

1 *Supplement of*

# 2 **Construction of homogenized daily surface air** 3 **temperature for Tianjin city during 1887-2019**

4 **Peng Si, Qingxiang Li, and Philip D Jones**

5 *Correspondence to:* Qingxiang Li (liqingx5@mail.sysu.edu.cn)

6 In this supplementary material, we report how to extend the daily maximum temperature time series  
7 for Tianjin from Jan 1 1887 to Aug 31 1890. From the information of reference data sources, only the  
8 Berkeley Earth daily maximum/minimum temperatures are available. So the maximum temperature  
9 Berkeley Earth-daily data corresponding to the site level Tianjin station is selected as the extension data.  
10 Two methods of standardized series and linear regression are both applied to the extension of the daily  
11 maximum series during the period of Jan 1 1887-Aug 31 1890 at Tianjin station (Si et al., 2017) and  
12 finally the reasonable extension series is determined by comparing the results of an error analysis  
13 between the two methods.

## 14 **S1 Standardized series method**

15 This method assumes that for all the stations in the same climate region, the meteorological  
16 information at a certain time point is similar to its anomaly from the multi-year average values. The  
17 method can be expressed as follows:

$$18 \quad Z_{Berkeley} = \frac{(x_i - \bar{x}_i)}{s_i} \quad (S1)$$

$$19 \quad y_{Tj} = Z_{Berkeley} s_{Tj} + \bar{y}_{Tj} \quad (S2)$$

20 Where,  $Z_{Berkeley}$  is the Berkeley Earth-daily standardized series,  $x_i$  is the Berkeley Earth-daily data  
21 for some date,  $\bar{x}_i$  and  $s_i$  are the mean and standard deviation of Berkeley Earth-daily data for some

22 date during 1961-1990.  $y_{Tj}$  is the fitted daily series at Tianjin station while  $\overline{y_{Tj}}$  and  $s_{Tj}$  are the mean  
 23 and standard deviation of observed Tianjin-daily data at some date covering 1961-1990.

## 24 S2 Linear regression method

25 Linear regression method is used to interpolate the candidate time series ( $\hat{y}_i$ ) by establishing a  
 26 quantitative linear statistical relationship between the candidate station and its neighbour. It can be  
 27 expressed as follows:

$$28 \quad \hat{y}_i = b_0 + bx_i \quad (S3)$$

29 In which, regression coefficients of  $b_0$  and  $b$  are obtained by the following formulas,

$$30 \quad \begin{cases} nb_0 + b \sum_{j=1}^n x_j = \sum_{j=1}^n y_j \\ b_0 \sum_{j=1}^n x_j + b \sum_{j=1}^n x_j^2 = \sum_{j=1}^n x_j y_j \end{cases} \quad (S4)$$

32 Where  $n = 30$ ,  $x_j$  is Berkeley Earth-daily data, and  $y_j$  is observed Tianjin-daily data at the same date.  
 33 Similar to standardized series method (described in the earlier subsection), a quantitative linear statistical  
 34 relationship is also established by multi-year average values during 1961-1990.

35 The significance test of the regression equation is carried out using the  $F$  statistic defined in Eq.  
 36 (S5), to investigate whether the established linear statistical relationship defined in Eq.(S3) is significant.  
 37 The  $F$  statistic follows an  $F$ -distribution with numerator degree of freedom of 1 and denominator degree  
 38 of freedom  $(n - 2)$ . If  $F > F_{0.05}$ , then the linear relation of regression equation is considered to be  
 39 significant at 5% significance level.

$$40 \quad F = \frac{\frac{U}{1}}{\frac{Q}{(n-2)}} \quad (S5)$$

$$41 \quad U = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2 \quad (S6)$$

$$42 \quad Q = \sum_{j=1}^n (y_j - \hat{y}_j)^2 \quad (S7)$$

43 Where  $n = 30$ . In Eq. (S6) - (S7),  $\hat{y}_j$ ,  $y_j$  and  $\bar{y}$  are the fitted daily data at Tianjin station, observed  
 44 Tianjin-daily data and the mean of observed Tianjin-daily data at some date during 1961-1990,  
 45 respectively.

