

Daily rainfall data infilling with a stochastic model

Huidong Jin^a, Quanxi Shao^b and Steven Crimp^c

^aCSIRO Data61, GPO Box 1700, Canberra ACT 2601 Australia

^bCSIRO Data61, Private Bag 5, Wembley Wa 6913, Australia

^cClimate Change Institute, Australian National University, Australia

Email: Warren.Jin@csiro.au

Abstract: Most models are premised on complete data without missing values, such as using a complete daily weather time series to simulate crop biomass accumulation and production, and predict pest risk, even assess climate change impacts. However, historical data series often have some missing values or have only the aggregated values over a period of time. For example, daily rainfall amount is normally recorded by hand during working days and the data during weekends and holidays are sometimes missing and reported as total during these periods. In addition, some data are still missing even nowadays with automated weather stations, due to instrument failure, power outages, operation interruption and so on. Daily rainfall time series data may suffer from several missing data problems. These include (1) individual missing days, (2) consecutive missing days (missing segments), (3) consecutive missing days with their aggregation available. Aggregation of daily data is most common following weekends or holidays. There are several methods to infill these missing values, such as distributed accumulated rainfall evenly over the accumulation period, spatial interpolation from records of surrounding stations, and climatology. These methods often under-estimate the dry day proportions, i.e., giving more wet days than normal, and smooth out extremely daily rainfall amount.

To infill these data gaps appropriately, we investigate a time-varying stochastic model to simulate daily rainfall time series based on true observations, including aggregated observations. A complete rainfall time series is constructed using a three-state Markov chain model to simulate the occurrence of dry, wet and extremely wet days, whilst rainfall amounts for wet and extremely wet days are modelled using a truncated Gamma distribution and an extended Burr XII distribution respectively. Smooth changes on state transition probabilities within a year are captured by time-varying model inputs. The proposed technique can infill the missing data with or without aggregated observations. Experiments on three Australian stations from different climatic zones illustrate its superior performance to a defacto operational approach in Australia and classic climatology method in terms of maintaining daily rainfall data characteristics such as dry day proportions, dry day spells, and rainfall amount distributions. For example, average dry day proportions for these three stations are around 80% based on truly daily rainfall records. They are around 50% for the daily data infilled by our proposed stochastic method. Because these missing daily data do not have substantially missing patterns, these proportions are more reasonable than around 20% for the infilled daily data from the defacto operational approach.

Keywords: *Markov Chain, extremes, temporal disaggregation, missing data*

1 INTRODUCTION

Various models are premised on complete and accurate input time series (Keating et al., 2003). However, real-world observations often have missing values due to instrument failure, power outages, operation interruptions etc. Rainfall time series data in Australia can suffer from several missing data problems, including (1) individual missing days, (2) consecutive missing days (missing segments), (3) consecutive missing days with their aggregation available (Fig. 1). These missing data may mislead model simulation or projection results (Viney and Bates, 2004; Crimp et al., 2019; Jin et al., 2015). For example, at the Australian Bureau of Meteorology’s weather station 040096 (Helidon post office, 27.55° S 152.12° E) in Queensland, the rainfall records start in the year of 1871, and have about 80 segments of missing data in the month of January, and 20 segments of missing data in the month of July or September. The station data also have 10 consecutive missing days but with their aggregated total. Despite becoming a fully automated weather station in April 1988, it still has missing data such as a 3-day-missing segment in July 2016.

The Queensland government, via the Science Delivery Division of Department of Environment and Science (DES), hosts Scientific Information for Land Owners (widely known as SILO) data. It is an enhanced climate database and contains Australian climate data from 1889 to yesterday in a number of ready-to-use formats for research and climate applications. All the station data, including the three stations used in this analysis are accessible at <https://silo.longpaddock.qld.gov.au/> under a Creative Commons Attribution 4.0 International license (CC BY 4.0). As detailed in (Jeffrey et al., 2001), data infilling for the operational SILO database is derived via a three-step process. 1) monthly rainfall for each station is normalised after a fractional power transformation, where the power parameter, mean, and variance are spatially interpolated; 2) ordinary kriging is applied to the normalised monthly rainfall; 3) the monthly values are disaggregated to daily values using the relative daily rainfall distributions generated by ordinary kriging of the observed daily values within each month. The SILO infilled data underestimate extreme rainfalls and have much lower probability of dry days (Jeffrey et al., 2001). Underestimation of dry day proportions is a common issue of the infilling process and some examples are given in Table 1.

