

# Modeling a drought index using a nonparametric approach

Ghamghami M.<sup>1</sup>, Hejabi S.<sup>1</sup>, Rahimi J.<sup>1</sup>, Bazrafshan J.<sup>1</sup>, Olya H.<sup>2,\*</sup>

<sup>1</sup>Meteorological Division, Dept. of Irrigation and Reclamation Engineering, University of Tehran, Karaj-Iran

<sup>2</sup>College of Hospitality & Tourism Management, Sejong University, 209, Neungdong-ro, Gwangjin-gu, Seoul 05006, Republic of Korea.

Received: 11/03/2016, Accepted: 17/12/2016, Available online: 16/02/2017

\*to whom all correspondence should be addressed:

e-mail: [Olya@sejong.ac.kr](mailto:Olya@sejong.ac.kr)

## Abstract

Alternatively, to other studies that used parametric distributions (e.g. Gamma) in the estimation of the Standardized Precipitation Index (SPI), this study aims to apply a nonparametric method based on Kernel Density Estimator (KDE) for calculating the SPI. Results of the proposed method were compared with the ones from the most widely used parametric distribution, using a long dataset of monthly precipitation of four meteorological stations in Iran (including Bushehr, Mashhad, Tehran and Esfahan) over a period of 107 water years (1895-2002). The capability of KDE-based SPI was compared with the Gamma-based SPI at four-time scales of 3, 6, 9 and 12 months. The frequencies of the drought classes of SPI were calculated and compared with corresponding expected frequencies. The results revealed that the KDE is more consistent with the expected values of the SPI drought/wet classes frequencies (especially in the extreme classes) at all stations as well as at the four-time scales, compared to the Gamma distribution. The greatest deviation from the expected frequencies for KDE and Gamma distribution were about 10% and 150%, respectively. This study proposes a new analytical approach in modeling SPI that provides more accurate results pertaining frequency of occurrences of extreme drought events. The output of the study can be used in many fields (e.g. tourism, agriculture, insurance, etc.) that are influenced by severe droughts.

**Keywords:** SPI, Gamma Distributions, Kernel Density Estimation, extreme drought events.

## 1. Introduction

The presence of an appropriate and accurate drought index to determine dry spells using a quantitative analysis is necessary and helpful in many disciplines (Efthimiou and Karavitis, 2016; Silva, 2003). Index-based assessment has been recommended for effective management of natural phenomena (Anđelković *et al.*, 2016; Olya and Alipour, 2015ab). Depending on the variety of the available data, different indices such as Palmer drought severity index (Palmer, 1965), crop moisture index (Palmer, 1968), surface water supply index (Shafer and Dezman, 1982), drought index (Bhalme and Mooley, 1980) standardized

precipitation index (McKee *et al.*, 1993), reconnaissance drought index (Tsakiris and Vangelis, 2005), stream-flow drought index (Nalbantis and Tsakiris, 2009) joint deficit index (Kao and Govindaraju, 2010) standardized precipitation evapotranspiration index (Vicente-Serrano *et al.*, 2010) have been used for quantitative assessment of the drought phenomena in different spatial and temporal scales. Various drought studies have shown that the contribution of precipitation in these indices is more important compared to other climatic variables, because precipitation is able to justify over 80% of the variability in these indices (Keyantash and Dracup, 2003; Moorhead *et al.*, 2015).

The Standardized Precipitation Index (SPI) is considered as an appropriate index due to its simple calculation, availability of its input data across the world, flexibility of time scale, and comparability of droughts by time and space (Hayes *et al.*, 1999; Guttman, 1999; Mishra and Singh, 2010; Moreira, 2015; Mundetia and Sharma, 2015; Ozelkan *et al.*, 2016). The SPI is calculated based on the precipitation probability in different time scales. Thus, fitting an appropriate probability distribution to the time series of precipitation at a certain time scale is the first step for calculating this index. According to the existing literature, it was determined that the parametric Gamma distribution demonstrates a suitable fitting to the monthly precipitation data series. Therefore, the SPI theory was firstly introduced on the basis of the Gamma distribution as a parametric function (Edwards and McKee, 1997; Guttman, 1999; Thom, 1996). However, parametric methods have some limitations in constructing drought indices; the empirical probability with distribution-free function can be used as a more reliable alternative for calculating a nonparametric standardized index (Huang *et al.*, 2015).

Calculation of parametric distributions is based on parameters that determine the properties of the probability curve such as shape, skewness and kurtosis. This can lead to misleading results, especially at a local scale. On the other hand, nonparametric statistical distributions (e.g. KDE) estimate the probability density function of observations, not only by taking into account

various parameters, but also by using all observations (depending on the Kernel Density Function being continuous or discrete). Compared with parametric methods, the main advantage of nonparametric methods is their functionality as a reliable tool for modeling natural phenomena in various fields, such as the generation of climatic data like temperature and precipitation (Lall *et al.*, 1996; Rajagopalan *et al.*, 1997b; Mehrotra *et al.*, 2006; Srikanthan *et al.*, 2004; Sharif and Burn, 2006), generation of hydrologic data such as streamflow (Sharma and O'neil, 2002) and plan growth modeling (Gommès, 2006). In drought-based studies, a nonparametric method was used by Cancelliere *et al.* (2006) to estimate transition probabilities of SPI classes corresponding to different severities of droughts. Similarly, Kim *et al.* (2003) used a nonparametric KDE approach to estimate bivariate return periods of drought in Mexico. In recent studies, nonparametric approaches in modeling drought indices have been frequently used and recommended by many researchers (Farahmand and AghaKouchak 2015; Hao and AghaKouchak, 2014; Huang *et al.*, 2015; Solakova *et al.*, 2013; Zhu *et al.*, 2015). For example, Solakova *et al.* (2013) compared parametric and nonparametric approaches for the calculation of two drought indices: the standardized precipitation index (SPI) and the standardized streamflow index (SSI). They utilised the Kolmogorov-Smirnov goodness-of-fit test for selecting a parametric distribution. For the nonparametric approach, the Weibull plotting position has been used to calculate the cumulative frequency. Results of time series analyses with both parametric and nonparametric approaches revealed how the differences between these two approaches are more evident in terms of severity and less in terms of duration, and inter-arrival time.

Hao and AghaKouchak (2014) noted drought monitoring based on a single variable may be insufficient for detecting drought conditions in a prompt and reliable manner, because of the complexity of drought phenomena in their causation and impact. They proposed a multivariate and multi-index drought monitoring framework, namely, the multivariate standardized drought index (MSDI), for describing droughts based on the states of precipitation and soil moisture. MSDI is calculated based on an empirical cumulative probability distribution, the Weibull plotting position formula, which compared to the Gamma-based SPI, provides more realistic and precise results. In this regards, Huang *et al.* (2015) and Zhu *et al.* (2015) introduced the Nonparametric Multivariate Standardized Drought Index (NMSDI), where precipitation and streamflow information were coupled to investigate the spatial and temporal characteristics of drought structure on the ground. They found that parametric methods have some limitations in constructing drought indices and compared to nonparametric techniques provide less accurate results.

