linear regularised regression method that exerts small amounts of bias by
287 adding two penalty components to the regressed line to decline the coefficients of independent
288 variables and thus, provides better long-term predictions. Given that these two penalty components
289 come from ridge regression and LASSO, the elastic net is considered as a hybrid model consists of
290 ridge and LASSO regressions, overcoming the limitations of both. The estimates from the ELN method
291 can be formulated as below (Zou and Hastie, 2005):

$$\hat{\beta}(\text{elastic net}) = \frac{(|\hat{\beta}(OLS)| - \lambda_1/2)}{1 + \lambda_2} \text{sgn}\{\hat{\beta}(OLS)\} \quad (2)$$

292

293 where $\hat{\beta}$ is the coefficient of each ELN independent variable, λ_1 and λ_2 are penalty coefficients of
294 LASSO and ridge regression respectively, $\hat{\beta}(OLS)$ is the coefficient of an independent variable
295 calculated based on ordinary least squares, and sgn stands for the sign function:



$$\text{sgn}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases} \quad (3)$$

296

297 The ELN regression is good at addressing situations when the training datasets have small samples
298 or when there are correlations between parameters. `sklearn.linear_model.ElasticNet` was used to
299 apply this method in Python, and the hyperparameters used were as follows: {'alpha': 0.01,
300 'fit_intercept': True, 'max_iter': 5000, 'normalize': False} based on grid search.

301

302 Panel data (PD)

303 Panel data is a multidimensional statistical method, mainly used in econometrics to analyse
304 datasets, which involve time series of observations amongst individual cross-sections (Baltagi, 1995)
305 usually based on ordinary least squares (OLS) or generalised least squares (GLS). A two-way panel
306 data model consists of two extra components above a CLR as follows (Baltagi, 1995; Hsiao et al., 2002;
307 Wooldridge, 2008):

$$y_{it} = \alpha + \beta X_{it} + u_{it} \quad i = 1, 2, \dots, N; \quad t = 1, 2, \dots, T \quad (4)$$

$$y_{it} = \alpha + \beta X_{it} + \mu_i + \lambda_t \quad (5)$$

308 where i and t denote the cross-section and time series dimension in a row, y is a dependent-variable
309 vector, X is an independent variable matrix, α is a scalar, β is the coefficient of the independent-
310 variable matrix, μ_i is the unobservable individual-specific effect, and λ_t is the unobservable time-
311 specific effect. Panel data abilities to provide a holistic analysis of different individuals, as well as
312 determining the specific impact of every single time caused its superiority over CLR.

313 Extreme Gradient Boost (XGB)

314 Extreme gradient boost is a reinforced method of Gradient Boost introduced in 1999 that works based
315 on parallel boosted decision trees and similar to RF can be used for a variety of data processing
316 purposes including classification and regression (Friedman, 2002; Jerome H. Friedman, 2001; Ye et al.,
317 2009). XGB method is resistive to over-fitting and provides a robust, portable and scalable algorithm
318 for large-scale boosting decision-trees-based techniques.
319 `sklearn.ensemble.GradientBoostingRegressor` was used to apply this method in Python, and its
320 hyperparameters were chosen based on grid search as follows: {'learning_rate': 0.001, 'max_depth': 8,
321 'reg_alpha': 0.1, 'subsample': 0.5}.

322

323 The Prophet Forecasting Model (FBP)

324 The Prophet Forecasting Model, also known as "prophet", is a time series forecasting model
325 developed by Facebook to manage the common features of business time series and designed to have



326 intuitive parameters that can be adjusted without knowing the details of underlying model (Taylor
327 and Letham, 2017). A decomposable time series model was used (Harvey and Peters, 1990) to develop
328 this model, with three main components: trend, seasonality, and holidays as the equation below
329 (Taylor and Letham, 2018):

$$y(t) = g(t) + s(t) + h(t) \quad (6)$$

330

331 where $g(t)$ is the trend function, which models non-periodic changes, $s(t)$ is a function to represent
332 periodic changes, e.g. seasonality, and $h(t)$ assesses the effects of potential anomalies which occur over
333 one or more days, e.g. holidays.

334

335 *2.3. The gap scenarios*

336 To find out the effect of gap size on the performance of our gap-filling algorithms, we trained
337 each of them using nine different window lengths (i.e. 1, 5, 10, 20, 30, 60, 90, 180 and 365 days). The
338 gap size for each trained algorithm was chosen as the same size of the corresponding training window,
339 e.g. the gap size for a 20-day training window was 20 days and so on. As such, in every scenario, the
340 entire data of 2013 were used step by step to test the performance of the algorithms as follows: at the
341 first step of each scenario, the gap began from 1 Jan 2013, while its corresponding training window
342 was the same size but came from the preceding period. For instance, for a 30-day gap, the first step
343 included training an algorithm based on the data of Dec 2012 and the testing period of the first month
344 of 2013. In the second step, the data of the first month of 2013 were used for training, while the data
345 of the second month of 2013 was considered as a gap, and this went to the end of 2013 consecutively.
346 As such, for the last step, the training window was the second last 30 days of 2013, and its
347 corresponding gap was the last 30 days of 2013. The only exception of the mentioned training strategy
348 was FBP as it needed a training dataset with at least a year to be developed. Therefore, here, the
349 training data for each gap was all data prior to that gap since the beginning of 2011. Overall, 18
350 variables, nine window lengths and eight gap-filling methods across five flux towers resulted in 6480
351 computations.

352 *2.4. Statistical performance measures*

353 Different statistical metrics were used to evaluate the performance of algorithms and enable
354 comparison between measured values from the flux towers with each gap-filling algorithm prediction.
355 These metrics included the coefficient of determination (R-squared) to measure the square of the
356 coefficient of multiple correlations (Devore, 1991), the variance of measured and modelled values (S^2)
357 to indicate how well algorithms could follow the variations of the recorded data, the root mean square
358 error (RMSE), the mean bias error (MBE) to capture distribution and bias of residuals, variance ratio
359 (VR) to compare the variance of estimated values with those of measured, and the Index of Agreement
360 to compare the sum of the squared error to the potential error (Bennett et al., 2013). Abbreviations and
361 formulas of these metrics are illustrated as follows (Bennett et al., 2013):



$$R^2 = \frac{[\sum(p_i - \bar{p})(o_i - \bar{o})]^2}{\sum(p_i - \bar{p})^2 \sum(o_i - \bar{o})^2} \quad (7)$$

362

$$S^2 = \frac{\sum(x_i - \bar{x})^2}{N - 1} \quad (8)$$

363

$$RMSE = \sqrt{\frac{\sum(p_i - o_i)^2}{N - 1}} \quad (9)$$

364

365

$$MBE = \frac{\sum o_i - p_i}{N - 1} \quad (10)$$

366

$$VR = \frac{\sigma_p^2}{\sigma_o^2} \quad (11)$$

367

$$IoAd = 1 - \frac{\sum_{i=1}^n (o_i - p_i)^2}{\sum_{i=1}^n (|p_i - \bar{p}| + |o_i - \bar{o}|)^2} \quad (12)$$

368

369 where o_i and p_i are individual measured and predicted values respectively, \bar{o} and \bar{p} are the means of
370 o and p , and σ^2 is the variance. S^2 is calculated separately for the observed and predicted values with
371 the respective values defined as x that represents every observed or predicted value. All of these
372 metrics were calculated for each of the gap scenarios, and then the results of different windows were
373 concatenated. Afterwards, the yearly metrics were calculated to avoid Simpson's paradox or any
374 relevant averaging issue as described by (Kock and Gaskins, 2016). Moreover, the average of daily
375 and seasonal differences between the estimated and measured values, as well as the associated
376 variance were calculated and plotted.

377 3. Results

378

379 3.1. Fluxes

380 3.1.1 Fc

381 Even though factors such as Fg and Fn are fluxes, we dealt with them as environmental drivers
382 since they drive the three major fluxes. The metrics used to evaluate the performance of the algorithms
383 (RMSE, R^2 , MBE, IoAd and VR) (Table 4) illustrated that overall, the performance of these algorithms,
384 particularly the ML ones, was similar. The algorithms, however, showed different levels of sensitivity
385 to training/testing window length, e.g. the ANNs showed less sensitivity, whereas the FBP showed



386 the most sensitivity (Figure 1). The XGB provided the lowest values of RMSE and one of the highest
387 R^2 , while the FBP and ELN had the lowest and highest values of RMSE and R^2 , respectively.