In order to produce weather-like daily rainfall time series, the infilling approach should produce disaggregated values, that when summed, match with the aggregated observation (if present); and reproduce other summary statistics such as serial (lagged) correlations, the distribution of extreme events and other station-specific features such as dry spells or wet spells; and be able to match variation in distributional statistics across months and seasons (Elshamy et al., 2006). We propose a method called Daily Rainfall Infilling using a Time-varying Hybrid Stochastic model (DRITHS), in order to adhere to the criteria listed above. To infill missing daily rainfall data, DRITHS uses a three-state Markov chain to model sequential occurrences of dry, wet and extremely wet days. Rainfall amounts for wet and extremely wet days are modelled respectively with different distributions. Monthly or seasonal changes are captured through the inclusions of time-varying model parameters. The proposed technique can impute missing data with or without aggregated records being present.

Besides the spatial interpolation based methods, e.g.(Jeffrey et al., 2001), there are mainly two types of infilling methods in the literature: time series techniques and temporal disaggregation techniques. Daily rainfall features like zero values and skewed distribution make classic time series techniques unappealing, as they are mean-based approaches and their simulations tend to miss extreme values (Shao et al., 2016). Most temporal disaggregation techniques, popular in hydroclimatology, are not suitable for daily rainfall due to different temporal granularity (Furrer and Katz, 2008; Burian and Durrans, 2002), except a few daily weather generators (Elshamy et al., 2006; Semenov et al., 2002). Our work extends a recent development of weather generators (Shao et al., 2016) for daily rainfall data infilling. Our work is also related with data imputation methods (Little and Rubin, 2014). These methods often assume the missing at random, and mainly focus on imputing tabulated data sets. For our application, the temporal dependency is better considered explicitly.

The paper is organised as follows. Section 2 presents DRITHS, including a Markov chain with different distribution types for different states that optimises the state-sequences for missing segments based on maximum likelihood. Section 3 gives its experimental results and comparison with existing techniques. Concluding discussions are given in Section 4.

2 INFILL DAILY RAINFALL DATA WITH A HYBRID STOCHASTIC MODEL

We present a rainfall time series infilling technique based on a time-varying hybrid stochastic model (DRITHS). As illustrated in Fig 1a, the infilling approach is designed to cope with missing across several consecutive days with the accumulated rainfall records, e.g. the total rainfall for d and $d + 1$ is 32mm, and

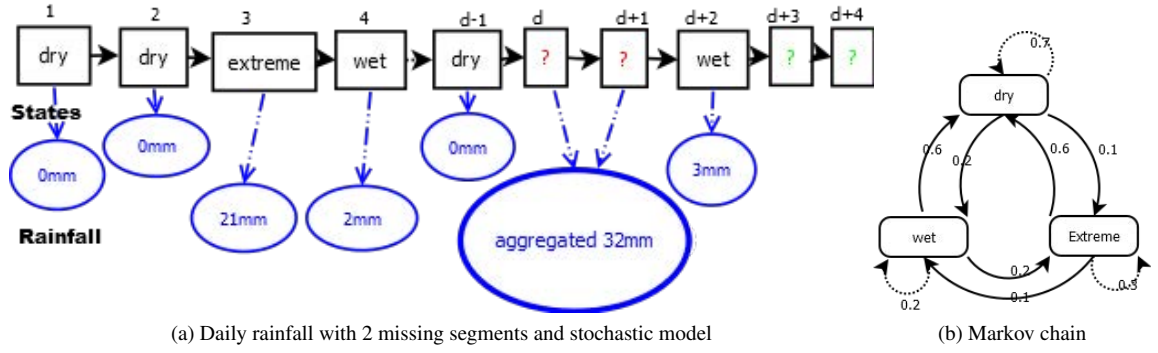


Figure 1. Overview of proposed time-varying hybrid stochastic rainfall model

for missing data without aggregated total, e.g. the segment of $d + 3$ and $d + 4$. We model the daily rainfall data with a three state Markov chain, as illustrated in Fig 1b. These three states, dry, wet, and extreme, are associated with different distribution types in order to address zero rainfall as well as extreme rainfall events. The underlying Markov chain models these state sequences, which transits from one to another state according to transition probabilities. The transition probabilities change smoothly with the day of the year. Fig 2e illustrates that daily observed precipitation amounts do not simply follow a single type of distribution such as a truncated Gaussian distribution or extreme value distribution, which made most weather generators (Elshamy et al., 2006; Furrer and Katz, 2008) less suitable. We present our hybrid stochastic model below.