Farahmand and AghaKouchak (2015) introduced the Standardized Drought Analysis Toolbox (SDAT) that offers a generalized framework for deriving nonparametric univariate and multivariate standardized indices. SDAT

functions based on a nonparametric framework can be applied to different climatic variables including precipitation, soil moisture, and relative humidity, without pre-assumption of representative parametric distributions. The most attractive feature of the framework is that it leads to statistically consistent drought indicators based on different variables. These studies demonstrated the superiority of nonparametric methods over parametric ones. Nevertheless, application of nonparametric continuous probability density functions in the SPI calculation procedure as an alternative method and their comparison with parametric approaches in terms of frequencies of extreme events have not been assessed. Sienz *et al.* (2012) compared several parametric methods in terms of frequencies of drought classes.

In this study, apart from drought monitoring, a nonparametric algorithm based on the KDE function to obtain the best results of a precipitation-based index such as the SPI is presented. This study contributes to the current knowledge of drought by: (I) modeling of drought using a large data set that provides more accurate estimation of the occurrences of severe drought events, (II) applying a nonparametric continuous probability density function with high flexibility and low complexity and (III) focusing on frequencies of drought classes, which are important in accurate calculation of return periods, especially in hydraulic studies.

## 2. Methodology

### 2.1 Preliminary Analysis of Data

In this study, long-term records of monthly precipitation data (1895-2002) are collected from four weather stations of Iran (Tehran, Bushehr, Esfahan, and Mashhad) to calculate the Standardized Precipitation Index (SPI). The length of precipitation records can affect the frequency of severe droughts. However, at least 30 years of precipitation data is sufficient for calculation of drought indices (McKee *et al.*, 1993), but a large sample size provides more accurate and reliable results for modeling of drought (Moreira, 2016). To calculate more reliable return periods of severe drought events, a long memory of the process is helpful, especially when a new approach is implemented and tested within a drought index. Thus, a 107-years period is selected as a dataset for modeling the SPI using a nonparametric technique. Before 1950, these datasets were collected and recorded by the World Weather Records and after that time by Iran Meteorological Organization (Khalili, 1996). Mean and standard deviation of precipitation and geographical properties of these four stations are shown in Table 1.

The time series help to statistically address data gaps caused by data unavailability in some years, especially during the world wars. Reconstruction of these gaps has been calculated using the auto-correlation method (Box *et al.*, 1994). The auto-correlation method is preferred over other methods (e.g. regression and nearest neighbor methods) mainly because the process of reconstructing data does not only depend on the availability of data from neighboring stations but also preserve the temporal

correlation of the data. To check the homogeneity of annual data, the datasets are divided into two periods

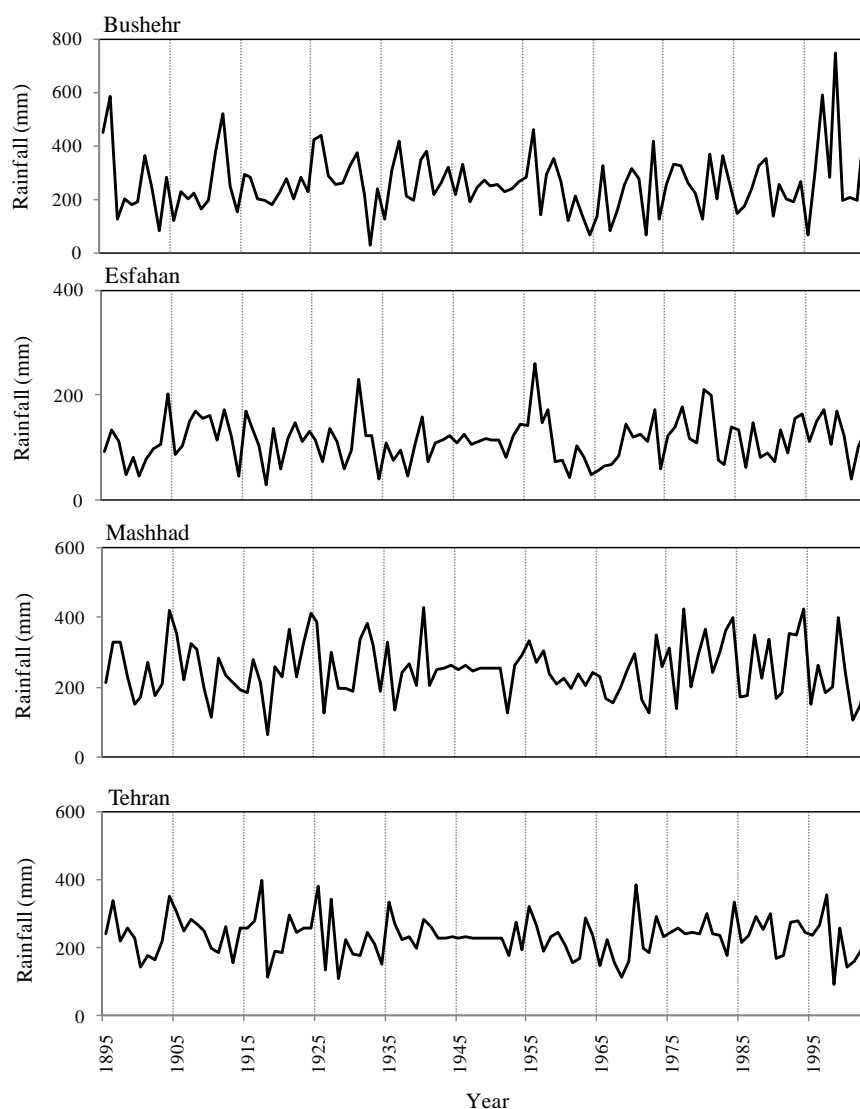
(before and after 1950) to be compared using the Run test (Kruger *et al.*, 2002).

**Table1.** The geographical properties and descriptive statistics of precipitation at four weather stations

Station name	Geographical coordinate		Height (m)	Precipitation Avg. (mm)	Precipitation Sd. (mm)
	X	Y			
Bushehr	50.83	28.95	19.6	257.7	115.5
Esfahan	51.76	32.64	1550.4	115.3	42.8
Mashhad	59.66	36.25	990	253.1	79.6
Tehran	51.35	35.66	1195.8	233.5	59.5

The result of the Run test shows that the datasets are homogeneous. In Fig. 1, time series of annual precipitation data are presented for all four stations. Trend analysis of annual and monthly data is very important due to the large data record available. Detrending is a very important pre-processing technique in case a negative or a positive trend

is detected in the data. As a trend can be caused by climate change, two methods namely the Mann-Kendall test and regression analysis based on Least Squares Errors are applied to evaluate significant trends for monthly and annual data.



