



Kwon, M., Kwon, H. H., & Han, D. (2018). A spatial downscaling of soil moisture from rainfall, temperature, and AMSR2 using a Gaussian-mixture nonstationary hidden Markov model. *Journal of Hydrology*, 564, 1194-1207. <https://doi.org/10.1016/j.jhydrol.2017.12.015>

Peer reviewed version

License (if available):
CC BY-NC-ND

Link to published version (if available):
[10.1016/j.jhydrol.2017.12.015](https://doi.org/10.1016/j.jhydrol.2017.12.015)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Elsevier at <https://www.sciencedirect.com/science/article/pii/S0022169417308338> . Please refer to any applicable terms of use of the publisher.

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/pure/about/ebr-terms>

1 A Spatial Downscaling of Soil Moisture from Rainfall, Temperature, and
2 AMSR2 Using a Gaussian-Mixture Nonstationary Hidden Markov Model
3

4
5 Moonhyuk Kwon¹, Hyun-Han Kwon^{2*} and Dawei Han¹
6
7
8
9

10
11 12/11/2017
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

27 *Corresponding Author: Hyun-Han Kwon, hkwon@jbnu.ac.kr
28

¹ Civil and Environmental Engineering, University of Bristol, United Kingdom

² Department of Civil Engineering, Chonbuk National University, Jeonju-si, Jeollabuk-do, South Korea

29 **Abstract**

30 A multivariate stochastic soil moisture estimation approach based on a Gaussian-mixture
31 nonstationary hidden Markov model (GM-NHMM) is introduced in this study to spatially
32 disaggregate the AMSR2 soil moisture data for multiple locations in the Yongdam dam
33 watershed in South Korea. Rainfall and air temperature are considered as additional
34 predictors in the proposed modeling framework. In GM-NHMM, a six-state model is
35 constructed with three predictors representing an unobserved state associated with soil
36 moisture. It is clearly seen that the rainfall predictor plays a substantial role in achieving the
37 overall predictability. Using weather variables (i.e., rainfall and temperature) can be effective
38 in picking up some of the predictability of local soil moisture that is not captured by the
39 AMSR2 data. On the other hand, larger scale dynamic features identified from the AMSR2
40 data seem to facilitate the identification of regional spatial patterns of soil moisture. The
41 efficiency of the proposed model is compared with that of an ordinary regression model
42 (OLR) using the same predictors. The mean correlation coefficient of the proposed model is
43 about 0.78, which is significantly greater than that of the OLR at about 0.49. The proposed
44 GM-NHMM method not only provides a better representation of the observed SM than the
45 OLR model but also preserves the spatial coherence across all stations reasonably well.

46

47 **Keywords:** Soil moisture, stochastic model, AMSR2, spatial downscaling, Gaussian mixture
48 model, and nonstationary hidden Markov model

49

50 **1. Introduction**

51
52 Soil moisture (SM) is a key hydrologic state variable for understanding hydrologic processes,
53 including runoff, infiltration, drought, crop growth, and many other phenomena closely
54 related to soil conditions (Albergel et al., 2008; Barrett and Petropoulos, 2013; Brocca et al.,
55 2011; Zhao and Li, 2013), even though the amount of water in the soil profile accounts for
56 less than 0.001 % of the total global water budget (Barrett and Petropoulos, 2013). Thus,
57 acquiring accurate SM information has been a priority in hydrology, meteorology, and
58 climatology. SM data can be obtained in several ways, including in-situ measurements,
59 remote sensing techniques, and soil moisture accounting models. However, each approach
60 has its own advantages and limitations, so different data sources are often integrated to
61 mitigate individual limitations. For more details, the reader is kindly referred to, e.g., Brocca
62 et al., (2017a), Owe et al., (2008), Parajka et al., (2006), and Zhuo and Han, (2016).

63 In-situ SM observations are generally regarded as the most reliable measurement to validate
64 remotely sensed soil moisture products. The reason for using in-situ SMs is their robustness
65 with respect to the SM retrieved through either remote sensing techniques or soil moisture
66 accounting models. However, in many parts of the world, it remains challenging to collect
67 spatially and temporally suitable ground-based soil moisture data (Brocca et al., 2017b; Peng
68 et al., 2017; Zhuo and Han, 2016). Another issue is that in-situ SM observations are rarely
69 representative of large-scale SM (Griesfeller et al., 2016; Merlin et al., 2012; Reichle et al.,
70 2007), and hydrological analysis is typically conducted on a catchment scale. Considering the
71 limitations of using point-based SM measurements, satellite remote sensing has become an
72 alternative way to monitor SM conditions on a regional scale (Brocca et al., 2011), providing
73 more comprehensive and coherent coverage both spatially and temporally to better
74 understand soil moisture variability in the context of water resource management (Zhao and
75 Li, 2013).

76 Satellite-based active and passive microwave sensors have the potential advantage of
77 estimating SM spatial fields. Specifically, microwave remote sensing techniques use a longer
78 wavelength than visible and infrared radiation, so they are less affected by cloud coverage,
79 haze, rainfall, and many other weather conditions (Barrett and Petropoulos, 2013; Zhao and
80 Li, 2013). SM data retrieved from various remote sensing sensors, such as Advanced
81 Microwave Scanning Radiometer 2 (AMSR2; JAXA, 2013) , the Soil Moisture Ocean
82 Salinity Satellite (SMOS; Kerr et al., 2012), Soil Moisture Active Passive (SMAP; Das et al.,
83 2011), and the Advanced Scatterometer (ASCAT; Albergel et al., 2008), have become widely
84 available in recent years, providing reasonable accuracy over a wide area with relatively high
85 spatial–temporal resolution. In the past few decades, many studies have explored the
86 accuracy of microwave sensors and improved their applicability to hydrology (Brocca et al.,
87 2017a; Cenci et al., 2016; Parajka et al., 2006; Zhuo and Han, 2016). The challenges
88 associated with these efforts have in turn led to the introduction of new methods to facilitate
89 the suitable use of satellite-based SM measurements with a reasonable degree of accuracy.
90 One major challenge in using satellite SM data for practical applications is their coarse spatial
91 resolution and uncertainties stemming from an inability to resolve sub-grid scale variability.
92 To overcome those limitations, various statistical approaches have used a downscaling
93 framework to achieve a higher spatial resolution for microwave SM data (Merlin et al., 2012;
94 Peng et al., 2016; Piles et al., 2014; Ranney et al., 2015; Zhao and Li, 2013). Those
95 techniques can be divided into two categories: statistical and dynamic downscaling
96 approaches. The downscaling methods also vary depending on the type of data being studied,
97 such as radar, optical/thermal, topography, or soil information data (Peng et al., 2017).
98 Optical/thermal sensor data (generally vegetation index, surface temperature, albedo, etc.)
99 have been widely used to disaggregate the original satellite SM products into fine-scale
100 estimates because they not only provide land surface parameters at higher spatial resolution

101 (Peng et al., 2016; Piles et al., 2011; Zhao and Li, 2013) but also have a significant
102 correlation with soil moisture (Fang and Lakshmi, 2014; Peng et al., 2015; Srivastava et al.,
103 2013). The basic idea behind these approaches is to build a statistical model (based on the
104 relationship between the satellite SM products and surface parameters) that can simulate SM
105 sequences using given surface parameters as predictors. The most frequently reported
106 practical limitation of this approach is that optical and thermal properties can be obtained
107 only under clear-sky conditions (Djamai et al., 2016; Park et al., 2017). Geo-information
108 data, such as topography, soil attributes, and vegetation, have also been used to disaggregate
109 coarse-scale SM values into fine-scale ones using a regression framework (Busch et al., 2012;
110 Ranney et al., 2015).

111 During the past few decades, machine learning techniques have been used to spatially
112 downscale satellite-based SM data for enhanced spatial resolution (Im et al., 2016; Park et al.,
113 2017; Srivastava et al., 2013; Xing et al., 2017). For example, Srivastava et al. (2013) tested
114 and compared several machine learning techniques, including an artificial neural network, a
115 support vector machine, and a relevance vector machine, to spatially downscale the SMOS
116 SM data sets. Specifically, they used Moderate Resolution Imaging Spectro-radiometer
117 (MODIS) land surface temperature as auxiliary information in disaggregating the SMOS SM
118 products. Park et al. (2017) developed a downscaling scheme based on a modified regression
119 tree model that combined multiple sensors (AMSR2 and ASCAT) with four other predictors:
120 MODIS land surface temperature, the normalized difference vegetation index, land cover,
121 and a digital elevation model.

122 However, the existing approaches all largely depend on a linear or nonlinear regression
123 model to spatially downscale the satellite SM products without considering the stochastic
124 nature of soil moisture dynamics. The spatiotemporal dynamics of soil moisture content
125 result from complicated and mutually related processes of hydro-meteorological elements,

126 such as subsurface flow, lateral flow, infiltration, precipitation, climate, and soil (Botter et al.,
127 2007; Ridolfi et al., 2003). The influence of spatiotemporal variability in precipitation and
128 temperature on the slow-varying behavior of basin-scale SM can be better represented within
129 a stochastic modeling framework (Botter et al., 2007). Recently, a stochastic downscaling
130 technique, a nonstationary Markov model with a gamma (or exponential) distribution, has
131 been widely used in both hydrology and meteorology (Cioffi et al., 2017; Khalil et al., 2010;
132 Mehrotra and Sharma, 2005; Robertson et al., 2004). The stochastic downscaling approaches
133 have been mainly used for rainfall simulation at multiple locations (Cioffi et al., 2017; Khalil
134 et al., 2010; Kwon et al., 2011, 2009, Mehrotra and Sharma, 2010, 2006; Robertson et al.,
135 2004; Stehlík and Bárdossy, 2002); they have rarely been applied to SM data by means of a
136 multivariate downscaling framework (no literature regarding SM has been found).

137 Given this background, we here investigate the following questions:

- 138 (1) Can daily soil moisture sequences conditional on intraseasonal variability in
139 climate be effectively clustered and discretized as a small set of states? In
140 addition, can the identified states of daily soil moisture and their transition
141 probability be explicitly considered to better characterize soil moisture
142 dynamics?
- 143 (2) Is it desirable to use a nonstationary stochastic model that considers climate
144 variables such as precipitation, temperature, and satellite-based soil moisture
145 products as predictors? Does a combination of climate variables and satellite-
146 based soil moisture better inform simulations?
- 147 (3) Can the proposed stochastic modeling framework be applied to simultaneously
148 simulate the daily sequences of soil moisture at multiple locations on a watershed
149 scale?

150 We here propose a multivariate Gaussian mixture nonstationary hidden Markov model (GM-
151 NHMM), which is primarily based on Hughes et al., (1999) and Yoo et al., (2015), to
152 investigate those questions, with the intention of providing a practical tool for the estimation
153 of daily soil moisture on the watershed scale for use in agricultural drought monitoring and
154 hydrologic modeling. In-situ SM observations at multiple stations are here used as a
155 dependent variable, and both air temperature and rainfall, as well as the AMSR2 data, are
156 considered as predictors. The proposed downscaling approach is applied to the Yongdam
157 dam watershed in South Korea. The performance of the proposed downscaling scheme is then
158 validated with 6 in-situ observations through a cross-validation procedure.

159

160 **2. Study Area and Data**

161 **2.1 Site description and observation data**

162 In this study, we apply the spatial downscaling approach to satellite SM measurements for
163 multiple stations in the Yongdam dam watershed in southwestern Korea (35.6° – 36.0° N
164 latitude and 127.3° – 127.7° E longitude). Most of the in-situ SM observation stations in this
165 catchment are in the forest, and the dominant soil type consists of sand (62.1 %), loam (20.7
166 %), and silt (17.0 %). The average annual precipitation and air temperature during the
167 investigation period (2014–2016) were 1,147 mm and 11.4° C, respectively. Figure 1 shows
168 the study area and six in-situ soil moisture stations where precipitation data were also
169 measured (<http://www.ydew.or.kr/kdrum/main/main.do>). Here, precipitation data are
170 averaged over the entire region. Additionally, air temperature (available for download from
171 <https://data.kma.go.kr/cmmn/main.do>) was measured at the Jangsu weather station operated
172 by the Korea Meteorological Administration (<https://web.kma.go.kr/eng/>). The soil moisture
173 observation network covers a drainage area of 930 km^2 with elevation ranging from 209 to
174 1,588 m a.s.l. The Korea Water Resources Corporation has continuously recorded in-situ SM

175 observations measured at half-hourly time interval since 2014 using a time domain
176 reflectometer (TDR; Topp et al., 1980). The specifications for the observation sites used in
177 this study are given in Table 1. Depth-averaged SM representing the mean soil moisture
178 content in the soil layer 0-60cm were used for subsequent study.

