

1 Evaluation often machine learning methods for estimating terrestrial evapotranspiration from
2 remote sensing

3

4 Corinne Carter and Shunlin Liang*

5

6 Department of Geographical Sciences, University of Maryland, College Park, MD USA

7 *corresponding author, sliang@umd.edu

8

9 **Abstract:** Remote sensing retrieval of evapotranspiration (ET), or surface latent heat exchange
10 (LE), is of great utility for many applications. Machine learning (ML) methods have been
11 extensively used in many disciplines, but so far little work has been performed systematically
12 comparing ML methods for ET retrieval. This paper provides an evaluation of ten ML
13 methods for estimating daily ET based on daily Global Land Surface Satellite (GLASS) radiation
14 data and high-level Moderate-Resolution Imaging Spectroradiometer (MODIS) data products
15 and ground measured ET data from 184 flux tower sites. Measurements of accuracy (RMSE, R^2 ,
16 and bias) and run time were made for each of ten ML methods with a smaller training data set (n
17 = 7910 data points) and a larger training data set ($n = 69,752$ data points). Inclusion of more input
18 variables improved algorithm performance but had little effect on run time. The best results were
19 obtained with the larger training data set using the bootstrap aggregation (bagging) regression
20 tree (validation RMSE = 19.91 W/m^2) and three hidden layer neural network (validation RMSE
21 = 20.94 W/m^2), although the less computationally demanding random kernel (RKS) algorithm
22 also produced good results (validation RMSE = 22.22 W/m^2). Comparison of results from sites
23 with different ecosystem types showed the best results for evergreen, shrub, and grassland sites,

24 and the weakest results for wetland sites. Generally, performance was not improved by training
25 with data from only the same ecosystem type.

26 **Introduction**

27 Evapotranspiration (ET), often expressed as an energy flux, the latent heat of evaporation (LE),
28 is an important linkage between the surface energy and water balances and an indicator of
29 vegetation health. Compared to the radiative elements of the surface energy balance, there is
30 more uncertainty in LE measurements. Ground-based measurements are made at small scales
31 with weighing lysimeters, and at scales of tens of meters to kilometers with flux towers
32 and scintillometers. However, these measurements are sparse outside the northern hemisphere
33 midlatitudes. Remote sensing data, reanalyses, and ground-based observations have been
34 combined in a variety of ways to retrieve LE. Reviews of methods for obtaining LE through
35 remote sensing are available in Zhang et al. (2016), Wang and Dickinson (2012), and Kalma et al.
36 (2008). Some of these methods (e. g. Wang and Liang 2008; Yao et al. 2011, 2013, 2015; Yebra
37 et al. 2013; Helman et al. 2015) use statistical regression techniques. Carter and Liang
38 (2018) evaluated a number of statistical regression formulas for obtaining LE.

39
40 Machine learning (ML) methods are means of extracting patterns from data sets with little prior
41 knowledge of those patterns. The best-known ML methods include neural networks (NN), tree
42 methods, and support vector machines (SVMs). The model tree ensemble technique has been
43 applied to the problem of determining global trends in LE by Jung et al. (2010). Multiple studies
44 have been conducted using machine learning techniques for downscaling LE (Ke et al. 2017,
45 2016; Kaheil et al. 2008) and drought detection and forecasting (Rhee and Im 2017; Park et al.
46 2016). There are also a number of studies comparing the performance of different ML techniques

47 for obtaining LE. In these studies (eg.Deo et al. 2016; Dou and Yang 2018)no single ML method
48 produced the best results. Most of these studies, with the notable exception of Jung et al.
49 (2010),involve training to measurements of LE from a relatively small number of locations (20
50 or fewer). In previous studies where ML method comparisons are performed four or fewer
51 methods are compared.

52
53 The goal of this study is to evaluate the utility of a range of ML algorithms for obtaining LE
54 from remote sensing data on a global basis, and to evaluate their performance for different
55 ecosystem types.

56
57

58 **Data**

59 The remote sensing data used in this study are Global Land Surface Satellite (GLASS) radiation
60 data and Moderate-Resolution Imaging Spectroradiometer (MODIS) high-level data products.
61 Ground-based Fluxnet tower site data were also used. The data variables, sources, and spatial
62 and time resolutions for each data set used are listed in Table 1.

63

64 Table 1: Input and validation data used in this study

65

Abbreviation	Variable	Source	Frequency	Spatial resolution
LE	Surface latent heat	Fluxnet	Half-hourly, averaged to daily	Flux tower footprint
R _n	Net radiation at surface	Fluxnet	Half-hourly, averaged to daily	Flux tower footprint
DSR	Downward surface radiation	GLASS	Daily	5 km

PAR	Photosynthetically active radiation	GLASS	Daily	5 km
NDVI	Normalized difference vegetation index	MODIS	16-day, interpolated to daily	250 meters
EVI	Enhanced vegetation index	MODIS	16-day, interpolated to daily	250 meters
LAI	Leaf area index	MODIS	8-day, interpolated to daily	500 meters
FPAR	Fraction of photosynthetically adjusted radiation	MODIS	8-day, interpolated to daily	500 meters
Albedo	Albedo	MODIS	8-day, interpolated to daily	500 meters
NBAR	Nadir BRDF-adjusted reflectance	MODIS	8-day, interpolated to daily	500 meters

66

67 The GLASS data set (Liang et al. 2013, 2014) consists of radiative and biophysical parameters
68 generated using data from multiple satellite sensors. The products used here are the downward
69 shortwave radiation (DSR) and photosynthetically active radiation (PAR).

70

71 Several parameters obtained from MODIS were also used in this analysis: Normalized-difference
72 vegetation index (NDVI) and enhanced vegetation index (EVI) (Didan 2015), leaf area index
73 (LAI) and fraction of photosynthetically active radiation absorbed (FPAR) (Myneni and
74 Knyazikhin 2015), surface albedo, and nadir BRDF-adjusted reflectance (NBAR) (Schaaf and
75 Wang 2015a, 2015b). Subsets of all MODIS products used were generated centered on the
76 coordinates of each flux tower site. All MODIS products were linearly interpolated to daily
77 frequency.