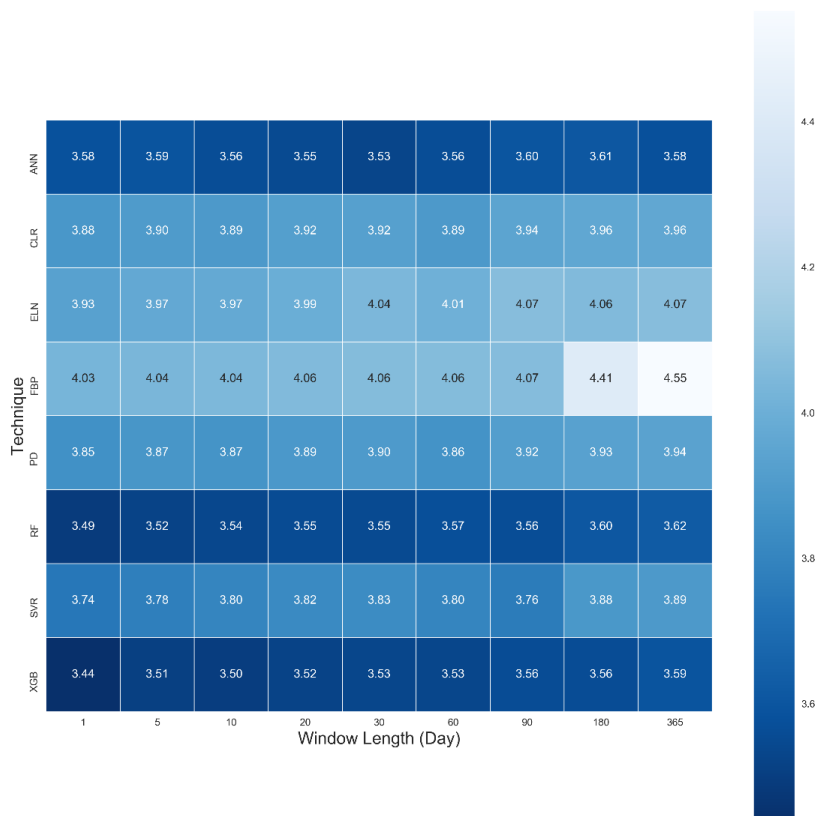
388 *Table 4. The average amounts of performance metrics for each gap-filling algorithm regarding F_c , which includes all window*
389 *lengths and sites, ranked by RMSE using the Tukey's HSD test at the level of 5 per cent.*

Algorithm	Mean RMSE	Mean R^2	Mean MBE	Mean IoAd	Mean VR
XGB	3.53 ^a	0.56	-0.44	0.89	0.59
RF	3.56 ^a	0.54	-0.38	0.90	0.70
ANNs	3.57 ^a	0.52	-0.34	0.89	0.68
SVR	3.81 ^b	0.47	-0.33	0.86	0.79
PD	3.89 ^b	0.45	-0.36	0.80	0.53
CLR	3.92 ^{b,c}	0.46	-0.37	0.80	0.54
ELN	4.01 ^c	0.40	-0.38	0.72	0.37
FBP	4.15 ^d	0.44	-0.06	0.77	0.68

390

391 These outcomes were expected for the XGB as it uses a more regularised model formalisation to
392 control over-fitting (Chen and Guestrin, 2016) that leads to better performance. The relatively poor
393 performance of FBP was also foreseen for unlike other algorithms, FBP did not use any feature to
394 estimate flux values, other than the previous time series of flux values. However, the weaker
395 performance of the ELN compared to CLR was unforeseen due to by adding two penalty components
396 to the regressed line, and the ELN is supposed to improve the long term prediction compared to the
397 traditional linear regression methods. Tukey's HSD (honestly significant difference) test at the level
398 of five per cent was applied to the results to find out whether the difference amongst the algorithms
399 was significant (Table 4). Where the null hypothesis was there is no significant difference between the
400 mean values of the RMSE. According to the results, there were significant differences between certain
401 algorithms, and the XGB, RF and ANNs were different from the rest, showing that these three
402 performed considerably better. Tukey's HSD test, however, did not reject the second error probability
403 between RF, XGB and ANNs meaning that the three algorithms were not significantly different from
404 each other. This result agrees with the results of (Falge et al., 2001) and (Moffat et al., 2007) in the sense
405 that ANNs are one of the best available algorithms, and there is no significant difference amongst the
406 appropriate algorithms. Nonetheless, it is worth mentioning that Tukey's HSD is well known as a
407 conservative test. That being said, despite no meaningful difference based on Tukey's HSD, XGB and
408 RF might have performed better than ANN, as the superiority of RF in gap-filling of methane flux has
409 recently been claimed by (Kim et al., 2020).

410



411

412 *Figure 1. A heat map of mean RMSE values of Fc across all sites based on 8 algorithms and 9 window lengths in 2013.*

413

414 To address the first objectives of this paper, finding out the sensitivity of each algorithm to the
 415 training and testing window length, we used the averaged RMSE, R^2 and MBE for each window length
 416 and gap size, using the output of all algorithms for all sites (Table 5). The outcome illustrates that the
 417 longer the window length got, the bigger the amounts of RMSE became. Yet, no such pattern was
 418 recognisable for the R^2 and MBE, particularly for the window lengths equal to or shorter than 60 days.
 419 As a result, based on our scenarios (using the same length for training window and gap size), choosing
 420 any training windows longer than 60 days, i.e. 90, 180 and 365 days, made the performance of the



421 algorithms worst. The phenomenon can be justified by the idea that longer windows do not let the
 422 algorithms to accommodate seasonal changes and therefore, different physiological behaviour of
 423 the canopy.

424 *Table 5. The average amounts of RMSE, R², and MBE for Fc gap-filling based on the window length including the outcome of all*
 425 *sites; the differences of RMSE values were tested using the Tukey's HSD test at the level of 5 per cent.*

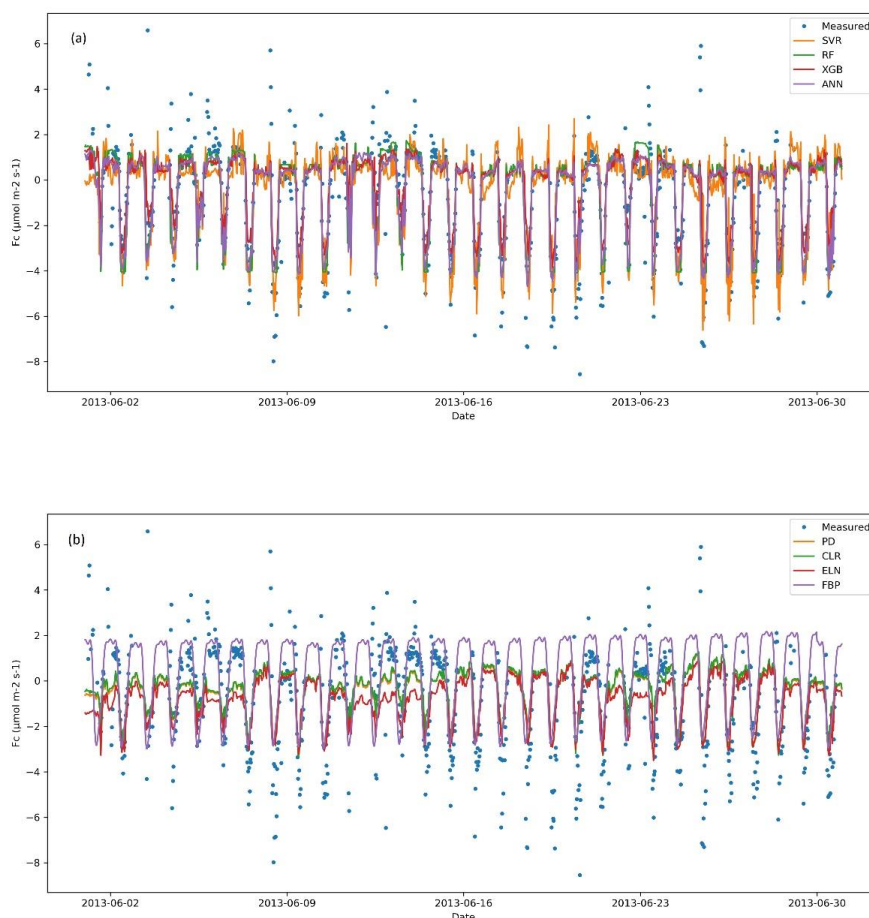
Window length	Mean RMSE	Mean R ²	Mean MBE
1-day	3.74 ^a	0.49	-0.26
5-days	3.77 ^a	0.49	-0.31
10-days	3.77 ^a	0.49	-0.29
20-days	3.79 ^a	0.48	-0.31
30-days	3.80 ^a	0.48	-0.31
60-days	3.79 ^a	0.48	-0.35
90-days	3.81 ^a	0.48	-0.39
180-days	3.88 ^a	0.47	-0.41
365-days	3.90 ^a	0.46	-0.37

426
 427 Besides, the metrics of the top three algorithms, XGB, RF and ANNs, did not show any sizeable
 428 difference for the window lengths shorter than 60 days. As such, finding the ideal window length, at
 429 least at this stage, is not distinctly noticeable and should rely on the local knowledge of each specific
 430 site. Nevertheless, as mentioned earlier, the ideal window generally cannot be longer than 60 days,
 431 unless for a monotonic ecosystem without a dramatic change during the year. According to the MBE
 432 values, mainly, all algorithms had negative amounts of MBE, showing overestimation of the Fc values.
 433 This bias varied from tower to tower and depended on the window lengths. For instance, absolute
 434 amounts of the MBE were bigger in Gingin and Tumberumba, while considerably smaller (closer to
 435 zero) at AliceSprings Mulga and Calperum (results not shown). The lower leaf area index of the two
 436 later sites, and thus their smaller amounts of photosynthesis is likely to be the reason that justifies the
 437 outcome. FBP, nonetheless, provided substantially lower mean bias (-0.06) compared to the other
 438 algorithms, which varied between -0.33 and -0.44.

439 Observations from the EC technique often include extremely low or high values, especially at
 440 night, when some of the theoretical assumptions might be violated. The nature of the EC technique
 441 associated with its practical challenges, often makes it difficult to distinguish between the good data
 442 and the noise (Aubinet et al., 2012; Burba and Anderson, 2010). This problem seems to affect the
 443 outcomes of the gap-filling algorithms in this research, as none of them performed ideally in capturing
 444 the observed variance (Figure 2). Even though RMSE, R² and IoAd showed the superiority of the XGB,
 445 RF and ANNs, the variance ratio between the estimated and measured values revealed different
 446 information (Table 4), which is also recognisable in Figure 2. The variance ratios (VR) showed that
 447 SVR captured the extreme values of Fc better than the other algorithms, 0.79 on average. The XGB, on



448 the other hand, provided smaller VR (0.59) compared with those of the RF (0.70) and ANNs (0.68),
449 especially for the window lengths longer than 10 days (not shown).