46 At the same time, the significance of the regression coefficient is also tested to further confirm the  
 47 credibility of established linear statistical relationship. The  $t$  statistic defined in Eq. (S8) satisfies a  
 48  $t$ -distribution with  $(n - 2)$  degrees of freedom. If  $|t| > t_{0.05}$ , then the linear relation is considered to  
 49 be significant at the 5% significance level.

$$50 \quad t = \frac{\frac{b}{\sqrt{c}}}{\frac{Q}{\sqrt{n-2}}} \quad (S8)$$

$$51 \quad c = \left[ \sum_{j=1}^n (x_j - \bar{x})^2 \right]^{-1} \quad (S9)$$

52 Where  $n = 30$ ,  $b$  is one of the regression coefficients in Eq. (S4) and  $Q$  is calculated using Eq. (S7).  
 53 In Eq. (S9),  $x_j$  is Berkeley Earth-daily data and  $\bar{x}$  represents the mean of Berkeley Earth-daily series  
 54 at some date during 1961 - 1990.

### 55 **S3 Error analysis**

56 In order to make the sample for error analysis large enough, the fitted daily maximum temperature  
 57 series formed by the above two methods are estimated by comparing errors between the fitted and the  
 58 actual observations over the same time period. The error assessment indicators include the Standard  
 59 Mean Square Error ( $SMSE$ ), Standard Error ( $SE$ ) and the proportion of the differences in the range of  
 60  $\pm 0.5$  °C between the fitted and observed data ( $P$ ). These indicators have been used in construction of the  
 61 monthly surface air temperature series for Baoding city during 1913 - 2014 (Si et al., 2017).

62 Standard Mean Square Error ( $SMSE$ ) is similar to Root Mean Square Error ( $RMSE$ ), which is used  
 63 to measure the deviation between fitted and true values. But in order to eliminate the differences from

64 units of different elements and climate change of the same element at different site location as well as  
 65 easy comparison of different elements, the  $SMSE$  is superior indicator compared to  $RMSE$  (Huang et al.,  
 66 2004).  $SMSE$  is shown as follows:

$$67 \quad SMSE = \sqrt{\sum_{i=1}^m \left( \frac{\hat{y}_i - y_i}{s} \right)^2} \quad (S10)$$

68 Where,  $m$  is the sample size,  $\hat{y}_i$  is the fitted daily data at Tianjin station and  $y_i$  is the actual  
 69 observed Tianjin-daily data at the same date, similarly hereinafter.  $s$  is the standard deviation of  
 70 observed Tianjin-daily data at the same date during 1961 - 1990.

71 Standard Error ( $SE$ ) is an indicator inferring statistical reliability. The smaller the values of  $SE$   
 72 implies the closeness of the sample statistics to the population parameter, thereby makes the sample as  
 73 ideal representative of the population and naturally therefore extrapolation of the population parameter  
 74 using sample statistics became more reliable. It can be expressed as follows:

$$75 \quad SE = \sqrt{\frac{1}{(m-1)} \sum_{i=1}^m (\varphi'_i - \overline{\varphi'})^2} \quad (S11)$$