179
180 [Insert Figure 1 and Table 1]

181

182 **2.2 Satellite data**

183 AMSR2 is on the GCOM-W1 satellite launched by the Japan Aerospace Exploration Agency
184 (JAXA) in May 2012. As a follow-on instrument to AMSR-E, which was operated from 2002
185 to 2012, the AMSR2 is a passive microwave sensor that measures the brightness temperature
186 at seven different frequencies between 6.9 GHz and 89.0 GHz (Imaoka et al., 2010). It is
187 widely acknowledged that microwaves measured from space are severely contaminated by
188 radio frequency interference (RFI) effects (Liu et al., 2011; Njoku et al., 2005; Zeng et al.,
189 2015). Therefore, a new 7.3-GHz channel was added to the AMSR2 to identify and address
190 RFI signals. Additionally, the AMSR2 has a larger antenna (2.0 m) than the AMSR-E (1.6 m)
191 to provide a higher spatial resolution. The AMSR2 provides geophysical products such as
192 integrated water vapor, integrated cloud liquid water, precipitation, sea surface temperature,
193 sea surface wind speed, sea ice concentration, snow depth, and soil moisture content (Imaoka
194 et al., 2010). For this study, we obtained the AMSR2 L3 SM products, derived from the
195 JAXA algorithm with 10 km spatial resolution, from the distributor's website ([https://gcom-
196 w1.jaxa.jp/auth.html](https://gcom-w1.jaxa.jp/auth.html)). Readers are referred to Koike (2013) for a detailed description of the
197 retrieval algorithm. The AMSR2 sensor provides volumetric SM content from 0 to 60 % with
198 1–2 day revisit frequency. The daily AMSR2 SM data are extracted by averaging the
199 ascending (1:30 pm) plus descending (1:30 am) overpasses over a three-year period (2014–
200 2016).

201 3. Methodology

202 3.1 Multivariate Gaussian-Mixture Nonstationary Hidden Markov Model

203
204 In this study, we propose a novel approach to stochastic modeling of soil moisture at multiple
205 locations that takes into account a set of exogenous variables: rainfall, temperature, and
206 satellite information. Here, we briefly present only the relevant details of a multivariate
207 hidden Markov model described elsewhere (Khalil et al., 2010; Kwon et al., 2011, 2009;
208 Robertson et al., 2004; Yoo et al., 2015) and primarily based on Hughes et al. (1999). Figure
209 2 shows schematically the procedure of this study.

210 [Insert Figure 2]

211 A hidden Markov model (HMM) describes a process in which part of the system dynamics is
212 hidden, and some other part of the system can be partially explained by other observations.

213 The HMM uses a Markovian process and a set of stochastic functions to generate plausible
214 sequences for a given time series based on stochastic sampling from probability distributions
215 conditioned on different hidden states (Daniel and Martin, 2017; Gharhramani, 2001).

216 Let \mathbf{SM}_t be an M -dimensional vector of in-situ soil moisture measurements corresponding
217 to M -stations at time t . Let $\mathbf{SM}_{1:T} = (\mathbf{SM}_1, \dots, \mathbf{SM}_T)$ denote a sequence of soil moisture with
218 length T . The sequence of observed soil moisture measurements $\mathbf{SM}_{1:T}$ is presumed to be
219 governed by a Markov property with the corresponding sequence $\mathbf{S}_{1:T} = (S_1, \dots, S_T)$ of a finite
220 number of hidden states, taking on values k in $\{1, K\}$. A joint distribution of $\mathbf{SM}_{1:T}$ and
221 $\mathbf{S}_{1:T}$ can be explicitly defined by taking the two conditional independence (CI) assumptions
222 (Bishop, 2006; Smyth et al., 1997), as formulated below.

223 First, assume that the sequence of hidden states $\mathbf{S}_{1:T}$ follows the stationary Markovian
224 process that relies only on the values of the previous k -th order states. Obviously, the

225 probability distribution for the current hidden state with a first-order model ($k = 1$) can be
 226 represented as equation (1) (Rabiner, 1989).

$$227 \quad p(S_1, \dots, S_T) = p(S_1) \prod_{t=2}^T p(S_t | S_{t-1}) \quad (1)$$

228 For a stationary HMM, $p(S_1)$ is the initial-state probability vector, and the state-transition
 229 probability matrix of a hidden state can be denoted as $p(S_i | S_{i-1}) = \{\gamma_{ij}\}$, $1 \leq i, j \leq K$.

230 Second, assume that individual in-situ observations \mathbf{SM}_t are conditionally independent of
 231 all other variables in the model given the current state S_t (Robertson et al., 2006; Smyth et
 232 al., 1997).

$$233 \quad p(\mathbf{SM}_{1:T} | \mathbf{S}_{1:T}) = \prod_{t=1}^T p(\mathbf{SM}_t | S_t) \quad (2)$$

234 The joint probability of the soil moisture data $\mathbf{SM}_{1:T}$ and the hidden states can then be
 235 formulated as equation (3) (Kwon et al., 2011, 2009; Robertson et al., 2006).

$$236 \quad p(\mathbf{SM}_{1:T}, \mathbf{S}_{1:T}) = \left[p(S_1) \prod_{t=2}^T p(S_t | S_{t-1}) \right] \left[\prod_{t=1}^T p(\mathbf{SM}_t | S_t) \right] \quad (3)$$

237 Soil moisture values, \mathbf{SM}_t^M , at time t for M stations are assumed to be conditionally
 238 independent of one another given the hidden state S_t . Here, spatial dependencies across
 239 multiple stations are indirectly modeled by the hidden state variable, as described in equation
 240 (4). Note that a more advanced approach to modeling the spatial structure of \mathbf{SM}_t across M
 241 sites could be of particular interest in situations with high spatial correlation. More
 242 specifically, the spatial coherence across stations is considered by assigning a state to each
 243 day, representing the spatial structure of soil moisture (Kwon et al., 2011, 2009; Robertson et
 244 al., 2006).

245
$$p(\mathbf{SM}_t^m | S_t) = \prod_{m=1}^M p(\mathbf{SM}_t^m | S_t) \quad (4)$$

246 The probability density function for the emission distribution at an individual soil moisture
 247 station \mathbf{SM}_t^m is assumed to be approximated by a Gaussian mixture function of C
 248 components for non-zero soil moisture, with $p_{i,m,c} \geq 0$ and $\sum_{c=1}^C p_{i,m,c} = 1$ for all
 249 $m = 1, \dots, M$ and $i = 1, \dots, K$, as follows:

250
$$p(\mathbf{SM}_t^m = r | S_t = i) = \sum_{c=1}^C p_{i,m,c} N(\mu_{i,m,c}, \sigma_{i,m,c}) \quad (5)$$

251 Here, μ and σ are the mean and variance of the Gaussian distribution, respectively, and
 252 the set of parameters associated with the transition matrix, the initial states, and the
 253 parameters of emission distribution are simultaneously estimated from the observed soil
 254 moisture data using the expectation-maximization (EM) algorithm in an optimization context.
 255 Gaussian mixture models are a statistical tool for multimodal density estimation (Bilmes,
 256 1998; Gauvain and Lee, 1994). Gaussian mixture models have been used for soil moisture
 257 modeling (Ryu and Famiglietti, 2005; Verhoest et al., 2015; Vilasa et al., 2017), and have
 258 also been used extensively in hydrologic field (Carreau et al., 2009; Lakshmanan and Kain,
 259 2010; Rings et al., 2012; Yoo et al., 2015). Unlike the HMM, the underlying assumption of
 260 the GM-NHMM is that soil moisture is generated in a stochastic process that sequentially
 261 depends on a set of predictors represented by rainfall, temperature, and the satellite product.
 262 Specifically, NHMMs can be constructed by imposing a non-stationarity assumption on the
 263 probability distribution of the response variables, which in turn depends on observed
 264 independent variables (Hughes et al., 1999; Hughes and Guttorp, 1994; Kwon et al., 2011).
 265 This soil moisture model can be substantially expanded by introducing a mixture model for
 266 soil moisture content into the existing HMM. In this study, we use a mixture of Gaussians to
 267 describe soil moisture at multiple stations in a stochastic framework to account for soil

268 moisture variability. Again, we use the EM algorithm to estimate the parameters (Dempster et
 269 al., 1977).

270 The concept of CI can be illustrated as edges in a directed acyclic graph of the GM-NHMM,
 271 as shown in Figure 3. Suppose $\mathbf{X}_{1:T} = (\mathbf{X}_1, \dots, \mathbf{X}_T)$ is a set of predictors representing soil
 272 moisture, such as rainfall, temperature, and AMSR2 soil moisture data. In a GM-NHMM, the
 273 state-transition matrix is assumed to be nonstationary, and therefore, the dynamic evolution
 274 of transition probability is a function of multivariate exogenous variables, $\mathbf{X}_{1:T}$. The GM-
 275 NHMM is then written as equation (6) (Khalil et al., 2010; Kirshner, 2005; Kwon et al., 2011,
 276 2009).

$$277 \quad p(\mathbf{SM}_{1:T}, \mathbf{S}_{1:T} | \mathbf{X}_{1:T}) = \left[p(S_1 | \mathbf{X}_1) \prod_{t=2}^T p(S_t | S_{t-1}, \mathbf{X}_t) \right] \left[\prod_{t=1}^T p(\mathbf{SM}_t | S_t) \right] \quad (6)$$

278 [Insert Figure 3]

279 In this study, we consider uniform priors, thus leading to the maximum likelihood approach
 280 to estimating a set of model parameters, $\arg \max_{\Theta} P(\mathbf{SM} | \mathbf{X}, \Theta)$. Again, note that the
 281 proposed model assumes that the observed soil moisture sequences from different years are
 282 conditionally independent. Under the GM-NHMM, the log-likelihood function $LL(\Theta)$ of
 283 the observed soil moisture data at multiple locations can be written as follows (Khalil et al.,
 284 2010):

$$285 \quad LL(\Theta) = \ln p(\mathbf{SM}_{1:T} | \mathbf{X}_{1:T}, \Theta) \\
 = \sum \ln \sum_{S_{1:T} \in \{1, \dots, K\}^T} \left[p(S_1 | \mathbf{X}_1, \Theta) \prod_{t=2}^T p(S_t | S_{t-1}, \mathbf{X}_t, \Theta) \right] \left[\prod_{t=1}^T p(\mathbf{SM}_t | S_t, \Theta) \right] \quad (7)$$

286 The parameter values cannot be obtained analytically, so we use the EM algorithm to
 287 estimate the value of the parameter vector Θ by maximizing equation (7). The EM
 288 algorithm is an iterative method for maximizing the likelihood function in a parameter space

289 Θ . Finally, the state evolutions over time in equation (6) are simulated by a multinomial
 290 logistic regression as follows (Kirshner, 2005; Kwon et al., 2011):

291 :

$$292 \quad p(S_t = \beta | S_{t-1} = \alpha, \mathbf{X}_t = \mathbf{x}) = \frac{\exp(\omega_{\alpha\beta} + \xi_{\beta}' \mathbf{x})}{\sum_{k=1}^K \exp(\omega_{\alpha\beta} + \xi_k' \mathbf{x})} \quad (8)$$

293 All the parameters ω are real, and ξ is a vector in a multi-dimensional parameter space.

294 Here, the prime denotes the transpose of the vector. Parameterization and prediction using

295 NHMM are well documented in the statistical literature and, thus, need not be elaborated

296 here. For more detailed description of the NHMM algorithm the reader is referred to Daniel

297 and Martin, (2017), Gharhramani, (2001), Rabiner, (1989), and Robertson et al., (2003).

298

299 **3.2 Ordinary Linear Regression (OLR)**

300 As a comparison to the GM-NHMM, we applied a linear regression model with the same

301 input variables used in the GM-NHMM to downscale the AMSR2 SM product for each

302 station m . Here, each parameter (β) is obtained from the least squares method. The linear

303 combination of predictors for estimating soil moisture can be written as follows:

304

$$305 \quad SM_t^m = (\beta_0^m + \beta_1^m \times R_t + \beta_2^m \times Tp_t + \beta_3^m \times ST_t) \quad (9)$$

306

307 where SM , R , and Tp are in-situ SM, rainfall, and temperature data, respectively, and ST is

308 10km AMSR2 SM data. Again note that predictor variables used here are averaged over the

309 entire region.

