



Daily soil moisture mapping at 1 km resolution based on SMAP data for areas affected by desertification in Northern China

Pinzeng Rao^{1,2}, Yicheng Wang², Fang Wang^{2*}, Yang Liu², Xiaoya Wang³, Zhu Wang²

¹State Key Laboratory of Hydrosience and Engineering, Department of Hydraulic Engineering, Tsinghua University, Beijing 100084, China.

²State Key Laboratory of Simulation and Regulation of Water Cycle in River Basin, China Institute of Water Resources and Hydropower Research, Beijing 100038, China.

³State Key Laboratory of Remote Sensing Science, Faculty of Geographical Science, Beijing Normal University, Beijing 100875, China.

*Correspondence to: Fang Wang (657563390@qq.com)

Abstract: Land surface soil moisture (SM) plays a critical role in hydrological processes and terrestrial ecosystems in areas affected by desertification. Passive microwave remote sensing products such as the Soil Moisture Active Passive (SMAP) have been shown to monitor surface soil water well. However, the coarse spatial resolution and lack of full coverage of these products greatly limit their application in areas undergoing desertification. In order to overcome these limitations, a combination of multiple machine learning methods, including multiple linear regression (MLR), support vector regression (SVR), artificial neural networks (ANN), random forest (RF) and extreme gradient boosting (XGB), have been applied to downscale the 36 km SMAP SM products and produce higher spatial-resolution SM data based on related surface variables, such as vegetation index and surface temperature. Areas affected by desertification in Northern China, which are very sensitive to SM, were selected as the study area, and the downscaled SM with a resolution of 1 km on a daily scale from 2015 to 2020 was produced. The results show a good performance compared with in situ observed SM data, with an average unbiased root mean square error value of 0.049 m³/m³. In addition, their time series are also consistent with precipitation and perform better than some common gridded SM products. The data can be used to assess soil drought and provide a reference for reversing desertification in the study area. This dataset is freely available at <https://doi.org/10.6084/M9.FIGSHARE.16430478.V5> (Rao et al., 2021).

Keywords: Soil moisture; SMAP; Multiple machine learning; Surface variables; Desertification.

1. Introduction

Surface soil moisture (SM) plays a very important role in water-energy cycle processes (Sandholt et al., 2002; De Santis et al., 2021) and is an important source of water for plants and soil microbes (Wang et al., 2007; Gu et al., 2008; Mallick et al., 2009). Large-scale areas of northern China are undergoing desertification because of scarce precipitation and insufficient SM. The accurate acquisition of SM is valuable to ecological conservation and revegetation in arid areas of Northern China.

In the past, SM data were mainly obtained through ground measurements or the assimilation of products based on land surface models such as the Global Land Data Assimilation System (GLDAS). Although most accurate SM data at different



33 soil depths can be obtained, field measurements and in situ observations are limited due to the high cost and labor intensity
 34 involved in their collection and are generally not representative of soil water status over larger areas (Rahimzadeh-Bajgiran et
 35 al., 2013; Zhao et al., 2018; Bai et al., 2019). With the development of remote sensing technologies, continuous SM estimates
 36 can be generated at regional and global scales (Peng et al., 2021). Compared to ground measurements, remote sensing products
 37 can provide good spatial and temporal coverage of SM with a relatively low cost to the user (Zeng et al., 2015; Zhao et al.,
 38 2018; Meng et al., 2020). Data assimilation products largely depend on the accuracy of the land surface model and the original
 39 data (Zawadzki and Kędzior, 2016). They generally have low accuracy in areas where ground measurements are scarce, which
 40 is a problem that can be overcome with remote sensing.

41 At present, there are many remotely sensed SM data, some of which are from microwave remote sensing satellites,
 42 including active and passive types. SM retrievals from active sensors like Synthetic Aperture Radar (SAR) are sensitive to
 43 scattering and greatly affected by the surface roughness and vegetation types (Lievens et al., 2011; Wagner et al., 2013). Unlike
 44 active sensors, passive microwave radiometers or sensors have almost no scattering and generate very stable SM products
 45 (Abbaszadeh et al., 2019). Common passive microwave SM products are listed in Table 1 below. Some studies have compared
 46 these products and found that SMAP SM products have higher accuracy and robustness than other remotely sensed SM
 47 products (Liu et al., 2019; Wang et al., 2021).

48 **Table 1: Information of five common passive microwave soil moisture (SM) products.**

SM Datasets (Abbreviation)	Name	Production source	Resolution	Temporal Coverage	Equator Crossing Time
AMSR-E/Aqua Daily L3	Advanced Microwave Scanning Radiometer- Earth Observing System	National Aeronautics and Space Administration (NASA) National Snow and Ice Data Center Distributed Active Archive Center (NSIDC)	25 km; Daily	2002- 2011	1:30 PM Ascending 1:30 AM Descending 6:00 PM
SMOS	Soil Moisture and Ocean Salinity	European Space Agency (ESA)	25 km; Daily	2010- present	Ascending 6:00 AM Descending 1:40 PM
FY3B	Fengyun-3B	National Satellite Meteorological Center	25 km; Daily	2011- present	Ascending 1:40 AM Descending 1:30 PM
GCOM- W1/AMSR2	Advanced Microwave Scanning Radiometer 2	Japan Aerospace Exploration Agency (JAXA)	0.25°/0.1°; Daily	2012- present	Ascending 1:30 AM Descending 6:00 PM
SMAP	Soil Moisture Active Passive	National Aeronautics and Space Administration (NASA)	36 km; Daily	2015- present	Ascending 6:00 AM Descending

49 Passive microwave SM products have been applied at watershed and national scale (Fang and Lakshmi, 2014; Meng et
 50 al., 2020). However, due to their coarse spatial resolution, microwave SM products have limited applicability to small-scale
 51 areas. Compared to microwave sensors, optical satellites such as MODIS and Landsat have a finer spatial resolution. Some
 52 observations generated from optical satellites provide good information about SM, such as vegetation index (VI) and land



53 surface temperature (LST) (Wang et al., 2007; Sun et al., 2012). Many experiments have tried to use these two parameters
54 from optical remote sensing to retrieve surface SM (Mallick et al., 2009; Fang et al., 2013). Based on the LST and VI triangle
55 space, Sandholt et al. (2002) proposed the temperature vegetation dryness index (TVDI) and used it to assess the SM status.
56 Despite their higher resolution, however, optical remote sensing data do not allow to directly retrieve true SM.

57 Some studies have tried to use surface variables from optical observations to improve the spatial resolution of passive
58 remotely sensed SM products (Peng et al., 2017). Zhao et al. (2017) used the triangle method and Landsat satellite observations
59 to disaggregate coarse-resolution SM data. Some studies have also shown that polynomial regression is effective in SM and
60 optical observations (Zhao and Li, 2013; Piles et al., 2016). However, these methods have some shortcomings in representing
61 the nonlinear relationship between SM and other surface variables (Zhao et al., 2018; Hu et al., 2020). Machine learning
62 methods can be applied to show the nonlinear relationships between SM and surface variables. Random forest (RF) and
63 artificial neural network (ANN) have been widely used in previous studies due to their high generalization ability and
64 robustness (Yao et al., 2017; Liu et al., 2020; Demarchi et al., 2020; Chen et al., 2021). Chen et al. (2021) developed the global
65 surface SM dataset covering 2003–2018 at 0.1° resolution with neural networks and some related variables. Im et al. (2016)
66 used machine learning approaches (RF, boosted regression trees, and Cubist) to downscale AMSR-E SM data in South Korea
67 and Australia and found RF to be superior to the other downscaling methods. Although these machine learning methods
68 perform well in constructing nonlinear regression models, there are still some shortcomings. For example, neural networks are
69 prone to overfitting when there are inefficient samples (Piotrowski and Napiorkowski, 2013) or variables that are weakly
70 correlated with the dependent variable (Elshorbagy and Parasuraman, 2008; Ågren et al., 2021). Extreme gradient boosting
71 (XGB), as a new ensemble learning method (Chen and Guestrin, 2016), performs well in some fields (Wang et al., 2020; Fan
72 et al., 2021; Ma et al., 2021), but it has rarely been used for soil moisture downscaling.

73 The selection of feature variables is critical for regression models. In addition to LST and VI mentioned above, variables
74 such as terrain and soil conditions also have a significant impact on SM. Abbaszadeh et al. (2019) downscaled SMAP
75 radiometer SM products over the continental United States using MODIS products (including NDVI and LST), precipitation
76 and topographic data, and also evaluated the influence of soil texture on SM. Zhao et al. (2018) added additional surface
77 variables, such as LST leaf area index (LAI), normalized difference water index (NDWI), surface albedo and the solar zenith
78 angle. Hu et al. (2020) added the normalized shortwave-infrared difference bare soil moisture index (NSDSI), horizontally
79 polarized Brightness Temperature (TBh) and vertically polarized Brightness Temperature (TBv) to the regression model. In
80 general, all these variables can be classified into vegetation, temperature, soil wetness, topography, and soil factors and sensors
81 conditions.

82 In recent years, the Chinese government has carried out afforestation activities in order to reverse desertification in the
83 North. Considering the role of SM in the ecological environment, it is urgent to obtain accurate SM with high temporal and
84 spatial resolution. This study aims to downscale SMAP SM products by constructing a nonlinear relationship between SM and
85 related surface variables by means of multiple machine learning methods and generate SM products with higher temporal and
86 spatial resolution in areas affected by desertification. The in situ observed SM data from the Maqu Monitoring Network and

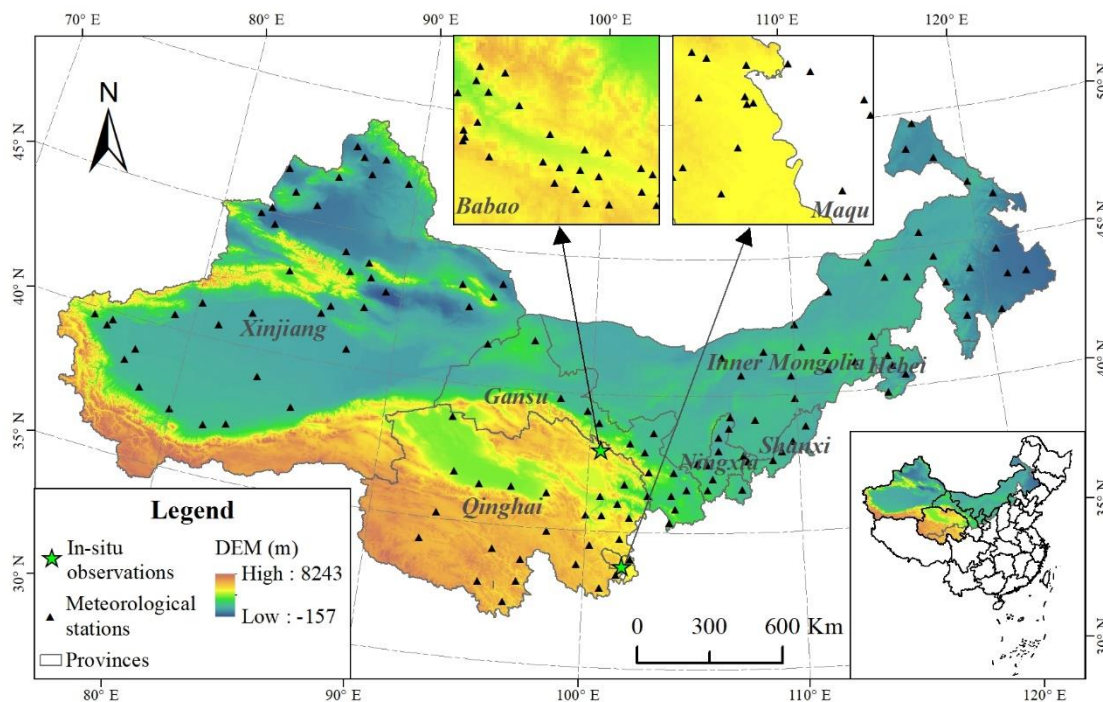


87 Babao Monitoring Network and precipitation and temperature data from 131 meteorological stations were used for validation
88 and analysis.