78

79 Flux tower data were used for validation of the ML algorithms, and also for testing the effects of
80 using remote sensing vs. ground-based radiation data as input. A total of 184 flux tower sites

81 were used, 119 from the Ameriflux network (<http://ameriflux.lbl.gov>) and 65 from the
82 Fluxnet2015 data set (<http://fluxnet.fluxdata.org/data/fluxnet2015-dataset/>). The half-hourly LE
83 and net radiation (R_n) variables from these data were pre-processed by removing all data days for
84 which there were not at least 40 of 48 possible observations present, then averaging the
85 remaining observations. A map of the site locations and information about their distribution
86 across ecosystem types is given in Carter and Liang (2018).

87

88 A total of 79098 site-days of data were used, randomly partitioned twice into 7,910/ 35,594/
89 35,594 and 63,278/ 7,910/ 7,910 site-days of training, validation, and test data,
90 respectively. These training, validation, and test data sets were used in every case except where
91 only data from individual ecosystem types was used. Each site-day includes a flux tower LE
92 value associated with the remote sensing parameters retrieved at the site for that day. For
93 purposes of this study, we treated the flux tower LE values as “ground truth”, although flux
94 tower footprints vary and may not always coincide closely with the pixel size of the remote
95 sensing data and also require adjustment to compensate for lack of energy balance closure.

96

97 **Methods**

98 In order to use the ML algorithms properly, it is necessary to adjust one or more tunable
99 parameters for each algorithm. This is done by training the algorithms with a training data set for
100 different parameter values and checking against a validation data set until the optimum
101 parameter values are found, for example by minimizing RMSE for the validation data set. Once
102 the optimization is performed, the optimized algorithm is checked against a test data set separate
103 from the training and validation data sets. Timing of a single iteration of training with each

104 training data set and checking with the validation data set was performed for each algorithm as a
 105 feasibility check, since it is necessary to repeat this process tens to hundreds of times to tune the
 106 algorithms. Timing was conducted on a server with 24 6-core 3.33GHz Intel Xenon X5680
 107 CPUs.

108

109 Fourteen ML algorithms were subjected to initial timing tests with the smaller training data set.
 110 Based on the results of this timing, 10 of the original algorithms were tuned with the smaller
 111 training data set. Of those 10, 8 were found to run with enough efficiency for systematic tuning
 112 with the larger training data set to be feasible. The 14 ML algorithms considered are listed in
 113 Table 2, with references to descriptions of each of the algorithms.

114

115 Table 2: Algorithms tested in this study.

116

Family	Full name	Abbreviation	Tuned with small training data set	Tuned with large training data set
Linear	Regularized linear regression*	RLR	Yes	Yes
	Least absolute shrinkage and selection operator regression*	LASSO	Yes	Yes
	Elastic net regularization*	ELASTIC	Yes	Yes
Kernel	Gaussian process regression*(Murphy 2012)	GPR	No	No
	Kernel ridge regression (Murphy 2012)	KRR	Yes	No
	Random kernel (Rahimi and Recht 2009, Perez-Suay et al. 2017)	RKS	Yes	Yes
	Variational	VHGPR	No	No

	heteroscedastic Gaussian process regression(Lazaro-Gredilla et al. 2014; Lazaro-Gredilla and Titsias 2011)			
Tree	Regression tree*	TREE	Yes	Yes
	Bootstrap aggregation (bagging) tree*	BAGTREE	Yes	Yes
	Boosted regression tree*	BOOST	Yes	Yes
Neural network	Standard neural network (1, 2, and 3 hidden layers)*	NN	Yes	Yes
	Extreme learning machine (Huang et al. 2006)	ELM	No	No
Support vector	Support vector regression* (Smola and Scholkopf 2004)	SVR	Yes	No
	Relevance vector machine(Thayananthan et al. 2006)	RVM	No	No

117 Algorithms marked with an asterisk (*) are described in Hastie et al. (2008).

118

119 Optimum values of the parameters are found by minimizing the root mean square error (RMSE)
120 of the algorithm when applied to the validation data set. The coefficient of determination R^2 and
121 bias were also used to characterize the correspondence of the modeled LE from different surface
122 types. The implementation in Matlab of all of these algorithms, with the exception of the random
123 kernel (RKS), was obtained from package “simpleR” (Lazaro-Gredilla et al. 2014). The RKS
124 algorithm code was obtained from <http://isp.uv.es/code/rks2017.html>(Pérez-Suay et al. 2017).

125

126 Initially, one training/ validation iteration was timed for each ML algorithm using a smaller
127 training data set. Algorithms that took more than ten minutes for one iteration were removed
128 from further consideration. The remaining algorithms were timed for one training/ validation

129 iteration with a larger training data set. Two of the algorithms (KRR and SVR) that were
130 tractable with the smaller training data set became too computationally demanding with the
131 larger training data set. These timing results are given in Results section 1. The linear regression,
132 boost tree, and RKS were used to test the effects on accuracy and computation time of using
133 different combinations of input variables, as discussed in Results section 2.

134

135 Once the most viable combinations of ML algorithms and input variables had been identified,
136 each algorithm was tuned by varying all parameters of each algorithm independently. The
137 optimal tuning parameters with respect to the validation data set (lowest validation RMSE) were
138 applied to the test data set. Variation in algorithm performance with tuning of parameters and
139 optimum algorithm performances are shown in Results section 3.

140

141 The RKS, BAGTREE, and 2 and 3 hidden layer NNs were applied to each of seven ecosystem
142 types' test data sets, first with the algorithms optimized using data from all of the sites, and then
143 optimized using data from sites of the same ecosystem type only. These results are given in
144 Results section 4.