310

311 **4. Results and Discussion**

312 **4.1 Quantile Mapping for Bias Correction**

313 The mismatch in spatial-temporal resolution between AMSR2 SM products and in-situ

314 observations causes inevitable systematic biases. Therefore, a statistical bias correction

315 approach is commonly applied to remove the systematic bias from the satellite SM data for
316 subsequent use in either downscaling or SM modeling (Kornelsen and Coulibaly, 2015). We
317 used a quantile mapping method in which the cumulative density function of the AMSR2
318 data is matched with that of the in-situ SM observations. In this study, t location-scale (eq.
319 (10)) and gamma (eq. (11)) distributions were selected to fit the AMSR2 and in-situ soil
320 moisture data, respectively, based on the Akaike information criterion (AIC) and the
321 Bayesian information criterion (BIC), respectively, as summarized in Table 2. As shown in
322 Figure 4, the bias-corrected AMSR2 SM data exhibit enhanced variability and match well
323 with the in-situ observations. We used these bias-corrected AMSR2 SM products for our
324 subsequent analyses.

$$325 \quad f(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sigma\sqrt{\nu\pi}\Gamma(\nu/2)} \left[\frac{\nu + (\frac{x-\mu}{\sigma})^2}{\nu} \right]^{-\frac{\nu+1}{2}} \quad (10)$$

$$326 \quad f(y) = \frac{y^{\theta-1}e^{-y/\tau}}{\tau^\theta\Gamma(\theta)} \quad (11)$$

327 where μ , σ , and ν are the location, scale, and shape parameters of the t location-
328 scale distribution, respectively, and $\Gamma(\cdot)$ is the gamma function. θ and τ are the
329 shape and scale parameters of the gamma distribution, respectively.
330

331 [Insert Figure 4 and Table 2]

332

333 4.2. Predictor Selection

334 It is important to identify a suitable set of predictors that consistently influences the response
335 variables. However, in a regression model, using several predictors can cause serious
336 overfitting, which results in unrealistic predictions (Khalil et al., 2010). For a parsimonious
337 model, we consider only three predictors, daily rainfall, air temperature, and AMSR2 data, and

338 we initially evaluate the cross-correlations for all lagged orders. The correlations are
339 statistically significant and strongly persistent, as illustrated in Figure 5. Note that here the
340 values are averaged over the entire watershed for a representation. The lag-1 correlation is high
341 for daily rainfall, and the correlations appear to be consistent with the lag in the temperature
342 and AMSR2 data. Therefore, we retained a set of 1 day time-lagged values for the three
343 predictors to simulate soil moisture content in the proposed GM-NHMM.

344
345 [Insert Figure 5]
346

347 **4.3 Stochastic Modeling of Soil Moisture Using GM-NHMM**

348 The performance of the GM-NHMM is greatly influenced by the number of hidden states
349 used to represent an unobserved SM state. In this study, we estimated the number of hidden
350 states by recursively maximizing the log-likelihood (or minimizing the BIC) in the context of
351 optimization. The maximized log-likelihoods for each state are shown in Figure 6, together
352 with the minimized BIC. As shown in Figure 6(a), the log-likelihoods gradually increase with
353 the number of hidden states, but we could not clearly identify an inflection point on the curve
354 to determine the optimal number of hidden states. On the other hand, the BIC decreases
355 rapidly at 4 states, and the degree of reduction beyond 6 hidden states is negligible.
356 Therefore, we used 6 hidden states to build our stochastic soil moisture model at multiple
357 locations.

358
359 [Insert Figure 6]
360

361 For the selected 6 hidden states, the most likely temporal sequences can be efficiently
362 determined using the Viterbi algorithm (Viterbi, 1967), which calculates the probability of
363 that a hidden state will occur as well as the probability that it will transition to another state at

364 a certain date. The estimated temporal sequences of observed SM are illustrated in Figure 7,
365 and considerable inter-annual and intraseasonal variability are clearly identified. The Viterbi
366 analysis is a useful tool not only to capture intra- and inter-annual variability but also to
367 quantify its intensity. More specifically, changes in the intra-annual sequence of observed SM
368 states are shown along a horizontal line, and inter-annual variability is represented by a
369 vertical line.

370

371 [Insert Figure 7]

372

373 The degree of soil wetness and the frequencies associated with hidden states are presented in
374 Figure 8. Figure 8 (a) shows boxplots representing station-averaged SM data corresponding
375 to each state in 2014–2016. Clearly, the lower states are closely related to drier soil
376 conditions, and vice versa. Moreover, the median SM value increases largely as a function of
377 the number of states (i.e., from 21% (state 1) to 29.3 % (state 6)). The percentage of days
378 falling into the 6 hidden states for SM data across 6 stations are 14.4, 14.8, 19.5, 19.8, 20.3,
379 and 11.1 %. States 3–5 occur dominantly during the entire period, accounting for 59.6 %,
380 whereas state 6, representing the wettest soil condition, has the lowest frequency, as shown in
381 Figure 8(b). The estimated transition probabilities of the NHMM are shown in Table 3. Note
382 that the state-transition in the GM-NHMM is assumed to be nonstationary and informed by
383 exogenous variables, such as rainfall and temperature. As expected, the self-transition
384 probability (more likely to stay in the current state than to transition to a new state) is
385 noticeably high, with state 1 being the most persistent (0.93) and state 6 being the least
386 persistent (0.70).

387

388 [Insert Figure 8 and Table 3]

389

390 The temporal patterns of the simulated SM and the in-situ observations at 6 stations are
391 illustrated in Figure 9. To verify the potential of the model to reproduce the variability observed
392 in the SM data, we conducted 100 simulations. The results show a fairly good agreement with
393 the in-situ observations. Here, the proposed GM-NHMM is illustrated across the entire period
394 (2014–2016), along with the OLR model, in Figure 10. The GM-NHMM comprises the vector
395 of observed SM data from 6 stations (as dependent variables) given a vector of observed
396 covariates (as independent variables). For comparison, we built an OLR model for each station
397 using the ordinary least square method for the best-fit model of SM data. Summary statistics
398 for the comparison between the GM-NHMM and OLR are presented in Table 4, and the GM-
399 NHMM outperforms the OLR model. More specifically, the SM data simulated through the
400 GM-NHMM agree well with the in-situ observations, with correlation coefficients (r) ranging
401 from 0.73 to 0.81 (mean: 0.78), and a root mean square error (RMSE) ranging from 1.47 % to
402 2.62 % (mean: 2.06 %), whereas the OLR has much lower performance (mean r : 0.49 and mean
403 RMSE: 2.58 %).

404

[Insert Figure 9-10 and Table 4]

406

407 To further ensure that the proposed modeling scheme can predict SM, we subdivided the SM
408 data into different groups and then validated the proposed GM-NHMM using a cross-
409 validation scheme. We partitioned a sample of SM data into three different subsets
410 corresponding to the year of interest, trained the model on one subset, and then validated the
411 model with the remaining data. In other words, a set of parameters for the GM-NHMM is
412 estimated in the training period, and the identified parameters are then used to simulate SM
413 for the validation. We performed 100 simulations for each cross-validation partition for both

414 the training and validation periods. As a representative case, the simulated SM values for 6
415 stations are compared with the values observed at those stations for the training period
416 (2014–2015) and the testing period (2016) in Figure 11. The SM data are reasonably well
417 reproduced by the proposed GM-NHMM for both the training and testing phases. The results
418 of the cross-validation using the GM-NHMM for the different partitions are summarized in
419 Table 5. We considered three goodness- of- fit measures, correlation coefficient (r), RMSE,
420 and bias, in evaluating the models. During the training periods, the 6-station averaged
421 correlation coefficient values range from 0.72 to 0.80, whereas during the validation period,
422 the r values show slightly lower correlations than during the training period. However, the
423 GM-NHMM can clearly generate the intraseasonal sequence of daily SM fairly well, and
424 other measures also show reasonable performance at multiple locations, leading to higher
425 correlations with the observed SM data. The RMSE and bias values are also generally better
426 for the training period than the validation period.

427

428 [Insert Figure 11 and Table 5]

429

430 For a multisite SM simulator, it is of particular importance to correctly reproduce the spatial
431 coherence of daily SM across multiple stations. Therefore, we estimated the spatial
432 correlations of the sequence of daily SM and compared them with the observed values. As
433 shown in Figure 12, the spatial correlations across the stations are reasonably well reproduced
434 by proposed GM-NHMM model.

435

436 [Insert Figure 12]

437

438 Table 6 shows the results of applying the GM-NHMM with different combinations of
439 predictors to examine the contribution of the AMSR2 SM data to the proposed model. The
440 use of rainfall and temperature without the AMSR2 data (case-1) led to a slightly lower
441 correlation coefficient of 0.73, compared to the results obtained with all three predictors
442 shown in Table 5. On the other hand, there was no significant change in the correlation
443 coefficient of 0.63 when we used rainfall alone as a predictor (case-2). Furthermore, we
444 found a similar trend in our cross-validation analysis. Therefore, the 1 day time-lagged
445 rainfall data might be the main factor in properly reproducing SM dynamics. Nonetheless,
446 combining rainfall with temperature and AMSR2 still yielded the highest correlation with the
447 in-situ observations.

448 [Insert Table 6]

449

450 **5. Concluding Remarks**

451 We have here presented a stochastic soil moisture estimation model based on a GM-NHMM
452 to spatially disaggregate AMSR2 SM data at multiple locations in the context of
453 downscaling. Given the close relationship with SM, we considered both rainfall and air
454 temperature as potential predictors in the proposed stochastic downscaling model. We used 1
455 day time-lagged values for the three predictors to simulate SM in the proposed GM-NHMM
456 model. Before applying the proposed downscaling scheme, we used the quantile mapping
457 approach to reduce the systematic bias in the AMSR2 SM products, and we then used those
458 bias-corrected AMSR2 SM products for subsequent analyses. In GM-NHMM terms, we
459 formulated a six-state model with three predictors representing an unobserved SM state based
460 on the BIC. The temporal sequences of unobserved hidden states and the dynamic evolution
461 of transition probability were estimated by the Viterbi algorithm. Consequently, the proposed
462 GM-NHMM was applied to simulate fine-resolution SM products in a multivariate

463 framework. We compared our results with in-situ observations from the Yongdam dam
464 watershed in South Korea. The key results obtained are summarized as follows.

- 465 1. The estimated small set of hidden states that most likely corresponds to localized soil
466 moisture dynamics is effectively captured and accounts for a certain fraction of the
467 soil moisture process, which improves understanding of the intraseasonal and inter-
468 annual variability of SM dynamics. Based on the identified state transition-
469 probability matrix, self-transitions are more significant than the probability of
470 transitioning to other states, indicating that the states seem to be persistent over time
471 due to the slow-varying behavior of basin-scale SM (Botter et al., 2007).
- 472 2. Given the relatively short length of the in-situ SM time series data, we considered a
473 cross-validation performance assessment of the simulations. The rainfall predictor
474 plays a substantial role in achieving overall predictability. Adding temperature and
475 AMSR2 data as predictors improves the fit to the SM data. Therefore, weather variables
476 (i.e., rainfall and temperature) could be effective in picking up some of the
477 predictability of local SM that is not captured by AMSR2 data. On the other hand,
478 large-scale dynamic features identified in remote-sensed SM data seem to facilitate the
479 identification of other SM states with well-defined regional spatial patterns. The results
480 presented here illustrate the potential of a stochastic model with a climate-predictor-
481 based forecast. However, the relatively small improvement in forecast skill that the
482 AMSR2 SM products offer in the model suggests that the AMSR2 data might not
483 sufficiently reflect the regional or seasonal characteristics of this study area.
- 484 3. We compared the efficiency of the proposed model with that of an ordinary regression
485 model using the same predictors. The mean correlation coefficient for the GM-NHMM
486 obtained by averaging over all the stations is about 0.78, which is significantly greater
487 than that of the OLR, about 0.23. The proposed model also yields a noticeable reduction

488 in RMSE. Moreover, the proposed GM-NHMM method not only provides a better
489 representation of the observed SM than the OLR model but also preserves spatial
490 coherence across all the stations, which is a fundamentally important property in
491 describing the spatial pattern of soil moisture and its association with runoff on a
492 catchment scale.

493 Our main contributions in this study are our insights into the soil moisture process and its
494 potential predictability, leading to the way for more applications in hydrologic studies. We
495 expect that future work will address this study's shortcomings with respect to the use of
496 satellite-based products and predictor selection and further investigate cross-validation
497 assessment of forecasts for different regions over a longer period of record, which are
498 required to support these applications.