89 2. Materials and methodology

90 2.1 Study area

91 Northern China is mostly arid with an annual precipitation of less than 400 mm. The region belongs to the temperate
92 continental monsoon climate and is subject to large-scale desertification. The desert areas of Northern China are susceptible
93 to climate and hydrological changes and have fragile ecosystems. Soil water is a key parameter of the water-vapor-ecosystem,
94 and its change greatly affects the survival of vegetation and agricultural production in areas affected by desertification. The
95 studied area used for this study covers 3.36 million km², encompassing seven provinces. The terrain is complex, and the
96 average elevation is approximately 1900 m, ranging from -192 m to 7439 m.



97
98

Figure 1: Location of the study area.

99 2.2 Observations for the production of soil moisture data

100 2.2.1 SMAP SM data

101 The SMAP satellite was launched on January 31, 2015. Its mission consists of an L-band radar and radiometer instrument
102 suite, which provides global measurements and monitoring of SM in the top 5 cm of soil. The Level-3 products are daily



103 composites of the Level-2 products and are the most commonly used for applications. The Level-3 products are available in
104 three spatial resolutions: 36 km passive, 9 km active-passive, and 3 km active (O'Neill et al., 2010). Following the
105 malfunctioning of its radar in 2015, SMAP radar data were replaced with those of Sentinel-1, limiting the application of active
106 and active-passive products.

107 The SMAP Level-3 passive daily SM product (L3_SM_P, Version 6) with a grid resolution of 36 km has been produced
108 since March 31, 2015. Zeng et al. (2015) showed that most of remotely sensed SM products were slightly better during daytime
109 than during nighttime, and the same conclusion for the SMAP SM product was confirmed by Zhao et al. (2018). Therefore,
110 the SMAP Level-3 SM product with the descending overpass time of 6:00 AM was used in this study. The data were
111 downloaded from NASA Earthdata (<https://search.earthdata.nasa.gov>).

112 2.2.2 MODIS products

113 MODIS provides continuous time-series predictors for important parameters, such as vegetation index and surface
114 temperature. This paper used MODIS products MOD09A1, MOD11A1, MOD13A2, MOD15A2H and MCD43D58 (Table 2).
115 The soil wetness related indexes, including NDWI, LSWI and NSDSI, were produced using bands of the MOD09A1 product.
116 Their formulas are:

$$117 \quad NDWI = (B_4 - B_2)/(B_4 + B_2) \quad (1)$$

$$118 \quad LSWI = (B_2 - B_6)/(B_2 + B_6) \quad (2)$$

$$119 \quad NSDSI = (B_6 - B_7)/B_6 \quad (3)$$

120 where B_2 , B_4 , B_6 and B_7 represent the MOD09A1 surface reflectance of the 2nd, 4th, 6th and 7th bands, respectively.

121 These MODIS products are available from NASA Earthdata (<https://search.earthdata.nasa.gov>), and all data were
122 obtained from 2015 to 2020.

123 2.2.3 Topographic data

124 Topographic factors are strongly related to SM, including elevation, slope and aspect. The Shuttle Radar Topography
125 Mission (SRTM) digital elevation model (DEM) was used as elevation. Slope and aspect can be generated based on the DEM.
126 These data were obtained from the Geospatial Data Cloud (<http://www.gscloud.cn/>), where slope and aspect have been
127 processed and provided directly.

128 2.2.4 Soil texture data

129 Soil texture, the proportions of sand, silt and clay particles, controls the water holding capacity of the soil. The soil data
130 used for this study used the China Soil Characteristics Dataset (CSCD) (Shangguan et al., 2012), obtained from National
131 Tibetan Plateau Data Center (<http://westdc.westgis.ac.cn/>).

132 2.2.5 In Situ SM observations



133 The in situ SM measurements were collected from the data provided by the Maqu Monitoring Network (Zhang et al.,
 134 2020) and the Babao Monitoring Network (Kang et al., 2017). The Maqu Monitoring Network covers 26 sites and provides
 135 SM for the surface layer (0-5 cm) at 15-minute intervals from 2009 to 2019; 19 of the available sites which have data after
 136 2015 were used in this study (Fig. 1). The Babao Monitoring Network covers 40 sites and provides hourly SM for the surface
 137 layer (4 cm, 10 cm and 20 cm) from 2013 to 2017; 29 of the available sites have data after 2015 and were used in this study
 138 (Fig. 1). To compare with the simulated results, they were all processed into daily time series.

139 2.2.6 Precipitation data

140 The precipitation data were acquired from 131 meteorological stations from the China Meteorological Data Service
 141 Centre (<http://data.cma.cn>). The spatial locations of these meteorological stations are shown in Fig. 1.

142 **Table 2: Main predictors used in the study and corresponding datasets**

Datasets	Predictors	Spatial resolution	Temporal resolution	Number of granules (Years×tiles)
SMAP	SM	~36 km	Daily	2064
MOD11A1	LST	1 km	Daily	17460
MOD13A2	NDVI; EVI	1 km	16-day	1104
MOD15A2H	LAI; FAPAR	500 m	8-day	2208
MOD09A1	NDWI; LSWI; NSDSI	500 m	8-day	2208
MCD43D58	Albedo	30 Arcsec	Daily	2192
SRTM	DEM; Slope; Aspect	90 m	-	32
CSCD	Sand; Silt; Clay	1 km	-	1

143 2.2.7 Other gridded SM datasets

144 Some other gridded SM data were used to compare the simulation results (Table 3). The SMAP Level-2 product
 145 (L2_SM_SP) merges SMAP radiometer and processed Sentinel-1A/1B SAR observations. It is available at 3 km and 1 km
 146 resolutions. The Global Change Observation Mission Water (GCOM-W1) AMSR2 product is produced by the Japan
 147 Aerospace Exploration Agency (JAXA), and SM data at a 0.1° spatial resolution were selected for this study. The Copernicus
 148 Climate Change Service (C3S) produces a global SM gridded dataset from 1978 to present from satellite sensors such as SMOS,
 149 AMSR2 and SMAP. It has a spatial resolution of 0.25 degrees and offers three types of products: active, passive and combined.
 150 The combined product that we used in this study is generated by merging the active and passive products. The fifth generation
 151 ECMWF reanalysis dataset (ERA5) provides several variables including volumetric soil water over several decades. In the
 152 dataset, the soil is divided into four layers and the depth of the top layer is 0-7 cm. In this study, we downloaded the hourly
 153 volumetric soil water data of the top layer and processed them as daily averages. Famine Early Warning Systems Network
 154 (FEWS NET) Land Data Assimilation System (FLDAS) provides daily SM at a 0.01° spatial resolution over the Central Asia
 155 region (30-100° E, 21-56° N), which covers part of our study area. The product consists of four layers of SM, and the SM at
 156 the top layer (0-10 cm) was selected for this study.

157 **Table 3: The gridded SM products used in this study**

Institu tion	Name	Soil layers	TYPES	Temporal resolution	Grid spacing	Data link
-----------------	------	-------------	-------	------------------------	-----------------	-----------



NASA	SMAP/ (L2_SM_SP)	Sentinel-1	One layer (0-5 cm)	Active microwave	1-2 days	1/3 km	https://cmr.earthdata.nasa.gov/search/concepts/C1931663473-NSIDC_ECS.html
JAXA	GCOM-W1/AMSR2		One layer (~)	Passive microwave	Daily	0.1°/0.25°	https://gportal.jaxa.jp/gpr/
ECM WF	C3S		One layer (~)	Passive, active and combined	Daily	0.25°	https://cds.climate.copernicus.eu/cdsapp#!/dataset/satellite-soil-moisture
ECM WF	ERA5		Four layers (0-7 cm, 7-28 cm, 28-100 cm, 100-289 cm)	Reanalysis	Hourly	0.1°	https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-land
NASA	FLDAS		Four layers (0-10 cm, 10-40 cm, 40-100 cm, 100-200 cm)	Reanalysis	Daily	0.01°	https://cmr.earthdata.nasa.gov/search/concepts/C2020764153-GES_DISC.html

158 2.3 Downscaling approach based on multi-machine learning

159 According to the selected variable indicators (mainly including topographic data, soil data and some MODIS products)
 160 and machine learning methods, we constructed a framework to downscale SMAP SM based on multiple machine learning
 161 methods (Fig. 2).

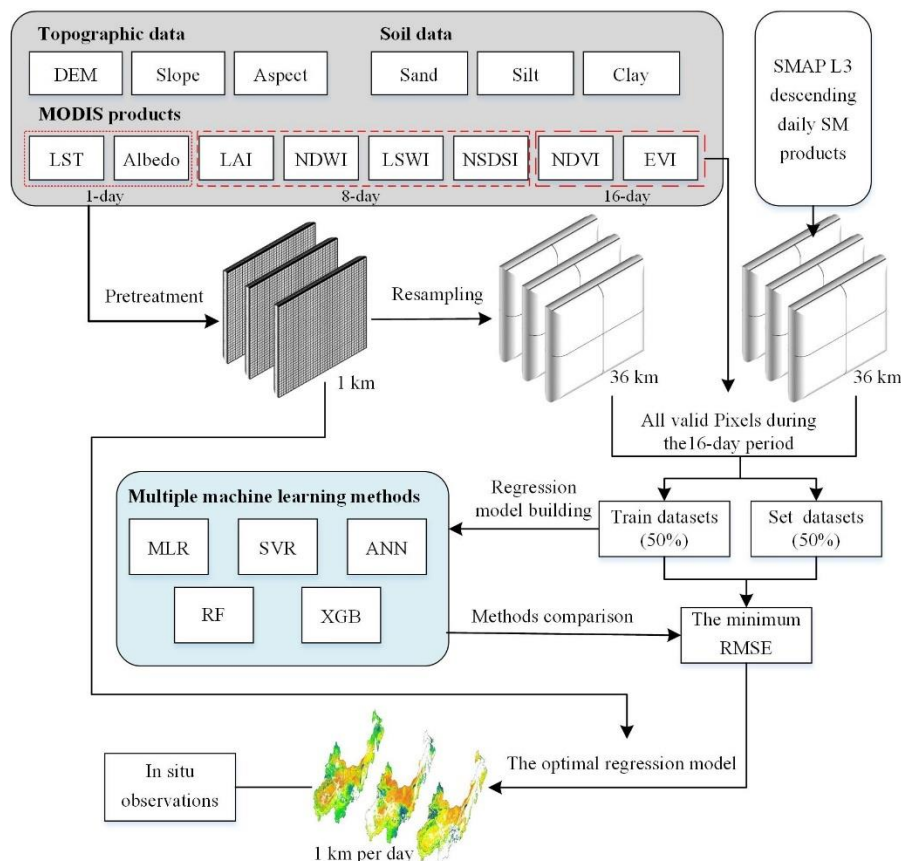


Figure 2: Schematic of the SMAP soil moisture downscaling framework

162
163

164 2.3.1 Machine learning methods

165 Machine learning methods are widely used in regression and classification. We selected machine learning methods that
 166 are currently widely used to build regression models for SM and its related variables. We studied five methods: Multiple linear
 167 regression (MLR), support vector regression (SVR), artificial neural networks (ANN), random forest (RF) and extreme
 168 gradient boosting (XGB). MLR and SVR have been widely used as regression methods in the past (Yu et al., 2012; Achieng,
 169 2019; Wang et al., 2019). ANN is currently one of the most popular machine learning methods and is used in many fields,
 170 including remote sensing of soil moisture inversion (Del Frate et al., 2003; Elshorbagy and Parasuraman, 2008; Yao et al.,
 171 2017; Chen et al., 2021).