145

146

147 **Results**

148 1. Initial time trials of ML algorithms

149 The time in seconds for each algorithm to run a single iteration of training and validation with
150 all input variables is shown in Table 3. If an algorithm took longer than 10 minutes to run a single
151 iteration, it is labeled “prohibitive” and no further testing was done for that algorithm.

152
 153
 154
 155
 156

Table 3: Time in seconds for one iteration of training and validation for each algorithm.

Algorithm	Small training data set (7,910 data points)	Large training data set (69,752 data points)
RLR	0.0014	0.0087
LASSO	19.3208	118.1654
ELASTIC	21.8523	109.0231
GPR	prohibitive	
KRR	219.6609	prohibitive
RKS, D = 100	0.0945	1.1625
D = 400	0.3950	4.7896
D = 1000	1.1794	13.026
D = 4000	9.1198	91.2732
TREE	20.2009	351.7759
BAGTREE	15.9202	114.309
BOOST (200 trees)	3.1619	9.1393
NN, 1 HL, 5 neurons	4.0422	102.3333
1 HL, 30 neurons	6.2271	207.9746
2 HL, 5 x 5 neurons	5.5793	108.4992
2 HL, 10 x 10 neurons	6.0898	131.7211
2 HL, 30 x 30 neurons	10.0308	436.8679
3 HL, 5 x 5 x 2 neurons	4.6199	128.9298
3 HL, 10 x 10 x 10 neurons	7.8482	153.7699
3 HL, 50 x 5 x 2 neurons	12.849	prohibitive
3 HL, 150 x 30 x 10 neurons	249.671	prohibitive
ELM	prohibitive	
SVR	41.6029	prohibitive
RVM	prohibitive	
VHGPR	prohibitive	

157 Note: For RKS, “D” represents the number of random functions used. The number of hidden
 158 layers (HL) and neurons in each of the NN trials is also indicated.

159

160

161

162 2. Combinations of input variables

163 In order to test the effects on speed and accuracy of using different combinations of input
 164 variables, trials of a single training and testing cycle were done with the linear regression, boost
 165 tree, and RKS methods using the small training data set. Boost tree RMSE was found after
 166 optimizing the number of trees, but timing trials were done for 100 and 1000 trees. Results are
 167 summarized in Table 4. Generally, including more input variables produced similar or more
 168 accurate results at little additional computational cost. Using radiation information from surface
 169 measurements produced results of similar accuracy to using the GLASS radiation variables.

170

171 Table 4: Accuracy and timing tests for different combinations of input variables using the
 172 smaller training data set.

173

Variables	Linear regression RMSE (W/m ²)	Linear regression timing (s)	BOOST RMSE (W/m ²)	BOOST timing (100 trees) (s)	BOOST timing (1000 trees) (s)	RKS RMSE (W/m ²)	RKS timing (s)
R _n + NDVI	34.68	6.74 x 10 ⁻⁴	32.52	2.71	28.25	31.71	0.17
R _n + NBAR	31.80	0.0032	29.18	3.15	31.11	28.10	0.20
R _n + NDVI + EVI + LAI + FPAR + NBAR + Albedo	31.78	0.0098	28.03	4.33	38.08	28.10	0.20

DSR + NDVI	33.83	8.69×10^{-4}	33.07	2.70	28.51	31.73	0.17
PAR + NDVI	33.67	6.72×10^{-4}	32.89	2.77	27.74	31.84	0.18
DSR + PAR + NDVI	33.66	0.0016	32.84	2.69	28.88	31.19	0.18
DSR + PAR + EVI	33.29	0.0013	32.23	2.94	28.96	30.41	0.18
DSR + PAR + NDVI + EVI	33.21	0.0014	31.62	3.04	28.27	29.78	0.18
DSR + PAR + FPAR	33.62	0.0014	32.75	2.99	27.83	31.09	0.17
DSR + PAR + LAI	34.47	0.0016	32.46	2.94	28.31	31.10	0.20
DSR + PAR + LAI + FPAR	33.63	0.0014	31.84	2.80	27.85	30.36	0.20
DSR + PAR + NDVI + EVI + LAI + FPAR	32.90	0.0030	30.53	3.28	31.55	29.14	0.20
DSR + PAR + Albedo	36.96	0.0015	35.85	2.70	26.92	35.49	0.17
DSR + PAR + NDVI + EVI + LAI + FPAR + Albedo	32.89	0.0028	30.26	3.26	31.31	28.95	0.18
DSR + PAR + NBAR	31.27	0.0044	29.44	3.46	33.81	28.44	0.19
DSR + PAR + NDVI + EVI + LAI + FPAR + NBAR	31.21	0.0058	28.53	3.74	35.94	28.32	0.17

174

175 The first set of three trials was made with the R_n taken from the ground-based measurements.

176 The second set of trials tested the effects of using DSR or PAR or both in combination with

177 NDVI and EVI. Using both radiation variables with NDVI produced better results than using

178 either of them separately. Using all four variables produced the lowest RMSEs at little additional

179 computational cost. For all subsequent trials, both DSR and PAR were included.

180

181 The third set of comparisons tested the use of LAI and FPAR as input variables. The fourth set of
182 trials included albedo as one of the input data variables. When albedo was the only input other
183 than DSR and PAR the highest RMSEs for any of the combinations of input variables resulted.
184 Including NDVI, EVI, LAI, and FPAR along with albedo improved the results to be similar to,
185 or slightly better than, the trial with those variables but without albedo.

186 The final set of trials included NBAR as an input and produced the lowest RMSEs of any
187 combination.

188

189 Based on the overall patterns in these results, further tuning of all algorithms was conducted
190 using all of the remote sensing input data variables: DSR, PAR, NDVI, EVI, LAI, FPAR, albedo,
191 and NBAR.

192

193 3. Tuning of ML algorithms

194 Each of the algorithms that ran sufficiently quickly to be iterated with each training data set was
195 optimized with that data set. The overall minimum RMSE results for the validation and test data
196 sets for all algorithms are shown in Table 5. Notable aspects of the tuning are then described
197 below for each “family” of algorithms.