499

500

501 **Appendix A**

List of Abbreviations	
AIC	Akaike information criterion
AMSR2	Advanced Microwave Scanning Radiometer 2
ASCAT	Advanced Scatterometer
BIC	Bayesian information criterion
CI	Conditional independence
EM	Expectation-maximization
GM-NHMM	Gaussian mixture nonstationary hidden Markov model
HMM	Hidden Markov model
JAXA	Japan Aerospace Exploration Agency
MODIS	Moderate Resolution Imaging Spectroradiometer
OLR	Ordinary regression model
r	Correlation coefficient
RMSE	Root mean square error
SM	Soil moisture
SMAP	Soil Moisture Active Passive
SMOS	Soil Moisture Ocean Salinity Satellite

502

503

504 **Appendix B**

List of Symbols	
T_p	Temperature
R	Rainfall
ST	AMSR2 SM data
$\Gamma(\cdot)$	Gamma function
θ	Shape parameter of the gamma distribution
ν	Shape parameter of the t location-scale distribution
τ	Scale parameter of the gamma distribution
\mathbf{SM}_t^M	M-dimensional vector of in-situ soil moisture measurements at time t.
$S_{\nu\tau}$	Finite number of hidden states
X	A set of predictors
$LL(\Theta)$	Log-likelihood function
μ	Location parameter of the t location-scale distribution
σ	Scale parameter of the t location-scale distribution

505

506

507

508 **Acknowledgement**

509 This research was supported by a grant (17AWMP-B121100-02) from the Water
510 Management Research Program funded by Ministry of Land, Infrastructure and Transport of
511 Korean government.

512

513

514 **References**

- 515 Albergel, C., Rüdiger, C., Carrer, D., Calvet, J.-C., Fritz, N., Naeimi, V., Bartalis, Z.,
516 Hasenauer, S., 2008. An evaluation of ASCAT surface soil moisture products with in-
517 situ observations in southwestern France. *Hydrol. Earth Syst. Sci. Discuss.* 5, 2221–
518 2250. doi:10.5194/hessd-5-2221-2008
- 519 Barrett, B., Petropoulos, G., 2013. Satellite Remote Sensing of Surface Soil Moisture.
520 *Remote Sens. Energy Fluxes Soil Moisture Content* 85–120. doi:doi:10.1201/b15610-6
- 521 Bilmes, J.A., 1998. A Gentle Tutorial of the EM Algorithm and its Application to Parameter
522 Estimation for Gaussian Mixture and Hidden Markov Models, International Computer
523 Science Institute. doi:10.1016/S0550-3213(97)00753-0
- 524 Bishop, C.M., 2006. *Pattern Recognition and Machine Learning*. Springer.
525 doi:10.1117/1.2819119
- 526 Botter, G., Porporato, A., Rodriguez-Iturbe, I., Rinaldo, A., 2007. Basin-scale soil moisture
527 dynamics and the probabilistic characterization of carrier hydrologic flows: Slow,
528 leaching-prone components of the hydrologic response. *Water Resour. Res.* 43, 1–14.
529 doi:10.1029/2006WR005043
- 530 Brocca, L., Ciabatta, L., Massari, C., Camici, S., Tarpanelli, A., 2017a. Soil Moisture for
531 Hydrological Applications: Open Questions and New Opportunities. *Water* 9, 140.
532 doi:10.3390/w9020140
- 533 Brocca, L., Crow, W.T., Ciabatta, L., Massari, C., De Rosnay, P., Enenkel, M., Hahn, S.,
534 Amarnath, G., Camici, S., Tarpanelli, A., Wagner, W., 2017b. A Review of the
535 Applications of ASCAT Soil Moisture Products. *IEEE J. Sel. Top. Appl. Earth Obs.*
536 *Remote Sens.* 10, 2285–2306. doi:10.1109/JSTARS.2017.2651140
- 537 Brocca, L., Hasenauer, S., Lacava, T., Melone, F., Moramarco, T., Wagner, W., Dorigo, W.,
538 Matgen, P., Martínez-Fernández, J., Llorens, P., Latron, J., Martin, C., Bittelli, M., 2011.
539 Soil moisture estimation through ASCAT and AMSR-E sensors: An intercomparison
540 and validation study across Europe. *Remote Sens. Environ.* 115, 3390–3408.
541 doi:10.1016/j.rse.2011.08.003
- 542 Busch, F.A., Niemann, J.D., Coleman, M., 2012. Evaluation of an empirical orthogonal
543 function-based method to downscale soil moisture patterns based on topographical
544 attributes. *Hydrol. Process.* 26, 2696–2709. doi:10.1002/hyp.8363
- 545 Carreau, J., Naveau, P., Sauquet, E., 2009. A statistical rainfall-runoff mixture model with
546 heavy-tailed components. *Water Resour. Res.* 45. doi:10.1029/2009WR007880

547 Cenci, L., Laiolo, P., Gabellani, S., Campo, L., Silvestro, F., Delogu, F., Boni, G., Rudari, R.,
548 2016. Assimilation of H-SAF Soil Moisture Products for Flash Flood Early Warning
549 Systems. Case Study: Mediterranean Catchments. *IEEE J. Sel. Top. Appl. Earth Obs.*
550 *Remote Sens.* PP, 5634–5646. doi:10.1109/JSTARS.2016.2598475

551 Cioffi, F., Conticello, F., Lall, U., Marotta, L., Telesca, V., 2017. Large scale climate and
552 rainfall seasonality in a Mediterranean Area: Insights from a non-homogeneous Markov
553 model applied to the Agro-Pontino plain. *Hydrol. Process.* 31, 668–686.
554 doi:10.1002/hyp.11061

555 Daniel, J., Martin, H., 2017. Speech and Language Processing, in: Ch 9. Stanford University.
556 doi:10.1007/978-1-61779-400-1_22

557 Das, N.N., Entekhabi, D., Njoku, E.G., 2011. An Algorithm for Merging SMAP Radiometer
558 and Radar Data for High Resolution Soil Moisture Retrieval. *IEEE Trans. Geosci.*
559 *Remote Sens.* in press, 1–9. doi:10.1109/TGRS.2010.2089526

560 Dempster, A.P., Laird, N.M., Rubin, D.B., 1977. Maximum likelihood from incomplete data
561 via the EM algorithm. *J. R. Stat. Soc. Ser. B Methodol.* 39, 1–38.
562 doi:http://dx.doi.org/10.2307/2984875

563 Djamai, N., Magagi, R., Goïta, K., Merlin, O., Kerr, Y., Roy, A., 2016. A combination of
564 DISPATCH downscaling algorithm with CLASS land surface scheme for soil moisture
565 estimation at fine scale during cloudy days. *Remote Sens. Environ.* 184, 1–14.
566 doi:10.1016/j.rse.2016.06.010

567 Fang, B., Lakshmi, V., 2014. Soil moisture at watershed scale: Remote sensing techniques. *J.*
568 *Hydrol.* 516, 258–272. doi:10.1016/j.jhydrol.2013.12.008

569 Gauvain, J.-L., Lee, C.-H., 1994. Maximum a posteriori estimation for multivariate Gaussian
570 mixture observations of Markov chains. *IEEE Trans. Speech Audio Process.* 2, 291–298.
571 doi:10.1109/89.279278

572 Gharhramani, Z., 2001. An Introduction to Hidden Markov Models and Bayesian Networks.
573 *J. Pattern Recognit. Artif. Intell.* 15, 9–42.

574 Griesfeller, A., Lahoz, W. a., Jeu, R.A.M. d., Dorigo, W., Haugen, L.E., Svendby, T.M.,
575 Wagner, W., 2016. Evaluation of satellite soil moisture products over Norway using
576 ground-based observations. *Int. J. Appl. Earth Obs. Geoinf.* 45, 155–164.
577 doi:10.1016/j.jag.2015.04.016

578 Hughes, J.P., Guttorp, P., 1994. A class of stochastic models for relating synoptic
579 atmospheric patterns to regional hydrologic phenomena. *Water Resour. Res.* 30, 1535–
580 1546. doi:10.1029/93WR02983

581 Hughes, J.P., Guttorp, P., Charles, S.P., 1999. A non-homogeneous hidden Markov model for
582 precipitation occurrence. *J. R. Stat. Soc. Ser. C (Applied Stat.* 48, 15–30.
583 doi:10.1111/1467-9876.00136

584 Im, J., Park, S., Rhee, J., Baik, J., Choi, M., 2016. Downscaling of AMSR-E soil moisture
585 with MODIS products using machine learning approaches. *Environ. Earth Sci.* 75, 1–19.
586 doi:10.1007/s12665-016-5917-6

587 Imaoka, K., Kachi, M., Fujii, H., Murakami, H., Hori, M., Ono, A., Igarashi, T., Nakagawa,
588 K., Oki, T., Honda, Y., Shimoda, H., 2010. Global change observation mission (GCOM)
589 for monitoring carbon, water cycles, and climate change. *Proc. IEEE* 98, 717–734.
590 doi:10.1109/JPROC.2009.2036869

591 Kerr, Y.H., Waldteufel, P., Richaume, P., Wigneron, J.P., Ferrazzoli, P., Mahmoodi, A.,
592 Bitar, A. Al, Cabot, F., Gruhier, C., Juglea, S.E., Leroux, D., Mialon, A., Delwart, S.,
593 2012. The SMOS Soil Moisture Retrieval Algorithm. *Geosci. Remote Sens.* 50, 1384–
594 1403. doi:10.1109/TGRS.2012.2184548

595 Khalil, A.F., Kwon, H.-H., Lall, U., Kaheil, Y.H., 2010. Predictive downscaling based on
596 non-homogeneous hidden Markov models. *Hydrol. Sci. J.* 55, 333–350.
597 doi:10.1080/02626661003780342

598 Kirshner, S., 2005. Modeling of multivariate time series using hidden Markov models.
599 University of California, Irvine.

600 Koike, T., 2013. Description of GCOM-W1 AMSR2 Soil Moisture Algorithm, in:
601 Descriptions of GCOM-W1 AMSR2 Level 1R and Level 2 Algorithms. Japan
602 Aerospace Exploration Agency Earth Observation Research Center, p. 8.1-8.13.

603 Kornelsen, K.C., Coulibaly, P., 2015. Reducing multiplicative bias of satellite soil moisture
604 retrievals. *Remote Sens. Environ.* 165, 109–122. doi:10.1016/j.rse.2015.04.031

605 Kwon, H.H., Lall, U., Obeysekera, J., 2009. Simulation of daily rainfall scenarios with
606 interannual and multidecadal climate cycles for South Florida. *Stoch. Environ. Res. Risk*
607 *Assess.* 23, 879–896. doi:10.1007/s00477-008-0270-2

608 Kwon, H.H., Sivakumar, B., Moon, Y. Il, Kim, B.S., 2011. Assessment of change in design
609 flood frequency under climate change using a multivariate downscaling model and a

610 precipitation-runoff model. *Stoch. Environ. Res. Risk Assess.* 25, 567–581.
611 doi:10.1007/s00477-010-0422-z

612 Lakshmanan, V., Kain, J.S., 2010. A Gaussian Mixture Model Approach to Forecast
613 Verification. *Weather Forecast.* 25, 908–920. doi:10.1175/2010WAF2222355.1

614 Liu, Y.Y., Parinussa, R.M., Dorigo, W.A., De Jeu, R.A.M., Wagner, W., M. Van Dijk, A.I.J.,
615 McCabe, M.F., Evans, J.P., 2011. Developing an improved soil moisture dataset by
616 blending passive and active microwave satellite-based retrievals. *Hydrol. Earth Syst.*
617 *Sci.* 15, 425–436. doi:10.5194/hess-15-425-2011

618 Mehrotra, R., Sharma, A., 2010. Development and application of a multisite rainfall
619 stochastic downscaling framework for climate change impact assessment. *Water Resour.*
620 *Res.* 46, 1–17. doi:10.1029/2009WR008423

621 Mehrotra, R., Sharma, A., 2006. A nonparametric stochastic downscaling framework for
622 daily rainfall at multiple locations. *J. Geophys. Res. Atmos.* 111, 1–16.
623 doi:10.1029/2005JD006637

624 Mehrotra, R., Sharma, A., 2005. A nonparametric nonhomogeneous hidden Markov model
625 for downscaling of multisite daily rainfall occurrences. *J. Geophys. Res. D Atmos.* 110,
626 1–13. doi:10.1029/2004JD005677

627 Merlin, O., Rüdiger, C., Al Bitar, A., Richaume, P., Walker, J.P., Kerr, Y.H., 2012.
628 Disaggregation of SMOS soil moisture in Southeastern Australia. *IEEE Trans. Geosci.*
629 *Remote Sens.* 50, 1556–1571. doi:10.1109/TGRS.2011.2175000

630 Njoku, E.G., Ashcroft, P., Chan, T.K., Li, L., 2005. Global survey and statistics of radio-
631 frequency interference in AMSR-E land observations. *IEEE Trans. Geosci. Remote*
632 *Sens.* 43, 938–946. doi:10.1109/TGRS.2004.837507

633 Owe, M., de Jeu, R., Holmes, T., 2008. Multisensor historical climatology of satellite-derived
634 global land surface moisture. *J. Geophys. Res. Earth Surf.* 113, 1–17.
635 doi:10.1029/2007JF000769

636 Parajka, J., Naeimi, V., Blöschl, G., Wagner, W., Merz, R., Scipal, K., 2006. Assimilating
637 scatterometer soil moisture data into conceptual hydrologic models at the regional scale.
638 *Hydrol. Earth Syst. Sci.* 10, 353–368. doi:10.5194/hessd-2-2739-2005

639 Park, S., Park, S., Im, J., Rhee, J., Shin, J., Park, J., 2017. Downscaling GLDAS Soil
640 Moisture Data in East Asia through Fusion of Multi-Sensors by Optimizing Modified
641 Regression Trees. *Water* 9, 332. doi:10.3390/w9050332

642 Peng, J., Loew, A., Merlin, O., Verhoest, N.E.C., 2017. A review of spatial downscaling of
643 satellite remotely sensed soil moisture. *Rev. Geophys.* 1–26.
644 doi:10.1002/2016RG000543

645 Peng, J., Loew, A., Zhang, S., Wang, J., 2016. Spatial downscaling of global satellite soil
646 moisture data using temperature vegetation dryness index. *IEEE Trans. Geosci. Remote*
647 *Sens.* 1, 558–566.