2.1 A time-varying hybrid stochastic model

The daily precipitation observations are denoted as z_t ($t = 1, 2, \dots, n$) with n being the number of daily observations; the corresponding aggregated precipitation amounts are y_t ($t = 1, 2, \dots, n$). In our application, y_t is the moving average of 31 neighbouring days centered around day t , and only daily rainfall records are used. Two predefined constants c_{dry} and c are used to differentiate precipitation amount of wet days from dry days and extremely wet days, respectively. c_{dry} is the minimum daily rainfall records for rainy days, and equals 0, and c is 4mm per day for our experiments in this work. A three-state process for precipitation occurrence is as follows.

$$x_t = \begin{cases} 2, & \text{if } z_t > c, \\ 1, & \text{if } c_{dry} < z_t \leq c, \\ 0, & \text{if } z_t \leq c_{dry}, \end{cases} \quad (1)$$

which are termed as extreme, wet and dry state respectively. This setting provides flexibility to model the distributions of the precipitation amount on the wet and the extremely wet days. Let $J_{t,i} = \begin{cases} 1, & \text{if the } i\text{-th state is observed,} \\ 0, & \text{otherwise at time } t, \end{cases}$ ($t = 1, 2, \dots$) represent the i -th state ($i = 1, \dots, m$). m indicates the number of states, where $m = 3$ in our applications (Fig 1a). Let U_t be the collection of covariates and possible lag variables governing the state process. We assume the state process depends on the previous state as well as a moving average of daily rainfall. Thus, $U_t = \left(1, J_{t,1}, \dots, J_{t,(m-1)}, y_t, \sin\left(2\pi \frac{d(t)}{365.25}\right), \cos\left(2\pi \frac{d(t)}{365.25}\right)\right)$. The first element is the intercept in the regression for simple notation in the equations below. $\sin\left(2\pi \frac{d(t)}{365.25}\right)$ and $\cos\left(2\pi \frac{d(t)}{365.25}\right)$ aim to capture possible seasonality with the Julian day of t , $d(t)$. The application of these two terms, and the 31 day-moving-average y_t , distinguish DRITHS from the model in Shao et al. (2016), which considers monthly effects but is not suitable for data infilling as the missing segments could occur crossing a month. The conditional probability is defined

$$\pi_{ti} = \Pr \{J_{ti} = 1|U_{t-1}\}. \quad (2)$$

From a modelling perspective, only $m - 1$ conditional probabilities need to be modelled because all the probabilities $(\pi_{t1}, \dots, \pi_{tm})$ add up to one. The multinomial logit model is frequently used to form the linear

relationship between the ratio π_{ti}/π_{tm} and covariates as

$$\log \{ \pi_{ti}/\pi_{tm} \} = \beta'_i U_{t-1} = \beta_{i0} + \sum_{l=1}^{m-1} \beta_{il} J_{(t-1),l} + \beta_{im} y_t + \beta_{i,m+1} \sin \left(\frac{2\pi d(t)}{365.25} \right) + \beta_{i,m+2} \cos \left(\frac{2\pi d(t)}{365.25} \right). \quad (3)$$

Here $\beta_i = (\beta_{i0}, \beta_{i1}, \dots, \beta_{i,m+2})'$ ($i = 1, \dots, m-1$) are the vectors of regression coefficients, and is obtained via maximum likelihood estimation based on training data. So, β_i are constants after fitting, but π_{ti} will vary smoothly with day t .

For wet days, the precipitation amounts are bounded between c_{dry} and c . To model them, we use a statistical distribution with a finite support between 0 and c . Following Shao et al. (2016), we employ a truncated Gamma distribution defined by a density function within $0 \leq x \leq c$ as

$$f_{tG}(x; \alpha_1, \alpha_2) = \frac{e^{-x/\alpha_1} x^{\alpha_2-1}}{K_1(\alpha_1, \alpha_2)} \quad (4)$$

with $\alpha_1 > 0$ and $\alpha_2 > 0$ being the scale and the shape parameters respectively, where $K_1(\alpha_1, \alpha_2) = \int_0^c e^{-x/\alpha_1} x^{\alpha_2-1} dx = \alpha_1^{\alpha_2} \gamma(\alpha_2, c/\alpha_1)$ acts as the normalisation factor to ensure that the definition in Eq.4 is a distribution, and $\gamma(a, t) = \int_0^t e^{-x} x^{a-1} dx$ is the truncated Gamma function.