198

199 Table 5: RMSE in W/m^2 for each optimized ML algorithm

200

Algorithm	Small training data set		Large training data set	
	Validation	Test	Validation	Test
RLR	30.55	29.84	31.22	30.23
LASSO	30.55	29.84	31.22	30.23
ELASTIC	30.55	29.84	31.22	30.23

KRR	23.85	23.41	Prohibitive	
RKS	25.35	25.52	22.22	22.10
TREE	29.19	28.71	25.14	25.45
BAGTREE	24.50	23.91	19.91	20.15
BOOST	28.86	28.33	28.21	27.86
NN, 1 HL	26.42	26.78	23.18	23.48
NN, 2 HL	25.76	25.20	21.58	22.69
NN, 3 HL	25.59	25.51	20.94	21.79
SVR	24.13	23.63	Prohibitive	

201 Note: Algorithms that are too computationally demanding for training with the large data set are
 202 labeled “Prohibitive”.

203 a. Linear regression variants

204 The linear regression variants demonstrated the weakest performance of all algorithm types. All
 205 linear regression variants (RLR, LASSO, and Elastic Net) show a pattern of optimum
 206 performance at zero or low regularization parameter values (<100 for RLR, <0.01 for LASSO),
 207 then worse performance or failure to converge as regularization parameters increase. Including
 208 the regularization parameters provided no advantage over a standard linear regression.

209

210 b. Kernel methods

211 When optimized, the kernel ridge regression performed better than any of the other algorithms
 212 with the small training data set, but it was too computationally demanding for use with the larger
 213 training data set. The RKS, which is in a sense a faster approximation of the KRR, did not
 214 perform as well with the smaller data set, but improved on that performance significantly with
 215 the larger training data set. Both KRR and RKS had more sensitivity to the kernel width
 216 parameter than to regularization, except for the RKS using a high number (> 1000) of random
 217 functions. Increasing the number of random functions usually produces better performance when
 218 optimized, but at the expense of more sensitivity to the other algorithm parameters.

219

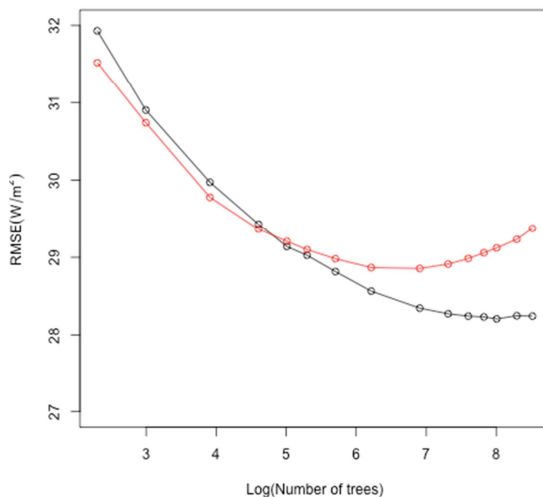
220 c. Tree methods

221 The simple tree method was not sensitive to degree of pruning or number of data points required
222 per partition. Therefore, these parameters were not adjusted in the trials with more complex tree
223 algorithms.

224

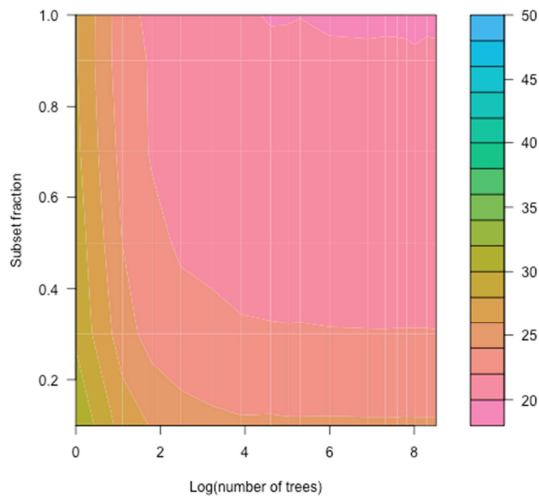
225 Performance of the boosting tree method improved with increasing number of trees up to about
226 500 trees, then saturated with the larger data set and showed evidence of overfitting with the
227 smaller data set with higher numbers of trees. (Figure 1). Boost tree algorithm performance was
228 generally weak overall. The bagging tree algorithm was the strongest performer out of all of the
229 algorithms with the large data set and shows improved performance with increasing number of
230 trees and fraction of input data used to construct each, although a saturation effect is evident
231 when the number of trees exceeds 100 (Figure 2).

232



233

234 Figure 1: Validation RMSE versus number of trees used in boost tree algorithm. Red: Small
235 training data set. Black: Large training data set.



236

237 Figure 2: Validation RMSE versus number of trees and fraction of data included in each bagging
 238 tree using large training data set

239

240 d. Neural networks

241 The most notable result of the neural network trials is that two and three hidden layer networks
 242 outperform the single-layer neural network, especially for the larger data set. Performance
 243 generally improves with number of neurons in each layer up to about 50 to 100 neurons in the
 244 first layer but is less sensitive to the number of neurons in the second or third layer if they are
 245 present. Some evidence of overfitting is also present in all neural network results, since RMSE
 246 with the test data set exceeds that with the validation data set, by about 1 W/m^2 in the case of the
 247 2 and 3 hidden layer NNs.

248

249 e. Support vector regression

250 The support vector regression method performed only modestly well with the smaller training
 251 data set, and tuning was computationally prohibitive with the larger training data set.

252

253 4. Trials with different ecosystem classes

254 The two and three hidden layer NN, RKS, and BAGTREE algorithms were used with the test
 255 data sets for each of seven ecosystem types. Initially, the algorithms were optimized using
 256 training and validation data from all sites. Then, each algorithm was tuned using the training and
 257 validation data sets for each ecosystem type, then tested using the test data set for the same type.
 258 Results with the full training data set are shown in Table 6, and results with the like-type-only
 259 training data sets are given in Table 7.