648 Peng, J., Niesel, J., Loew, A., 2015. Evaluation of soil moisture downscaling using a simple
649 thermal-based proxy-the REMEDHUS network (Spain) example. *Hydrol. Earth Syst.*
650 *Sci.* 19, 4765–4782. doi:10.5194/hess-19-4765-2015

651 Piles, M., Camps, A., Vall-llossera, M., Corbella, I., Panciera, R., Rüdiger, C., Kerr, Y.H.,
652 Walker, J., 2011. Downscaling SMOS-Derived Soil Moisture Using MODIS Visible /
653 Infrared Data. *IEEE Trans. Geosci. Remote Sens.* 49, 3156–3166.

654 Piles, M., Sánchez, N., Vall-Llossera, M., Camps, A., Martínez-Fernandez, J., Martínez, J.,
655 Gonzalez-Gambau, V., 2014. A downscaling approach for SMOS land observations:
656 Evaluation of high-resolution soil moisture maps over the Iberian peninsula. *IEEE J. Sel.*
657 *Top. Appl. Earth Obs. Remote Sens.* 7, 3845–3857. doi:10.1109/JSTARS.2014.2325398

658 Rabiner, L.R., 1989. A Tutorial on Hidden Markov Models and Selected Applications in
659 Speech Recognition. *Proc. IEEE.* doi:10.1109/5.18626

660 Ranney, K.J., Niemann, J.D., Lehman, B.M., Green, T.R., Jones, A.S., 2015. A method to
661 downscale soil moisture to fine resolutions using topographic, vegetation, and soil data.
662 *Adv. Water Resour.* 76, 81–96. doi:10.1016/j.advwatres.2014.12.003

663 Reichle, R.H., Koster, R.D., Liu, P., Mahanama, S.P.P., Njoku, E.G., Owe, M., 2007.
664 Comparison and assimilation of global soil moisture retrievals from the Advanced
665 Microwave Scanning Radiometer for the Earth Observing System (AMSR-E) and the
666 Scanning Multichannel Microwave Radiometer (SMMR). *J. Geophys. Res. Atmos.* 112,
667 1–14. doi:10.1029/2006JD008033

668 Ridolfi, L., D’Odorico, P., Porporato, A., Rodriguez-Iturbe, I., 2003. Stochastic soil moisture
669 dynamics along a hillslope. *J. Hydrol.* 272, 264–275. doi:10.1016/S0022-
670 1694(02)00270-6

671 Rings, J., Vrugt, J.A., Schoups, G., Huisman, J.A., Vereecken, H., 2012. Bayesian model
672 averaging using particle filtering and Gaussian mixture modeling: Theory, concepts, and
673 simulation experiments. *Water Resour. Res.* 48. doi:10.1029/2011WR011607

674 Robertson, A.W., Kirshner, S., Smyth, P., 2004. Downscaling of daily rainfall occurrence
675 over Northeast Brazil using a hidden Markov model. *J. Clim.* 17, 4407–4424.
676 doi:10.1175/JCLI-3216.1

677 Robertson, A.W., Kirshner, S., Smyth, P., Charles, S.P., Bates, B.C., 2006. Subseasonal-to-
678 interdecadal variability of the Australian monsoon over North Queensland. *Q. J. R.*
679 *Meteorol. Soc.* 132, 519–542. doi:10.1256/qj.05.75

680 Robertson, A.W., Sergey, K., Padhraic, S., 2003. Hidden Markov models for modeling daily
681 rainfall occurrence over Brazil, Technical Report UCI-ICS 03-27. Information and
682 Computer Science University of California, Irvine.

683 Ryu, D., Famiglietti, J.S., 2005. Characterization of footprint-scale surface soil moisture
684 variability using Gaussian and beta distribution functions during the Southern Great
685 Plains 1997 (SGP97) hydrology experiment. *Water Resour. Res.* 41, 1–13.
686 doi:10.1029/2004WR003835

687 Smyth, P., Heckerman, D., Jordan, M.I., 1997. Probabilistic Independence Networks for
688 Hidden Markov Probability Models. *Neural Comput.* 9, 227–269.
689 doi:10.1162/neco.1997.9.2.227

690 Srivastava, P.K., Han, D., Ramirez, M.R., Islam, T., 2013. Machine Learning Techniques for
691 Downscaling SMOS Satellite Soil Moisture Using MODIS Land Surface Temperature
692 for Hydrological Application. *Water Resour. Manag.* 27, 3127–3144.
693 doi:10.1007/s11269-013-0337-9

694 Stehlík, J., Bárdossy, A., 2002. Multivariate stochastic downscaling model for generating
695 daily precipitation series based on atmospheric circulation. *J. Hydrol.* 256, 120–141.
696 doi:10.1016/S0022-1694(01)00529-7

697 Topp, G.C., Davis, J.L., Annan, A.P., 1980. Electromagnetic Determination of Soil Water
698 Content: Measurements in Coaxial Transmission Lines. *Water Resour. Res.* 16, 574–
699 582. doi:10.1029/WR016i003p00574

700 Verhoest, N.E.C., Van Den Berg, M.J., Martens, B., Lievens, H., Wood, E.F., Pan, M., Kerr,
701 Y.H., Al Bitar, A., Tomer, S.K., Drusch, M., Vernieuwe, H., De Baets, B., Walker, J.P.,
702 Dumedah, G., Pauwels, V.R.N., 2015. Copula-based downscaling of coarse-scale soil
703 moisture observations with implicit bias correction. *IEEE Trans. Geosci. Remote Sens.*
704 53, 3507–3521. doi:10.1109/TGRS.2014.2378913

705 Vilasa, L., Miralles, D.G., de Jeu, R.A.M., Dolman, A.J., 2017. Global soil moisture
706 bimodality in satellite observations and climate models. *J. Geophys. Res. Atmos.* 122,
707 4299–4311. doi:10.1002/2016JD026099

708 Viterbi, A., 1967. Error bounds for convolutional codes and an asymptotically optimum
709 decoding algorithm. *IEEE Trans. Inf. Theory* 13, 260–269.
710 doi:10.1109/TIT.1967.1054010

711 Xing, C., Chen, N., Zhang, X., Gong, J., 2017. A Machine Learning Based Reconstruction
712 Method for Satellite Remote Sensing of Soil Moisture Images with In Situ Observations.
713 *Remote Sens.* 9, 484. doi:10.3390/rs9050484

714 Yoo, J., Kwon, H.H., So, B.J., Rajagopalan, B., Kim, T.W., 2015. Identifying the role of
715 typhoons as drought busters in South Korea based on hidden Markov chain models.
716 *Geophys. Res. Lett.* 42, 2797–2804. doi:10.1002/2015GL063753

717 Zeng, J., Li, Z., Chen, Q., Bi, H., Qiu, J., Zou, P., 2015. Evaluation of remotely sensed and
718 reanalysis soil moisture products over the Tibetan Plateau using in-situ observations.
719 *Remote Sens. Environ.* 163, 91–110. doi:10.1016/j.rse.2015.03.008

720 Zhao, W., Li, A., 2013. A downscaling method for improving the spatial resolution of
721 AMSR-E derived soil moisture product based on MSG-SEVIRI data. *Remote Sens.* 5,
722 6790–6811. doi:10.3390/rs5126790

723 Zhuo, L., Han, D., 2016. Could operational hydrological models be made compatible with
724 satellite soil moisture observations? *Hydrol. Process.* 30, 1637–1648.
725 doi:10.1002/hyp.10804
726

Table 1. Specification and characteristics of soil observation sites in the Yongdam dam watershed. Site	Elevation	Longitude	Latitude	Annual rainfall	Observation	Land Cover
	(m a.s.l)	(°)	(°)	(mm/yr)	depth (cm)	
SM & Rainfall						
Station 1	313	127.55	35.87	1,107	10, 20, 40, 60	Forest
Station 2	330	127.43	35.97	1,224	10, 20, 40, 60	Forest
Station 3	396	127.4	35.86	1,191	10, 20, 40, 60	Forest
Station 4	334	127.49	35.8	1,120	10, 20, 40, 60	Agriculture
Station 5	453	127.63	35.81	1,049	10, 20, 40, 60	Agriculture
Station 6	409	127.51	35.68	1,193	10, 20, 40, 60	Forest
Temperature						
Jangsu	406	127.52	35.66	-	-	-

727

728 Table 2. BIC and AIC scores with respect to distribution models.

In-situ			AMSR2		
Distribution	BIC	AIC	Distribution	BIC	AIC
Gamma	44,677	44,663	t-location scale	31,445	31,425
Log-logistic	45,051	45,037	Log-logistic	32,316	32,303
Normal	45,128	45,114	Gamma	36,550	36,536
t-location scale	45,137	45,116	Weibull	38,680	38,666
Weibull	45,259	45,246	Normal	43,660	43,646

729

730 Table 3. Transition probability matrix of 6 hidden states for soil moisture at 6 stations in the
 731 Yongdam watershed.

Site	Station 1	Station 2	Station 3	Station 4	Station 5	Station 6
Station 1	0.93	0.01	0.05	0.00	0.00	0.01
Station 2	0.02	0.90	0.01	0.04	0.00	0.04
Station 3	0.04	0.03	0.88	0.00	0.03	0.02
Station 4	0.00	0.04	0.00	0.92	0.02	0.02
Station 5	0.00	0.00	0.07	0.05	0.79	0.09
Station 6	0.00	0.00	0.00	0.00	0.30	0.70

732

733

734 Table 4. Comparison between in-situ and simulated SM.

Site	BC AMSR2		GM-NHMM		OLR	
	<i>r</i>	RMSE (%)	<i>r</i>	RMSE (%)	<i>r</i>	RMSE (%)
Station 1	0.34	4.55	0.79	2.62	0.49	3.36
Station 2	0.10	4.07	0.78	2.02	0.55	2.42
Station 3	0.31	2.55	0.73	1.52	0.49	1.83
Station 4	0.38	2.54	0.81	1.47	0.54	1.86
Station 5	0.17	4.34	0.79	2.22	0.41	2.95
Station 6	0.10	4.93	0.79	2.50	0.48	3.06
Average	0.23	3.83	0.78	2.06	0.49	2.58

735

736

737

738

739

740 Table 5. Comparison between in-situ and simulated SM.

741

Site	Training (2014–2015)			Validation (2016)			Training (2015–2016)			Validation (2014)			Training (2014, 2016)			Validation (2015)		
	<i>r</i>	RMSE (%)	Bias	<i>r</i>	RMSE (%)	Bias	<i>r</i>	RMSE (%)	Bias	<i>r</i>	RMSE (%)	Bias	<i>r</i>	RMSE (%)	Bias	<i>r</i>	RMSE (%)	Bias
Station 1	0.79	2.69	0.43	0.80	2.47	0.32	0.83	2.14	0.32	0.62	3.34	0.26	0.77	2.79	0.72	0.68	3.14	1.53
Station 2	0.79	2.10	0.66	0.75	1.85	0.19	0.86	1.65	0.37	0.63	2.57	1.36	0.73	2.25	0.83	0.86	2.26	0.69
Station 3	0.76	1.49	0.28	0.67	1.57	0.10	0.75	1.45	0.25	0.69	1.65	0.00	0.68	1.70	0.44	0.74	1.80	1.02
Station 4	0.80	1.57	0.24	0.83	1.24	0.06	0.73	1.44	0.20	0.74	1.86	0.06	0.76	1.57	0.35	0.59	1.91	0.99
Station 5	0.79	2.37	0.67	0.78	1.89	0.23	0.76	2.18	0.48	0.60	2.63	0.19	0.65	2.54	0.87	0.66	3.47	2.11
Station 6	0.83	2.23	0.68	0.73	2.96	1.04	0.88	1.88	0.41	0.68	2.43	0.30	0.71	2.80	1.01	0.86	2.60	1.27
Average	0.79	2.08	0.49	0.76	2.00	0.33	0.80	1.79	0.34	0.66	2.41	0.36	0.72	2.28	0.70	0.73	2.53	1.27