450

451 *Figure 2. Measured vs estimated values of Fc for Calperum based on the 30-day window during June (Austral winter) 2013*

452 This substantial smaller VR calls into question the ability of XGB to provide a solid gap-filling for long
453 gaps. The linear algorithms, CLR, PD, and ELN, performed worse with respect to the VR compared
454 to the ML algorithms. The estimated versus measured values of Fc for Calperum during June 2013
455 (Figure 2) confirms the information achieved by the VR. Based on the figure, and VR of 0.79, the SVR
456 captured the extreme values of Fc the best, whereas the ELN, as expected, performed the worst (0.37).
457 Although the XGB (VR of 0.59) provided relatively well while estimating the maximum values
458 (respiration), it was not capable of assessing the minimum values, thereby provided a constant
459 overestimation of NEE during the day. The RF (VR of 0.70), in contrast, captured both negative and
460 positive extremes better than the XGB, while the performance of the ANNs (VR of 0.68) was



461 somewhere in between. The rest of the algorithms performed poorly, particularly during the night,
462 except the FBP. It is noteworthy that CLR, PD, and ELN frequently predicted nocturnal
463 photosynthesis.

464 Apart from the objectives of this paper, tracing the performance of gap-filling algorithms
465 based on the hourly time step and seasonality has been as of the research interests. Thus, as an aside,
466 the differences between the average of estimations and measured values, as well as the difference
467 between the variances for the top three algorithms (XGB, RF and ANNs) were calculated for the 24-h
468 and seasonal ranges. These algorithms showed different anomalies in different towers and hours of
469 the day, except for Tumbarumba, where the patterns of anomalies were almost similar (Figure 3). The
470 average variance of differences was slightly lower during the night while the largest values of
471 anomalies usually occurred around noon, as expected due to the bigger variations of carbon uptake
472 caused by photosynthesis. Here, the RF and XGB performed better than ANNs, with the curves closer
473 to the basis, except for Tumbarumba. According to the seasonal anomalies, however, the algorithms
474 showed more similarities and closer outcomes, particularly for Tumbarumba, where all three
475 overestimated F_c values for the whole year (Figure 4). Similar to the 24-hour scale, the anomaly values
476 varied from site to site based on the season and the algorithm. Although the performance of the
477 algorithms was less here as against the daily scale, it seems that the XGB and RF still show superiority
478 over the ANNs. Apart from Tumbarumba, XGB, RF and ANNs showed a significant bias during
479 spring (July, August and September) in Howard Springs, when the site receives lower precipitation
480 due to the dry season.

481

482

483

484

485

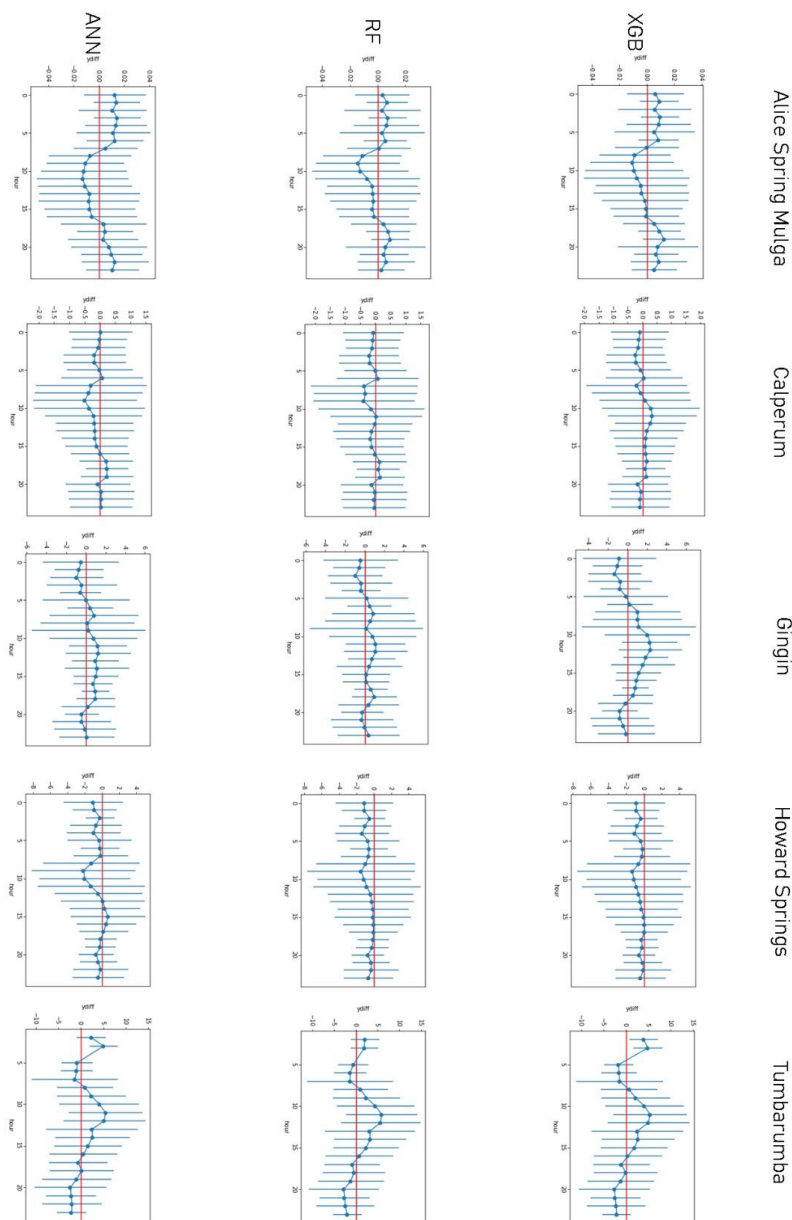
486

487

488

489

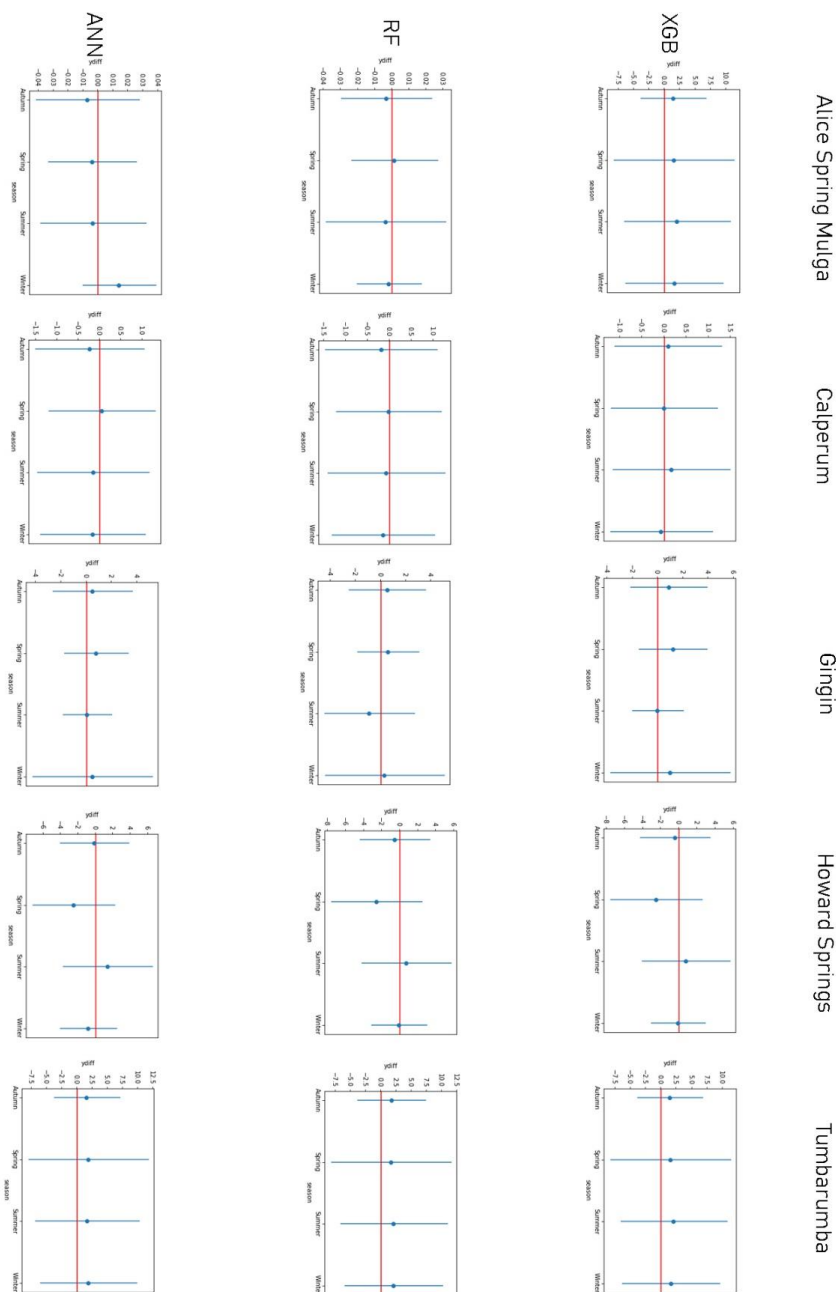
490



491

492 *Figure 3. 24-h Anomalies of XGB, RF and ANNs based on the Fc average of differences, and associated variances between the*
493 *estimated and measured values for all towers during 2013.*