**Figure 1.** Annual precipitation series at four selected stations

## 2.2 Standardized Precipitation Index (SPI)

The SPI was introduced by McKee *et al.* (1993) to monitor and classify the drought/wet events based on the precipitation time series at a given point. This index is considered as a meteorological drought indicator calculated based on precipitation. The flexibility of SPI

enables researchers to calculate drought severity at different time scales, which can be used for various types of droughts such as meteorological, agricultural and hydrological droughts. Meteorological and agricultural droughts occur at short time scales, and hydrological drought occurs at long time ranges (Smakhtin and Hughes,

2004; Keyanthash and Dracup, 2004; Barua, 2010). The SPI algorithm is based on the probability transformation of the Cumulative Distribution Functions (CDF) of aggregated precipitation at a moving time window to obtain a standard normal variable or SPI. Standardizing of the index is necessary because it results in a dimensionless index making it comparable in different times and locations (Guttman, 1999). In this study, the SPI was calculated at several time scales such as 3, 6, 9, and 12 months. The 3-month time scale indicates seasonality of the precipitation process (Ji and Peters, 2003) and the 12-month time scale corresponds to a medium-term trend in the precipitation

pattern (Potop *et al.*, 2012) and may provide an annual estimation of hydrological condition.

Since the SPI is a normalized variate, it is expected that the occurrence probabilities of different wet and drought classes follow the Normal distribution. These probabilities are represented in Table 2 as the *expected probability for each class of SPI*. The extreme, severe and moderate drought classes are represented by 2.3%, 4.4% and 9.2% of the SPI values, respectively. The remaining values are allocated to the wet and normal classes (Table 2).

**Table 2.** The SPI classes and their expected probabilities based on the standard normal distribution.

SPI class	Class symbol	Description	Expected probability (%)
Larger than 2.0	Ew	extreme wet	2.3
1.5 to 2.0	Sw	severe wet	4.4
1.0 to 1.5	Mw	moderate wet	9.2
-1.0 to 1.0	N	normal	68.2
-1.5 to -1.0	md	moderate drought	9.2
-2.0 to -1.5	Sd	severe drought	4.4
Less than -2.0	Ed	extreme drought	2.3

The deviation of the aggregated precipitation CDF from the standard Normal distribution leads to overestimation or underestimation of the frequencies of estimated classes with respect to the expected ones (Table 2) (Sienz *et al.*, 2012). This is highly dependent on well-fitting of the CDF to the data. For example, if the frequency of the extreme drought class is estimated equal to 2.4% (2.2%) (from a given CDF fitted to precipitation data), it will be about 4% overestimated (underestimated) compared to the expected frequency of the same class in Table 2.

### 2.3 Nonparametric distribution of Kernel Density Estimator (KDE)

The KDE method estimates the density function without assuming that the data must conform to a particular parametric distribution. Hence, checking the goodness of fit does not require to perform statistical tests (Lall, 1995). In contrast, to calculate the density function using a parametric distribution requires that the distribution fits the data well. Sharma *et al.* (1998) developed the KDE method. In this method, to estimate the PDF of the data, a given observation ( $x$ ) is selected among other observed values and the contribution of each observation is determined by a kernel density function. The effective parameter in this function is the *bandwidth*. The function value of variable  $x$  is obtained from equation (1):

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) \quad (1)$$

Where  $K\left(\frac{x-x_i}{h}\right)$  is equal to 1 for  $x_i$  within the *bandwidth* range,  $h$ , and equal to zero for  $x_i$  out of this region. In fact, KDE allows all observations to participate in the estimation of the PDF value of a certain observation ( $x$ ). Different forms of kernel functions have been used in different studies, the most popular being the *standard normal*

*function*. It has been analytically proven that the form of the kernel function has no major role in the performance of the method (Rajagopalan *et al.*, 1997a; Dinardo and Tobis, 2001). The value of the density function is estimated from equation (2) and based on the standard normal function:

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n \frac{1}{(2\pi)^{1/2}} \exp\left[-\frac{(x-x_i)^2}{2h^2}\right] \quad (2)$$

The value for the exponential component of this equation is between 0 and 1, and it determines the participation level of independent data in estimating the value of the PDF for each  $x$ . Observations closer to the  $x$  observation have a higher contribution in estimating the value of the PDF. Hence, to estimate the density of any observation, a standard normal Kernel function, centered on the given observation, is fitted to the data. In other words, the  $n$  normal kernel functions are fitted to the  $n$  independent observations to determine the probability density function. In addition, determination of a suitable value for  $h$  is important. The large values of the *bandwidth* cause an excessively even estimation and its small values result in an uneven estimation with additive variance. There are different methods to estimate the optimized *bandwidth*. However, Silverman (1986) introduced an analytical equation (3) to estimate the optimized *bandwidth*:

$$h = 1.06\sigma n^{-1/5} \quad (3)$$

The estimated *bandwidth* by equation (3) is referred to as the *global bandwidth* in references. The *local bandwidth* was used to estimate the probability curve with a high level of accuracy (Sharma, 1996). In this case, a separate  $h$  parameter was defined to estimate every observed PDF. Abramson (1982) introduced the relation between the *local* and *global bandwidth* as:

$$h_i = h \left[ \hat{f}(x_i) / g \right]^{-1/2} \quad (4)$$

Where  $\hat{f}(x_i)$  is the probability density function based on the *global bandwidth* obtained from equation (3) and  $g$  is the geometric mean of  $\hat{f}(x_i)$ . According to equation (2), the optimized *bandwidth* is the only required parameter by the aforementioned nonparametric method and is obtained by means of simplified formulas.

### 3. Results and Discussion

In this section, results of the SPI based on KDE and Gamma distribution are evaluated and compared with estimations of expected frequency of SPI classes using long records at

weather stations of Iran. The SPI is calculated at four scales 3, 6, 9 and 12 months aiming at the evaluation of the flexibility of the PDFs (KDE and Gamma).

#### 3.1 Trend Analysis Results

The results obtained from the Mann-Kendall method, on monthly and annual scales are presented in Table 3. According to the results, no significant trend was found for the annual scale, at all stations considered. On a monthly scale, a negative trend was observed in 37.5% of the series, while the rest of them did not follow any specific linear/nonlinear trends. To re-examine the significant trends, an analysis based on linear regression for the annual time series was applied and the statistical significance of the trends was evaluated using a t-student test.

**Table 3.** Trend analysis using Mann-Kendall method at the studied stations, \* Statistically significant at a 5% significant level.