172 RF and XGB are tree based ensemble algorithms, which have prediction accuracy and good generalization ability, and
 173 are not prone to overfitting (Rao et al., 2018; Abbaszadeh et al., 2019). RF is a multiple-tree algorithm improved by Bootstrap
 174 to reduce decision tree bias in determining the splits (Mohana et al., 2021). Many studies have used RF to build regression
 175 models of remotely sensed SM and related variables, and almost all achieved better results compared to other regression
 176 methods (Zhao et al., 2018; Qu et al., 2019; Hu et al., 2020). In contrast, the application of XGB, which applies a regularized



177 gradient boosting framework, is still very limited. However, XGB has incomparable advantages in generalization performance
178 and accuracy (Wang et al., 2020). Compared with RF and other some methods, XGB has significantly faster calculation speed
179 (Fan et al., 2018; Shi et al., 2021). Some studies have shown that XGB is a better regression and classification algorithm than
180 RF and other machine learning methods (Ågren et al., 2021; Fan et al., 2021).

181 2.3.2 Downscaling process

182 The downscaling process is shown in Fig. 2. First, due to the difference in spatial resolution and data format, all required
183 data were pre processed. All selected variables, including LST, Albedo, LAI, NDWI, LSWI, NSDSI, NDVI, EVI, DEM, slope,
184 aspect, sand, silt and clay, were aggregated into a resolution of 1 km with a geotiff format. These variables were further
185 resampled to the spatial resolution of the SMAP SM data (36 km) using the nearest neighbor interpolation method. The
186 regression model was then defined according to the selected machine learning method:

$$187 \quad SM = f(LST, Albedo, LAI, NDWI, LSWI, NSDSI, NDVI, EVI, \\ 188 \quad \quad \quad DEM, slope, aspect, sand, silt \text{ and } clay) \quad (4)$$

189 where f represents the regression function of the machine learning method (MLR, SVR, ANN, RF or XGB).

190 Then, the regression model based on multiple machine learning was built. The MODIS products and SMAP SM data
191 have different temporal resolutions. Since it is severely affected by noise (such as clouds), MOD11A1 only provides daily
192 valid clear-sky LST values onto grids. The variables from MOD13A2 and MOD15A2H are the best composite within 16 days
193 and 8 days, respectively. In addition, because each SMAP image has a narrow coverage, there may be few or no valid samples
194 if only the data of a certain day are selected to build the regression. To overcome the limitation, we chose to build regression
195 models within 16 day periods (the lowest temporal resolution from these dynamic variables). All valid data (including training
196 and test datasets) within 16 days were used as the samples in the regression model. For instance, for NDVI and EVI on January
197 1, 2020, which are composite results from January 1 to January 15, the valid data during the period were used as samples. The
198 number of valid samples for surface variables and SMAP SM for each period in 2015-2020 is shown in Fig. 3. Since there are
199 fewer available SMAP SM grid data in the cold season, there may be few valid samples we can obtain for the period.
200 Considering the number of samples is critical to the accuracy of the regression model, we only selected periods with more than
201 100 samples to build the model and DOY of 2016017, 2018017, 2018353, 2019001 and 2019177 were excluded. The valid
202 samples were divided into a training set and a test set, each accounting for 50% of the total number of samples. Then, we used
203 these samples and multiple machine learning methods (MLR, SVR, ANN, RF and XGB) to build a regression model for each
204 16-day period.

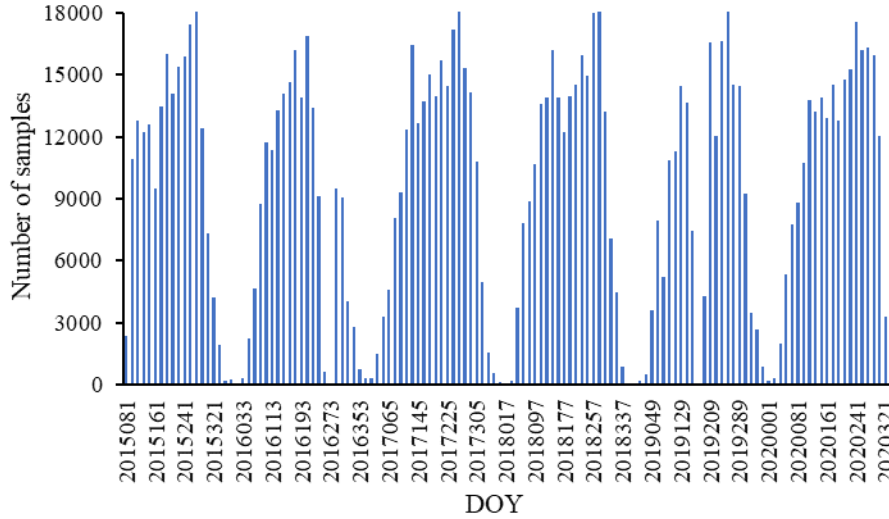


Figure 3: The number of valid samples for a 16-day period in 2015-2020

2.3.3 Evaluation method

The correlation coefficient (R) and the root mean square error (RMSE) were used to evaluate the accuracy of the regression model based on these machine learning methods (MLR, SVR, ANN, RF and XGB). They are calculated as:

$$R = \frac{Cov(SM_I, SM_P)}{\sqrt{Var(SM_I)Var(SM_P)}} \quad (5)$$

$$RMSE = \sqrt{\frac{1}{n}(SM_P - SM_I)^2} \quad (6)$$

where SM_I is the SMAP SM, SM_P is the corresponding SM predicted by the regression model, Cov represents the covariance function, Var is the variance, and n is the number of valid samples for SM_I or SM_P .

The regression model with the smallest average RMSE of training and test datasets was selected as the optimal model.

$$\overline{RMSE} = \frac{RMSE_{Training} + RMSE_{Test}}{2} \quad (7)$$

where $RMSE_{Training}$ and $RMSE_{Test}$ are the RMSE of the training set and test set for these models, respectively.

We used the selected optimal model with these surface variables with a resolution of 1 km within 16 days to simulate SM at 1 km resolution on the corresponding date. Taking 16 days as a period, we predicted all daily SM data with a spatial resolution of 1 km from 2015 to 2020. In addition, to obtain a more complete time series of SM data, we used the model of the previous period when the number of valid samples was less than 100.

The in situ SM measurements were used to validate the downscaled results. In addition to R and RMSE, unbiased RMSE (ubRMSE) and bias were also calculated according to:

$$ubRMSE = \sqrt{\frac{1}{n}((SM_{In} - \overline{SM_{In}}) - (SM_D - \overline{SM_D}))^2} \quad (8)$$



224
$$bias = \overline{SM_{In}} - \overline{SM_D} \quad (9)$$

 225 where SM_{In} is the in situ observed SM, SM_d is the downscaled SM of the corresponding grid, and n is the number of valid
 226 samples for SM_{In} or SM_D .

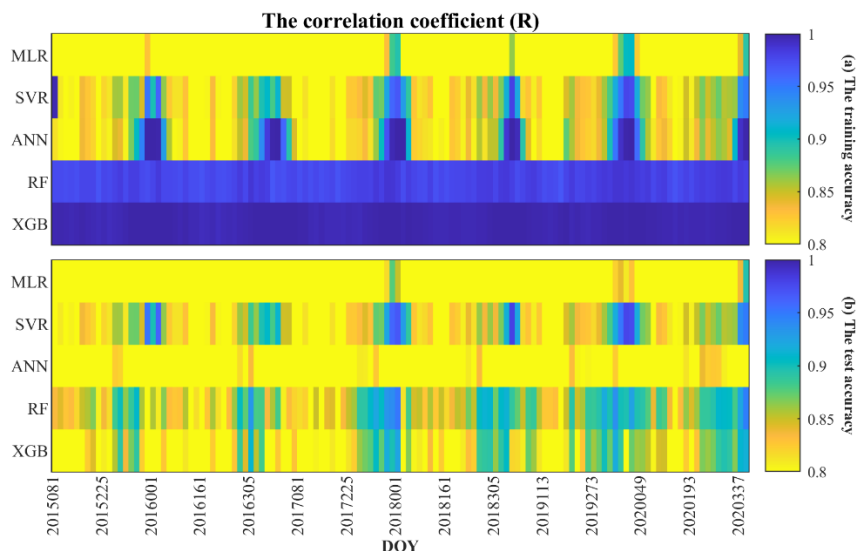
227 3. Results

228 3.1 Model comparison

229 The daily SM from DOY 81 in 2015 to DOY 366 in 2020 were simulated producing 128 regression results every 16 days.
 230 According to Equation 3, among the 128 regression results, there were 123 from the XGB model, and 5 from RF (including
 231 DOY 2015241, 2016161, 2016209, 2017241 and 2017257).

232 The correlation coefficient (R) and the root mean square error (RMSE) of each regression result for the training set and
 233 the test set are shown in Fig. 4 and Fig. 5, respectively. For all models, R is greater than 0.8 and RMSE is less than 0.1 both
 234 for the training and the test set. For the training set using XGB, Rs are all above 0.96, generally higher than for other methods;
 235 Similarly, the RMSEs of XGB are all lower than 0.02, generally lower than those of other methods. The R of RF is second
 236 only to that of XGB, and for some periods it is higher than for XGB; the RMSEs of RF are also generally lower than 0.02 and
 237 are lower than those of XGB in some periods. SVR and ANN perform better in the cold season, and worse in other seasons.
 238 In general, their results are inferior to those of XGB and RF. The simulation results of MLR are relatively poor both in terms
 239 of RMSE and R.

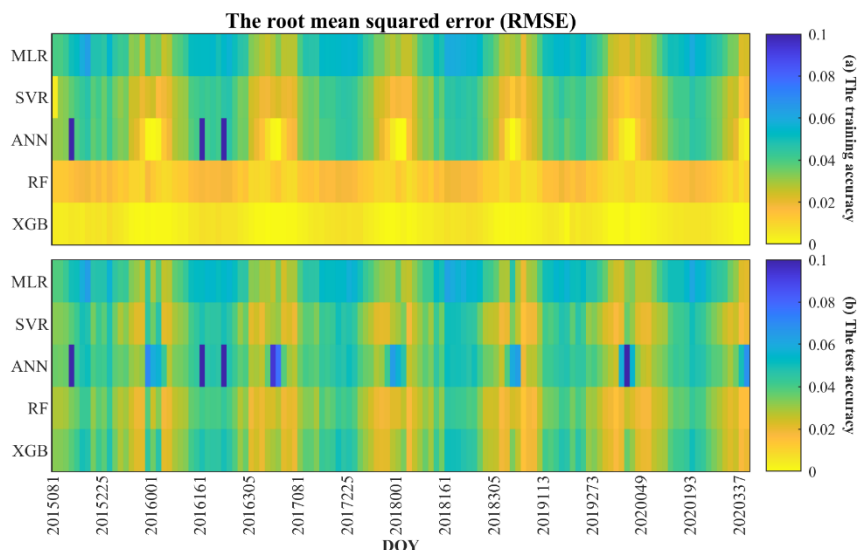
240 The results of the test set show that XGB, RF and SVR perform better than ANN and MLR, and are better in the cold
 241 season. Table 4 shows the average RMSE and R values of the training and test sets over all periods, and the performance order
 242 of the model can be obtained as XGB>RF>SVR >ANN >MLR.