494



495

496 *Figure 4. Seasonal anomalies of XGB, RF and ANNs based on the Fc averages of differences, and the associated variances*
 497 *between the estimated and measured values for all towers during 2013 (Jan, Feb and Mar as Summer, Apr, May and Jun as*
 498 *Autumn, Jul, Aug and Sep as winter).*



499

3.1.2 Fe

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

The performance of algorithms for Fe was similar to that for Fc regarding RMSE, MBE and R², as shown in *Table 6*. This similarity was not surprising since these processes are partially coupled via stomatal conductance (Scanlon and Kustas, 2010; Scanlon and Sahu, 2008). Again, the top three ML algorithms performed better, with a significant superiority of the XGB and RF, as shown by the Tukey's HSD (*Table 6*). Besides, the null hypothesis was not rejected while comparing FBP and SVR, whereas the better performance of the other algorithms was confirmed. As a result, the FBP and SVR provided the most unsatisfactory results in estimating Fe, according to the average values of the RMSE. No significant improvement in RMSE occurred when the window lengths of training and testing became shorter than 90 days, meaning that the performance of the algorithms did not vary considerably from a 60-day to a one-day window. The results of CLR and PD were very similar to those for Fc, showed lower RMSE and higher R² values as against ELN, but the ELN led to slight lower MBE. The MBE values also showed moderately high values for the SVR, meaning that there was an absolute bias in its outcome, which might be related to overfitting. The source of the bias shown by the SVR algorithm (*Figure 5*), was because it could not capture the minimum values appropriately, resulting in a considerable overestimation. A common issue in estimating Fe values, which had affected all algorithms other than the FBP, was not assessing the negative values. In contrast to Fc results, the ANNs did not perform as solid as the XGB and RF, which could be due to not being able to capture the maximum values as satisfying as its rivals were.

518

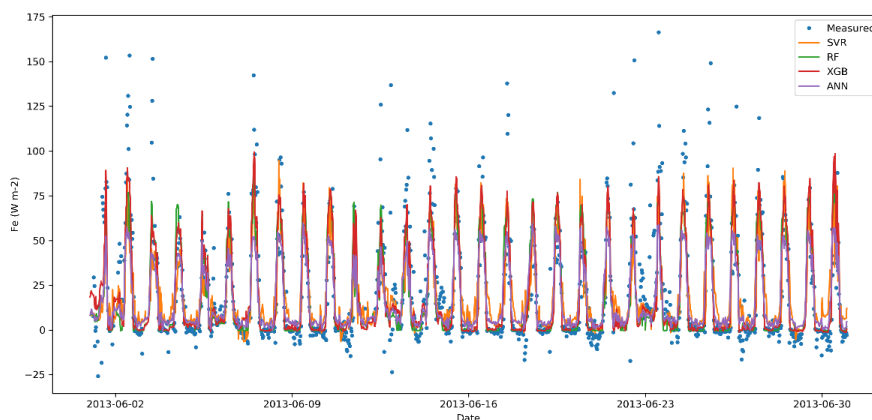
519

Table 6. The average of metrics for Fe gap-filling based on the algorithms, ranked by RMSE using the Tukey's HSD test at the level of 5 per cent.

Algorithm (Fe)	Mean RMSE	Mean R ²	Mean MBE
XGB	37.27 ^a	0.69	-3.19
RF	37.98 ^a	0.68	-3.00
ANNs	40.62 ^b	0.61	-3.48
PD	42.45 ^{b,c}	0.58	-5.50
CLR	42.67 ^{b,c}	0.58	-5.95
ELn	43.48 ^c	0.53	-5.00
SVR	48.42 ^d	0.53	-21.08
FBP	49.46 ^d	0.44	2.03

520

521



522

523 *Figure 5. Measured vs estimated values of Fe for Calperum based on a 60-day window during June 2013*

524

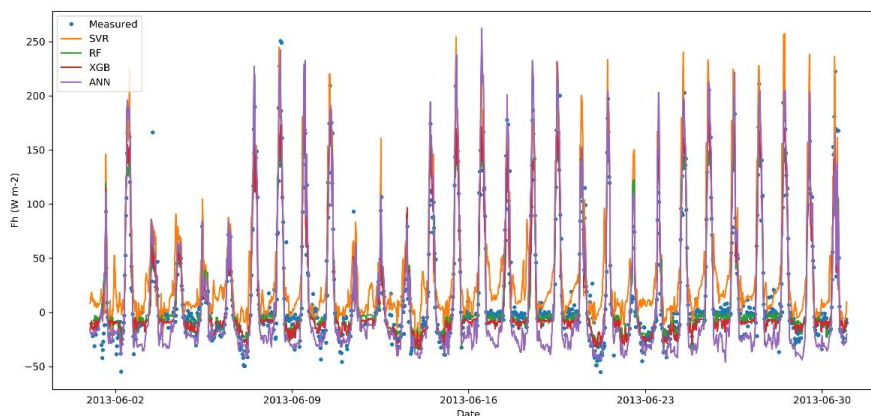
3.1.3 Fh

525

526 As with the other flux results, the metrics (RMSE, R^2 and MBE) showed slight superiority of
 527 XGB and RF, as well as the inferiority of the SVR and FBP over the other algorithms (Table 7).
 528 Likewise, the SVR provided relatively large negative values of MBE, showing considerable
 529 overestimation. The Tukey's HSD test of the average RMSE values confirmed that the performance of
 530 the FBP was significantly different from the rest at the level of 5 per cent, making FBP the weakest
 531 performer for Fh. On the other hand, even though there was no significant difference amongst the
 532 XGB, ANNs and RF, the first one was considerably superior over the other five algorithms as regards
 533 the Tukey's HSD test. Like Fe, estimated values of Fh using SVR had a negative bias (Figure 6) because
 534 it was not able to provide appropriate estimations of Fh negative values. In contrast, the ANNs
 535 performed the best in capturing the minimum values, while the XGB and RF performed relatively
 536 well, close to each other. Despite this superiority in assessing the minimum amounts, ANNs did not
 537 carry out as solid as XGB and RF concerning the overall values, resulted in higher RMSE. Finally,
 538 similar to the other fluxes, the PD performed slightly better than the CLR and ELN.

539 *Table 7. The average metrics for Fh gap-filling based on the algorithms, ranked by RMSE using the Tukey's HSD test at the level*
 540 *of 5 per cent.*

Algorithm (Fh)	Mean RMSE	Mean R^2	Mean MBE
XGB	37.26 ^a	0.93	-1.00
RF	38.08 ^{ab}	0.93	-1.35
ANNs	40.48 ^{ab,c}	0.92	-0.41
PD	41.83 ^{bc}	0.92	-0.27
CLR	42.14 ^{bc}	0.92	-0.05
Eln	42.28 ^{bc}	0.92	0.04
SVR	43.98 ^c	0.91	-8.28
FBP	67.19 ^d	0.74	1.25



541

542 *Figure 6. Measured vs estimated values of Fh for Calperum based on a 60-day window during June 2013*

543

544 *3.2. Meteorological and Environmental Drivers*

545 Since meteorological and environmental drivers are needed to fill the gaps of the three
546 substantial fluxes, F_c , F_e and F_h , the eight algorithms were used to fill the gaps of these drivers. The
547 metrics of R^2 , RMSE, and MBE were calculated for all five towers and nine window lengths (16
548 meteorological and environmental drivers and three fluxes). Overall, for most meteorological drivers,
549 the linear algorithms, especially the CLR and PD, performed slightly better than the ML algorithms
550 such as the XGB, RF, ANNs and SVR, except for A_h , F_g and F_n . This unexpected superiority can be
551 explained based on the two following reasons. Firstly, unlike the fluxes, the input and output features
552 were the same here, e.g. T_a for T_a , which led to strong correlations (e.g. up to 0.99 for atmospheric
553 pressure - p_s) as well as strong linear relationships between the independent and dependent features.
554 These strong correlations helped the linear algorithms to perform well, while nullified the ability of
555 ML algorithms to capture non-linear behaviour of complicated problems. Second, the slight inferiority
556 of ML algorithms could be due to data noise where simple linear algorithms such as the CLR are
557 usually less sensitive to the noise relatively. Therefore, over-fitting is not an issue for them when the
558 number of observations is big enough (i.e. at least 10 to 20 observations per parameter (Harrell, 2014)).
559 The exceptions were A_h , F_n and F_g , for which values were estimated more accurately by the XGB,
560 ANNs and RF, especially the latest one (the RMSE of 30.23 versus 35.24 provided by the RF and CLR
561 for F_g , respectively). Tukey's HSD test for the mean RMSE values of F_g confirmed that The XGB,
562 ANNs and RF provided better results at the level of 5 per cent, while, like all other fluxes and drivers,
563 the FBP confirmed to be the worst algorithm (Table 8). Yet, according to the same test for the other
564 drivers, there was not any significant difference between the algorithms, other than the FBP, which
565 provided the most significant mean values of the RMSE (results not shown). Importantly, though,
566 none of the algorithms offered adequate estimations for soil moisture (Sws), particularly in drier
567 regions. This weak performance happened because Sws changes dramatically during rainfall in a



568 pulsed manner often from zero to saturation in short space of time, whereas, the algorithms had been
569 trained based on the datasets mostly reflecting non-rainy periods. These datasets, consequently, could
570 not fit the algorithms in a way that they could estimate Sws accurately when precipitation occurs and
571 the soil moisture increases dramatically. For instance, in a wet region like Tumarumba, where the
572 soil faces rainy days frequently, the time series are much less spikey. Thus, the overall performance
573 was better in these regions compared with the drier ones, e.g. R^2 of 0.43 and 0.25 on average for
574 Tumarumba and Calperum, respectively. Besides, the dataset used to gap-fill the soil moisture was
575 a model derivation from gridded data or regional reanalysis and therefore, can be not close to reality.
576 Another challenge of estimating soil moisture comes from the low spatial coherence of soil moisture
577 is that it can be extremely different just a couple of hundred metres away, due to storms, topography,
578 soil structure heterogeneity, etc. (Reichle et al., 2004; Sahoo et al., 2008).