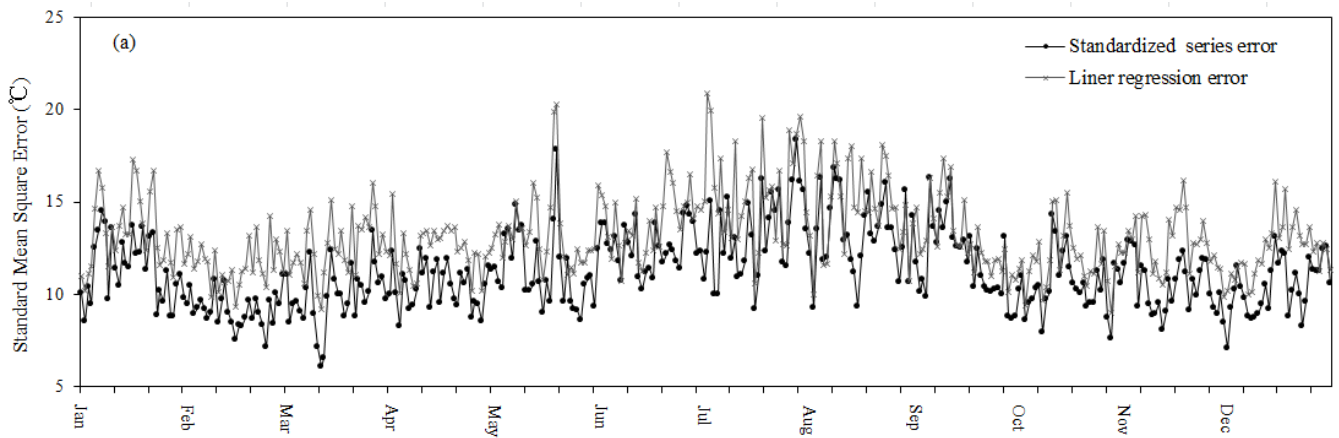
76 In Eq. (S11),  $\varphi'_i = \varphi_i - \overline{\varphi_i}$ ,  $\overline{\varphi'} = \frac{1}{m} \sum_{i=1}^m \varphi'_i$ ,  $\varphi_i = \hat{y}_i - y_i$

$$77 \quad P = \frac{m_p}{m} \times 100\% \quad (S12)$$

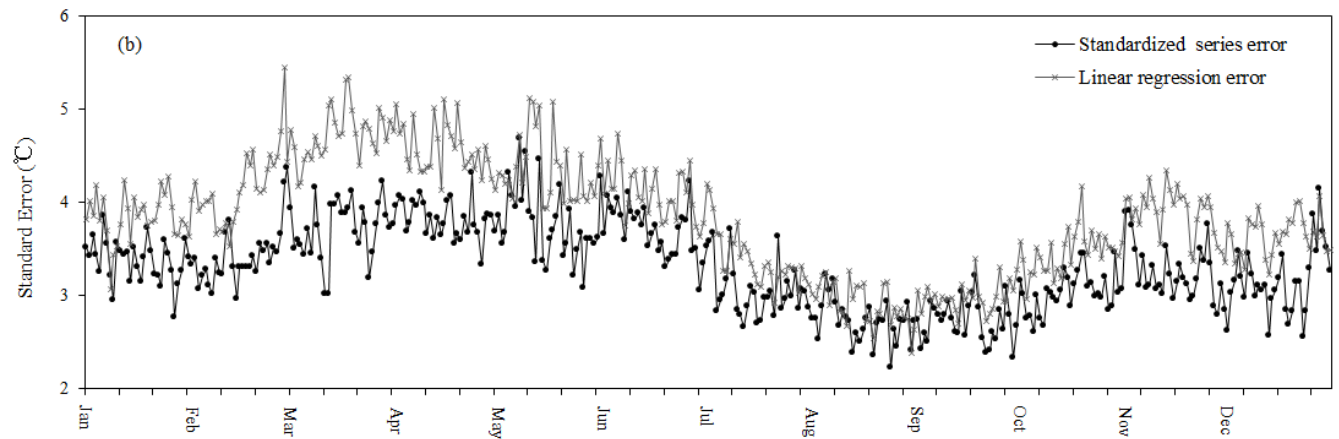
78 In Eq. (S12),  $m_p$  is the number of differences between fitted data and the actual observations within  
 79  $\pm 0.5$  °C.

80 Figure S1 displays the errors between fitted and actual observations of daily maximum temperature  
 81 series calculated using Eq. (S10) - (S12) in the period of Jan 1 1891 - Dec 31 2018 at Tianjin station.