243



244 **Figure 4: The correlation coefficient (R) of the models (MLR, SVR, ANN, RF and XGB) on different periods: (a) The training**
 245 **accuracy; (b) The test accuracy.**



246 **Figure 5: The root mean square error (RMSE) of the models (MLR, SVR, ANN, RF and XGB) for different periods: (a) The training**
 247 **accuracy; (b) The test accuracy.**

248 **Table 4: Accuracy of the models based on correlation coefficient (R) and root mean square error (RMSE)**

Model		MLR	SVR	ANN	RF	XGB
Training set	RMSE	0.688	0.843	0.864	0.979	0.985
	R	0.042	0.032	0.028	0.013	0.010
Test set	RMSE	0.677	0.843	0.660	0.857	0.861
	R	0.043	0.029	0.047	0.030	0.029

250 3.2 Comparison with the in situ data and precipitation

251 To evaluate the performance of the downscaling approach, the downscaled 1 km gridded SM were compared with the in
 252 situ SM observations. The SM before and after downscaling were both compared with the in situ SM data of the Maqu Network
 253 and Babao Network (Fig. 6). Due to the difference in sensors, soil depth and measurement scale (point observation in case of
 254 the in situ measured SM and 1 km grid for the downscaled SM), there is a certain deviation between in situ observation data
 255 and the downscaled gridded SM data. The downscaled SM of most sites at the Maqu Network are highly correlated with the
 256 in situ measured SM ($R > 0.6$). The ubRMSEs with an average of $0.049 \text{ m}^3/\text{m}^3$ are all less than $0.073 \text{ m}^3/\text{m}^3$, and the bias ranges
 257 from -0.10 to $0.15 \text{ m}^3/\text{m}^3$. The comparative results of the Babao Network are not as good as that of the Maqu Network. The
 258 SM data of most sites at the Babao Network have larger ubRMSE and bias, and the correlation coefficients between in situ
 259 observed SM and the downscaled SM are generally lower. That may be mainly because the measured soil depth at the Babao
 260 Network is 4 cm, which means that there could be a systematic error between the datasets.

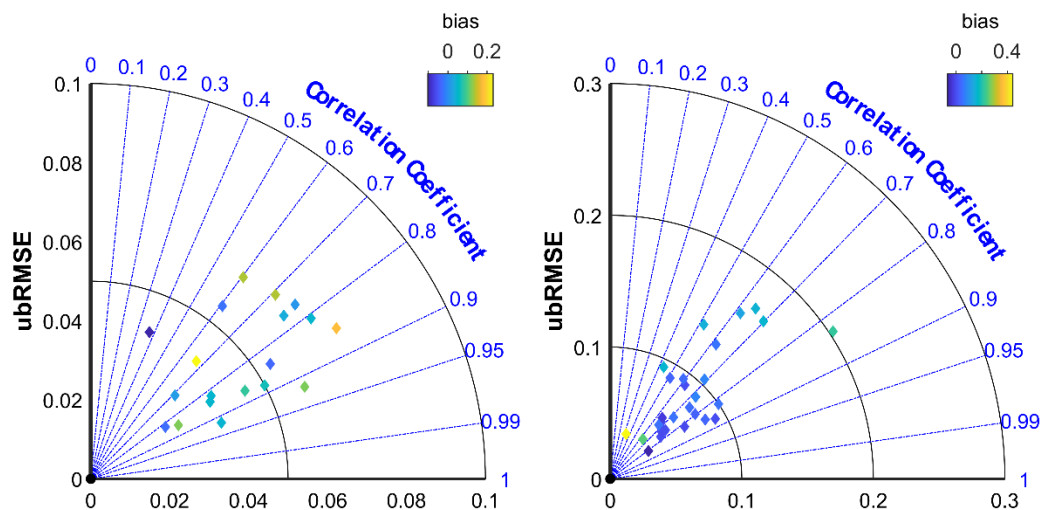


Figure 6: The relationships between in situ SM and downscaled SM. (a) Maqu Network; (b) Babao Network.

To better understand the reason for these poor results, the scatter plots comparing the two sets of data were drawn. Figure 7 shows the results of the 19 sites of the Maqu Network. All four statistical metrics, namely, R, RMSE, ubRMSE and bias were calculated, and their fitting line of the scatter was also plotted. Not surprisingly, the relationship is generally improved where there are more points. The same conclusion can be drawn according to Fig. S1, which shows the comparative results of 29 sites at the Babao Network.

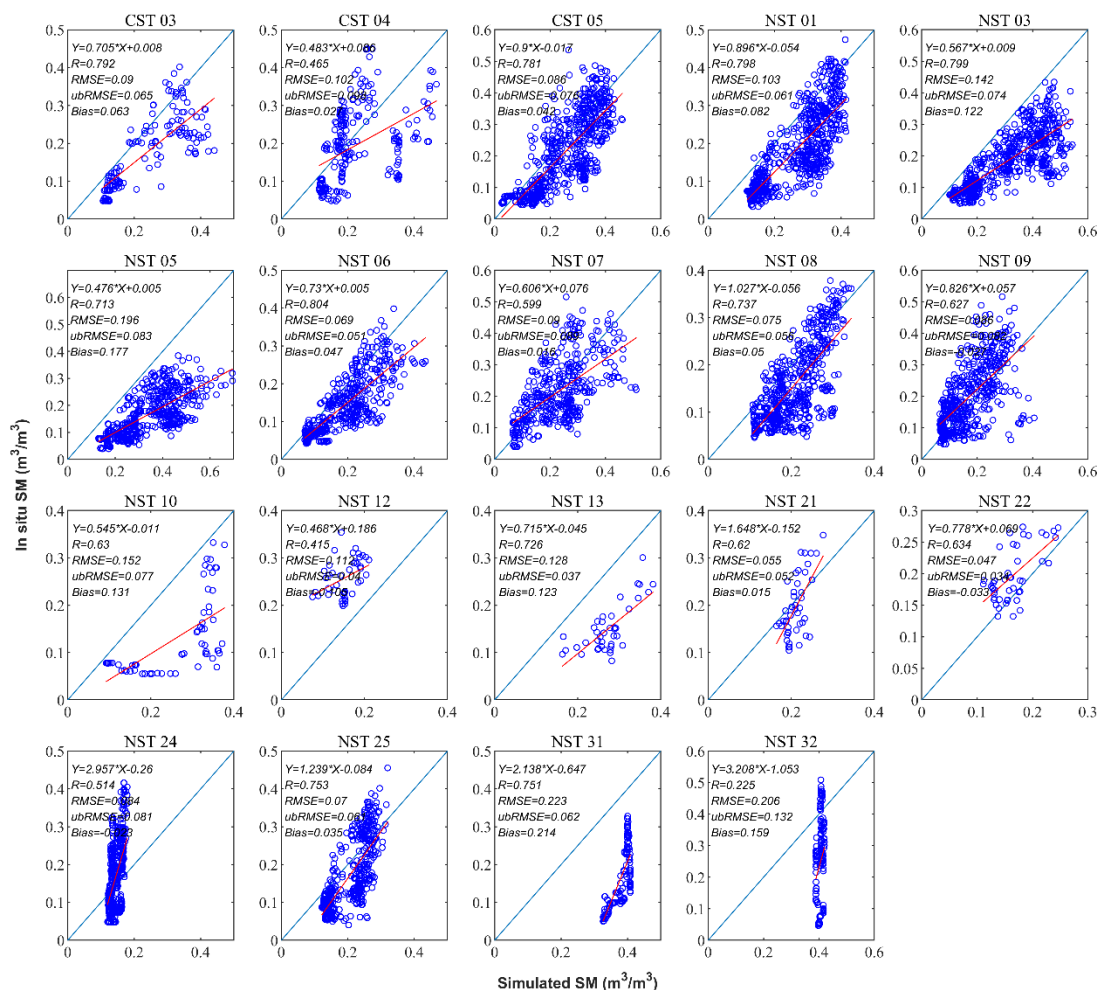


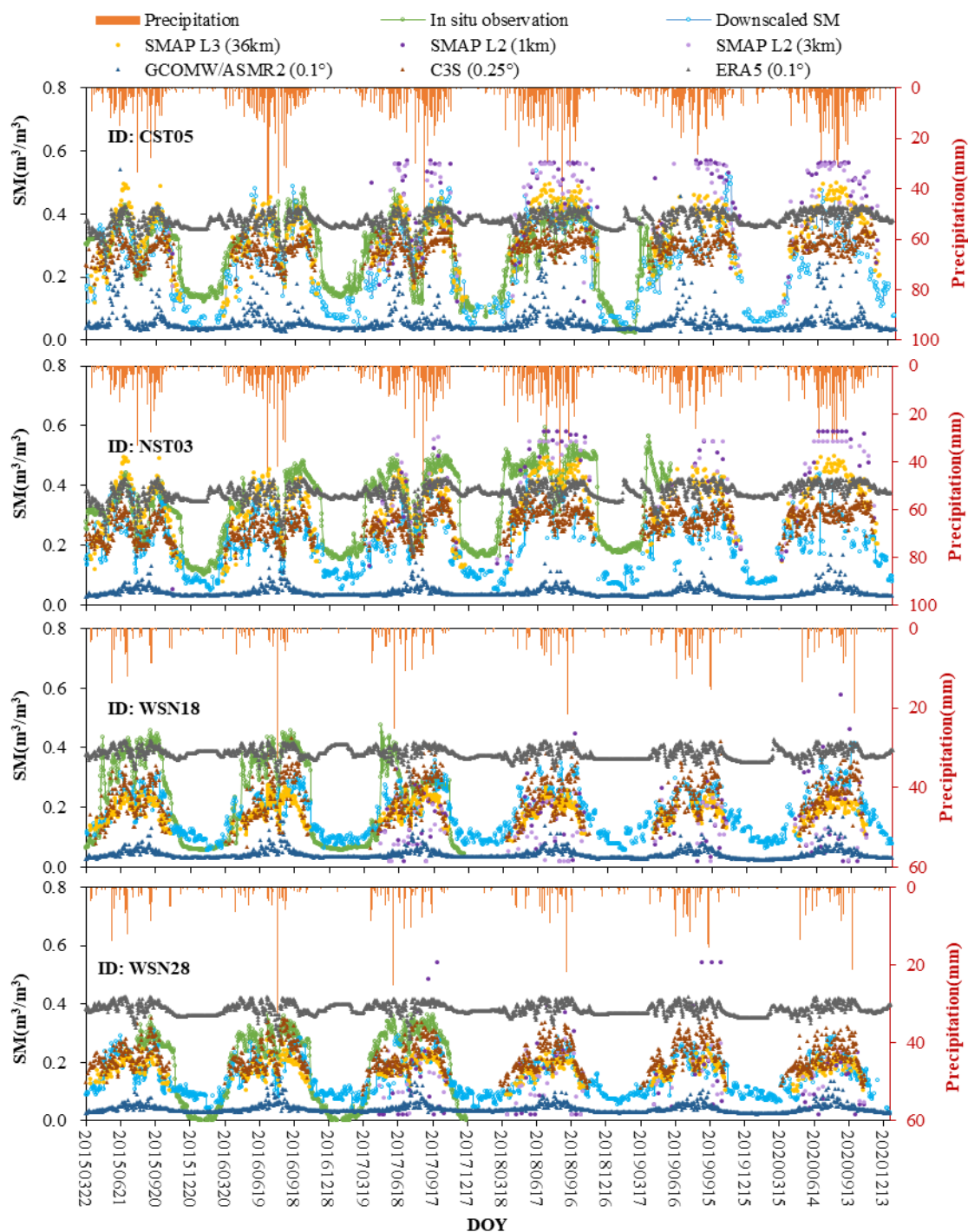
Figure 7: Comparison between the downscaled SM and in situ SM of the Maqu Network.