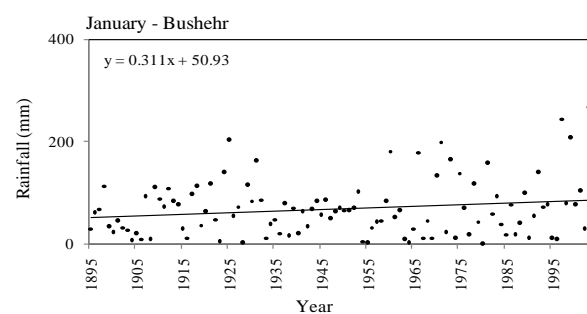
Time series	Bushehr	Esfahan	Mashhad	Tehran
Oct.	-7.798*	-3.891*	-0.56	-1.864
Nov.	0.037	-0.738	-0.147	0.492
Dec.	0.1	-1.016	-1.629	-0.77
Jan.	-0.702	-0.765	-1.587	0.388
Feb.	-0.649	-1.194	-0.765	0.676
Mar.	-1.220	-1.278	0.885	0.885
Apr.	-1.922*	0.529	1.121	1.713
May	-6.525*	-1.393	1.618	-0.608
Jun.	-15.292*	-6.808	-0.697	-2.634*
Jul.	-15.308*	-9.505*	-7.646*	-4.19*
Aug.	-12.433	-9.961	-8.584	-5.284
Sep.	-15.413	-9.285	-7.693	-5.117*
Water year	0.047	-0.911	0.1	0.304

The negative and positive slopes represent the descending and ascending trends, respectively. The only statistically significant trend on a monthly scale was confirmed at the Bushehr station in January (Fig. 2). Based on these results, there is no overall statistically significant trend (negative or positive) over all time series. With regards to the difference between the findings of the two methods, it is not possible to identify the existence of significant trends in precipitation time series. The Mann-Kendall test did not perform well for the summer months. When zero values increase in the series, the Mann-Kendall statistic tends to higher negative values and it can come up with artificially statistically significant results. The results of the regression analysis confirmed the lack of linear trends in precipitation series, except for January (Fig. 2). The findings are consistent with Khalili's (1996) study in which none of the tests detected significant climate change effects on precipitation for the data of the four stations.

#### 3.2 KDE parameter

To estimate the *global bandwidth*, the value of  $\hat{f}(x_i)$  was calculated on a monthly scale. Estimating the geometric mean of  $\hat{f}(x_i)$ , the *local bandwidth* for any event was exclusively determined, thus avoiding excessively smooth

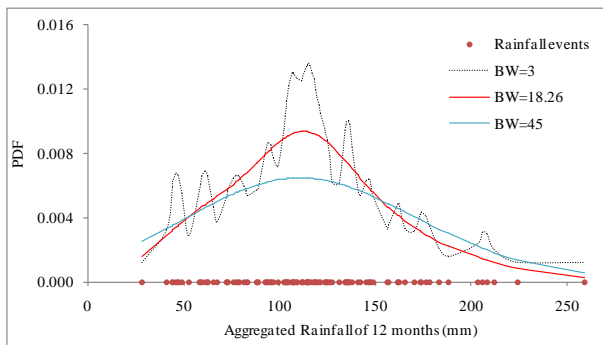
and uneven estimations in the vicinity of distribution tails and modes, respectively.



**Figure 2.** A linear regression model fitted to precipitation data of January at Bushehr station. The t-student statistics is equal to 2.656 identifying a statistically significant trend

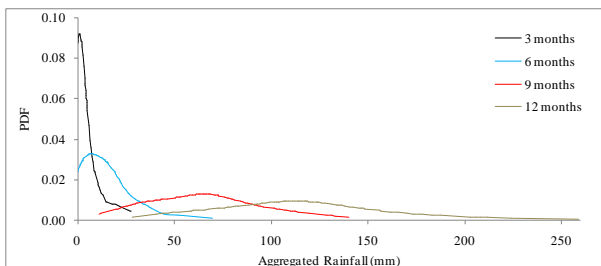
According to equation (4), ranges with low densities (such as tails) have a large *bandwidth* and ranges with high densities (such as modes) have a low *bandwidth*. To determine the *local bandwidth*, however, the accurate determination of the *global bandwidth* is important. The importance of the primary choice of the *bandwidth* is illustrated in Fig. 3. The graph shows PDF curves for the aggregated series of precipitation of Esfahan station at a

12-month time scale (from February to January). The points on the horizontal axis represent the 107 precipitation events. The KDE function is fitted to the aforementioned data with three different *global bandwidths*. The red curve indicates the KDE function considering the one that equals 18.26, determined by equation (3) as *referenced bandwidth*. The dotted and blue curves represent *bandwidths* lower and higher than the one of the red curve, respectively. Setting the bandwidth smaller and larger than the *referenced bandwidth*, the PDF curves appear strongly uneven and even, respectively.



**Figure 3.** Effect of the bandwidth value (BW) on the smoothness of the PDF curves related to the 12-month aggregated series of precipitation (i.e. sum of precipitation from February of last year to January of current year) at the Esfahan station

The KDE function can be properly fitted to positively or negatively skewed data as well as to symmetric ones without the need for other parameters. This is a distinct advantage of the KDE compared to parametric methods such as the Gamma distribution. The PDF curves based on the KDE for four-time scales are presented in Figure 4.



**Figure 4.** PDF curves based on the KDE for four-time scales related to the aggregated series of precipitation of January recorded at Esfahan station

The curves demonstrate the aggregated series of precipitation of January recorded at Esfahan station. It can be clearly seen that, as the time scale of SPI increases, the frequency of zero values would decrease and the PDF curve would be closer to the one of the Normal distribution. Accordingly, the 3-month curve is positively skewed and the 12-month curve is almost symmetric. Thus, the KDE function successfully generates the variabilities in the PDF curves of the various ranges of observations (Fig. 3).

### 3.3 Comparisons

In this subsection, the performance of two PDFs is compared in terms of the frequencies of drought SPI

classes. According to the different structure nature of the two methods, performing the conventional statistical comparisons such as applying goodness-of-fit tests would not be possible. Therefore, comparisons of the frequencies of drought SPI classes have been considered. As previously stated, any deviation from the probabilities mentioned in Table 2 results in underestimation or overestimation of the frequency of a drought class, which is important in drought risk assessment. The consequence of such abnormalities affects the theoretical computations of drought return periods.

Bar graphs of different percentages between expected frequencies and the ones resulting from PDFs are provided in Fig. 5 to 8. Each graph shows four bar graphs for the various time scales- 3, 6, 9, and 12 months. The black and white bars indicate the KDE and Gamma PDFs, respectively. Bars above the horizontal axis (positive deviations) represent overestimation and bars under the horizontal axis (negative deviations) represent underestimation. The reason for these deviations is directly connected to the appropriateness of the PDF. This is particularly important in the lower tail of the precipitation distribution representing extreme droughts, as any overestimation (underestimation) of the probabilities in this area leads to underestimation (overestimation) of the frequency of extreme events. Since SPI was used as the drought index, an analysis of the drought classes was conducted.