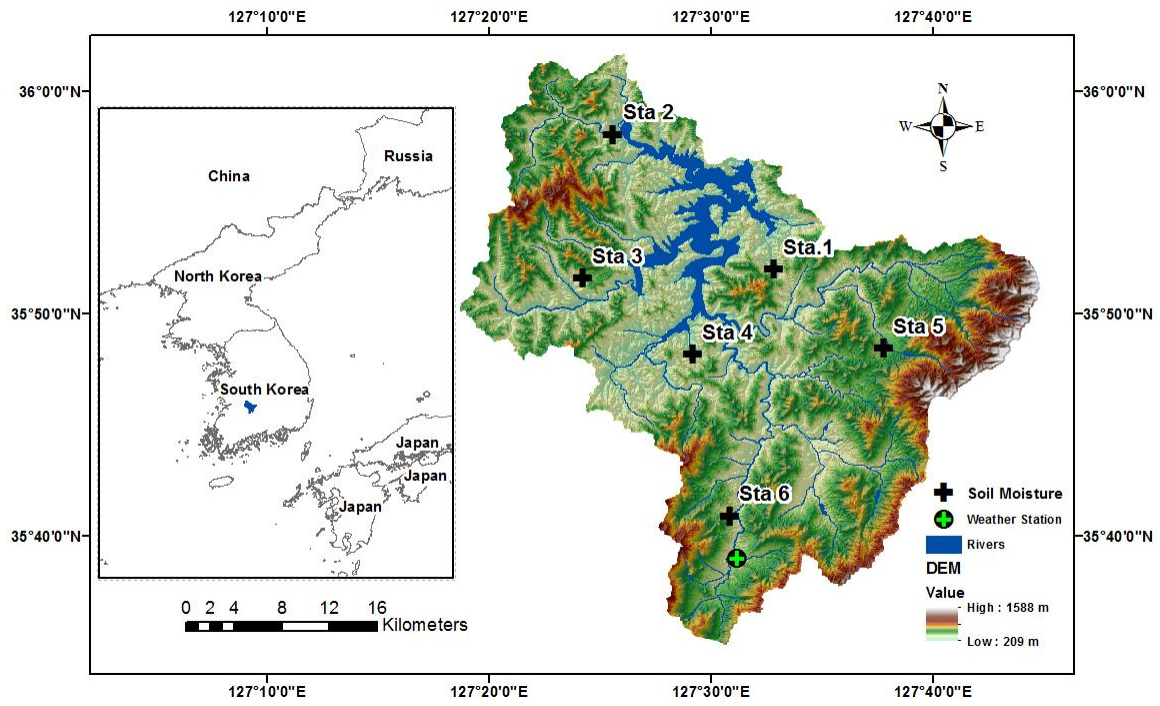
742 Table 6. Comparison of r values with respect to different combinations of predictors.

Sta. No	Modeling	Cross Validation					
	Entire period (2014–2016)	Training (2014– 2015)	Validation (2016)	Training (2015– 2016)	Validation (2014)	Training (2014, 2016)	Validation (2015)
(Case 1) Predictors: Rainfall, Temperature							
Station 1	0.75	0.76	0.72	0.76	0.64	0.71	0.59
Station 2	0.73	0.74	0.70	0.84	0.60	0.56	0.74
Station 3	0.63	0.70	0.48	0.69	0.66	0.52	0.65
Station 4	0.78	0.78	0.79	0.70	0.71	0.81	0.66
Station 5	0.73	0.75	0.68	0.70	0.54	0.64	0.42
Station 6	0.75	0.77	0.72	0.87	0.60	0.55	0.66
Average	0.73	0.75	0.68	0.76	0.63	0.63	0.62
(Case 2) Predictor: Rainfall							
Station 1	0.78	0.78	0.79	0.72	0.81	0.80	0.70
Station 2	0.39	0.45	0.22	0.23	0.51	0.38	0.47
Station 3	0.62	0.66	0.53	0.57	0.67	0.56	0.62
Station 4	0.81	0.80	0.84	0.79	0.83	0.84	0.70
Station 5	0.62	0.61	0.64	0.49	0.75	0.67	0.53
Station 6	0.57	0.61	0.50	0.49	0.67	0.59	0.63
Average	0.63	0.65	0.58	0.55	0.71	0.64	0.61

743

744

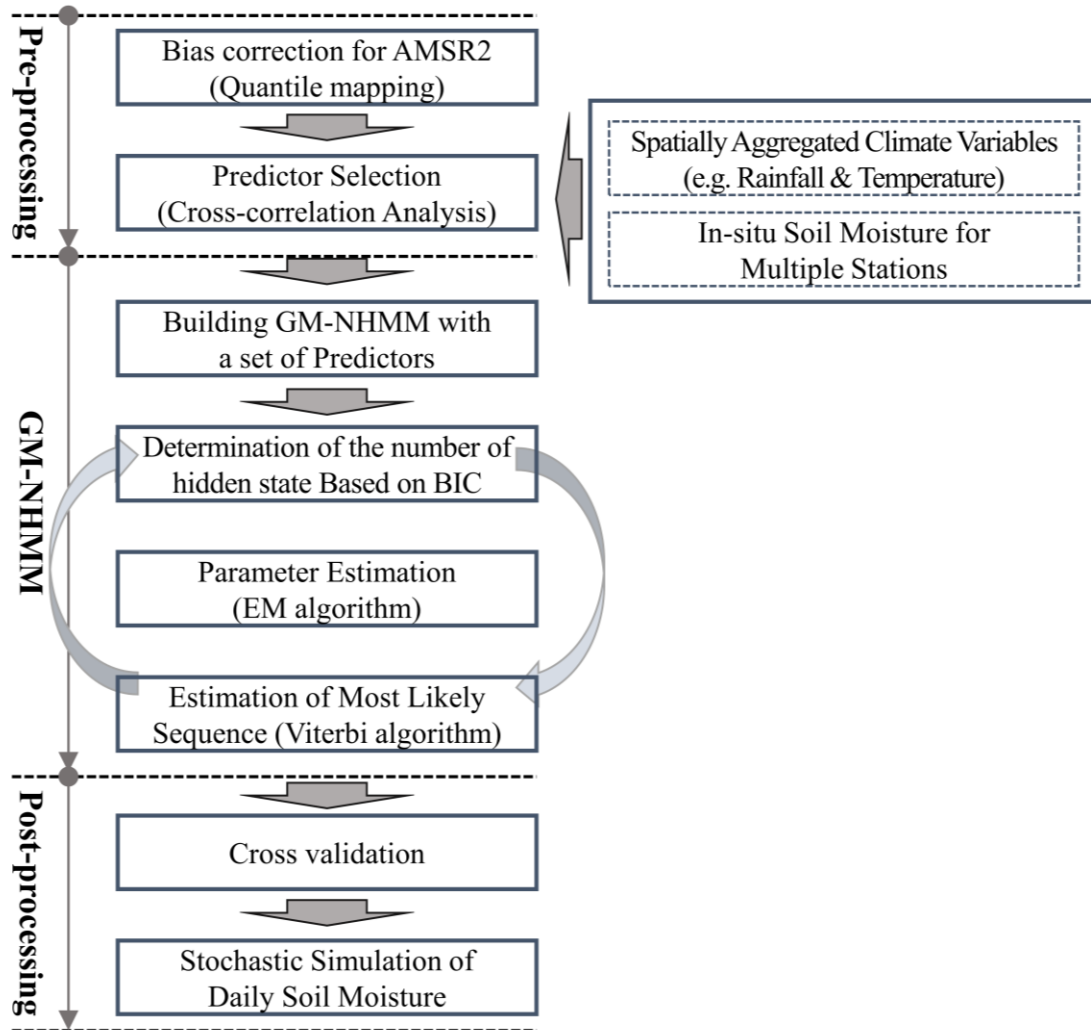
745



746

747 Figure 1. The study site with topography and observation stations.

748

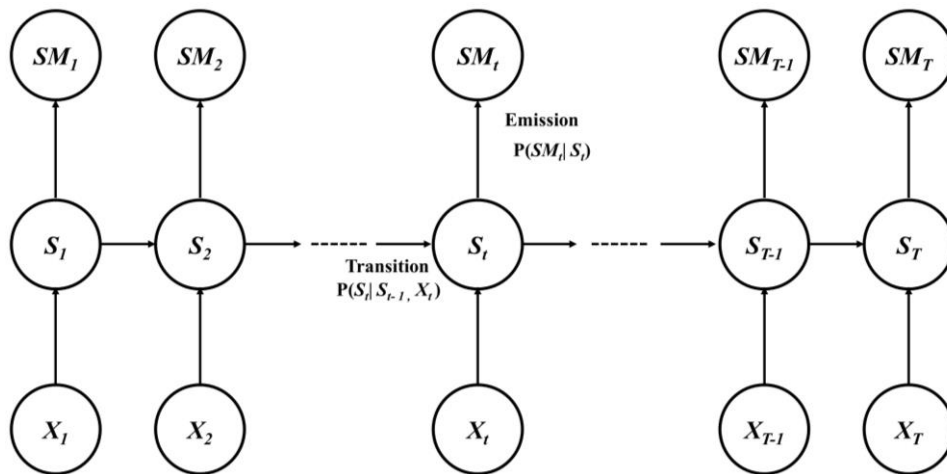


749

750 Figure 2. Schematic diagram representing the processing steps.

751

752



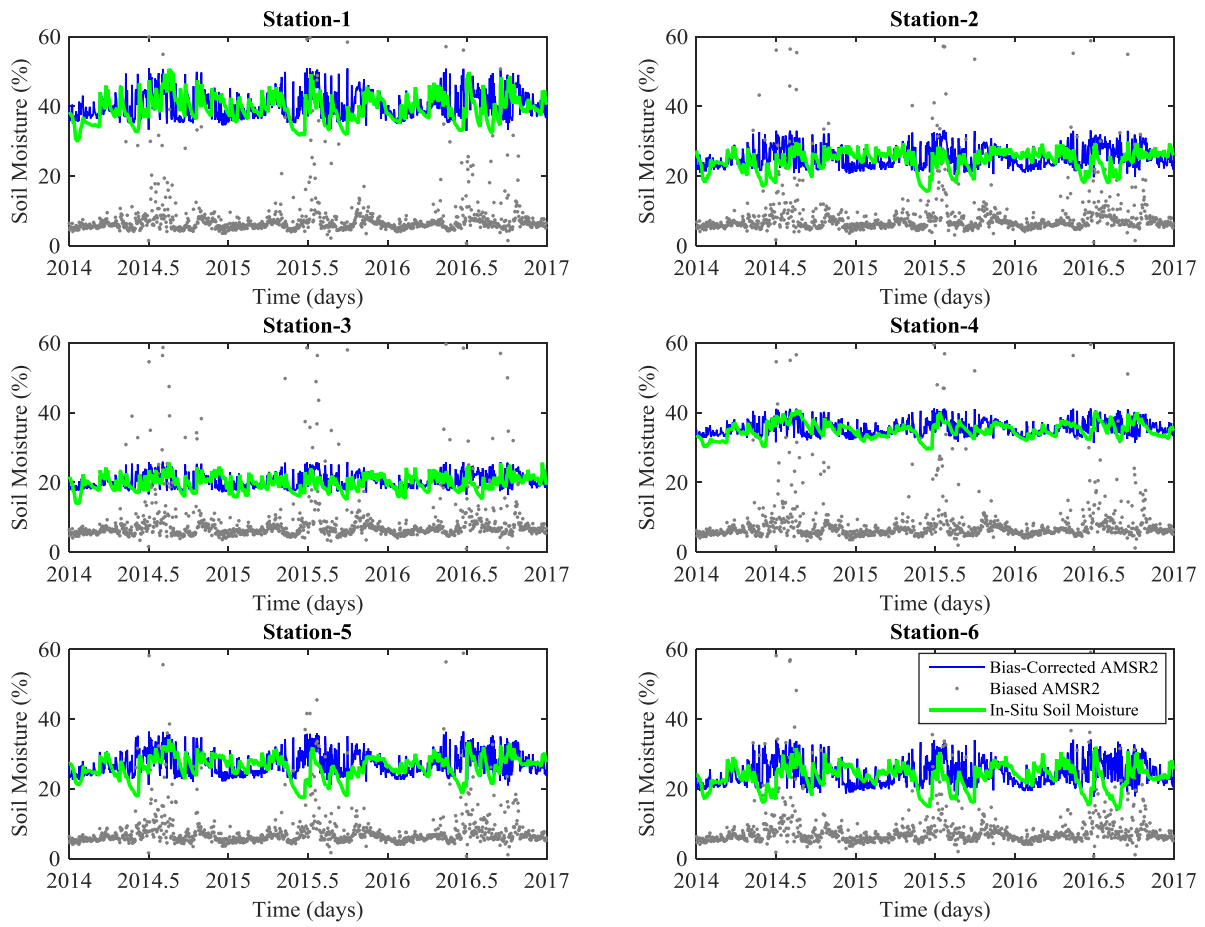
753

754 Figure 3. Graphical model representation of nonhomogeneous hidden Markov model. Here,
 755 SM , S , X indicate soil moisture, hidden state and exogenous variable (i.e., rainfall,
 756 temperature, and AMSR2), respectively.