Extreme wet days, with possible long-tail properties in precipitation amount, require a distribution with flexible tail properties. Extreme value distributions or its generalised versions could be used. Following Shao et al. (2004), we use the censored Burr XII (EBXII) distribution with density function defined over $x > c$ as

$$f_{cE}(x; \alpha_3, \alpha_4, \alpha_5) = \begin{cases} K_2^{-1} \left[\frac{\alpha_4}{\alpha_3} \left(\frac{x}{\alpha_3} \right)^{\alpha_4-1} \left\{ 1 - \alpha_5 \left(\frac{x}{\alpha_3} \right)^{\alpha_4} \right\}^{\frac{1}{\alpha_5-1}} \right], & \alpha_5 \neq 0, \\ K_2^{-1} \left[\frac{\alpha_4}{\alpha_3} \left(\frac{x}{\alpha_3} \right)^{\alpha_4-1} \exp \left\{ - \left(\frac{x}{\alpha_3} \right)^{\alpha_4} \right\} \right], & \alpha_5 = 0, \end{cases} \quad (5)$$

where $K_2(\alpha_3, \alpha_4, \alpha_5) = \begin{cases} \{1 - \alpha_5 (c/\alpha_3)^{\alpha_4}\}^{1/\alpha_5}, & \alpha_5 \neq 0, \\ \exp \{ - (c/\alpha_3)^{\alpha_4} \}, & \alpha_5 = 0, \end{cases}$ acts as the normalisation factor to ensure

that the definition in Eq.5 is a distribution. For robust overall fitting, we use the quantile-based least squares method (Shao et al., 2016) for estimating the distributional parameters in Eqs 4 and 5.

2.2 Infilling daily rainfall time series with a trained hybrid stochastic model

With a trained model from daily observations, we infill any missing value in a given daily rainfall time series. Our infilling procedure has three steps.

Step 1 is to estimate the best state sequence in terms of maximum likelihood for each segment of consecutive missing data. For a segment of l days, there are totally m^l possible state sequences. For each of them $\{s_{d_1}, s_{d_2}, \dots, s_{d_l}\}$, we can calculate its likelihood as we have the state transition probability matrix $\pi_{t,i,j}$ for each day t from state i to state j . Note we may have the states before and after the segment, that is s_{d_0} and $s_{d_{l+1}}$, which are the constraints for the likelihood estimation

$$likelihood(s_{d_0}, s_{d_1}, s_{d_2}, \dots, s_{d_l}, s_{d_{l+1}}) = \prod_{k=1}^{l+1} \pi_{d_{k-1}, s_{d_{k-1}}, s_{d_k}}. \quad (6)$$

When s_{d_0} (or $s_{d_{l+1}}$) is missing, for easy computation, we force $\pi_{d_0,i,j} = 1$ ($\pi_{d_l,i,j} = 1$) for any i or j . When the segment length l is large, we use the aggregated rainfall amount y_a to restrict state sequences, because at most $\lfloor \frac{y_a}{c} \rfloor$ days are in extreme state as c is the minimum amount for extremely wet days.

Step 2 is to simulate daily rainfall amount within the time segment. Dry days have zero rainfall. The rainfall amount for wet days is simulated with Eq 4, and extreme state rainfall is simulated with Eq 5.

Step 3 is post-simulation adjustment to meet the aggregated rainfall amount if it is available. Adjusting procedures are needed to restore consistency with aggregated rainfall amount (such as 32mm for d and $d + 1$

Table 1. State frequency at three stations. ‘Obser’ is for truly observed daily rainfall, and others are for infilled data by three methods. ‘Clima’ indicates the climatology method. 1212 indicates all 12 months