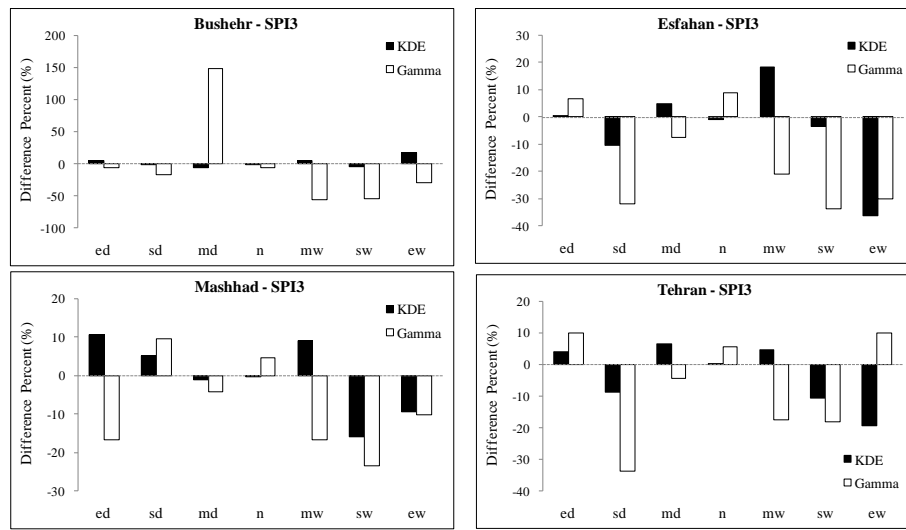
Fig. 5 represents the SPI on a 3 months' time scale. In Bushehr, the frequency of extreme droughts has been about 4% overestimated by the KDE and about 7% underestimated by the Gamma distribution. The KDE seems to perform relatively better, however, both methods do not show significant deviations. The frequency of severe droughts has been underestimated by the Gamma about 16%, while the Gamma distribution is judged extremely unacceptable for moderate droughts (overestimation up to 150%). In both extreme and severe classes, the KDE provides relatively satisfactory outcomes.

In Esfahan, the frequency of extreme droughts produced by the KDE is close to the expected value, whereas the Gamma distribution overestimates it about 7%. The frequency of severe droughts has been underestimated by both KDE and Gamma methods, but KDE is relatively more efficient. This condition also holds for moderate droughts (Fig. 5). KDE (Gamma) overestimated (underestimated) up to 11% (up to 17%) the frequency of extreme drought in Mashhad. Findings for the KDE overperform the ones for the Gamma in all other drought classes, too. Tehran's findings are the same as the ones for Esfahan, excluding the moderate droughts where the Gamma distribution provides better results. According to an overall assessment of the aforementioned results, the KDE is judged more suitable than the Gamma to estimate the frequency of drought classes on the time scale of 3 months. In some positions, however, the Gamma distribution provides acceptable results because it fits the data well.

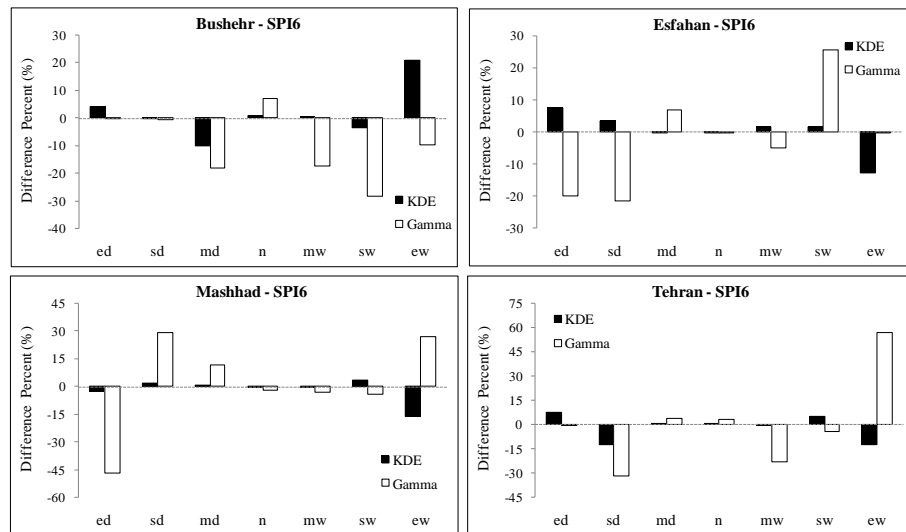
Fig. 6 illustrates differences (in %) between the estimated frequencies of each class of the 6 months SPI based on KDE and Gamma distribution and the corresponding expected

frequencies mentioned in Table 2. In an overall view, the differences for KDE are less than 10% and for the Gamma distribution are larger than 45%. To provide more information, in Bushehr KDE has a poorer performance than the Gamma distribution (overestimates about 5%,

while for the Gamma this percentage is close to 0). Both methods estimated the frequency of severe droughts close to the expected value (percentages close to 0). In the case of the moderate class, both methods underestimate the expected PDF, indicating that KDE is better.



**Figure 5.** Percentages of the 3 months SPI classes based on KDE and Gamma distribution for the stations of interest. The class symbols (x-axis) were introduced in Table 2



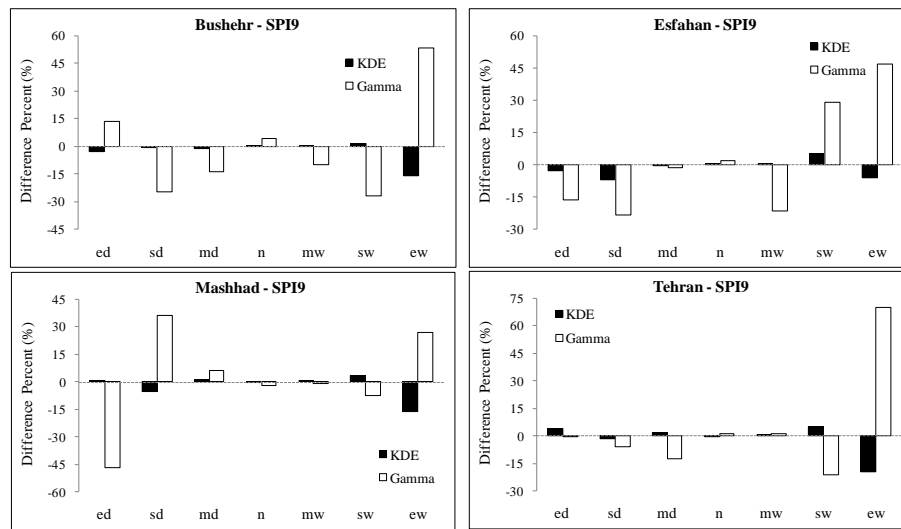
**Figure 6.** Percentages of the 6 months SPI classes based on KDE and Gamma distribution for the stations of interest

In Esfahan, the Gamma has underestimated the frequencies of extreme and severe drought classes up to 20% while KDE has overestimated them below 10%. Also, the better result is found for KDE in the case of the moderate drought class. The graph for Mashhad highlights considerable underestimation about 45% and overestimation about 30% by the Gamma distribution in extreme and severe classes, respectively. Results by the KDE do not show considerable deviations for both classes. Such a result is more gently observed for the moderate class. In Tehran, the Gamma distribution produces satisfactory results for estimating extreme drought class of SPI as well as KDE for other classes of SPI. To sum up, the KDE method fitted to the aggregated series of precipitation

at a time scale of 6 months, as for the case of the 3 months' time scale, performs better than the Gamma distribution.