260

261 Table 6: RMSE, R^2 , and bias for different ecosystem types when ML algorithms are trained with
 262 data from all sites.

Agricultural	RMSE(W/m²)	R²	Bias(W/m²)
2 hidden layer NN	32.3557	0.6680	-0.6399
3 hidden layer NN	23.9863	0.8035	-1.0610
RKS	26.2128	0.7637	-1.1795
BAGTREE	17.8557	0.8950	-0.8671
Deciduous			
2 hidden layer NN	20.2389	0.7416	3.4893
3 hidden layer NN	18.9362	0.7741	2.0059
RKS	20.3107	0.7399	4.2604
BAGTREE	13.4676	0.8918	3.2371
Evergreen			
2 hidden layer NN	19.0075	0.6874	-0.5384

3 hidden layer NN	18.0615	0.7186	-0.4556
RKS	19.3665	0.6751	-0.0450
BAGTREE	12.4909	0.8745	-0.0113
Grassland			
2 hidden layer NN	18.4500	0.7853	-0.1038
3 hidden layer NN	17.9347	0.7981	-0.1971
RKS	18.7163	0.7791	-0.2180
BAGTREE	12.0934	0.9140	-0.0882
Savannah			
2 hidden layer NN	16.1591	0.8025	-0.8624
3 hidden layer NN	14.7330	0.8351	-0.5861
RKS	16.3950	0.8006	-1.3287
BAGTREE	16.3950	0.8006	-1.3287
Shrub			
2 hidden layer NN	33.8698	0.3823	-0.3441
3 hidden layer NN	16.6110	0.7777	-0.1439
RKS	18.0345	0.7395	-0.8539
BAGTREE	11.4718	0.9025	-0.6610
Wetland			
2 hidden layer NN	29.5516	0.8038	-1.3708
3 hidden layer NN	28.9601	0.8122	-1.8941
RKS	35.2762	0.7212	-4.4138

BAGTREE	22.9066	0.8878	-4.3023
---------	---------	--------	---------

263

264

265

266 Table 7: RMSE, R^2 , and bias for different ecosystem types with ML algorithms optimized with
 267 training and validation data from the same ecosystem type.

Agricultural	RMSE (W/m²)	R²	Bias (W/m²)
2 hidden layer NN	27.7259	0.7354	-0.3191
3 hidden layer NN	29.6665	0.6971	-1.3922
RKS	27.2730	0.7440	-0.5729
BAGTREE	18.5799	0.8851	-0.7525
Deciduous			
2 hidden layer NN	28.3526	0.4983	-2.0546
3 hidden layer NN	33.7258	0.2953	-15.1065
RKS	28.7253	0.5318	6.5990
BAGTREE	25.8699	0.5878	0.1247
Evergreen			
2 hidden layer NN	24.2847	0.4945	-7.8675
3 hidden layer NN	24.8439	0.4796	-6.3154
RKS	23.6216	0.5208	-6.5283
BAGTREE	23.7131	0.5169	-3.6958
Grassland			

2 hidden layer NN	19.6553	0.7563	-0.3862
3 hidden layer NN	32.7546	0.3608	-0.3290
RKS	19.7727	0.7533	0.0463
BAGTREE	14.6462	0.8689	0.5006
Savannah			
2 hidden layer NN	16.9383	0.7830	0.6077
3 hidden layer NN	15.1067	0.8267	-0.2368
RKS	15.3770	0.8203	-0.3738
BAGTREE	12.9098	0.8779	-0.3341
Shrub			
2 hidden layer NN	17.0708	0.7674	0.5069
3 hidden layer NN	17.9521	0.7403	0.7350
RKS	17.5338	0.7522	0.0286
BAGTREE	12.0077	0.8851	0.3511
Wetland			
2 hidden layer NN	29.4605	0.8055	-0.3557
3 hidden layer NN	29.4891	0.8055	-0.3742
RKS	48.4014	0.4812	-0.1796
BAGTREE	24.9083	0.8606	1.5604

268

269 The results shown in Tables 6 and 7 show that the ML algorithms performed best for evergreen,

270 grassland, and shrub sites. Performance was usually worst for wetland sites. The BAGTREE

271 algorithm was the best performer in most cases, except for the savannah sites when the

272 algorithms were trained with all data and the evergreen sites when training was done with data
273 from the same site type only. Training with data of the same type led to improved algorithm
274 performance only in the case of savannah sites. This is probably related to the fact that the
275 optimized algorithm parameters for the smaller individual site type training data sets reached less
276 complexity (fewer neurons in neural networks, fewer random functions in RKS, and fewer trees
277 used by BAGTREE) before overfitting became an issue than for the larger all site training data
278 set.

279

280 **Discussion**

281 Here we systematically compared several machine learning methods for obtaining LE from a
282 smaller or larger remote sensing only input data set. The best results for the small training data set
283 were with the kernel ridge regression (KRR), which was not viable with the large training set.
284 Three of the other algorithms (RKS, BAGTREE, and multi-layer neural networks) produced a
285 lower RMSE with the large training data set than the lowest RMSE attained with the small
286 training data set. The cloud-detection example given in Pérez-Suay et al. (2017) also
287 demonstrated this dynamic between the KRR and RKS methods. Here we also had good
288 performance with the RKS, but even better performance with the bagging tree and multi-layer
289 neural network.

290

291 Other than weaker performance by the linear regression variants, no family of
292 methods outperformed the rest. Regularization of the linear regression variants did not produce
293 any improvement to the algorithm results over a standard linear regression.