(a) Helidon					(b) Bencubin					(c) Yenda				
State freq(%) for Helidon					State freq(%) for Bencubin					State freq(%) for Yenda				
	Month	Dry	Wet	Ext		Month	Dry	Wet	Ext		Month	Dry	Wet	Ext
Obser	1	73.2	10.5	16.3	Obser	1	89.9	6.2	3.9	Obser	1	85.7	7.2	7.2
DRITHS	1	52.9	9.8	37.3	DRITHS	1	66.7	0	33.3	DRITHS	1	50	0	50
SILO	1	19.6	33.3	47.1	SILO	1	0	33.3	66.7	SILO	1	0	0	100
Clima	1	0	52.9	47.1	Clima	1	0	0	100	Clima	1	0	0	100
Obser	7	86.4	7.6	6	Obser	7	56.4	33.4	10.3	Obser	7	72.5	18.3	9.2
DRITHS	7	57.7	19.2	23.1	DRITHS	7	0	100	0	DRITHS	7	64.3	14.3	21.4
SILO	7	30.8	46.2	23.1	SILO	7	50	0	50	SILO	7	14.3	57.1	28.6
Clima	7	0	69.2	30.8	Clima	7	0	100	0	Clima	7	0	85.7	14.3
Obser	1212	80.7	8.8	10.6	Obser	1212	77.3	16.5	6.2	Obser	1212	81.8	10.2	7.9
DRITHS	1212	57	10.2	32.8	DRITHS	1212	44.4	33.3	22.2	DRITHS	1212	58.6	9.2	32.2
SILO	1212	23.5	36.9	39.5	SILO	1212	11.1	50	38.9	SILO	1212	19.5	42.5	37.9
Clima	1212	0	53.1	46.9	Clima	1212	0	61.1	38.9	Clima	1212	0	59.8	40.2

in Fig 1a). Following (Koutsoyiannis and Manetas, 1996; Shao et al., 2016), we simply scale up the rainfall amount simulated. Provided that y_a is the aggregated amount for the segment, and the synthetic series \tilde{z}_s ($s = 1, 2, \dots, l$) has been generated by our hybrid stochastic model, the simple proportional adjustment procedure is according to

$$\hat{z}_s = \left(\frac{y_a}{\sum_{j=1}^l \tilde{z}_j} \right) \tilde{z}_s \quad (s = 1, 2, \dots, l) \tag{7}$$

To ensure daily rainfall amount for extreme states is larger than the threshold c , c is deducted from these days (as well as y_a accordingly) before scaling, and then added back afterwards.

3 EXPERIMENTS AND COMPARISONS

To illustrate if missing daily values are infilled appropriately like truly observed daily weather data (Kokic et al., 2013), we use several metrics such as the state (dry/wet/extreme) proportions, state spells, state transition probabilities, and the rainfall amount density distributions. These infilled rainfall data are classified into different states according to Eq.1. The state proportions, such as dry day proportions, as well as the serial dependence, are important factors affecting the performance of cropping or ecosystem modelling (Shao et al., 2016). The latter is characterised by spell lengths (i.e. the number of consecutive days having a given weather state) and state transition probabilities. The rainfall distributions are to check whether the rainfall amount is distributed like true weather observations. As we do not have ground truth for these missing values, we will use the metrics from the truly observed daily rainfall data as the target metrics for these infilled data.

Besides the Helidon station, we also use Station 75079 (Yenda – Henry Street, 34.25° S 146.20° E, NSW), and Station 10007 (Bencubbin – 30.80° S 117.86° E, Western Australian). These weather stations reside in three different climatological zones and are thus very different in terms of rainfall statistical characteristics. Besides SILO, a widely used data repository in Australia, we also compare DRITHS with a widely used climatology method. It infills data using the long-term average data for a given day in years (Kokic et al., 2013).

In Table 1, the dry day proportion in January for Helidon, based on truly observed daily rainfall data, is 73.2%. That of the infilled data from DRITHS is 52.9%, which is much closer than those of the other two methods. Note the dry day frequency of SILO infilled data is only 19.6%. Similarly, DRITHS is better able to simulate the wet state and extreme state probabilities than its two counterparts. Quite similar comparison results could

be found for July data, for all the 12 months data (indicated by 1212), as well as for the other two stations (with an exception on Bencubin for July where only a few days are missing).