As shown in Fig. 7, results of the KDE have a better overall efficacy compared to the Gamma at the time scale of 9 months. Estimating the frequency of extreme droughts, the largest percentage deviation of KDE is about 4% (Tehran) and of the Gamma distribution it is about 46% (Mashhad).

Better results for the Gamma distribution are obtained in the extreme drought class in Tehran in which different percentages are close to 0. Similar results are found for the time scale of 12 months (Fig. 8). Therefore, KDE offered significantly better performance in comparison with the Gamma distribution in terms of accuracy of SPI results and of its integrity in all time scales.

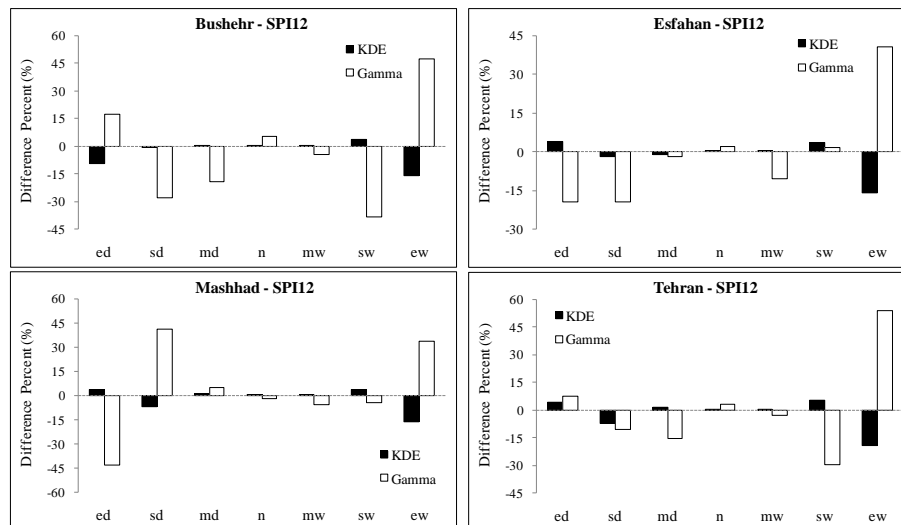


**Figure 7.** Percentages of the 9 months SPI classes based on KDE and Gamma distribution for the stations of interest

According to all bar graphs (Fig. 1-8), in both KDE and Gamma methods, the percentage of the normal class of SPI (representative of the middle part of SPI distribution) is slightly different (the difference being less than 10%) than the corresponding expected percentage of Table 2. The maximum differences for the normal class of SPI based on KDE and Gamma were 1% and 9%, respectively, which occurred in Esfahan and at the time scale of 3 months.

According to the results, the KDE is more reliable than the Gamma in the middle of the distribution (i.e. normal class) as well as in the tail (i.e. drought classes). The Gamma has created deviations more than 100% in the tail of the distributions, while its efficacy is better in the middle.

To demonstrate the effects of the deviations on historical series of SPI based on KDE, Gamma, and empirical distribution, a given period (1960-1976), with extreme events, has been shown in Fig. 9.



**Figure 8.** Percentages of the 12 months SPI classes based on KDE and Gamma distribution for the stations of interest

The segment of SPI's time series given in Fig. 9 belongs to Bushehr station. In this graph, the calculated SPI9 based on the theoretical functions is shown as a curve graph while the bar graph represents SPI9 based on the empirical function. The graph covers the period of 1960-1975. According to Fig. 7, the Gamma function overestimated the class of extreme droughts up to 15%. The same outcome is clearly perceived from the Fig 9.

#### 4. Discussions

The reason for the superiority of a given distribution compared to other distributions may be related to the statistical features of the available data and flexibility of the

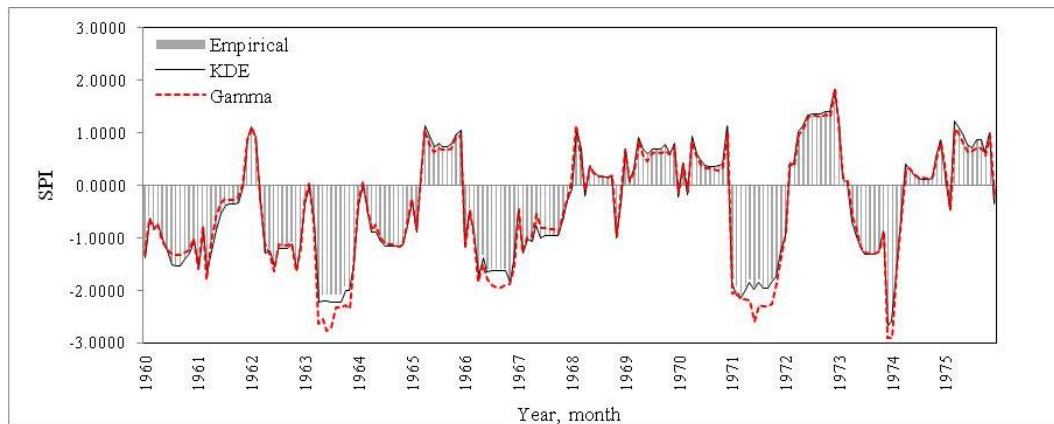
given distribution. For example, the Normal distribution is not suitable for positively or negatively skewed data. The Gamma distribution may be a suitable option when data are positively skewed due to the existence of the shape parameter in its structure, moreover, there is a higher probability that it can be selected as the best distribution when data are skewed (Guttman, 1999).

In spite of the flexibility of the Gamma to fit various data, the quite bad fitting to the distribution tails, median and mode of the data results in considerable deviations in frequencies of drought classes. These results are important for post-monitoring processes such as calculating return periods. An underestimation (overestimation) in the



frequencies of SPI drought classes can affect the return period of droughts. For example, if the return period of a drought class is equal to 20 years, the underestimation by a theoretical distribution equals 50% resulting in return periods of that drought class changing from 20 to 40 years. Thus, the wrong selection of a distribution provides

misleading results. Since the concept of return period plays a key role in studies of hydrology and water resource management, reduction of the estimation error of drought frequencies is necessary for efficient management and making right decisions under drought conditions (Sienz *et al.*, 2012).