579

580 *Table 8. The average amounts of RMSE for Fg gap-filling based on the algorithms, using the Tukey's HSD test at the level of 5*
581 *per cent.*

Algorithm (Fg)	Mean RMSE
RF ^a	30.17
XGB ^{a,b}	30.70
ANNs ^{b,c}	30.86
SVR ^c	32.77
CLR ^d	34.93
PD ^d	34.94
ELN ^d	34.94
FBP ^e	39.10

582

583 4. Discussion

584 All algorithms performed similarly in estimating the meteorological and environmental drivers
585 (turbulent fluxes included) across all stations, except the FBP, which performed poorly for it did not
586 use any ancillary data. The best results were achieved using training/gap windows of 60 days or
587 shorter, while the worst results obtained for the most extended window, 365 days. Although most of
588 the algorithms performed almost equally well in estimating of meteorological and environmental
589 drivers, the linear algorithms, the CLR, ELN and PD, performed slightly better (not significant using
590 a Tukey's HSD test, though). The only clear exception was Fg, which the RF provided more accurate
591 and robust estimations. The ML algorithms, on the other hand, showed their superiority over the
592 linear algorithms while estimating the main fluxes, Fc, Fe and Fh. For Fc, the XGB, RF and ANNs
593 performed significantly better than the SVR, FBP and all linear algorithms, i.e. the CLR, PD and ELN.
594 The superiority of the ML algorithms, as well as their close performance, agreed with the results of
595 previous researches, e.g. (Falge et al., 2001; Moffat et al., 2007), that showed the superiority of non-
596 linear algorithms and no significant difference amongst the top algorithms in estimating Fc. Besides,
597 the slight superiorities of RF over ANNs, mainly unnoticeable by a conservative test like Tukey's HSD,
598 confirms RF performs better regarding the EC flux gap-filling (Kim et al., 2020).



599 The XGB was the most novel ML algorithm used in this research and based on the most
600 performance metrics provided comparatively robust results in estimating the fluxes. However, the
601 XGB failed to capture the minimum values of F_c as against the SVR and RF, and thus, provided biased
602 results, while assessing the maximum values of F_c well. In estimating the meteorological drivers, the
603 XGB did not show any superiority over the other algorithms, especially the linear ones. Hence, we do
604 not recommend the XGB as an alternative to the current alternative algorithms, especially for long
605 gaps. Nevertheless, because of its local superiorities, this algorithm is suitable to use in an ensemble
606 model alongside the algorithms with different weakness points.

607 The RF was the best all-around algorithm amongst the eight algorithms used in this study,
608 providing the best estimates of the fluxes (similar to XGB) but also captured both minimum and
609 maximum values of F_c . Unlike the RF, all other algorithms generally struggled with estimating either
610 minimum or maximum values of major fluxes, comparatively. It also provided the best results for F_g ,
611 where the linear algorithms did not perform well. Another advantage of the RF over the XGB was that
612 it required less time (approximately a quarter) for training, which was a challenge while using the
613 XGB.

614 The ANNs estimated the fluxes better than the linear algorithms, notably for F_c , yet not as
615 robust as the XGB and RF in general. For F_c and F_h , the ANNs provided bias, mainly due to
616 overestimation of minimum values when the window lengths were 90 days or longer. However, since
617 the superiority of the XGB and RF was not considerable, it is difficult at this point to suggest using
618 XGB or RF as better alternatives. That is because ANNs have been checking out for a long time in
619 different locations and considered as one of the most reliable algorithms in the field for more than a
620 decade (Aubinet et al., 2012; Hagen et al., 2006; Kunwor et al., 2017; Moffat et al., 2007). Furthermore,
621 there are a wide variety of different ANNs algorithms used in the field (Beringer et al., 2016b; Hagen
622 et al., 2006; Isaac et al., 2017; Kunwor et al., 2017; Moffat et al., 2007), and this minor superiority of RF
623 and XGB cannot be generalised without enough additional proves. As such, we suggest other
624 researches to use the RF, especially regarding F_h and F_c alongside the ANNs to find out which one
625 performs better in the challenging scenarios, e.g. when the gaps are long. Another option is to develop
626 ensemble models using since, according to the literature, there is no room to improve the results
627 substantially based on a single algorithm (Moffat et al., 2007).

628 On the other hand, a model with a higher level of flexibility is required in the field (Hagen et
629 al., 2006; Kunwor et al., 2017; Richardson and Hollinger, 2007). Finally, according to the environmental
630 drivers, The ANNs, like the other ML algorithms, could not show a consistent superiority over the
631 linear algorithms. Therefore, we do not recommend using ML algorithms in such scenarios, except for
632 F_g , for which RF seems to be a better option.

633 The SVR showed consistent inferiority over the other ML algorithms and did not fulfilled our
634 expectations, neither for the meteorological drivers nor for the major fluxes. The only strength of the
635 SVR was that it captured the extreme values better than any other algorithm, as revealed in the plots
636 of F_c . However, according to its larger RMSE amounts, the mentioned advantage seems to be achieved
637 suspiciously and might have occurred due to over-fitting. This dubious performance shows the SVR
638 is more vulnerable to the over-fitting issues regarding these types of data. Hence, we suggest the SVR



639 not to be used in any kind of environmental modelling related to the reviewed drivers and fluxes,
640 whatsoever.

641 The CLR, the simplest algorithm used in this research, provided a comparatively acceptable
642 performance in estimating the meteorological drivers, except for Fg. This algorithm, however, could
643 not perform well in assessing the fluxes, especially Fc, mainly because of its inability to capture the
644 extreme values caused by the non-linear nature of Fc. Overall, considering the CLR simplicity,
645 resource-saving and robust performance for drivers, this algorithm seems to be the most suitable way
646 to fill the gaps of meteorological parameters in similar scenarios, where the same ancillary dataset is
647 available.

648 The PD performed slightly better than the CLR, yet it could not fulfil the expectations to show
649 a significant superiority over the other linear algorithms used in the research. This unforeseen weak
650 performance can be explained due to a couple of reasons. First, one of the assumptions of using the
651 PD is that the behaviour of the cross-sections, here towers, is similarly under the similar conditions
652 (the independent variables), and the only thing leads to the difference is the specific characteristics of
653 each individual cross-section. Contrariwise, it seems that the five towers selected in this research
654 violated this assumption due to their absolute different ecosystems. Based on the previous studies in
655 which the PD performed satisfying (Izady et al., 2013, 2016; Mahabbati et al., 2017), (Izady et al., 2016)
656 and (Mahabbati et al., 2017), it appears that a decent level of homogeneity is vital for the PD to perform
657 satisfactorily. As in all previous cases, the ecosystem of the cross-sections had significant similarities,
658 and the distance between them were tens to hundreds of kilometres, not thousands. Therefore, the
659 characteristics of cross-sections, such as radiation, climate, rainfall, etc. had considerable more
660 similarity and homogeneity compared with the towers used in this research. Finally, it is worth
661 mentioning that PD has been commonly used to analyse the time series with a time resolution of
662 weekly or longer, with some exceptional daily-scale cases. In this research, the resolution of data was
663 half-hourly instead, which dramatically increased the computational demands of the algorithm, led
664 to days of processing for a single run. This demand happened because the algorithm creates a dummy
665 variable for each time step and the relevant matrix of variables becomes too large to compute by a
666 regular PC. Considering the expenses of this algorithm, we recommend other researches not to use
667 PD when the time resolution is shorter than daily. Despite the limitation, we still encourage further
668 using of PD whenever there is a decent level of homogeneity amongst the cross-sections and the time
669 resolution is daily or longer (ideally weekly or monthly).

670 The ELN, as a hybrid linear model, did not show any superiority over the CLR, despite its
671 modifications to provide more accurate estimations. Even though ELN performed well in estimating
672 the drivers with slight supremacy in some occasions, e.g. Fld, the CLR is a more proper algorithm to
673 choose for gap-filling the drivers due to its simplicity and less calculation requirement.

674 The FBP was a unique algorithm used in this research, as it did not use any independent
675 variables to estimate the values of drivers and fluxes. The FBP performance was significantly more
676 unsatisfactory than the other algorithms. Therefore FBP cannot be considered as a reliable alternative
677 for current algorithms to fill the gaps, especially the long ones.