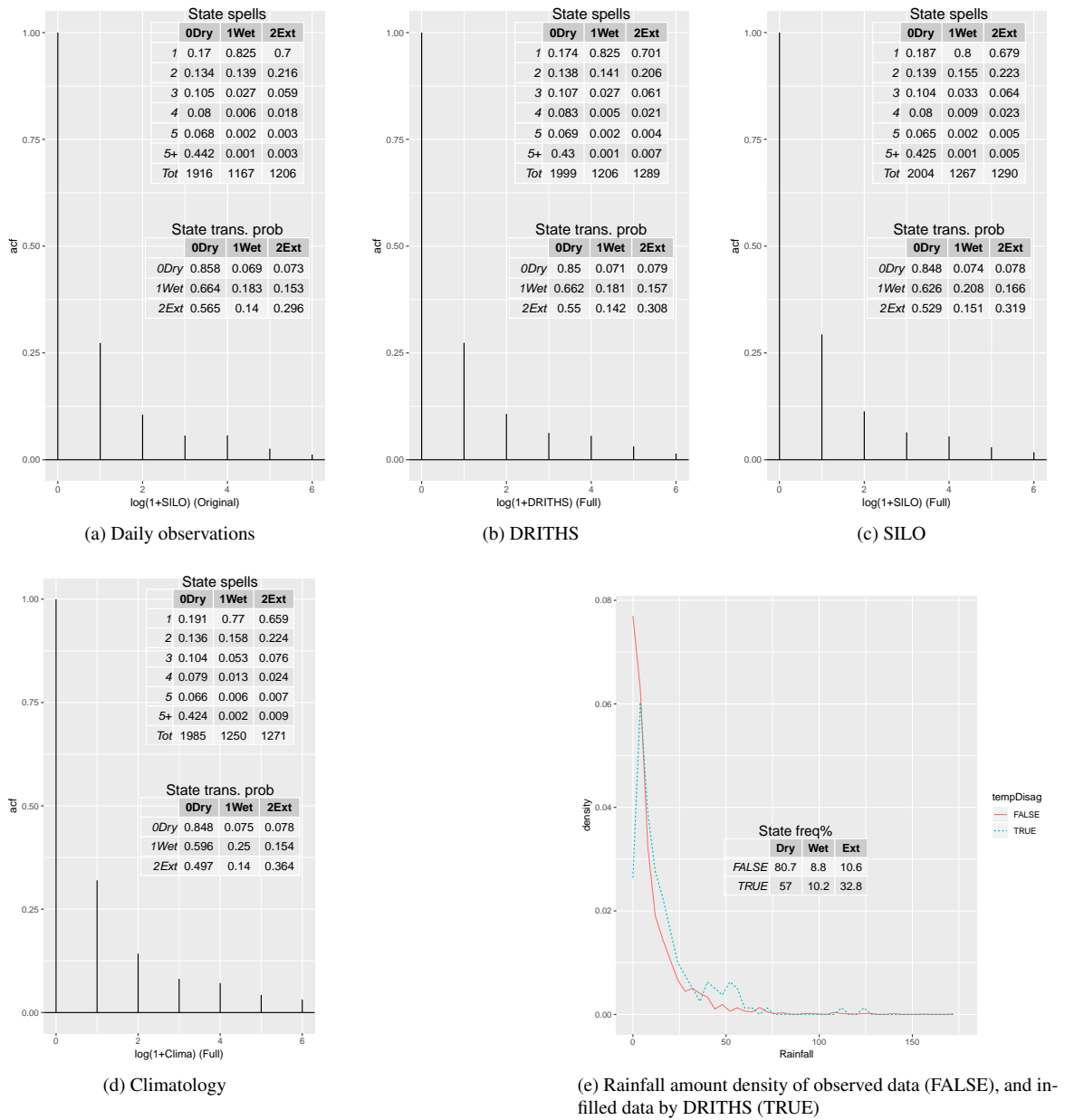


Figure 2. Auto-correlation function (ACF), state spells, and averaged state transition probabilities for Helidon. States are indicated by ‘0Dry’, ‘1Wet’, and ‘2Ext’ respectively.

Fig 2 illustrates different metrics for the Helidon station. Figs 2b to 2d present the state spells (y-axis is spell lengths in days), average state transition probabilities, and Auto-Correlation Function (ACF) with lag 1 to 5+ based on the complete rainfall time series infilled by DRITHS, SILO, and climatology respectively for Helidon. Comparing with those based on the daily observations (Fig 2a), results of DRITHS have small difference. For 5-day dry spells, e.g., the probability is 6.8%, 6.9%, 6.5% and 6.6% respectively for daily observations only, DRITHS, SILO and climatology. DRITHS has the closest result to the observed. The lag-1 autocorrelation on $\log(1+\text{rainfall})$ from DRITHS (Fig 2b) is very close to that based on observations in Fig 2a. The other two approaches produce higher lag-1 auto-correlations than DRITHS. Fig 2e shows the rainfall

amount density from daily observations only (solid line), and from infilled data using DRITHS (dotted line). DRITHS captures the long-tail of the distribution (i.e., extreme rainfall) quite well, although does not match perfectly near zero rainfall amount, while its two counterparts perform considerably worse.