The observed SM of sites with a greater number of observed data were compared with these gridded SM data at different resolutions and precipitation. Figure 8 shows the temporal variations of these SM at four sites. The relationship between in situ observed SM and precipitation at all four sites is very consistent, showing annual fluctuation. The greater SM corresponds to more precipitation during the hot season, and the smaller SM corresponds to less precipitation during the cold season.

Except for GCOMW/ASMR2 SM, the variation trends of these acquired gridded SM and the downscaled SM are basically the same despite the large difference in spatial resolution. GCOMW/ASMR2 significantly underestimates SM compared to other products. Both the SMAP L2 SM at 1 km and 3 km are overestimated compared with in situ observations. Moreover, SMAP L2 SM has some valid data mainly on hot days and almost no valid data during cold seasons. The peak values of the ERA5 SM are close to those of the in situ observations, but the low values are overestimated. The C3S SM is similar to the 36 km SMAP SM, and its peak values are simulated more accurately, while the minimum values have little valid data. Compared with the original data (36 km SMAP L3), the downscaled SM has a more complete time series, especially during the cold



281 season. The downscaled SM data almost all match well with the in situ measured SM data, and all of them are also consistent
282 with the precipitation. The difference between the downscaled SM and the in situ measured SM is mainly reflected in the
283 magnitude of the variation, which is probably due to the difference in spatial resolution.



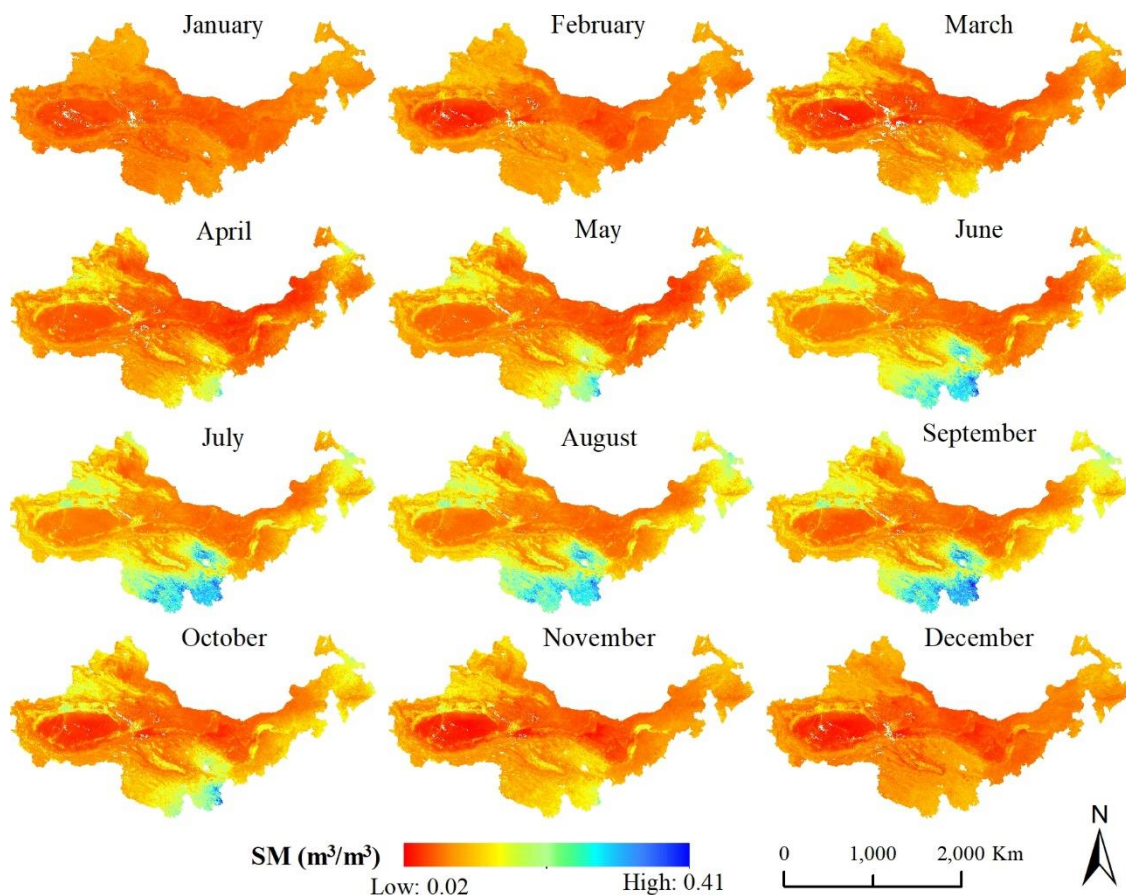
284
 285
 286

Figure 8: Time series of the in situ observed SM, the downscaled SM, the acquired gridded SM products and daily precipitation at the four selected SM sites (From Maqu Network and Babao Network, respectively).



287 3.3 Mapping of the downscaled SM

288 SM varies greatly in different months in desertified areas. Figure 9 shows the average SM in each month in the study area.
289 The SM shows a monthly change pattern, and the values from June to September are bigger than in other months, especially
290 in southern Qinghai Province, eastern Inner Mongolia Province, and western Xinjiang Province, which is consistent with the
291 process of vegetation growth. The SM in some areas is low throughout the year, such as in the Tarim Basin of Xinjiang
292 Province, western Inner Mongolia Province and most of Gansu Province.



293

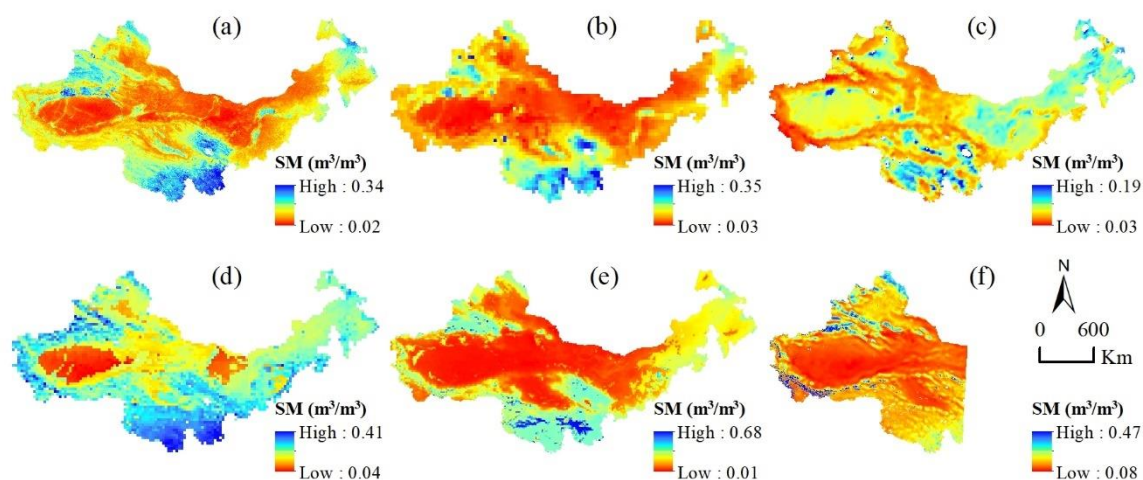
294

Figure 9: Monthly average SM in the study area.

295 The annual average SM was also calculated (Fig. S2). Overall, there is little variation in SM in different years. Further,
296 we compared the spatial patterns of the downscaled SM with the gridded SM products with different resolutions. Figure 10
297 shows the daily average SM of these products from 2015 to 2020. The spatial patterns of the downscaled SM and 36 km SMAP
298 SM are basically consistent, but the downscaled data show better details in some areas such as near rivers. The overall values
299 of GCOMW SM are relatively small, and exhibit some obvious errors in some areas. For example, SM in the Tarim Basin is
300 higher than in the surrounding area, which is completely inconsistent with other SM data. The spatial pattern of the C3S SM



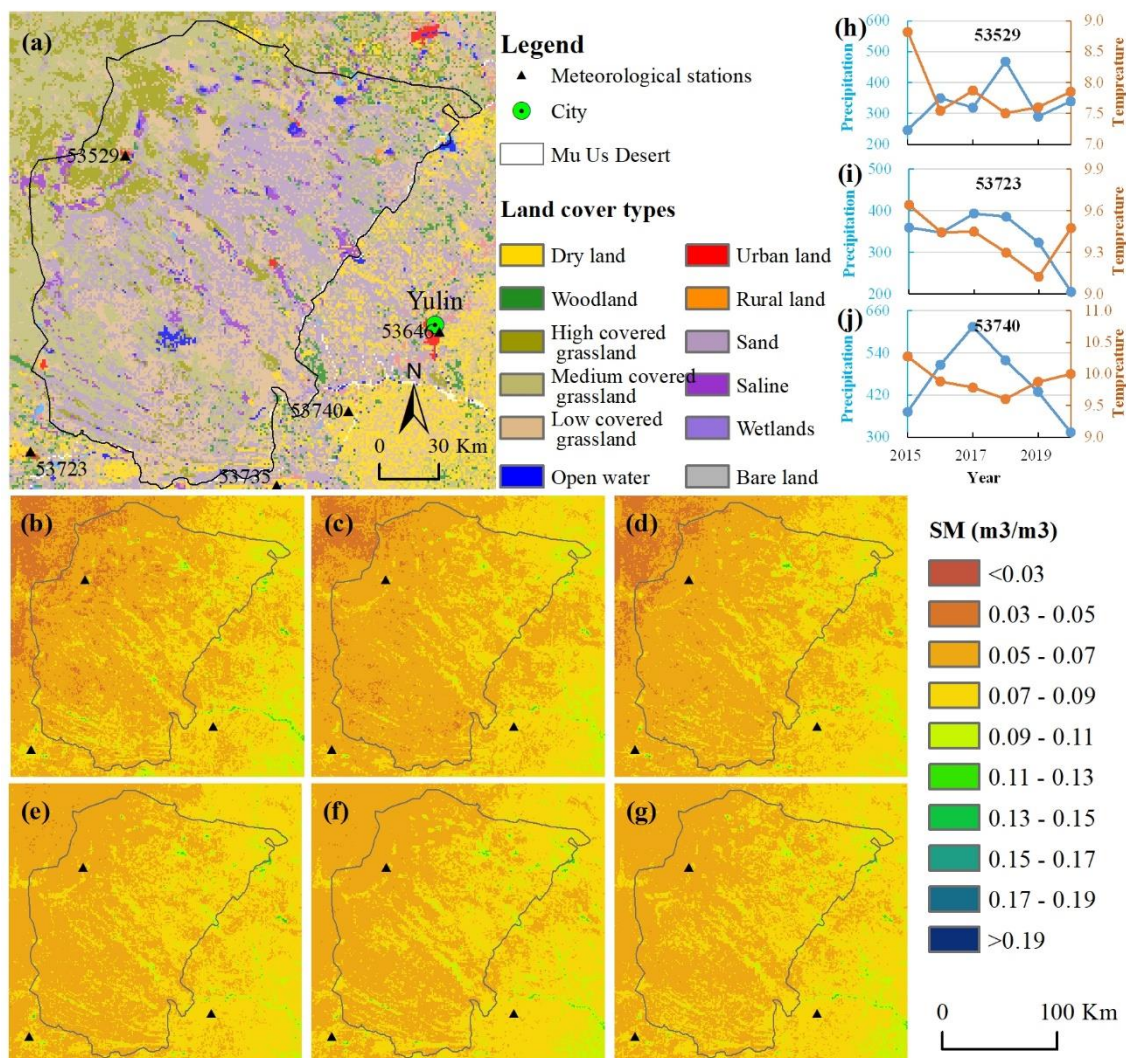
301 is close to the downscaled SM and the 36 km SMAP SM, but some details are not presented. For example, SM in the Hetao
302 Plain along the Yellow River is much higher than that in its surrounding area, which can be found in the downscaled SM and
303 the SMAP SM, but not in the C3S SM. The average SM of the ERA5 products is polarised. In some areas the values are very
304 large, and in some small areas they are very small. The FLDAS SM has high resolution, and its overall spatial pattern is
305 relatively consistent with the downscaled SM and 36 km SMAP SM. The difference is that the FLDAS SM is significantly
306 larger in higher elevation areas of the west than in other regions, which is quite different from other products. This suggests
307 that the FLDAS SM may be overestimated in these regions. In addition, FLDAS SM does not show wetter soil along the river.