**Figure 9.** Comparison of the historical values of SPI in Bushehr and at a 9 months' time scale

Utilization of KDE in SPI structure has two advantages: 1. The KDE only has one parameter easy to calculate mathematically. Furthermore, KDE does not require the inclusion of other parameters to assess skewness and kurtosis of data; 2. The KDE presents more efficient results than the Gamma distribution which is revealed by comparisons of percentage differences at each class of SPI. Since KDE was not embedded in SPI in the past studies, it is unlikely to compare the results of this study with findings of other studies. Nevertheless, comparison of the findings of the study with past studies that used other nonparametric approaches, as presented in Section 3.3, indicated that the results are in accordance with findings of Solakova *et al.* (2013), Hao and AghaKouchak (2014), Farahmand and AghaKouchak (2015). Moreover, the underperformance of the Gamma distribution function is acknowledged by Sienz *et al.* (2012).

## 5. Conclusion

A suitable CDF provides more accurate results, especially for extreme drought events. Small differences in smoothing distribution tails can create significant errors in calculating frequencies of drought classes as well as return periods of hydrological and meteorological drought events. Application of a nonparametric approach in calculation of a drought index (i.e. probability transformation) provide more precise estimation of that index.

To achieve the main aim of this study, the efficiency of a parametric (Gamma) and a nonparametric (KDE) probability density function are compared to estimate frequencies of drought classes from SPI at 3, 6, 9 and 12-month time scales at the selected weather stations. The main features of this study are: 1) this study used a long precipitation dataset to calculate SPI. This helps to compare the obtained frequencies to the expected values, and 2) this study is the first one that applied the KDE nonparametric distribution in the skeleton of the SPI

algorithm. Due to the nonparametric nature of the KDE, its functionality has been also approved by others studies.

As discussed in this study, introducing one fixed parametric distribution might not be possible as it would be a pre-assumption for all months and stations. If one distribution function among parametric methods is considered suitable in a region, it does not necessarily produce the best results for other regions and months, and should be examined for this purpose. In contrast, the KDE nonparametric distribution could be introduced as an overall probability distribution for all months of the year at the chosen weather stations.

Although the parametric methods have quite a high accuracy in the middle part of the data distribution, these methods usually overestimate or underestimate the extremes. Thus, it is necessary to apply a distribution function that is not just suitable for any region but also it provides more accurate estimation of the occurrence of extremes as well as of the observations in the central part of the distribution. The KDE nonparametric method presents both these two advantages.

As aforementioned, the reason of the overestimation/underestimation of the SPI is related to the underestimation/overestimation of the calculated CDF by theoretical functions. However, these drawbacks were considerably improved by using the KDE method. The expected frequencies of the SPI drought classes are very close to the observed values for the KDE nonparametric method, compared to the parametric methods. This is helpful in the drought risk studies where the estimation of return periods is a very important procedure.

According to the results of the present paper, it is recommended that a nonparametric probability distribution based on kernel functions can be used in the SPI computation algorithm. This approach increases the accuracy of drought return period that improves the

process of knowledge-based management of various fields (e.g., agriculture, tourism, insurance, etc.). Generalizing the univariate to the bivariate index enables us to use joint distributions based on KDE with parameters estimated using the Maximum Likelihood or other similar methods. In addition, it is worthy to note the scrutiny needed to select a proper distribution of drought characteristics such as severity and duration. Hence, the three following pathways are recommended for further improvement of the drought research: 1) Rerunning the model with data collected from many stations, 2) monitoring and analysis of drought characteristics in GIS with KDE-based SPI at regional, global, and continental levels, 3) developing KDE-based indices using more than one variable such as temperature and evapotranspiration.

### Acknowledgement

The authors would like to thank the editor and the anonymous reviewers of the Global NEST journal for their constructive insights and helpful comments and suggestions that have helped to improve the quality of the manuscript in terms of content, conceptualization, style, and readability.

### References

- Abramson I.S. (1982), On bandwidth variation in kernel estimates – a square root law, *Ann. Stat.*, **10**, 1217–1223.
- Andelković G., Pavlović S., Đurđić S., Belij M. and Stojković S. (2016), Tourism Climate Comfort Index (TCCI) – An attempt to evaluate the climate comfort for tourism purposes: The example of Serbia, *Global NEST Journal*, **18**(3), 482–493.
- Barua S. (2010), *Drought Assessment and forecasting using a nonlinear aggregated drought index*. PHD Thesis, School of Engineering and Science Faculty of Health, Engineering and Science Victoria University, Australia.
- Bhalme H.N. and Mooley D.A. (1980), Large-scale drought/floods and monsoon circulation, *Monthly Weather Review*, **108**, 1197–1211.
- Box G.E.P., Jenkins G.M. and Reinsel G.C. (1994), *Time Series Analysis: Forecasting and Control* (3rd ed.). Upper Saddle River, NJ: Prentice–Hall.
- Cancelliere A., Bonaccorso B. and Di maruo G. (2006), A non parametric approach for drought forecasting through the Standardized Precipitation Index, *Proceedings of XXXI IAHR Congress on Water Engineering for the future: Choice and Challenges*, Seoul, Korea, 11–16 September 2006, 3252–3260.
- Dinardo J. and Tobis L. (2001), Nonparametric density and regression estimation, *J. Econ. Perspect.*, **15**(4), 11–28.
- Edwards D.C. and McKee T.B. (1997), Characteristics of 20th century drought in the United States at multiple time scales, *Atmospheric Science Paper*, **634**, 1–30.
- Efthimiou N. and Karavitis C. (2016), Towards improving the applicability of the RUSLE model at mountainous Mediterranean catchments, *Global NEST Journal*, **18**(4), 778–793.
- Farahmand A. and AghaKouchak A. (2015), A generalized framework for deriving nonparametric standardized drought indicators, *Advances in Water Resources*, **76**, 140–145.
- Gommes R. (2006), Non-parametric crop yield forecasting. A didactic case study for Zimbabwe. In remote sensing support to crop yield forecast and area estimates. ISPRS Archives XXXVI-8/W48 Workshop proceedings, 79–84.
- Guttman N.B. (1999), Accepting the standardized precipitation index: a calculation algorithm, *Journal of the American Water Resources Association*, **35**, 311–322.
- Hao Z. and AghaKouchak A. (2014), A nonparametric multivariate multi-index drought monitoring framework, *Journal of Hydrometeorology*, **15**(1), 89–101.
- Hayes M.J., Svoboda M.D., Wilhite D.A. and Vanyarkho O.V. (1999), Monitoring the 1996 drought using the Standardized precipitation Index, *Bulletin of the American Meteorological Society*, **80**(3), 429–437.
- Huang S., Huang Q., Chang J., Zhu Y., Leng G. and Xing L. (2015), Drought structure based on a nonparametric multivariate standardized drought index across the Yellow River basin, China, *Journal of Hydrology*, **530**, 127–136.
- Ji L. and Peters A.J. (2003), Assessing vegetation response to drought in the Northern Great Plains using vegetation and drought indices, *Remote Sens. Environ.*, **87**, 85–98.
- Kao S.C. and Govindaraju R.S. (2010), A copula-based joint deficit index for droughts, *J Hydrol.*, **380**, 121–134.  
<http://dx.doi.org/10.1016/j.jhydrol.2009.10.029>.
- Keyantash J. and Dracup J.A. (2003), The quantification of drought: an evaluation of drought indices. *American Meteorological Society*, **83**(8), 1167–1180.
- Keyantash J.A. and Dracup J.A., (2004), An aggregate drought index: Assessing drought severity based on fluctuations in the hydrologic cycle and surface water storage, *Water Resources Research*, **40** (W09304), doi:10.1029/2003WR002610.
- Khalili A. (1996), Statistical evaluation of climate change and PDF, based on secular data in five Iranian old stations. *Proc. Summary of the First Regional Conference of Climate Change*. Tehran, 14–25.
- Kim T., Valdés J.B. and Yoo C. (2003), A nonparametric approach for estimating return periods of droughts in arid regions. Accepted for publication by ASCE, *Journal of Hydrologic Engineering*.
- Kruger A.C., Makamo L.B. and Shongwe S. (2002), An analysis of Skukuza climate data, *Koedoe*, **45**(1), 87–92.
- Lall U. (1995), Nonparametric function estimation. *Recent Hydrologic Applications*, Us National Reports, 1991–1994.
- Lall U., Rajagopalan B. and Tarboton D.G. (1996), A nonparametric wet/dry spell model for re-sampling daily precipitation, *Water Resources*, **32**(9), 2803–2823.
- McKee T.B., Doesken N.J. and Kleist J. (1993), The relationship of drought frequency and duration to time scales. *Proceedings of the 8th Conference on Applied Climatology*, Anaheim, CA, USA, 179–184.
- Mehrotra R., Srikanthan R. and Sharma A. (2006), A comparison of three stochastic multi-site precipitation occurrence generators, *Journal of Hydrology*, **331**, 280–292.
- Mishra A.K. and Singh V.P. (2010), A review of drought concepts, *Journal of Hydrology*, **391**, 202–216.
- Moreira E.E. (2015), SPI drought class prediction using log-linear models applied to wet and dry seasons, *Physics and Chemistry of the Earth, Parts A/B/C*. doi:10.1016/j.pce.2015.10.019
- Moorhead J.E., Gowda P.H., Singh V.P., Porter D.O., Marek T.H., Howell T.A. and Stewart B.A. (2015), Identifying and Evaluating a Suitable Index for Agricultural Drought Monitoring in the Texas High Plains, *JAWRA Journal of the American Water Resources Association*, **51** 3), 807–820.