678 5. Conclusions

679 Eight different gap-filling algorithms for estimating 16 meteorological drivers as well as the
680 three key ecosystem turbulent fluxes (sensible heat flux (Fh), latent heat flux (Fe), and net carbon flux
681 (Fc)) were investigated and their performance evaluated based on the datasets of five towers in
682 Australia. Overall, three ML algorithms, XGB, RF and ANNs, performed nearly equally well and
683 significantly better than their linear rivals (the CLR, PD, and ELN) in estimating the flux values.
684 However, the linear algorithms performed almost as equally well as the ML algorithms in assessing
685 the meteorological drivers. Amongst these eight algorithms, the RF showed the highest level of
686 robustness and reliability in estimating the Fc, as its closest rival, the XGB, could not capture the
687 minimum values equally well, despite providing slightly better RMSE and R². The PD was expected
688 to perform better than the linear methods and hoped to compete with the ML algorithms in estimating
689 the fluxes, but it failed to do so. The SVR was the only ML algorithm that did not perform at the same
690 level as the rest ML algorithms and was suspected of enduring over-fitting issues. Considering the
691 outcomes of the other researches undertaken in the OzFlux Network, e.g. (Cleverly et al., 2013; Isaac
692 et al., 2017), none of the ML algorithms used in this research was proven to provide substantially
693 better flux estimations compared with the standard method (ANNs). Nonetheless, amongst the
694 algorithms tested in this research, the RF showed some potential capabilities as an alternative due to
695 its more consistent performance regarding the long gaps. Eventually, we recommend suggestions
696 below to improve the results for similar prospective researches, as well as the QC and gap-filling
697 procedure of OzFlux Network:

698 1) Since the RF remained more consistent compared to its competitors -including the ANNs-, It is a
699 good idea to use RF alongside the commonly used algorithms in the challenging scenarios, such as
700 long gaps, to figure out whether this superiority can be generalised.

701 2) It appears that, even after three levels of quality control process done by the PyFluxPro platform,
702 the data are still noisy. This noisy data are an essential source of both uncertainty and inaccuracy of
703 the outcome, regardless of the algorithm used to gap-fill the data. As a result, another level of quality
704 control methods, such as Wavelets or Matrix Factorialisation, in addition to the current classical ones
705 used by the PyFluxPro and other similar platforms, can probably improve the data quality and thereby
706 improve the final imputation results.

707 3) For future researches, using recurrent neural networks (RNNs) instead of feedforward neural
708 networks (FFNN) could improve the predictions. That is likely because RNNs help the model to
709 consider temporal dynamic behaviour of time series, as unlike FFNN, wherein the activations flow
710 only from the input layer to the output layer, RNNs also have neuron connections pointing backwards
711 (Géron, 2019). This demand to an algorithm capable of considering time has been mentioned in
712 previous researches as one of the reasons why testing the new algorithms is needed (Richardson and
713 Hollinger, 2007).

714 3) Developing ensemble models using algorithms with different weaknesses and strengths may also
715 enhance the results where a single algorithm shows performance deficiency.



716 4) Given that some of the environmental drivers affect the F_c differently during the day versus night,
717 separating the diurnal and nocturnal datasets to train the algorithms possibly entails an improvement
718 in the outcome. Mainly because of the u^* threshold filtering and other problems associated with the
719 nocturnal period, the portion of diurnal data is generally, by far, outweighs the nocturnal data portion,
720 which potentially leads to a bias in the algorithm.

721 5) The same solution as number 4 is suggested for soil moisture estimation, as the behaviour of the
722 system on sunny days is utterly different from its conduct during the rainy periods. Moreover, the
723 system memory and the antecedent condition are undeniable features associated with soil moisture
724 (Ogle et al., 2015). Therefore, using the models that are capable of addressing these considerations are
725 more likely to improve the estimations.

726

727 6. Data availability

728 The data were used in this research are available through the following sources: The L3 and L4
729 data are accessible from the OzFlux data portal (<http://data.ozflux.org.au/portal>). Current ACCESS-R
730 and data are available from the BoM OPeNDAP server (<https://www.opendap.org/>). Likewise, the
731 data coming from the BoM AWS are accessible from (<http://www.bom.gov.au/climate/data>). Lastly,
732 the BIOS2 data are accessible from the ECMWF datasets portal
733 (<https://www.ecmwf.int/en/forecasts/datasets>). All data used in this research are available in this
734 repository address: ([https://research-repository.uwa.edu.au/en/datasets/a-comparison-of-gap-filling-
735 algorithms-for-eddy-covariance-fluxes](https://research-repository.uwa.edu.au/en/datasets/a-comparison-of-gap-filling-algorithms-for-eddy-covariance-fluxes)); DOI: [10.26182/5f292ee80a0c0](https://doi.org/10.26182/5f292ee80a0c0).

736

737 *Author contributions.* The ideas for this study originated in discussions with A. Mahabbati, J. Beringer,
738 and M. Leopold. A. Mahabbati carried out the analysis, supported by I. McHugh and P. Isaac. The
739 paper was prepared with contributions from all authors.

740

741 *Competing interests.* The authors declare that they have no conflict of interest.

742

743 *Acknowledgements.* The authors would like to acknowledge Terrestrial Ecosystems Research Network
744 (TERN) (www.tern.gov.au) and the OzFlux Network as a part of TERN for supporting the grants and
745 providing the required data, respectively. A. Mahabbati also personally thanks Prajwal Kalfe, Caroline
746 Johnson and Cacilia Ewenz for their support as regards Python programming, English academic
747 writing and PyFluxPro technical issues.

748

749

750 References

751 Allison, P. D.: Multiple Imputation for Missing Data: A Cautionary Tale, *Sociol. Methods Res.*, 28(3), 301–309,
752 doi:10.1177/0049124100028003003, 2000.

753 Altman, D. G. and Bland, J. M.: Missing data, *Br. Med. J.*, 334(7590), 424, doi:10.1136/bmj.38977.682025.2C, 2007.



- 754 Aubinet, M., Grelle, A., Ibrom, A., Rannik, Ü., Moncrieff, J., Foken, T., Kowalski, A. S., Martin, P. H., Berbigier, P., Bernhofer, C.,
755 Clement, R., Elbers, J., Granier, A., Grünwald, T., Morgenstern, K., Pilegaard, K., Rebmann, C., Snijders, W., Valentini, R. and
756 Vesala, T.: Estimates of the Annual Net Carbon and Water Exchange of Forests: The EUROFLUX Methodology, *Adv. Ecol. Res.*, 30,
757 113–175, doi:10.1016/S0065-2504(08)60018-5, 1999.
- 758 Aubinet, M., Vesala, T. and Papale, D.: *Eddy Covariance: A Practical Guide to Measurement and Data Analysis.*, 2012.
- 759 Baltagi, B.: *Econometric analysis of panel data*, [online] Available from: [http://www.sidalc.net/cgi-](http://www.sidalc.net/cgi-bin/wxis.exe/?IisScript=book2.xis&method=post&formato=2&cantidad=1&expresion=mfn=001143)
760 [bin/wxis.exe/?IisScript=book2.xis&method=post&formato=2&cantidad=1&expresion=mfn=001143](http://www.sidalc.net/cgi-bin/wxis.exe/?IisScript=book2.xis&method=post&formato=2&cantidad=1&expresion=mfn=001143) (Accessed 13 March 2018), 1995.
- 761 Barr, A. G., Black, T. A., Hogg, E. H., Kljun, N., Morgenstern, K. and Nestic, Z.: Inter-annual variability in the leaf area index of a boreal
762 aspen-hazelnut forest in relation to net ecosystem production, *Agric. For. Meteorol.*, 126(3–4), 237–255,
763 doi:10.1016/J.AGRFORMET.2004.06.011, 2004.
- 764 Barr, A. G., Richardson, A. D., Hollinger, D. Y., Papale, D., Arain, M. A., Black, T. A., Bohrer, G., Dragoni, D., Fischer, M. L., Gu, L.,
765 Law, B. E., Margolis, H. A., Mccaughey, J. H., Munger, J. W., Oechel, W. and Schaeffer, K.: Use of change-point detection for friction-
766 velocity threshold evaluation in eddy-covariance studies, *Agric. For. Meteorol.*, 171–172, 31–45, doi:10.1016/j.agrformet.2012.11.023,
767 2013.
- 768 Bennett, N. D., Croke, B. F. W., Guariso, G., Guillaume, J. H. A., Hamilton, S. H., Jakeman, A. J., Marsili-Libelli, S., Newham, L. T. H.,
769 Norton, J. P., Perrin, C., Pierce, S. A., Robson, B., Seppelt, R., Voinov, A. A., Fath, B. D. and Andreassian, V.: Characterising
770 performance of environmental models, *Environ. Model. Softw.*, 40, 1–20, doi:10.1016/j.envsoft.2012.09.011, 2013.
- 771 Beringer, J., Hutley, L. B., McHugh, I., Arndt, S. K., Campbell, D., Cleugh, H. A., Cleverly, J., De Dios, V. R., Eamus, D., Evans, B.,
772 Ewenz, C., Grace, P., Griebel, A., Haverd, V., Hinko-Najera, N., Huete, A., Isaac, P., Kanniah, K., Leuning, R., Liddell, M. J.,
773 MacFarlane, C., Meyer, W., Moore, C., Pendall, E., Phillips, A., Phillips, R. L., Prober, S. M., Restrepo-Coupe, N., Rutledge, S.,
774 Schroder, I., Silberstein, R., Southall, P., Sun Yee, M., Tapper, N. J., Van Gorsel, E., Vote, C., Walker, J. and Wardlaw, T.: An
775 introduction to the Australian and New Zealand flux tower network - OzFlux, *Biogeosciences*, 13(21), 5895–5916, doi:10.5194/bg-13-
776 5895-2016, 2016a.
- 777 Beringer, J., McHugh, I. and KLJUN, N.: Dynamic INtegrated Gap filling and partitioning for Ozflux (DINGO), *Biogeosciences*
778 *Discuss.*, *OzFlux spe(In prep)*, 1457–1460, doi:doi:10.5194/bg-2016-188, 2016b.
- 779 Beringer, J., McHugh, I., Hutley, L. B., Isaac, P. and Kljun, N.: Technical note: Dynamic INtegrated Gap-filling and partitioning for
780 OzFlux (DINGO), *Biogeosciences*, 14(6), 1457–1460, doi:10.5194/bg-14-1457-2017, 2017.
- 781 Burba, G. and Anderson, D.: A brief practical guide to eddy covariance flux measurements: principles and workflow examples for
782 scientific and industrial applications. [online] Available from:
783 https://books.google.com/books?hl=en&lr=&id=mCs11_8GdriC&oi=fnd&pg=PA6&dq=A+Brief+Practical+Guide+to+Eddy+Covarianc
784 [e+Flux+Measurements&ots=TKTg25Yq5X&sig=eBYc819N7Jh3gNhJnfEL1e40eM](https://books.google.com/books?hl=en&lr=&id=mCs11_8GdriC&oi=fnd&pg=PA6&dq=A+Brief+Practical+Guide+to+Eddy+Covarianc) (Accessed 11 February 2020), 2010.
- 785 Chen, T. and Guestrin, C.: XGBoost: A scalable tree boosting system, *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 13-17-
786 Augu, 785–794, doi:10.1145/2939672.2939785, 2016.
- 787 Cleverly, J., Boulain, N., Villalobos-Vega, R., Grant, N., Faux, R., Wood, C., Cook, P. G., Yu, Q., Leigh, A. and Eamus, D.: Dynamics of
788 component carbon fluxes in a semi-arid *Acacia* woodland, central Australia, *J. Geophys. Res. Biogeosciences*, 118(3), 1168–1185,
789 doi:10.1002/jgrg.20101, 2013.
- 790 Devore, J. L.: *Probability and Statistics for Engineering and the Sciences.*, *Biometrics*, 47(4), 1638, doi:10.2307/2532427, 1991.
- 791 Dragoni, D., Schmid, H. P., Grimmond, C. S. B. and Loescher, H. W.: Uncertainty of annual net ecosystem productivity estimated
792 using eddy covariance flux measurements, *J. Geophys. Res.*, 112(D17), D17102, doi:10.1029/2006JD008149, 2007.
- 793 Dreyfus, S. E.: Artificial neural networks, back propagation, and the kelley-bryson gradient procedure, *J. Guid. Control. Dyn.*, 13(5),
794 926–928, doi:10.2514/3.25422, 1990.
- 795 Drucker, H., Surges, C. J. C., Kaufman, L., Smola, A. and Vapnik, V.: Support vector regression machines, in *Advances in Neural*
796 *Information Processing Systems*, vol. 1, pp. 155–161., 1997.
- 797 Falge, E., Baldocchi, D., Olson, R., Anthoni, P., Aubinet, M., Bernhofer, C., Burba, G., Ceulemans, R., Clement, R., Dolman, H., Granier,
798 A., Gross, P., Grünwald, T., Hollinger, D., Jensen, N. O., Katul, G., Keronen, P., Kowalski, A., Lai, C. T., Law, B. E., Meyers, T.,
799 Moncrieff, J., Moors, E., Munger, J. W., Pilegaard, K., Rannik, Ü., Rebmann, C., Suyker, A., Tenhunen, J., Tu, K., Verma, S., Vesala,
800 T., Wilson, K. and Wofsy, S.: Gap filling strategies for defensible annual sums of net ecosystem exchange, *Agric. For. Meteorol.*,
801 107(1), 43–69, doi:10.1016/S0168-1923(00)00225-2, 2001.