308
309 **Figure 10: Daily average SM from 2015-2020 in the study area. (a)-(f) are the downscaled SM (1 km), SMAP L3 SM (36 km),**
310 **GCOMW/ASMR2 SM (0.1°), C3S SM (0.25°), ERA5 SM (0.1°) and FLDAS SM (0.1°), respectively.**

311 To better demonstrate the differences in SM, a case of the Mu Us Desert was selected (Fig. 11). The Mu Us Desert is
312 located in a semi-arid area with annual average precipitation of generally less than 400 mm, decreasing gradually from
313 southeast to northwest. The main types of land cover are grassland and sandy land, and the salinization is serious in a few
314 areas. Desertification has been severe for a long time in the past but has been significantly reversed with artificial afforestation
315 in recent years.

316 SM shows an overall trend of gradual decrease from the southeast to the northwest (Fig. 11 (b)-(g)), which is consistent
317 with the distribution of precipitation. The average SM of the same location changes little from year to year. Overall, it is
318 relatively large in 2018 and relatively small in 2015, which is also roughly consistent with annual precipitation patterns. Land
319 cover types also have a certain influence on the spatial difference of SM. The northwestern portion of the Mu Us Desert is
320 mainly grassland, which is strongly dependent on precipitation (Fig. 11 (h)). The southeastern area is mainly cultivated land
321 and is less affected by precipitation as it relies on pumping groundwater rather than natural precipitation (Fig. 11 (j)).



322
 323 **Figure 11: Soil moisture estimated for the Mu Us Desert. (a) Land cover distribution over the study area; (b)-(g) annual average SM**
 324 **from 2015-2020; (h)-(j) annual precipitation and annual average temperature of three sites (53529, 53723 and 53740), whose**
 325 **surroundings are mainly grassland, cultivated land, and cultivated land, respectively.**

326 **4. Discussion**

327 **4.1 Regression variable importance**

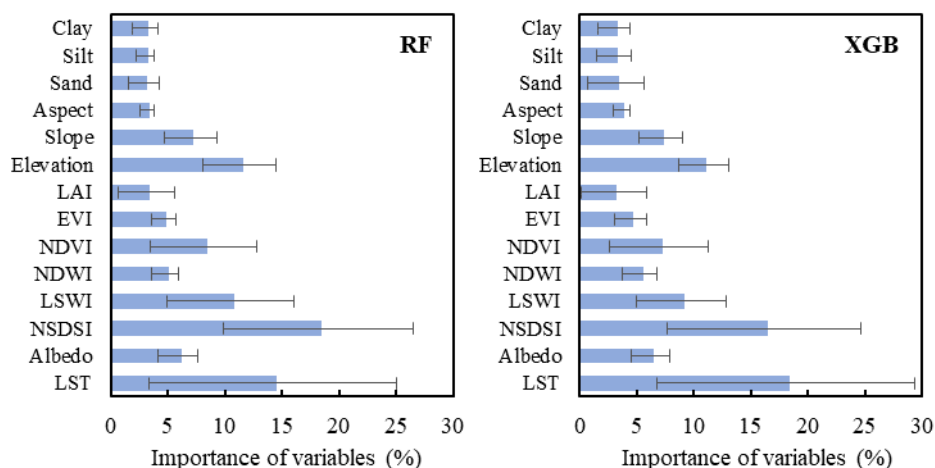
328 The selection of variables is an important step of a nonlinear regression model. The importance analysis of the variables
 329 carried out for this research found that a larger number of variables can improve the regression effect of these models to some
 330 extent. Figure 12 shows the average importance scores of each variable for the RF and XGB models across all available days.



331 The importance scores of different variables in the RF based model and the XGB based model are similar. LST and surface
 332 albedo both affect surface energy exchange and partition. LST is a very important variable in both models, which is consistent
 333 with the study of Zhao et al. (2018). NSDSI is the most sensitive soil moisture index compared to LSWI and NDWI, which
 334 was demonstrated in Yue et al. (2019). Topographical factors also exhibit importance on SM, especially elevation. The
 335 influence of soil texture (sand, silt and clay) is relatively weak, but it cannot be completely ignored.

336 The standard deviation of the importance scores of each variable is shown with error bars in Fig. 12. Its changes are
 337 mainly affected by the samples used in the regression model and the temporal variations in surface variables. For static
 338 variables such as soil structure and topographic factors, the changes in their importance scores mainly depend on the number
 339 and the location of the samples. Figure 12 also shows that their standard deviation is relatively small. Compared with static
 340 variables, the standard deviation of the importance scores of dynamic variables is significantly larger, especially for LST and
 341 LAI. This indicates that it is not reliable to construct a single regression model for a long time series.

342 In general, the variable importance analysis suggests that the selected variables are suitable for the construction of the
 343 regression model. Moreover, choosing 16 days as a time period to build a regression model benefits from obtaining a sufficient
 344 number of samples, especially since the surface variables were found still unchanged during these intervals.



345
 346 **Figure 12: The average importance scores of variables for the RF based approach and XGB based approach. Note: The importance**
 347 **scores are presented by IncNodePurity where the sum value is normalized for the RF model; The XGB model uses Gain to reflect**
 348 **the weight of variables.**

349 4.2 Advantages of model combination

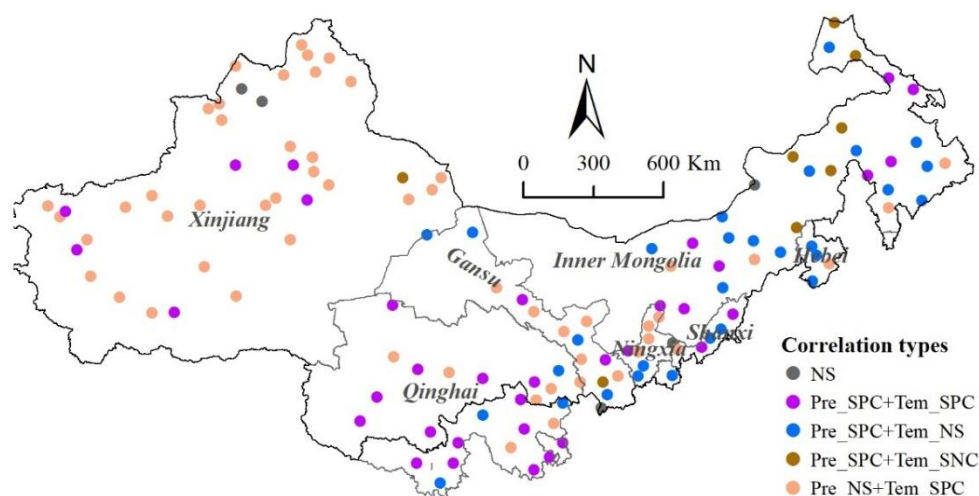
350 The simulation results of long time series will inevitably suffer the interference of various noises. A combination of
 351 multiple methods can reduce overfitting and uncertainties (Zanotti et al., 2019; Yu et al., 2021). The five methods (MLR, SVR,
 352 ANN, RF and XGB) in this study have indicated the potential flaws of a single model. Although XGB generally perform better
 353 than other models, it still has some shortcomings. As it can be seen from Figs. 4 and 5, compared with the training accuracy,



354 the test accuracy of the XGB model is significantly reduced in several periods. This means that the simulation results of the
355 XGB model is likely to have a certain degree of overfitting. In contrast, the difference between training accuracy and test
356 accuracy of the RF model is even smaller. It showed better stability than XGB at some periods (Figs. 4 and 5). The training
357 accuracy of MLR and SVR has a small difference from the test accuracy, but the overall accuracy is obviously lower, which
358 is not suitable for remote sensing SM prediction (Table 4). Some studies have also proved that SVR may also perform better
359 than some ensemble algorithms (Yu et al., 2012; Fan et al., 2018). The fitting effect of ANN varies greatly in different periods,
360 indicating that its generalization is lower than other models (Piotrowski and Napiorkowski, 2013). In general, the XGB and
361 RF models provide the best combination of prediction accuracy and stability.

362 4.3 Analysis of the relationship with precipitation and temperature

363 Precipitation and temperature are important factors affecting SM. To evaluate the impact of precipitation and temperature
364 on SM, we performed a partial correlation analysis on the data of all meteorological stations. Figure 13 shows that SM is
365 mainly positively correlated with precipitation and temperature, and a few regions are significantly negatively correlated with
366 temperature. In terms of spatial distribution, SM of the sites in the eastern region (including Inner Mongolia Province, Hebei
367 Province and Shanxi Province) is mainly significantly affected by precipitation. Due to the influence of glaciers and snowmelt,
368 the SM of the sites in the western region (Xinjiang Province and Gansu Province) is more affected by temperature. In addition,
369 the number of sites with significant positive correlation with precipitation and temperature is the largest in Qinghai Province.
370 This indicates that precipitation and temperature in the eastern part of the Tibetan Plateau both have a great influence on SM.



371
372 **Figure 13: Partial correlation between monthly downscaled SM and precipitation and temperature (Pre: precipitation; Tem:**
373 **temperature; NS: Not significance; SPC: Significantly positive correlation; SNC: Significantly negative correlation).**

374 4.4 Uncertainty and Prospects



375 While this study greatly improved the spatial resolution of SM data from 2015-2020 in the desertifying areas of North
376 China by downscaling SMAP SM products, it still presents some shortcomings. For example, due to the image quality and
377 coverage of SMAP and the impact of noise from clouds on the MODIS products, the number of valid samples for a 16-day
378 period may still be less than 100 points. This study replaced the periods with less than 100 samples with the model of the
379 previous periods. Due to the limited number of available samples, the simulation in the cold season is relatively poor (Fig. 8).
380 In addition, the upscaling (from 1 km to 36 km resolution) of surface variables also has a certain impact on the accuracy of the
381 model.