294

295 It has been shown that, while some of the ML algorithms perform well in terms of both accuracy
296 and computational demand, there is also some tradeoff between training efficiency and
297 performance. This is seen most clearly in the results with the large training data set, where the
298 BAGTREE algorithm produced the lowest RMSE but required more run time than the RKS,
299 boost tree, or smaller neural networks. The RKS algorithm is appealing due to its computational
300 efficiency and low test RMSE. Increasing the number of random functions in the RKS generally
301 reduces the optimized error, but also renders the algorithm more sensitive to its other parameters.
302 It is notable that the “deeper” 2- and 3- layer neural networks tested in this study performed
303 better than the single-layer neural network, since most studies in which neural networks are
304 applied to the LE problem only make use of single hidden layer neural networks. The multi-layer
305 neural networks only performed at their best if there were at least 50 neurons in the first layer but
306 showed less sensitivity to the numbers of neurons in deeper layers. The neural networks showed
307 more evidence of overfitting than any of the other algorithms, although the difference between
308 validation and test data set RMSEs was only about 1 W/m^2 .

309
310 Comparison of ML algorithm performance when trained with data from individual ecosystem
311 types instead of data from all sites usually showed worse performance, except for savannah sites.
312 This contrasts with the modest improvement found by Carter and Liang (2018) when non-ML
313 LE algorithms were tuned using data from individual ecosystem types. However, the poor
314 performance of the algorithms for wetland sites is consistent with Carter and Liang (2018). It
315 appears that the ability of the ML algorithms to extract more complex patterns from larger data
316 sets usually outweighs any advantage gained by restricting training data to one site type only.
317

318 **Conclusions**

319 A comparison of ten ML methods for obtaining LE from a combination of remote sensing data
320 (GLASS and MODIS) was performed in terms of accuracy and speed. The results showed wide
321 variation in algorithm efficiency. Including more input variables improved the results with little
322 or no additional computational cost. Use of GLASS radiation products produced results
323 comparable to using ground-based net radiation measurements. Inclusion of NBAR as one of the
324 parameters produced the best results.

325
326 The best performance with a smaller training data set was obtained using the kernel ridge
327 regression (KRR), which was too computationally demanding for use with the larger data set.
328 The best performance with the larger data set was achieved by the bootstrap aggregation tree
329 (BAGTREE) method, followed by the random kernel (RKS) and multiple hidden layer neural
330 network (NN) methods. The BAGTREE, neural network, and RKS algorithm performance could
331 be improved modestly for some ecosystem types by using training data from that ecosystem type
332 only.

333
334 Since the machine learning techniques evaluated here can be applied to any combination of input
335 variables, it should be possible to use them to generate global, long-term records of LE. The
336 GLASS data sets (Liang et al. 2013, 2014), which include albedo (Qu et al. 2014; Liu et al.
337 2013), leaf area index (Xiao et al. 2016, 2017a), and NDVI (Xiao et al. 2017b) in addition to
338 radiation variables, are based on the AVHRR and MODIS records, and therefore provide the
339 opportunity to examine global LE trends over decades.

340

341 Acknowledgements: This work was partially supported by funding from the United States
342 National Aeronautics and Space Administration (NASA) and National Oceanic and Atmospheric
343 Administration (NOAA).Funding for AmeriFlux data resources was provided by the U.S.
344 Department of Energy’s Office of Science.

345

346

347 References:

- 348 Carter, Corinne, and Shunlin Liang. 2018. “Comprehensive Evaluation of Empirical Algorithms
349 for Estimating Land Surface Evapotranspiration.” *Agricultural and Forest Meteorology*
350 256–257: 334–45.
- 351 Deo, Ravinesh C., Pijush Samui, and Dookie Kim. 2016. “Estimation of Monthly Evaporative
352 Loss Using Relevance Vector Machine, Extreme Learning Machine and Multivariate
353 Adaptive Regression Spline Models.” *Stochastic Environmental Research and Risk*
354 *Assessment* 30 (6): 1769–84. <https://doi.org/10.1007/s00477-015-1153-y>.
- 355 Didan, K. 2015. “MOD13A1 MODIS/Terra Vegetation Indices 16-Day L3 Global 500m SIN
356 Grid V006.” NASA EOSDIS Land Processes DAAC.
357 <https://doi.org/10.5067/MODIS/MOD13A1.006>.
- 358 Dou, Xianming, and Yongguo Yang. 2018. “Evapotranspiration Estimation Using Four Different
359 Machine Learning Approaches in Different Terrestrial Ecosystems.” *Computers and*
360 *Electronics in Agriculture* 148 (May): 95–106.
361 <https://doi.org/10.1016/j.compag.2018.03.010>.
- 362 Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical*
363 *Learning: Data Mining, Inference, and Prediction*. Springer Science + Business Media,
364 LLC.
- 365 Helman, D., A. Givati, and I. M. Lensky. 2015. “Annual Evapotranspiration Retrieved from
366 Satellite Vegetation Indices for the Eastern Mediterranean at 250 m Spatial Resolution.”
367 *Atmospheric Chemistry and Physics* 15 (21): 12567–79. [https://doi.org/10.5194/acp-15-](https://doi.org/10.5194/acp-15-12567-2015)
368 [12567-2015](https://doi.org/10.5194/acp-15-12567-2015).
- 369 Jung, Martin, Markus Reichstein, Philippe Ciais, Sonia I. Seneviratne, Justin Sheffield, Michael
370 L. Goulden, Gordon Bonan, et al. 2010. “Recent Decline in the Global Land
371 Evapotranspiration Trend Due to Limited Moisture Supply.” *Nature* 467 (7318): 951–54.
372 <https://doi.org/10.1038/nature09396>.
- 373 Kaheil, Y.H., E. Rosero, M.K. Gill, M. McKee, and L.A. Bastidas. 2008. “Downscaling and
374 Forecasting of Evapotranspiration Using a Synthetic Model of Wavelets and Support
375 Vector Machines.” *IEEE Transactions on Geoscience and Remote Sensing* 46 (9): 2692–
376 2707. <https://doi.org/10.1109/TGRS.2008.919819>.
- 377 Kalma, Jetse D., Tim R. McVicar, and Matthew F. McCabe. 2008. “Estimating Land Surface
378 Evaporation: A Review of Methods Using Remotely Sensed Surface Temperature Data.”
379 *Surveys in Geophysics* 29 (4–5): 421–69. <https://doi.org/10.1007/s10712-008-9037-z>.