- 802 Farley, B. G. and Clark, W. A.: Simulation of self-organizing systems by digital computer, *IRE Prof. Gr. Inf. Theory*, 4(4), 76–84,
803 doi:10.1109/TIT.1954.1057468, 1954.
- 804 Freedman, D. A.: *Statistical Models: Theory and Practice*. Cambridge University Press - 2nd edition. [online] Available from:
805 <https://www.cambridge.org/au/academic/subjects/statistics-probability/statistical-theory-and-methods/statistical-models-theory-and-practice-2nd-edition?format=PB> (Accessed 21 March 2020), 2009.
806
- 807 Friedman, J. H.: Stochastic gradient boosting, *Comput. Stat. Data Anal.*, 38(4), 367–378, doi:10.1016/S0167-9473(01)00065-2, 2002.
- 808 Gani, A., Mohammadi, K., Shamsirband, S., Altameem, T. A., Petković, D. and Ch, S.: A combined method to estimate wind speed
809 distribution based on integrating the support vector machine with firefly algorithm, *Environ. Prog. Sustain. Energy*, 35(3), 867–875,
810 doi:10.1002/ep.12262, 2016.
- 811 Géron, A.: *Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems*.
812 [online] Available from: <https://books.google.com.au/books?hl=en&lr=&id=HHetDwAAQBAJ&oi=fnd&pg=PP1&dq=hands-on+machine+learning+with+&ots=0KvfZqlgOo&sig=5tH2IHRsUaTMTy6CfQ6lw3UDKa4> (Accessed 7 February 2020), 2019.
813
- 814 Hagen, S. C., Braswell, B. H., Linder, E., Frolking, S., Richardson, A. D. and Hollinger, D. Y.: Statistical uncertainty of eddy flux - Based
815 estimates of gross ecosystem carbon exchange at Howland Forest, Maine, *J. Geophys. Res. Atmos.*, 111(8), 1–12,
816 doi:10.1029/2005JD006154, 2006.
- 817 Harrell, F. E.: *Regression Modeling Strategies: With Applications to Linear Models, Logistic, in books.google.nl*. [online] Available
818 from:
819 <https://books.google.com.au/books?hl=en&lr=&id=94RgCgAAQBAJ&oi=fnd&pg=PR7&dq=regression+modeling+strategies+frank+h>
820 [arrell&ots=ZAAt4Rsa51r&sig=mikE1s9G4IXzqZKEie-iVA9GTV0&redir_esc=y#v=onepage&q=regression+modeling+strategies+frank](https://books.google.com.au/books?hl=en&lr=&id=94RgCgAAQBAJ&oi=fnd&pg=PR7&dq=regression+modeling+strategies+frank+h)
821 [harrell&f=false](https://books.google.com.au/books?hl=en&lr=&id=94RgCgAAQBAJ&oi=fnd&pg=PR7&dq=regression+modeling+strategies+frank+h) (Accessed 11 February 2020), 2014.
- 822 Harvey, A. C. and Peters, S.: Estimation procedures for structural time series models, *J. Forecast.*, 9(2), 89–108,
823 doi:10.1002/for.3980090203, 1990.
- 824 Haverd, V., Briggs, P., Trudinger, C., Nieradzik, L. and Canadell, P.: *BIOS2 – Frontier Modelling of the Australian Carbon and Water*
825 *Cycles*, 2015.
- 826 Ho, T. K.: Random decision forests, *Proc. Int. Conf. Doc. Anal. Recognition, ICDAR*, 1, 278–282, doi:10.1109/ICDAR.1995.598994, 1995.
- 827 Ho, T. K.: 00709601.Pdf, , 20(8), 832–844, 1998.
- 828 Hollinger, D. Y., Goltz, S. M., Davidson, E. A., Lee, J. T., Tu, K. and Valentine, H. T.: Seasonal patterns and environmental control of
829 carbon dioxide and water vapour exchange in an ecotonal boreal forest, *Glob. Chang. Biol.*, 5(8), 891–902, doi:10.1046/j.1365-
830 2486.1999.00281.x, 1999.
- 831 Hsiao, C., Hashem Pesaran, M. and Kamil Tahmiscioglu, A.: Maximum likelihood estimation of fixed effects dynamic panel data
832 models covering short time periods, *J. Econom.*, 109(1), 107–150, doi:10.1016/S0304-4076(01)00143-9, 2002.
- 833 Hui, D., Wan, S., Su, B., Katul, G., Monson, R. and Luo, Y.: Gap-filling missing data in eddy covariance measurements using multiple
834 imputation (MI) for annual estimations, *Agric. For. Meteorol.*, 121(1–2), 93–111, doi:10.1016/S0168-1923(03)00158-8, 2004.
- 835 Hutley, L. B., Leuning, R., Beringer, J. and Cleugh, H. a: The utility of the eddy covariance technique as a tool in carbon accounting:
836 tropical savanna as a case study, *Aust. J. Bot.*, 53, 663–675, 2005.
- 837 Isaac, P., Cleverly, J., McHugh, I., Van Gorsel, E., Ewenz, C. and Beringer, J.: OzFlux data: Network integration from collection to
838 curation, *Biogeosciences*, 14(12), 2903–2928, doi:10.5194/bg-14-2903-2017, 2017.
- 839 Izady, A., Davary, K., Alizadeh, A., Moghaddam Nia, A., Ziaei, A. N. and Hashemina, S. M.: Application of NN-ARX Model to
840 Predict Groundwater Levels in the Neishaboor Plain, Iran, *Water Resour. Manag.*, 27(14), 4773–4794, doi:10.1007/s11269-013-0432-y,
841 2013.
- 842 Izady, A., Abdalla, O. and Mahabbati, A.: Dynamic panel-data-based groundwater level prediction and decomposition in an arid
843 hardrock–alluvium aquifer, *Environ. Earth Sci.*, 75(18), 1–13, doi:10.1007/s12665-016-6059-6, 2016.
- 844 Jerome H. Friedman: *Greedy Function Approximation: A Gradient Boosting Machine* on JSTOR, *Ann. Stat.*, 29, 1189–1232 [online]
845 Available from: https://www.jstor.org/stable/2699986?seq=1#metadata_info_tab_contents (Accessed 27 August 2019), 2001.
- 846 Kang, H.: The prevention and handling of the missing data, *Korean J. Anesthesiol.*, 64(5), 402–406, doi:10.4097/kjae.2013.64.5.402, 2013.