757

758

759

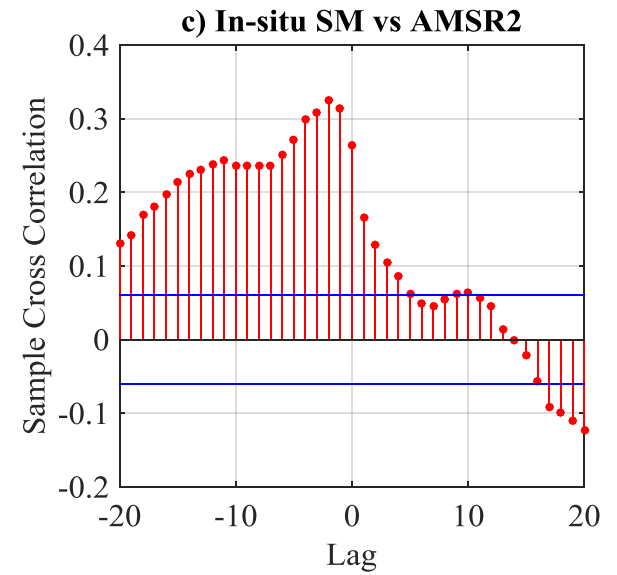
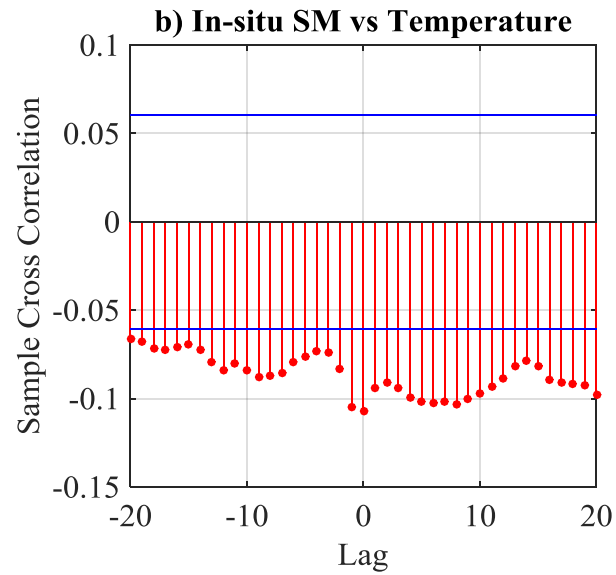
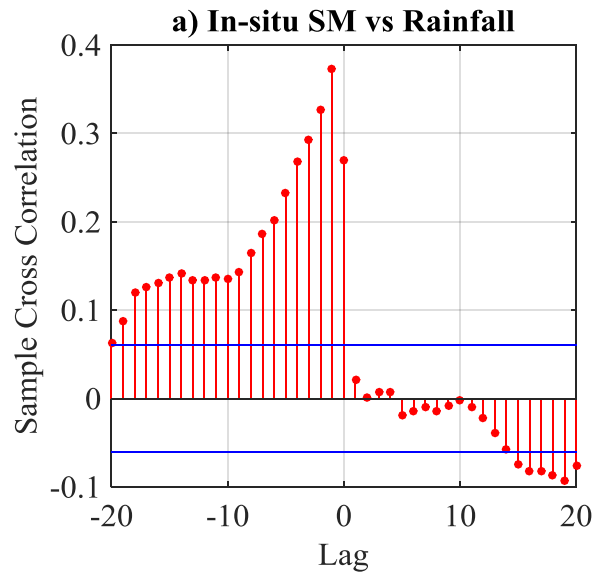


760

761 Figure 4. Bias-uncorrected and bias-corrected AMSR2 SM time series data with in-situ
 762 observations during the study period, 2014–2016.

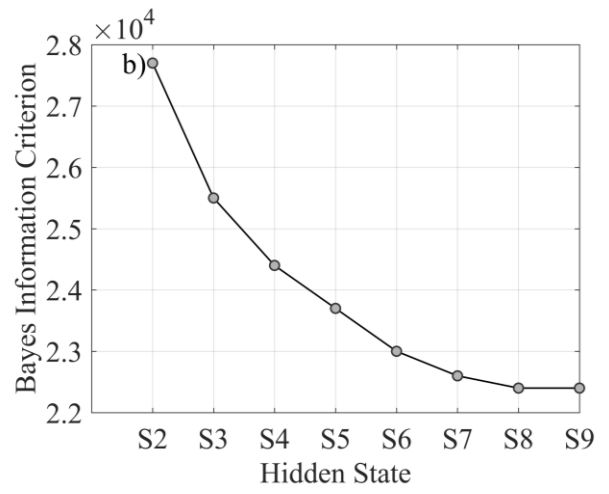
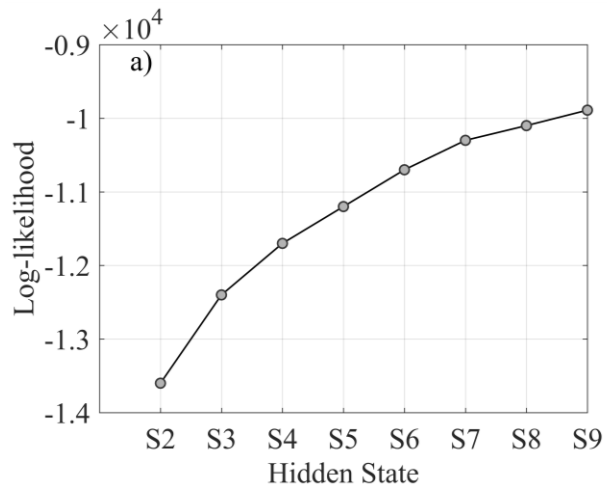
763

764



765

766 Figure 5 Sample cross correlation between the in-situ soil moisture and a set of predictors: a) rainfall, b) temperature, and c) AMSR2 soil
 767 moisture data. All values are averaged over the entire watershed.
 768

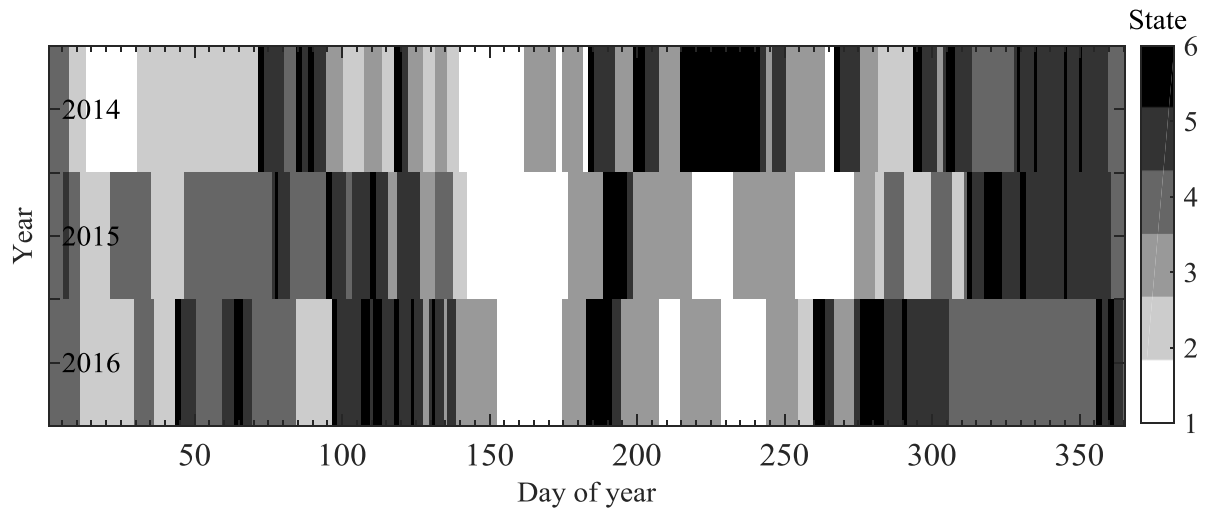


769

770 Figure 6. Log-likelihood and BIC values in terms of hidden states.

771

772



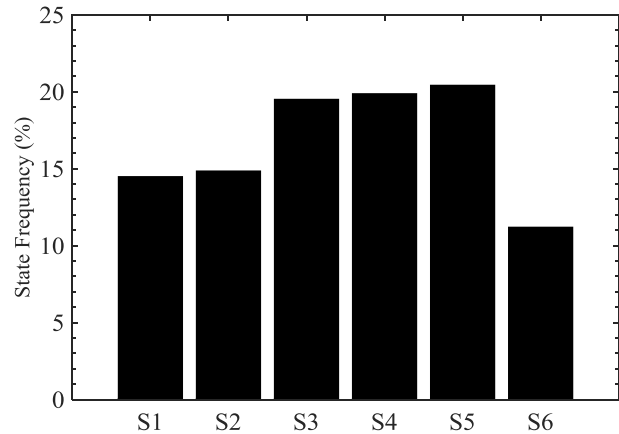
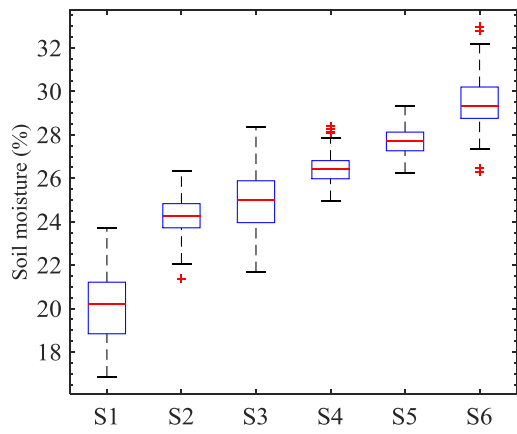
773