Quite similar comparison results could be found for the other two stations. For 5-day dry spells, e.g., the probability is 7.5%, 7.6%, 7.8% and 7.9% respectively for daily observations only, DRITHS, SILO and climatology for Yenda. DRITHS again has the best result.

4 CONCLUSIONS AND DISCUSSIONS

We have proposed the DRITHS approach to infill missing values in daily rainfall time series based on a time varying three-state Markov chain. It uses different rainfall amount distribution types for the three states so as to model extreme rainfalls better. Experimental results demonstrate that DRITHS could infill daily rainfall data better than two existing techniques in terms of maintaining weather data characteristics, such as dry day proportions, state spells, and daily rainfall distributions.

More comprehensive, such as leave-N-out cross-validation experiments, will be carried out. The 4mm per day rainfall threshold selection is left for future efforts. It is worth investigating whether extending the lag-one Markov assumption to lag-two will help. The proposed technique has not considered the possible dependency from nearby stations or other climate variables, though our preliminary results suggest the dependency on daily maximum or minimum temperature may change rapidly from one day to another.

Acknowledgements

The work is funded by CSIRO Digiscape Future Science Platform.

REFERENCES

- Burian, S. and S. Durrans (2002). Evaluation of an artificial neural network rainfall disaggregation model. *Water science and technology* 45(2), 99–104.
- Crimp, S., H. Jin, P. Kokic, S. Bakar, and N. Nicholls (2019). Possible future changes in south east Australian frost frequency: an inter-comparison of statistical downscaling approaches. *Climate dynamics* 52(1-2), 1247–1262.
- Elshamy, M. E., H. S. Wheeler, N. Gedney, and C. Huntingford (2006). Evaluation of the rainfall component of a weather generator for climate impact studies. *Journal of Hydrology* 326(1), 1–24.
- Furrer, E. M. and R. W. Katz (2008). Improving the simulation of extreme precipitation events by stochastic weather generators. *Water Resources Research* 44(12), W12439:1–13.
- Jeffrey, S. J., J. O. Carter, K. B. Moodie, and A. R. Beswick (2001). Using spatial interpolation to construct a comprehensive archive of Australian climate data. *Environmental Modelling & Software* 16(4), 309–330.
- Jin, H., J. Kokic, G. Hopwood, J. Ricketts, and S. Crimp (2015). A new quantile projection method for producing representative future daily climate based on mixed effect state-space model and observations. In *MODSIM2015*, pp. 1544–1550.
- Keating, B. A. et al. (2003). An overview of APSIM, a model designed for farming systems simulation. *European Journal of Agronomy* 18(3), 267–288.
- Kokic, P., H. Jin, and S. Crimp (2013). Improved point scale climate projections using a block bootstrap simulation and quantile matching method. *Climate dynamics* 41(3-4), 853–866.
- Koutsoyiannis, D. and A. Manetas (1996, July). Simple disaggregation by accurate adjusting procedures. *Water Resources Research* 30(7), 2015–2117.
- Little, R. J. and D. B. Rubin (2014). *Statistical analysis with missing data*, Volume 333. John Wiley & Sons.
- Semenov, M. A., E. M. Barrow, and A. Lars-Wg (2002, Aug). *LAWS-WG: A stochastic weather generator for use in climate impact studies* (Version 3.0 ed.). Harpenden, Hertfordshire, AL5 2JQ, UK: Rothamsted Research.
- Shao, Q., H. Wong, J. Xia, and W.-C. Ip (2004). Models for extremes using the extended three-parameter Burr XII system with application to flood frequency analysis. *Hydrological Sciences Journal* 49(4), 685–702.
- Shao, Q., L. Zhang, and Q. Wang (2016). A hybrid stochastic-weather-generation method for temporal disaggregation of precipitation with consideration of seasonality and within-month variations. *Stochastic Environmental Research and Risk Assessment* 30, 1705–1724.
- Viney, N. R. and B. C. Bates (2004). It never rains on Sunday: the prevalence and implications of untagged multi-day rainfall accumulations in the Australian high quality data set. *International Journal of Climatology: A Journal of the Royal Meteorological Society* 24(9), 1171–1192.