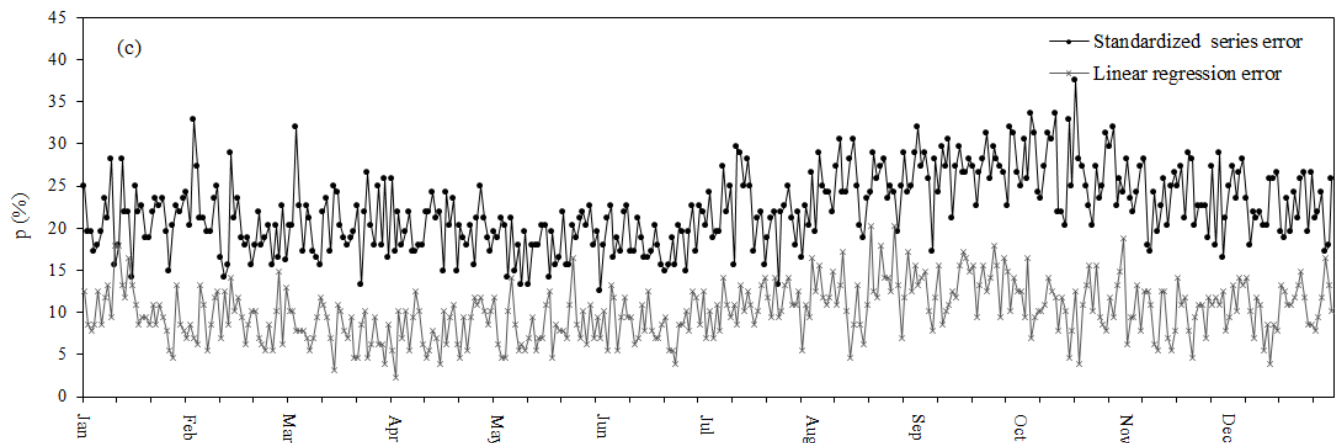
82



83



84



85 **Figure S1.** The indices of Standard Mean Square Error (*SMSE*) (a), Standard Error (*SE*) (b), and *P* (c) for the two fitted  
 86 daily maximum temperatures at Tianjin station from Jan 1 1891 to Dec 31 2018.

87 **Table S1.** The statistics of *SMSE*, *SE*, and *P* for the two fitted data during Jan 1 1891-Dec 31 2018.

| <i>SMSE</i> (°C) |        | <i>SE</i> (°C) |        | <i>P</i> (%) |        |
|------------------|--------|----------------|--------|--------------|--------|
| standardized     | linear | standardized   | linear | standardized | linear |

|         | series | regression | series | regression | series | regression |
|---------|--------|------------|--------|------------|--------|------------|
| Maximum | 18.4   | 20.9       | 4.7    | 5.4        | 37.5   | 20.3       |
| Minimum | 6.1    | 9.0        | 2.2    | 2.4        | 12.5   | 2.3        |
| Median  | 11.0   | 12.8       | 3.3    | 3.8        | 21.9   | 10.2       |

88 As shown in Fig. S1, it is obvious that errors are different in the case of the two fitted daily data at  
89 Tianjin station. Both of them are above 93% of *SMSE* (Fig. S1a) and *SE* (Fig. S1b) using the  
90 standardized series method which appears to be smaller than that from the linear regression method  
91 covering 1891-2018. Index *P* in both the cases is over 99% of difference ratio within  $\pm 0.5$  °C from  
92 standardized series method which is higher than that from the linear regression method (Fig. S1c). As  
93 shown in Table S1, standardized series exhibits smaller error which agrees well with the statistics *SMSE*,  
94 *SE* and *P* in Fig. S1. The maximum, minimum and median values of *SMSE* and *SE* from the standardized  
95 series method are all smaller than those from the linear regression method. For index *P*, the median of  
96 difference ratio within  $\pm 0.5$  °C between the standardized series and the actual observation is 21.9%, but it  
97 is just 10.2% for the linear regression fitted series. The minimum value of index *P* from the linear  
98 regression method is only 2.3%, while it is 12.5% from the standardized series method. Hence based on  
99 error analysis, the fitted daily maximum temperature series by the standardized series method is selected  
100 for the extension during the periods from Jan 1 1887 to Aug 31 1890 at Tianjin station. However, the  
101 daily minimum series still begins with the date on Sep 1 1890 due to scarcity of reference data sources.

102