- 847 Kim, Y., Johnson, M. S., Knox, S. H., Black, T. A., Dalmagro, H. J., Kang, M., Kim, J. and Baldocchi, D.: Gap-filling approaches for eddy
848 covariance methane fluxes: A comparison of three machine learning algorithms and a traditional method with principal component
849 analysis, *Glob. Chang. Biol.*, 26(3), 1499–1518, doi:10.1111/gcb.14845, 2020.
- 850 Kock, N. and Gaskins, L.: Simpson’s paradox, moderation and the emergence of quadratic relationships in path models: an
851 information systems illustration, *Int. J. Appl. Nonlinear Sci.*, 2(3), 200, doi:10.1504/ijans.2016.077025, 2016.
- 852 Kunwor, S., Starr, G., Loescher, H. W. and Staudhammer, C. L.: Preserving the variance in imputed eddy-covariance measurements:
853 Alternative methods for defensible gap filling, *Agric. For. Meteorol.*, 232, 635–649, doi:10.1016/j.agrformet.2016.10.018, 2017.
- 854 Law, B. E., Falge, E., Gu, L., Baldocchi, D. D., Bakwin, P., Berbigier, P., Davis, K., Dolman, A. J., Falk, M., Fuentes, J. D., Goldstein, A.,
855 Granier, A., Grelle, A., Hollinger, D., Janssens, I. A., Jarvis, P., Jensen, N. O., Katul, G., Mahli, Y., Matteucci, G., Meyers, T., Monson,
856 R., Munger, W., Oechel, W., Olson, R., Pilegaard, K., Paw U H, K. T., Thorgeirsson, H., Valentini, R., Verma, S., Vesala, T., Wilson, K.
857 and Wofsy, S.: Jourassess2, *Agric. For. Meteorol.*, 113(113), 97–120, 2002.
- 858 Lee, X., Fuentes, J. D., Staebler, R. M. and Neumann, H. H.: Long-term observation of the atmospheric exchange of CO₂ with a
859 temperate deciduous forest in southern Ontario, Canada, *J. Geophys. Res. Atmos.*, 104(D13), 15975–15984,
860 doi:10.1029/1999JD900227, 1999.
- 861 Little, R. J. A.: *Statistical analysis with missing data*, 2nd ed., edited by D. B. Rubin, Wiley, Hoboken, N.J., 2002.
- 862 Mahabhati, A., Izady, A., Mousavi Baygi, M., Davary, K. and Hasheminia, S. M.: Daily soil temperature modeling using ‘panel-data’
863 concept, *J. Appl. Stat.*, 44(8), 1385–1401, doi:10.1080/02664763.2016.1214240, 2017.
- 864 Moffat, A. M., Papale, D., Reichstein, M., Hollinger, D. Y., Richardson, A. D., Barr, A. G., Beckstein, C., Braswell, B. H., Churkina, G.,
865 Desai, A. R., Falge, E., Gove, J. H., Heimann, M., Hui, D., Jarvis, A. J., Kattge, J., Noormets, A. and Stauch, V. J.: Comprehensive
866 comparison of gap-filling techniques for eddy covariance net carbon fluxes, *Agric. For. Meteorol.*, 147(3–4), 209–232,
867 doi:10.1016/j.agrformet.2007.08.011, 2007.
- 868 Molenberghs, G., Fitzmaurice, G., Kenward, M. G., Tsiatis, A., Verbeke, G., Fitzmaurice, G., Kenward, M. G., Tsiatis, A. and Verbeke,
869 G.: *Handbook of Missing Data Methodology*, Chapman and Hall/CRC., 2014.
- 870 Ogle, K., Barber, J. J., Barron-Gafford, G. A., Bentley, L. P., Young, J. M., Huxman, T. E., Loik, M. E. and Tissue, D. T.: Quantifying
871 ecological memory in plant and ecosystem processes, *Ecol. Lett.*, 18(3), 221–235, doi:10.1111/ele.12399, 2015.
- 872 Papale, D. and Valentini, R.: A new assessment of European forests carbon exchanges by eddy fluxes and artificial neural network
873 spatialization, *Glob. Chang. Biol.*, 9(4), 525–535, doi:10.1046/j.1365-2486.2003.00609.x, 2003.
- 874 Pilegaard, K., Hummelshøj, P., Jensen, N. O. and Chen, Z.: Two years of continuous CO₂ eddy-flux measurements over a Danish
875 beech forest, *Agric. For. Meteorol.*, 107(1), 29–41, doi:10.1016/S0168-1923(00)00227-6, 2001.
- 876 Reichle, R. H., Koster, R. D., Dong, J. and Berg, A. A.: Global soil moisture from satellite observations, land surface models, and
877 ground data: Implications for data assimilation, *J. Hydrometeorol.*, 5(3), 430–442, doi:10.1175/1525-
878 7541(2004)005<0430:GSMFSO>2.0.CO;2, 2004.
- 879 Richardson, A. D. and Hollinger, D. Y.: A method to estimate the additional uncertainty in gap-filled NEE resulting from long gaps in
880 the CO₂ flux record, *Agric. For. Meteorol.*, 147(3–4), 199–208, doi:10.1016/j.agrformet.2007.06.004, 2007.
- 881 Richardson, A. D., Braswell, B. H., Hollinger, D. Y., Burman, P., Davidson, E. A., Evans, R. S., Flanagan, L. B., Munger, J. W., Savage,
882 K., Urbanski, S. P. and Wofsy, S. C.: Comparing simple respiration models for eddy flux and dynamic chamber data, *Agric. For.*
883 *Meteorol.*, 141(2–4), 219–234, doi:10.1016/J.AGRFORMET.2006.10.010, 2006.
- 884 Richardson, A. D., Aubinet, M., Barr, A. G., Hollinger, D. Y., Ibrom, A., Lasslop, G. and Reichstein, M.: Uncertainty Quantification, in
885 *Eddy Covariance*, pp. 173–209., 2012.
- 886 Sahoo, A. K., Dirmeyer, P. A., Houser, P. R. and Kafatos, M.: A study of land surface processes using land surface models over the
887 Little River Experimental Watershed, Georgia, *J. Geophys. Res. Atmos.*, 113(20), doi:10.1029/2007JD009671, 2008.
- 888 Scanlon, T. M. and Kustas, W. P.: Partitioning carbon dioxide and water vapor fluxes using correlation analysis, *Agric. For. Meteorol.*,
889 150(1), 89–99, doi:10.1016/j.agrformet.2009.09.005, 2010.
- 890 Scanlon, T. M. and Sahu, P.: On the correlation structure of water vapor and carbon dioxide in the atmospheric surface layer: A basis
891 for flux partitioning, *Water Resour. Res.*, 44(10), doi:10.1029/2008WR006932, 2008.



- 892 Staebler, M.: Long-term observation of the atmospheric exchange of CO₂ with a temperate deciduous forest in southern Ontario,
893 Canada ecosystem (net ecosystem production turbulence is turbulent, *Data Process.*, 104, 975–984, 1999.
- 894 Tannenbaum, C. E.: The empirical nature and statistical treatment of missing data., *Diss. Abstr. Int. Sect. A Humanit. Soc. Sci.*, 70(10-
895 A), 3825 [online] Available from: http://gateway.proquest.com/openurl?url_ver=Z39.88-2004&rft_val_fmt=info:ofi/fmt:kev:mtx:dissertation&res_dat=xri:pqdiss&rft_dat=xri:pqdiss:3381876%5Cnhttp://ovidsp.ovid.com/ovidweb.cgi?T=JS&PAGE=reference&D=psyc7&NEWS=N&AN=2010-99071-044, 2010.
- 898 Taylor, S. J. and Letham, B.: Business Time Series Forecasting at Scale, , doi:10.7287/peerj.preprints.3190v2, 2017.
- 899 Taylor, S. J. and Letham, B.: Forecasting at Scale, *Am. Stat.*, 72(1), 37–45, doi:10.1080/00031305.2017.1380080, 2018.
- 900 Tenhunen, J. D., Valentini, R., Köstner, B., Zimmermann, R. and Granier, A.: Variation in forest gas exchange at landscape to
901 continental scales, *Ann. des Sci. For.*, 55(1–2), 1–11, doi:10.1051/forest:19980101, 1998.
- 902 Wooldridge, J. M.: *Econometric Analysis of Cross Section and Panel Data.*, 2008.
- 903 Ye, J., Chow, J.-H., Chen, J. and Zheng, Z.: Stochastic gradient boosted distributed decision trees, in *Proceeding of the 18th ACM*
904 *conference on Information and knowledge management - CIKM '09*, p. 2061, ACM Press, New York, New York, USA., 2009.
- 905 Zhao, X. and Huang, Y.: A comparison of three gap filling techniques for eddy covariance net carbon fluxes in short vegetation
906 ecosystems, *Adv. Meteorol.*, 2015, 1–12, doi:10.1155/2015/260580, 2015.
- 907 Zou, H. and Hastie, T.: Regularization and variable selection via the elastic net. [online] Available from:
908 <https://citeseerx.ist.psu.edu/viewdoc/download?jsessionid=22250F01CC77D55C54B6BAFF4512C9E3?doi=10.1.1.124.4696&rep=rep1&type=pdf> (Accessed 28 August 2019), 2005.
- 909
- 910