382 The Chinese government focuses on desertification reduction through afforestation and the establishment of grasslands.
383 SM data with high temporal and spatial resolution can provide a reference for the next steps of revegetation.

384 **5 Code and data availability**

385 The codes mainly used in this paper mainly includes sample selection, the building of the optimal regression model and
386 the result prediction. These codes based on the R language can be found in the supplementary documents. The downsampled
387 daily SM dataset at 1 km spatial resolution is available at <https://doi.org/10.6084/M9.FIGSHARE.16430478.V5> (Rao et al.,
388 2021). The data maps are all provided in Geotiff format, and the value has expanded 10, 000 times to make them easier to
389 store. The filenames reflect the production date in Julian Day format.

390 **6 Conclusions**

391 In this study, a framework was proposed for downscaling 36 km SMAP SM products using MODIS optical products and
392 other surface variables (mainly topographic data and soil data) based on multiple machine learning methods. Overall, the
393 regression performance of the five methods is, in order: XGB>RF>SVR>ANN>MLR. Compared with MLR, SVR and ANN,
394 XGB and RF have much better regression accuracy, and they were used in combination to produce daily 1 km downsampled SM
395 in a period of 16 days. The validation shows that the downsampled SM are highly related to most in situ measured SM. The
396 ubRMSE with an average of $0.049 \text{ m}^3/\text{m}^3$ is generally less than $0.073 \text{ m}^3/\text{m}^3$ at the Maqu Network. Time series of SM data
397 from in situ observation sites are also compared. The results show that the downsampled SMs are highly related to SMAP SMs,
398 and provide a more complete time series and match better with the in situ measured SM. Compared with some commonly used
399 gridded SM products such as SMAP L2 (1 km or 3 km), GCOMW/ASMR2, C3S, ERA5 and FLDAS SMs, the downsampled
400 SM data not only have higher spatial resolution, but also have a more reliable accuracy whether in time series or spatial
401 distribution.

402 The maps of downsampled SM show larger values from June to September, which coincides with the vegetation growing
403 season. The difference in annual mean SM is small. Spatially, SM is relatively large in Qinghai Province and in northeastern
404 Inner Mongolia, especially in summer. In arid areas such as the Tarim Basin, SM is relatively small throughout the year.



405 Moreover, precipitation and temperature both have a great influence on SM in the study area. Precipitation has a greater impact
406 on SM in the eastern part of the study area, while the effect of temperature appears to be more pronounced in the west.

407 This approach makes it possible to more accurately assess the soil moisture status in the study area. The results can support
408 regional agricultural planting and revegetation efforts and can be applied to limit desertification in other areas in the future.

409

410 **Author contributions.** FW and PR designed the research, developed the methodology, performed the analysis, and wrote the
411 paper; YW, YL, XW, and ZW edited and revised the paper.

412

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

414

415 **Acknowledgements.** This work was supported by the National Key Research and Development Program of China
416 (2018YFC0408103), the National Pilot Project for Ecological Protection and Restoration of Mountains, Rivers, Forests,
417 Farmlands, Lakes and Grasslands (Grant No. WR0203A552018), and the Desertification Monitoring Project of National
418 Forestry and Grass Administration (Grant No. 2020062012). We thank all data providers and the anonymous reviewers for
419 their detailed and constructive comments.

420 References

421 Abbaszadeh, P., Moradkhani, H., and Zhan, X.: Downscaling SMAP Radiometer Soil Moisture Over the CONUS Using an Ensemble
422 Learning Method, *Water Resour. Res.*, 55, 324–344, <https://doi.org/10.1029/2018WR023354>, 2019.

423 Achieng, K. O.: Modelling of soil moisture retention curve using machine learning techniques: Artificial and deep neural networks vs support
424 vector regression models, *Computers & Geosciences*, 133, 104320, <https://doi.org/10.1016/j.cageo.2019.104320>, 2019.

425 Ågren, A. M., Larson, J., Paul, S. S., Laudon, H., and Lidberg, W.: Use of multiple LIDAR-derived digital terrain indices and machine
426 learning for high-resolution national-scale soil moisture mapping of the Swedish forest landscape, *Geoderma*, 404, 115280,
427 <https://doi.org/10.1016/j.geoderma.2021.115280>, 2021.

428 Bai, J., Cui, Q., Zhang, W., and Meng, L.: An Approach for Downscaling SMAP Soil Moisture by Combining Sentinel-1 SAR and MODIS
429 Data, *Remote Sensing*, 11, 2736, <https://doi.org/10.3390/rs11232736>, 2019.

430 Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System, in: Proceedings of the 22nd ACM SIGKDD International Conference
431 on Knowledge Discovery and Data Mining, KDD '16: The 22nd ACM SIGKDD International Conference on Knowledge Discovery
432 and Data Mining, San Francisco California USA, 785–794, <https://doi.org/10.1145/2939672.2939785>, 2016.

433 Chen, Y., Feng, X., and Fu, B.: An improved global remote-sensing-based surface soil moisture (RSSSM) dataset covering 2003–2018,
434 *Earth Syst. Sci. Data*, 13, 1–31, <https://doi.org/10.5194/essd-13-1-2021>, 2021.



- 435 De Santis, D., Biondi, D., Crow, W. T., Camici, S., Modanesi, S., Brocca, L., and Massari, C.: Assimilation of Satellite Soil Moisture
436 Products for River Flow Prediction: An Extensive Experiment in Over 700 Catchments Throughout Europe, *Water Res*, 57,
437 <https://doi.org/10.1029/2021WR029643>, 2021.
- 438 Del Frate, F., Ferrazzoli, P., and Schiavon, G.: Retrieving soil moisture and agricultural variables by microwave radiometry using neural
439 networks, *Remote Sensing of Environment*, 84, 174–183, [https://doi.org/10.1016/S0034-4257\(02\)00105-0](https://doi.org/10.1016/S0034-4257(02)00105-0), 2003.
- 440 Demarchi, L., Kania, A., Ciężkowski, W., Piórkowski, H., Oświecimska-Piasko, Z., and Chormański, J.: Recursive Feature Elimination and
441 Random Forest Classification of Natura 2000 Grasslands in Lowland River Valleys of Poland Based on Airborne Hyperspectral and
442 LiDAR Data Fusion, *Remote Sensing*, 12, 1842, <https://doi.org/10.3390/rs12111842>, 2020.
- 443 Elshorbagy, A. and Parasuraman, K.: On the relevance of using artificial neural networks for estimating soil moisture content, *Journal of*
444 *Hydrology*, 362, 1–18, <https://doi.org/10.1016/j.jhydrol.2008.08.012>, 2008.
- 445 Fan, J., Yue, W., Wu, L., Zhang, F., Cai, H., Wang, X., Lu, X., and Xiang, Y.: Evaluation of SVM, ELM and four tree-based ensemble
446 models for predicting daily reference evapotranspiration using limited meteorological data in different climates of China, *Agricultural*
447 *and Forest Meteorology*, 263, 225–241, <https://doi.org/10.1016/j.agrformet.2018.08.019>, 2018.
- 448 Fan, J., Zheng, J., Wu, L., and Zhang, F.: Estimation of daily maize transpiration using support vector machines, extreme gradient boosting,
449 artificial and deep neural networks models, *Agricultural Water Management*, 245, 106547, <https://doi.org/10.1016/j.agwat.2020.106547>,
450 2021.
- 451 Fang, B. and Lakshmi, V.: Soil moisture at watershed scale: Remote sensing techniques, *Journal of Hydrology*, 516, 258–272,
452 <https://doi.org/10.1016/j.jhydrol.2013.12.008>, 2014.
- 453 Fang, B., Lakshmi, V., Bindlish, R., Jackson, T. J., Cosh, M., and Basara, J.: Passive Microwave Soil Moisture Downscaling Using
454 Vegetation Index and Skin Surface Temperature, *Vadose Zone Journal*, 12, vzj2013.05.0089er,
455 <https://doi.org/10.2136/vzj2013.05.0089er>, 2013.
- 456 Gu, Y., Hunt, E., Wardlow, B., Basara, J. B., Brown, J. F., and Verdin, J. P.: Evaluation of MODIS NDVI and NDWI for vegetation drought
457 monitoring using Oklahoma Mesonet soil moisture data, *Geophys. Res. Lett.*, 35, L22401, <https://doi.org/10.1029/2008GL035772>,
458 2008.
- 459 Hu, F., Wei, Z., Zhang, W., Dorjee, D., and Meng, L.: A spatial downscaling method for SMAP soil moisture through visible and shortwave-
460 infrared remote sensing data, *Journal of Hydrology*, 590, 125360, <https://doi.org/10.1016/j.jhydrol.2020.125360>, 2020.
- 461 Im, J., Park, S., Rhee, J., Baik, J., and Choi, M.: Downscaling of AMSR-E soil moisture with MODIS products using machine learning
462 approaches, *Environ Earth Sci*, 75, 1120, <https://doi.org/10.1007/s12665-016-5917-6>, 2016.



- 463 Kang, J., Jin, R., Li, X., Ma, C., Qin, J., and Zhang, Y.: High spatio-temporal resolution mapping of soil moisture by integrating wireless
464 sensor network observations and MODIS apparent thermal inertia in the Babao River Basin, China, *Remote Sensing of Environment*,
465 191, 232–245, <https://doi.org/10.1016/j.rse.2017.01.027>, 2017.
- 466 Lievens, H., Verhoest, N. E. C., De Keyser, E., Vernieuwe, H., Matgen, P., Álvarez-Mozos, J., and De Baets, B.: Effective roughness
467 modelling as a tool for soil moisture retrieval from C- and L-band SAR, *Hydrol. Earth Syst. Sci.*, 15, 151–162,
468 <https://doi.org/10.5194/hess-15-151-2011>, 2011.
- 469 Liu, J., Chai, L., Lu, Z., Liu, S., Qu, Y., Geng, D., Song, Y., Guan, Y., Guo, Z., Wang, J., and Zhu, Z.: Evaluation of SMAP, SMOS-IC,
470 FY3B, JAXA, and LPRM Soil Moisture Products over the Qinghai-Tibet Plateau and Its Surrounding Areas, *Remote Sensing*, 11, 792,
471 <https://doi.org/10.3390/rs11070792>, 2019.
- 472 Liu, Y., Yao, L., Jing, W., Di, L., Yang, J., and Li, Y.: Comparison of two satellite-based soil moisture reconstruction algorithms: A case
473 study in the state of Oklahoma, USA, *Journal of Hydrology*, 590, 125406, <https://doi.org/10.1016/j.jhydrol.2020.125406>, 2020.
- 474 Ma, M., Zhao, G., He, B., Li, Q., Dong, H., Wang, S., and Wang, Z.: XGBoost-based method for flash flood risk assessment, *Journal of*
475 *Hydrology*, 598, 126382, <https://doi.org/10.1016/j.jhydrol.2021.126382>, 2021.
- 476 Mallick, K., Bhattacharya, B. K., and Patel, N. K.: Estimating volumetric surface moisture content for cropped soils using a soil wetness
477 index based on surface temperature and NDVI, *Agricultural and Forest Meteorology*, 149, 1327–1342,
478 <https://doi.org/10.1016/j.agrformet.2009.03.004>, 2009.
- 479 Meng, X., Mao, K., Meng, F., Shi, J., Zeng, J., Shen, X., Cui, Y., Jiang, L., and Guo, Z.: A fine-resolution soil moisture dataset for China in
480 2002–2018, *Geosciences – Geophysics*, <https://doi.org/10.5194/essd-2020-292>, 2020.
- 481 Mohana, R. M., Reddy, C. K. K., Anisha, P. R., and Murthy, B. V. R.: Random forest algorithms for the classification of tree-based ensemble,
482 *Materials Today: Proceedings*, S2214785321008853, <https://doi.org/10.1016/j.matpr.2021.01.788>, 2021.
- 483 O’Neill, P., Entekhabi, D., Njoku, E., and Kellogg, K.: The NASA Soil Moisture Active Passive (SMAP) mission: Overview, in: 2010 IEEE
484 International Geoscience and Remote Sensing Symposium, IGARSS 2010 - 2010 IEEE International Geoscience and Remote Sensing
485 Symposium, Honolulu, HI, USA, 3236–3239, <https://doi.org/10.1109/IGARSS.2010.5652291>, 2010.
- 486 Peng, J., Loew, A., Merlin, O., and Verhoest, N. E. C.: A review of spatial downscaling of satellite remotely sensed soil moisture: Downscale
487 Satellite-Based Soil Moisture, *Rev. Geophys.*, 55, 341–366, <https://doi.org/10.1002/2016RG000543>, 2017.
- 488 Peng, J., Albergel, C., Balenzano, A., Brocca, L., Cartus, O., Cosh, M. H., Crow, W. T., Dabrowska-Zielinska, K., Dadson, S., Davidson,
489 M. W. J., de Rosnay, P., Dorigo, W., Gruber, A., Hagemann, S., Hirschi, M., Kerr, Y. H., Lovergine, F., Mahecha, M. D., Marzahn, P.,
490 Mattia, F., Musial, J. P., Preuschmann, S., Reichle, R. H., Satalino, G., Silgram, M., van Bodegom, P. M., Verhoest, N. E. C., Wagner,
491 W., Walker, J. P., Wegmüller, U., and Loew, A.: A roadmap for high-resolution satellite soil moisture applications – confronting product
492 characteristics with user requirements, 15, 2021.