- Mundetia N. and Sharma D. (2015), Analysis of rainfall and drought in Rajasthan State, India, *Global NEST Journal*, **17**(1), 12-21.
- Nalbantis I. and Tsakiris G. (2009). Assessment of hydrological drought revisited, *Water Resour Manag*, **23**, 881-897.
- Olya H. and Alipour H. (2015a), Modeling tourism climate indices through fuzzy logic, *Climate Research*, **66**(1), 49-63.
- Olya H. and Alipour H. (2015b), Risk assessment of precipitation and the tourism climate index, *Tourism Management*, **50**, 73-80.
- Ozelkan E., Chen G. and Ustundag B.B. (2016), Multiscale object-based drought monitoring and comparison in rainfed and irrigated agriculture from Landsat 8 OLI imagery, *International Journal of Applied Earth Observation and Geoinformation*, **44**, 159-170.
- Palmer W.C. (1965). Meteorological Drought. Research Paper No. 45, U.S. Department of Commerce Weather Bureau, Washington, D.C.
- Palmer W.C. (1968), Keeping track of crop moisture conditions, nationwide: the new Crop Moisture Index, *Weather wise*, **21**, 156-161.
- Potop V., Mozny M. and Soukup J. (2012), Drought evolution at various time scales in the lowland regions and their impact on vegetable crops in the Czech Republic, *Agric. Forest Meteorol.*, **156**, 121-133.
- Rajagopalan B., Lall U. and Tarboton D.G. (1997a), Evaluation of kernel density estimation methods for daily precipitation re-sampling, *Stochastic Hydrology And Hydraulics*, **11**, 523-547.
- Rajagopalan B., Lall U., Tarboton D.G. and Bowles D.S. (1997b), Multivariate nonparametric re-sampling scheme for generation of daily weather variable, *Stochastic Hydrology And Hydraulics*, **11**, 65-93.
- Shafer B.A. and Dezman L.E (1982), Development of a Surface Water Supply Index (SWSI) to assess the severity of drought conditions in snow pack runoff areas. IN Proceedings of the (50th) Annual Western Snow Conference, 164-175. Fort Collins, CO: Colorado State University.
- Sharif M. and Burn D.H. (2006), Simulating climate change scenarios using an improved k nearest neighbor model, *Journal of Hydrology*, **325**, 179-196.
- Sharma A. and O'Neil R. (2002), A nonparametric approach for representing inter-annual dependence in monthly streamflow, *Water Res.*, **138**(7), 5-15-10.
- Sharma A., Lall U. and Tarboton D.G. (1998), Kernel bandwidth selection for a first order nonparametric streamflow simulation model, *Stoch. Hydrol. Hydraul.*, **12**, 35-52.
- Sienz F., Bothe O. and Fraedrich K. (2012), Monitoring and quantifying future climate projections of dryness and wetness extremes: SPI bias, *Hydrology and Earth System Sciences*, **16**(7), 2143-2157.
- Silva V.D.P.R. (2004), On climate variability in Northeast of Brazil, *Journal of Arid Environments*, **58**(4), 575-596.
- Silverman B.W. (1986), Density estimation for statistics and data analysis. Chapman and Hall, New York.
- Smakhtin V.U. and Hughes D.A. (2004), Review, automated estimation and analysis of drought indices in South Asia. Working Paper 83, International Water Management Institute, Sri Lanka.
- Soláková T., De Michele C. and Vezzoli R. (2013), Comparison between parametric and nonparametric approaches for the calculation of two drought indices: SPI and SSI, *Journal of Hydrologic Engineering*, **19**(9), 04014010.
- Srikanthan R., Sharma A. and McMahon T.A. (2004), Stochastic generation of monthly rainfall data using a nonparametric approach. Technical report 02/8, CRC for Catchment Hydrology, Monash University, Clayton, 40p.
- Thom H.C.S. (1966), Some methods of climatological analysis. WMO N. 199. Technical Note N. 81, Ginevra, 53 pp.
- Tsakiris G and Vangelis H. (2005), Establishing a drought index incorporating evapotranspiration, *Eur Water*, **9**(10), 1-9.
- Vicente-Serrano S.M., Beguería S. and López-Moreno J.I. (2010), A multiscalar drought index sensitive to global warming: the standardized precipitation evapotranspiration index, *J Climate*, **23**, 1696-1718.
- Zhu Y., Chang J., Huang S. and Huang Q. (2016), Characteristics of integrated droughts based on a nonparametric standardized drought index in the Yellow River Basin, China, *Hydrology Research*, **47**(2), 454-467.