380 Ke, Yinghai, Jungho Im, Seonyoung Park, and Huili Gong. 2016. “Downscaling of MODIS One
381 Kilometer Evapotranspiration Using Landsat-8 Data and Machine Learning Approaches.”
382 *Remote Sensing* 8 (3): 215. <https://doi.org/10.3390/rs8030215>.

383 ———. 2017. “Spatiotemporal Downscaling Approaches for Monitoring 8-Day 30 m Actual
384 Evapotranspiration.” *ISPRS Journal of Photogrammetry and Remote Sensing* 126 (April):
385 79–93. <https://doi.org/10.1016/j.isprsjprs.2017.02.006>.

386 Lazaro-Gredilla, Miguel, and Michalis K. Titsias. 2011. “Variation Heteroschedastic Gaussian
387 Process Regression.” *28th International Conference on Machine Learning (ICML 2011)*.

388 Lazaro-Gredilla, Miguel, Michalis K. Titsias, Jochem Verrelst, and Gustavo Camps-Valls. 2014.
389 “Retrieval of Biophysical Parameters With Heteroscedastic Gaussian Processes.” *IEEE*
390 *Geoscience and Remote Sensing Letters* 11 (4): 838–42.
391 <https://doi.org/10.1109/LGRS.2013.2279695>.

392 Liang, Shunlin, Xiaotong Zhang, Zhiqiang Xiao, Jie Cheng, Qiang Liu, and Xiang Zhao. 2014.
393 *Global LAnd Surface Satellite (GLASS) Products Algorithms, Validation and Analysis*.

394 Liang, Shunlin, Xiang Zhao, Suhong Liu, Wenping Yuan, Xiao Cheng, Zhiqiang Xiao, Xiaotong
395 Zhang, et al. 2013. “A Long-Term Global LAnd Surface Satellite (GLASS) Data-Set for
396 Environmental Studies.” *International Journal of Digital Earth* 6 (sup1): 5–33.
397 <https://doi.org/10.1080/17538947.2013.805262>.

398 Liu, Qiang, Lizhao Wang, Ying Qu, Nanfeng Liu, Suhong Liu, Hairong Tang, and Shunlin
399 Liang. 2013. “Preliminary Evaluation of the Long-Term GLASS Albedo Product.”
400 *International Journal of Digital Earth* 6 (sup1): 69–95.
401 <https://doi.org/10.1080/17538947.2013.804601>.

402 Murphy, Kevin P. 2012. *Machine Learning: A Probabilistic Perspective*. MIT Press.

403 Myneni, R., Y. Knyazikhin. 2015. “MOD15A2H MODIS/Terra Leaf Area Index/FPAR 8-Day
404 L4 Global 500m SIN Grid V006.” NASA EOSDIS Land Processes DAAC.
405 <https://doi.org/10.5067/MODIS/MOD15A2H.006>.

406 Park, Seonyoung, Jungho Im, Eunna Jang, and Jinyoung Rhee. 2016. “Drought Assessment and
407 Monitoring through Blending of Multi-Sensor Indices Using Machine Learning
408 Approaches for Different Climate Regions.” *Agricultural and Forest Meteorology* 216
409 (January): 157–69. <https://doi.org/10.1016/j.agrformet.2015.10.011>.

410 Pérez-Suay, Adrián, Julia Amorós-López, Luis Gómez-Chova, Valero Laparra, Jordi Muñoz-
411 Marí, and Gustau Camps-Valls. 2017. “Randomized Kernels for Large Scale Earth
412 Observation Applications.” *Remote Sensing of Environment* 202 (December): 54–63.
413 <https://doi.org/10.1016/j.rse.2017.02.009>.

414 Qu, Ying, Qiang Liu, Shunlin Liang, Lizhao Wang, Nanfeng Liu, and Suhong Liu. 2014.
415 “Direct-Estimation Algorithm for Mapping Daily Land-Surface Broadband Albedo From
416 MODIS Data.” *IEEE Transactions on Geoscience and Remote Sensing* 52 (2): 907–19.
417 <https://doi.org/10.1109/TGRS.2013.2245670>.

418 Rhee, Jinyoung, and Jungho Im. 2017. “Meteorological Drought Forecasting for Ungauged
419 Areas Based on Machine Learning: Using Long-Range Climate Forecast and Remote
420 Sensing Data.” *Agricultural and Forest Meteorology* 237–238 (May): 105–22.
421 <https://doi.org/10.1016/j.agrformet.2017.02.011>.

422 Schaaf, C. Z. Wang. 2015a. “MCD43A3 MODIS/Terra+Aqua BRDF/Albedo Daily L3 Global -
423 500m V006.” NASA EOSDIS Land Processes DAAC.
424 <https://doi.org/10.5067/MODIS/MCD43A3.006>.