- 493 Piles, M., Petropoulos, G. P., Sánchez, N., González-Zamora, Á., and Ireland, G.: Towards improved spatio-temporal resolution soil moisture
494 retrievals from the synergy of SMOS and MSG SEVIRI spaceborne observations, *Remote Sensing of Environment*, 180, 403–417,
495 <https://doi.org/10.1016/j.rse.2016.02.048>, 2016.
- 496 Piotrowski, A. P. and Napiorkowski, J. J.: A comparison of methods to avoid overfitting in neural networks training in the case of catchment
497 runoff modelling, *Journal of Hydrology*, 476, 97–111, <https://doi.org/10.1016/j.jhydrol.2012.10.019>, 2013.
- 498 Qu, Y., Zhu, Z., Chai, L., Liu, S., Montzka, C., Liu, J., Yang, X., Lu, Z., Jin, R., Li, X., Guo, Z., and Zheng, J.: Rebuilding a Microwave
499 Soil Moisture Product Using Random Forest Adopting AMSR-E/AMSR2 Brightness Temperature and SMAP over the Qinghai–Tibet
500 Plateau, China, *Remote Sensing*, 11, 683, <https://doi.org/10.3390/rs11060683>, 2019.
- 501 Rahimzadeh-Bajgirani, P., Berg, A. A., Champagne, C., and Omasa, K.: Estimation of soil moisture using optical/thermal infrared remote
502 sensing in the Canadian Prairies, *ISPRS Journal of Photogrammetry and Remote Sensing*, 83, 94–103,
503 <https://doi.org/10.1016/j.isprsjprs.2013.06.004>, 2013.
- 504 Rao, P., Jiang, W., Hou, Y., Chen, Z., and Jia, K.: Dynamic Change Analysis of Surface Water in the Yangtze River Basin Based on MODIS
505 Products, *Remote Sensing*, 10, 1025, <https://doi.org/10.3390/rs10071025>, 2018.
- 506 Rao, P., Wang, Y., Wang, F., Liu, Y., Wang, X., and Wang, Z.: Daily soil moisture mapping at 1-km resolution based on SMAP data for
507 areas affected by desertification in Northern China, <https://doi.org/10.6084/M9.FIGSHARE.16430478.V3>, 2021.
- 508 Sandholt, I., Rasmussen, K., and Andersen, J.: A simple interpretation of the surface temperature/vegetation index space for assessment of
509 surface moisture status, *Remote Sensing of Environment*, 79, 213–224, [https://doi.org/10.1016/S0034-4257\(01\)00274-7](https://doi.org/10.1016/S0034-4257(01)00274-7), 2002.
- 510 Shangguan, W., Dai, Y., Liu, B., Ye, A., and Yuan, H.: A soil particle-size distribution dataset for regional land and climate modelling in
511 China, *Geoderma*, 171–172, 85–91, <https://doi.org/10.1016/j.geoderma.2011.01.013>, 2012.
- 512 Shi, R., Xu, X., Li, J., and Li, Y.: Prediction and analysis of train arrival delay based on XGBoost and Bayesian optimization, *Applied Soft
513 Computing*, 109, 107538, <https://doi.org/10.1016/j.asoc.2021.107538>, 2021.
- 514 Sun, L., Sun, R., Li, X., Liang, S., and Zhang, R.: Monitoring surface soil moisture status based on remotely sensed surface temperature and
515 vegetation index information, *Agricultural and Forest Meteorology*, 166–167, 175–187,
516 <https://doi.org/10.1016/j.agrformet.2012.07.015>, 2012.
- 517 Wagner, W., Hahn, S., Kidd, R., Melzer, T., Bartalis, Z., Hasenauer, S., Figa-Saldaña, J., de Rosnay, P., Jann, A., Schneider, S., Komma, J.,
518 Kubu, G., Brugger, K., Aubrecht, C., Züger, J., Gangkofner, U., Kienberger, S., Brocca, L., Wang, Y., Blöschl, G., Eitzinger, J., and
519 Steinnocher, K.: The ASCAT Soil Moisture Product: A Review of its Specifications, Validation Results, and Emerging Applications,
520 *metz*, 22, 5–33, <https://doi.org/10.1127/0941-2948/2013/0399>, 2013.



- 521 Wang, G., Zhang, X., Yinglan, A., Duan, L., Xue, B., and Liu, T.: A spatio-temporal cross comparison framework for the accuracies of
522 remotely sensed soil moisture products in a climate-sensitive grassland region, *Journal of Hydrology*, 597, 126089,
523 <https://doi.org/10.1016/j.jhydrol.2021.126089>, 2021.
- 524 Wang, S., Liu, S., Zhang, J., Che, X., Yuan, Y., Wang, Z., and Kong, D.: A new method of diesel fuel brands identification: SMOTE
525 oversampling combined with XGBoost ensemble learning, *Fuel*, 282, 118848, <https://doi.org/10.1016/j.fuel.2020.118848>, 2020.
- 526 Wang, T., Yang, D., Fang, B., Yang, W., Qin, Y., and Wang, Y.: Data-driven mapping of the spatial distribution and potential changes of
527 frozen ground over the Tibetan Plateau, *Science of The Total Environment*, 649, 515–525,
528 <https://doi.org/10.1016/j.scitotenv.2018.08.369>, 2019.
- 529 Wang, X., Xie, H., Guan, H., and Zhou, X.: Different responses of MODIS-derived NDVI to root-zone soil moisture in semi-arid and humid
530 regions, *Journal of Hydrology*, 340, 12–24, <https://doi.org/10.1016/j.jhydrol.2007.03.022>, 2007.
- 531 Yao, P., Shi, J., Zhao, T., Lu, H., and Al-Yaari, A.: Rebuilding Long Time Series Global Soil Moisture Products Using the Neural Network
532 Adopting the Microwave Vegetation Index, *Remote Sensing*, 9, 35, <https://doi.org/10.3390/rs9010035>, 2017.
- 533 Yu, H., Wu, Y., Niu, L., Chai, Y., Feng, Q., Wang, W., and Liang, T.: A method to avoid spatial overfitting in estimation of grassland above-
534 ground biomass on the Tibetan Plateau, *Ecological Indicators*, 125, 107450, <https://doi.org/10.1016/j.ecolind.2021.107450>, 2021.
- 535 Yu, Z., Liu, D., Lü, H., Fu, X., Xiang, L., and Zhu, Y.: A multi-layer soil moisture data assimilation using support vector machines and
536 ensemble particle filter, *Journal of Hydrology*, 475, 53–64, <https://doi.org/10.1016/j.jhydrol.2012.08.034>, 2012.
- 537 Yue, J., Tian, J., Tian, Q., Xu, K., and Xu, N.: Development of soil moisture indices from differences in water absorption between shortwave-
538 infrared bands, *ISPRS Journal of Photogrammetry and Remote Sensing*, 154, 216–230, <https://doi.org/10.1016/j.isprsjprs.2019.06.012>,
539 2019.
- 540 Zanotti, C., Rotiroti, M., Sterlacchini, S., Cappellini, G., Fumagalli, L., Stefania, G. A., Nannucci, M. S., Leoni, B., and Bonomi, T.:
541 Choosing between linear and nonlinear models and avoiding overfitting for short and long term groundwater level forecasting in a
542 linear system, *Journal of Hydrology*, 578, 124015, <https://doi.org/10.1016/j.jhydrol.2019.124015>, 2019.
- 543 Zawadzki, J. and Kędzior, M.: Soil moisture variability over Odra watershed: Comparison between SMOS and GLDAS data, *International*
544 *Journal of Applied Earth Observation and Geoinformation*, 45, 110–124, <https://doi.org/10.1016/j.jag.2015.03.005>, 2016.
- 545 Zeng, J., Li, Z., Chen, Q., Bi, H., Qiu, J., and Zou, P.: Evaluation of remotely sensed and reanalysis soil moisture products over the Tibetan
546 Plateau using in-situ observations, *Remote Sensing of Environment*, 163, 91–110, <https://doi.org/10.1016/j.rse.2015.03.008>, 2015.
- 547 Zhang, P., Zheng, D., van der Velde, R., Wen, J., Zeng, Y., Wang, X., Wang, Z., Chen, J., and Su, Z.: Status of the Tibetan Plateau
548 observatory (Tibet-Obs) and a 10-year (2009–2019) surface soil moisture dataset, *Hydrology and Soil Science – Hydrology*,
549 <https://doi.org/10.5194/essd-2020-209>, 2020.



- 550 Zhao, W. and Li, A.: A Downscaling Method for Improving the Spatial Resolution of AMSR-E Derived Soil Moisture Product Based on
551 MSG-SEVIRI Data, *Remote Sensing*, 5, 6790–6811, <https://doi.org/10.3390/rs5126790>, 2013.
- 552 Zhao, W., Li, A., and Zhao, T.: Potential of Estimating Surface Soil Moisture With the Triangle-Based Empirical Relationship Model, *IEEE*
553 *Trans. Geosci. Remote Sensing*, 55, 6494–6504, <https://doi.org/10.1109/TGRS.2017.2728815>, 2017.
- 554 Zhao, W., Sánchez, N., Lu, H., and Li, A.: A spatial downscaling approach for the SMAP passive surface soil moisture product using random
555 forest regression, *Journal of Hydrology*, 563, 1009–1024, <https://doi.org/10.1016/j.jhydrol.2018.06.081>, 2018.

556