774 Figure 7. Estimated hidden state sequence for a 3-year period (2014–2016).

775

776

777

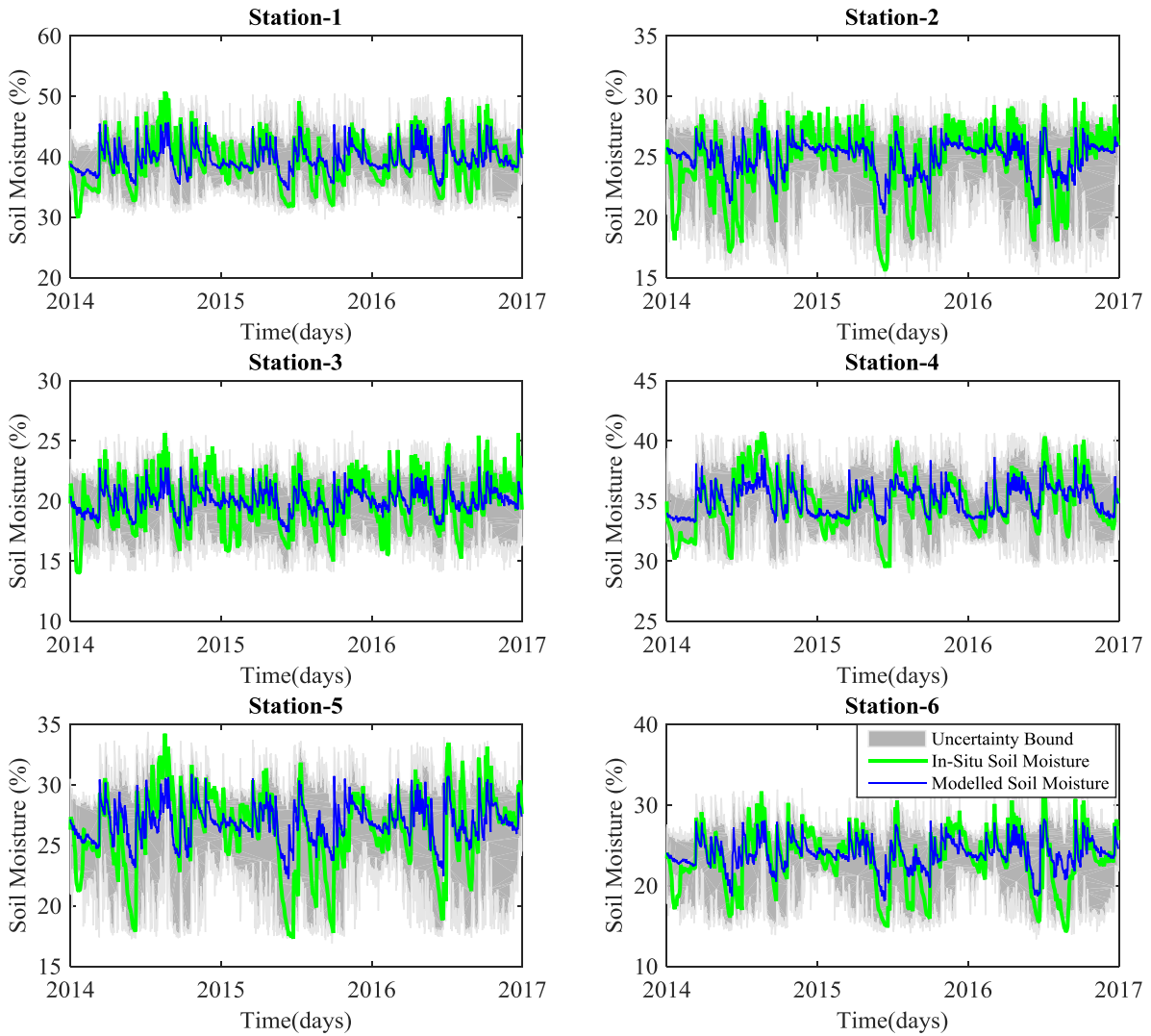


778
779

(a) SM per state

(b) Frequency of SM per state

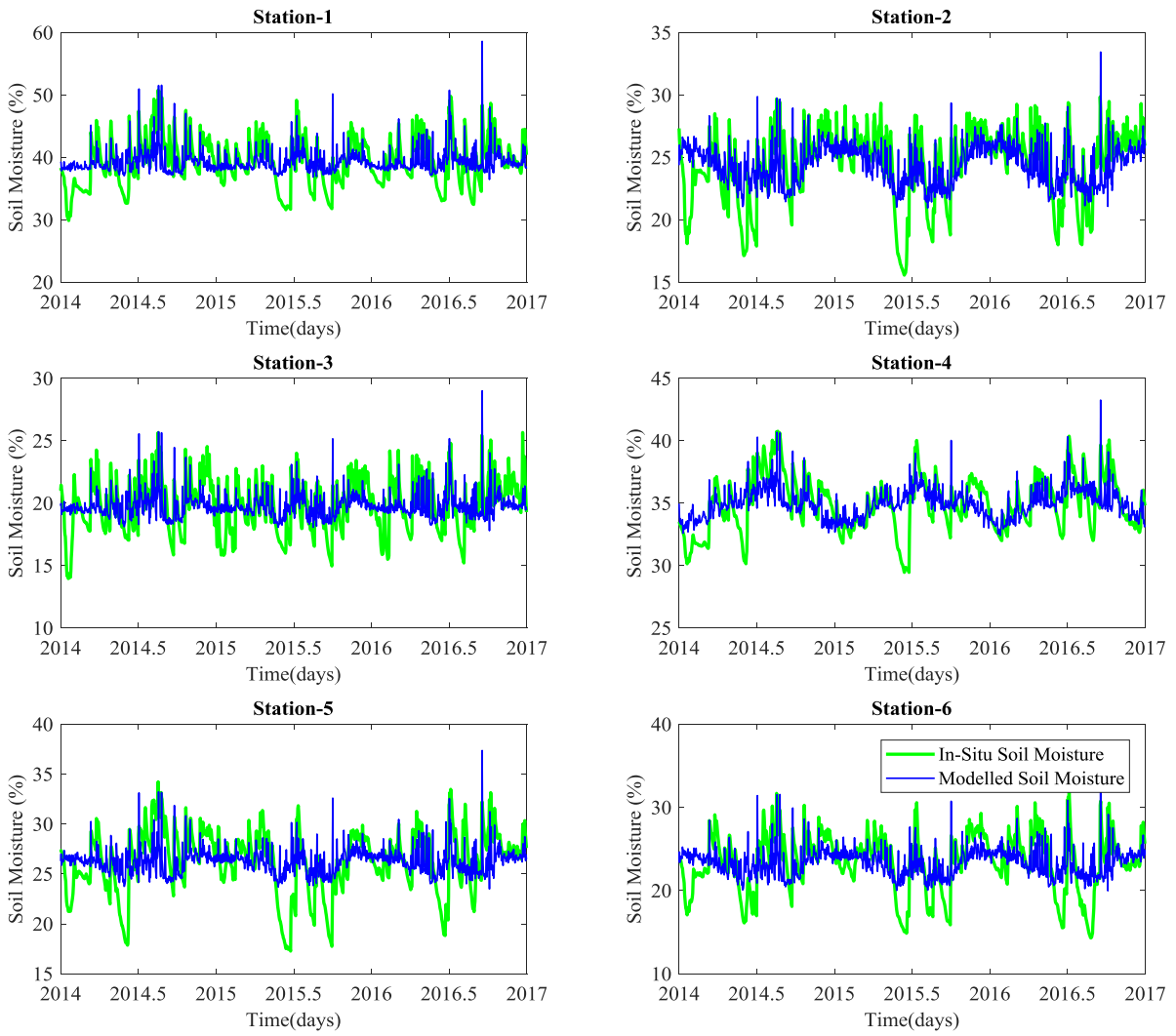
780 Figure 8. The estimated distribution and frequency of soil moisture in each state.



781

782 Figure 9. A comparison of time series data between the in-situ and GM-NHMM-simulated
 783 SM data for 2014–2016: the green line indicates the in-situ observations, and the blue line
 784 represents the median of 100 simulations. The shaded area represents the uncertainty bound
 785 of simulations (between 2.5% and 97.5%).

786

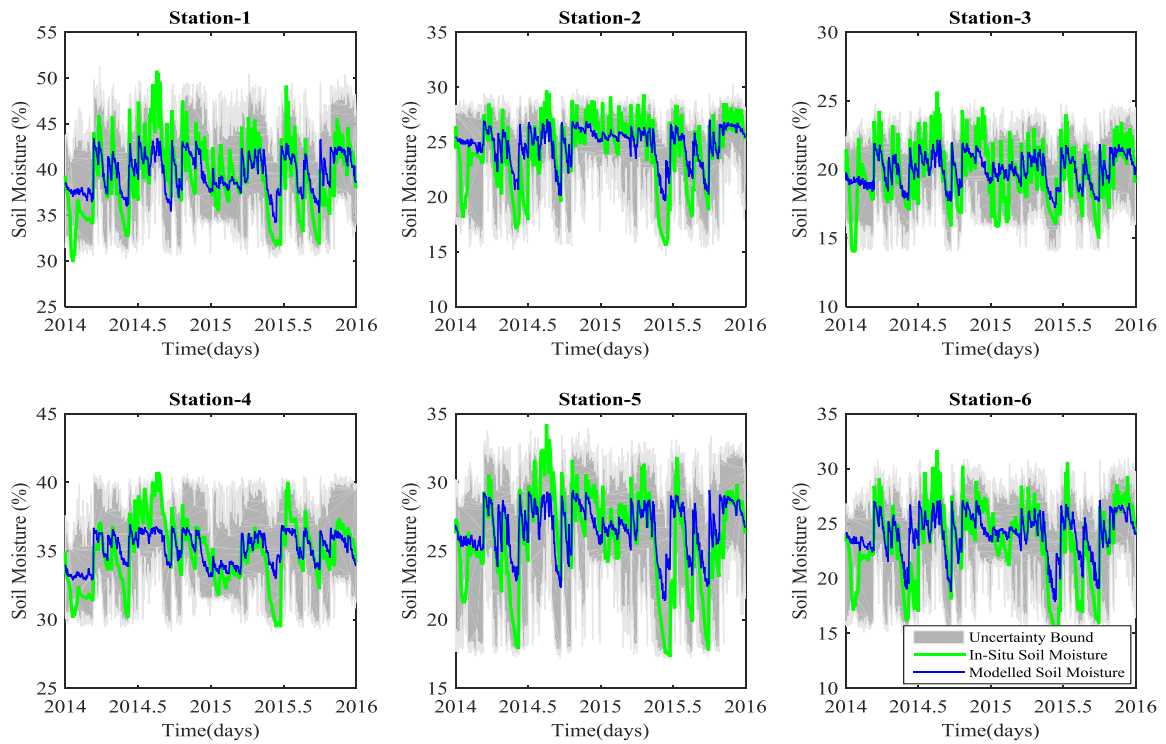


787

788 Figure 10. A comparison of time series data between the in-situ and OLR-simulated SM
 789 products for 2014–2016: the green line indicates in-situ observations, and the blue line
 790 represents OLR-simulated SM.

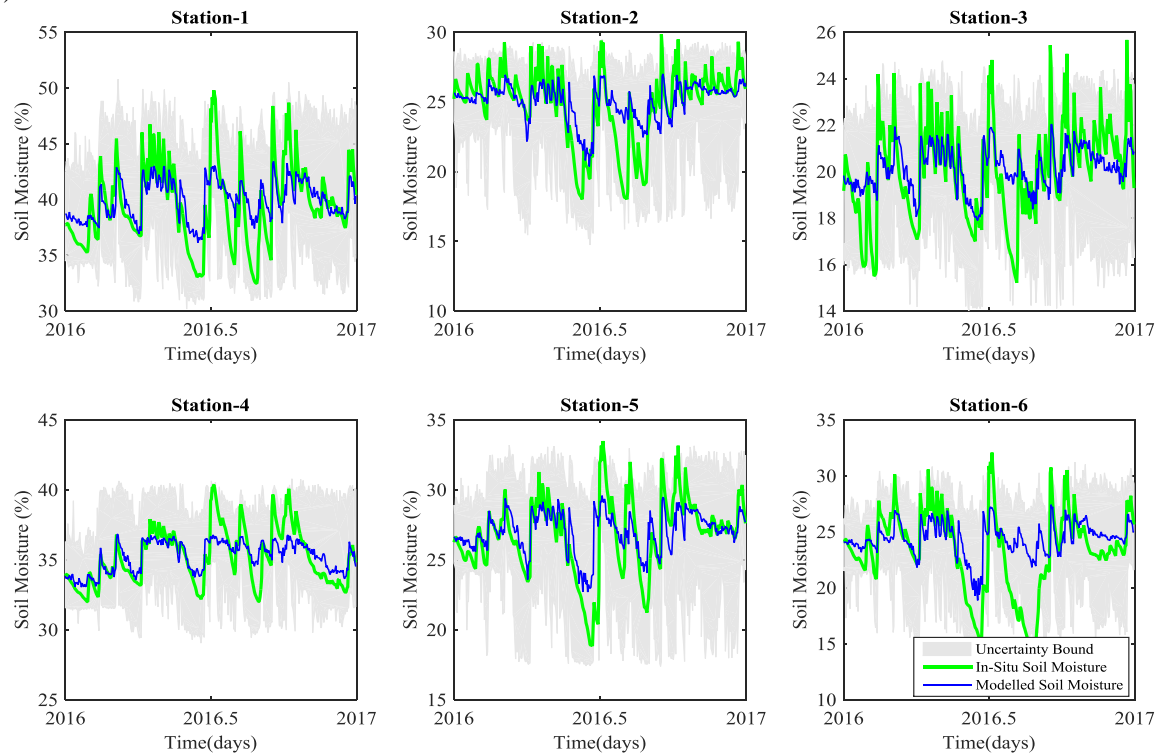
791

792 (a)



793

794 (b)



795

796 Figure 11. Comparisons between the sequences of simulated soil moisture and that observed
797 at multiple locations in the Yongdam watershed for a) the training period (2014–2015) and b)
798 the validation period (2016).