- 425 ———. 2015b. “MCD43A4 MODIS/Terra+Aqua BRDF/Albedo Nadir BRDF Adjusted
426 RefDaily L3 Global - 500m V006.” NASA EOSDIS Land Processes DAAC.
427 <https://doi.org/10.5067/MODIS/MCD43A4.006>.
- 428 Smola, A. J., and B. Scholkopf. n.d. “A Tutorial on Support Vector Regression.” *Statistics and*
429 *Computing* 14 (3): 199–222.
- 430 Thayananthan, A, R Navaratnam, B. Stenger, P. H. S. Torr, and R Cipolla. 2006. “Multivariate
431 Relevance Vector Machines for Tracking.” Edited by A Leonardis and A Pinz. *Computer*
432 *Vision- ECCV 2006, Pt 3, Proceedings, Lecture Notes in Computer Science*, 3953: 124–
433 38.
- 434 Wang, Kaicun, and Robert E. Dickinson. 2012. “A Review of Global Terrestrial
435 Evapotranspiration: Observation, Modeling, Climatology, and Climatic Variability:
436 GLOBAL TERRESTRIAL EVAPOTRANSPIRATION.” *Reviews of Geophysics* 50 (2).
437 <https://doi.org/10.1029/2011RG000373>.
- 438 Wang, Kaicun, and Shunlin Liang. 2008. “An Improved Method for Estimating Global
439 Evapotranspiration Based on Satellite Determination of Surface Net Radiation,
440 Vegetation Index, Temperature, and Soil Moisture.” *Journal of Hydrometeorology* 9 (4):
441 712–27. <https://doi.org/10.1175/2007JHM911.1>.
- 442 Xiao, Zhiqiang, Shunlin Liang, and Bo Jiang. 2017. “Evaluation of Four Long Time-Series
443 Global Leaf Area Index Products.” *Agricultural and Forest Meteorology* 246
444 (November): 218–30. <https://doi.org/10.1016/j.agrformet.2017.06.016>.
- 445 Xiao, Zhiqiang, Shunlin Liang, Xiaodan Tian, Kun Jia, Yunjun Yao, and Bo Jiang. 2017.
446 “Reconstruction of Long-Term Temporally Continuous NDVI and Surface Reflectance
447 From AVHRR Data.” *IEEE Journal of Selected Topics in Applied Earth Observations*
448 *and Remote Sensing* 10 (12): 5551–68. <https://doi.org/10.1109/JSTARS.2017.2744979>.
- 449 Xiao, Zhiqiang, Shunlin Liang, Jindi Wang, Yang Xiang, Xiang Zhao, and Jinling Song. 2016.
450 “Long-Time-Series Global Land Surface Satellite Leaf Area Index Product Derived From
451 MODIS and AVHRR Surface Reflectance.” *IEEE Transactions on Geoscience and*
452 *Remote Sensing* 54 (9): 5301–18. <https://doi.org/10.1109/TGRS.2016.2560522>.
- 453 Yao, Yunjun, Shunlin Liang, Jie Cheng, Shaomin Liu, Joshua B. Fisher, Xudong Zhang, Kun
454 Jia, et al. 2013. “MODIS-Driven Estimation of Terrestrial Latent Heat Flux in China
455 Based on a Modified Priestley-Taylor Algorithm.” *Agricultural and Forest Meteorology*
456 171–172 (April): 187–202. <https://doi.org/10.1016/j.agrformet.2012.11.016>.
- 457 Yao, Yunjun, Shunlin Liang, Xianglan Li, Jiquan Chen, Kaicun Wang, Kun Jia, Jie Cheng, et al.
458 2015. “A Satellite-Based Hybrid Algorithm to Determine the Priestley-Taylor Parameter
459 for Global Terrestrial Latent Heat Flux Estimation across Multiple Biomes.” *Remote*
460 *Sensing of Environment* 165 (August): 216–33. <https://doi.org/10.1016/j.rse.2015.05.013>.
- 461 Yao, Yunjun, Shunlin Liang, Qiming Qin, Kaicun Wang, and Shaohua Zhao. 2011. “Monitoring
462 Global Land Surface Drought Based on a Hybrid Evapotranspiration Model.”
463 *International Journal of Applied Earth Observation and Geoinformation* 13 (3): 447–57.
464 <https://doi.org/10.1016/j.jag.2010.09.009>.
- 465 Yebra, Marta, Albert Van Dijk, Ray Leuning, Alfredo Huete, and Juan Pablo Guerschman. 2013.
466 “Evaluation of Optical Remote Sensing to Estimate Actual Evapotranspiration and
467 Canopy Conductance.” *Remote Sensing of Environment* 129 (February): 250–61.
468 <https://doi.org/10.1016/j.rse.2012.11.004>.
- 469 Zhang, Ke, John S. Kimball, and Steven W. Running. 2016. “A Review of Remote Sensing
470 Based Actual Evapotranspiration Estimation: A Review of Remote Sensing

471 Evapotranspiration.” *Wiley Interdisciplinary Reviews: Water* 3 (6): 834–53.
472 <https://doi.org/10.1002/wat2.1168>.
473

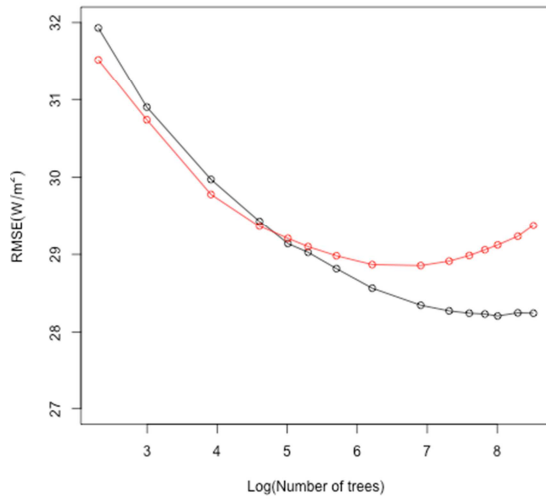


Figure 1: Validation RMSE versus number of trees used in boost tree algorithm. Red: Small training data set. Black: Large training data set.

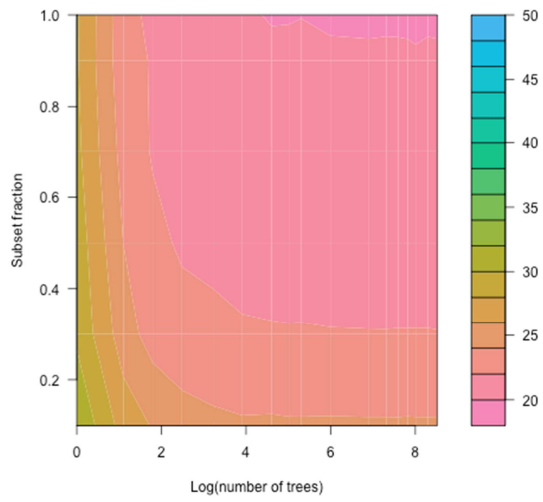


Figure 2: Validation RMSE versus number of trees and fraction of data included in each bagging tree using large training data set