



Drought frequency analysis using cluster analysis and bivariate probability distribution

Jiyoung Yoo^a, Hyun-Han Kwon^b, Tae-Woong Kim^{c,*}, Jae-Hyun Ahn^d

^a Department of Civil and Environmental Engineering, Hanyang University, Seoul 133-791, Republic of Korea

^b Department of Civil Engineering, Chonbuk National University, Jeonju 561-756, Republic of Korea

^c Department of Civil and Environmental Engineering, Hanyang University, Ansan 426-791, Republic of Korea

^d Department of Civil Engineering, Seokyeong University, Seoul 136-704, Republic of Korea

ARTICLE INFO

Article history:

Received 8 May 2011

Received in revised form 1 November 2011

Accepted 24 November 2011

Available online 1 December 2011

This manuscript was handled by Andras Bardossy, Editor-in-Chief, with the assistance of Fi-John Chang, Associate Editor

Keywords:

Drought
Frequency analysis
Clustering analysis
Bivariate distribution
Drought risk

SUMMARY

Analyses of drought frequency require long-term historical data to ensure reliable quantile estimates. Estimation of quantiles is difficult, because drought extremes are rare by definition, and the durations of extremes are often too short for reliable point frequency analysis. Regional frequency analysis provides a solution for these problems by using data from multiple sites, provided the sites are homogeneous, and this type of analysis yields appropriate estimates of quantiles at sites of interest. This study aims to develop a practical drought frequency analysis method based on a bivariate distribution by incorporating regional drought attributes that are associated with drought frequency (e.g., duration and severity). This study employed a kernel density function to describe joint probabilistic behavior of drought. Given the proposed approach, we estimated return periods according to the most severe drought events on record at each site, and ultimately assess the risks for occurrence of droughts exceeding the most severe droughts over the next 10, 50, 100, and 150 years.

© 2011 Elsevier B.V. All rights reserved.

1. Introduction

Global warming has accelerated since the 1980s, and the frequency of meteorological disasters has dramatically increased worldwide. According to precipitation outlook studies based on changes in rainfall patterns, increasing concentrations of green house gases are likely to trigger very different patterns of heavy rain, extreme drought, and heavy snow in some regions. Indeed, South Korea has experienced much more frequent and extreme rainfall since the late 1990s (Choi et al., 2008). Boo et al. (2004) calculated and analyzed the Palmer Drought Severity Index (PDSI) using the A2 climate change scenario based on a Regional Climate Model, RegCM3, covering the Korean Peninsula. They concluded that drought risk is likely to increase, despite concurrent increases in general precipitation, over the course of the twenty-first century. Meteorological disasters associated with enhanced droughts are expected to become more severe.

Various indices have been used to measure different drought characteristics depending on research objectives. The indices which represent the relationship between precipitation and streamflow are widely used to identify hydrological droughts. Pre-

cipitation-based drought indices with predetermined thresholds are effective because the main cause of drought is rainfall deficit. Precipitation-derived time series are used to investigate drought duration, severity, and frequency (Yevjevich, 1967).

Currently, one of the main limitations of drought analysis is the lack of ability to identify spatial characteristics. This ability is increasingly important because the effects of drought accumulate slowly over a considerable period of time, and move slowly to adjacent positions (Salas et al., 2005). Although an assumption that drought severity and duration are independent is commonly used in practice to make the problem simple, it is not true in real world. The multidimensional characteristics of drought make it difficult for univariate analysis to reveal significant relationships among drought properties (Kim et al., 2003). For this reason, multivariate distributions are supposed to be widely used for characterizing drought properties.

Numerous studies have been conducted to resolve the aforementioned limitations. Rossi et al. (1992) focused on spatial aspects of drought by examining all drought properties. In order to do this, the study quantified droughts at different sites using various types of hydrologic data (e.g., rainfall, streamflow, and reservoir levels) from each of the observation sites within the area of study. Taking into account the better understanding of spatial patterns, a forecasting model was developed and applied. Clause and

* Corresponding author. Tel.: +82 31 400 5184.

E-mail address: twkim72@hanyang.ac.kr (T.-W. Kim).

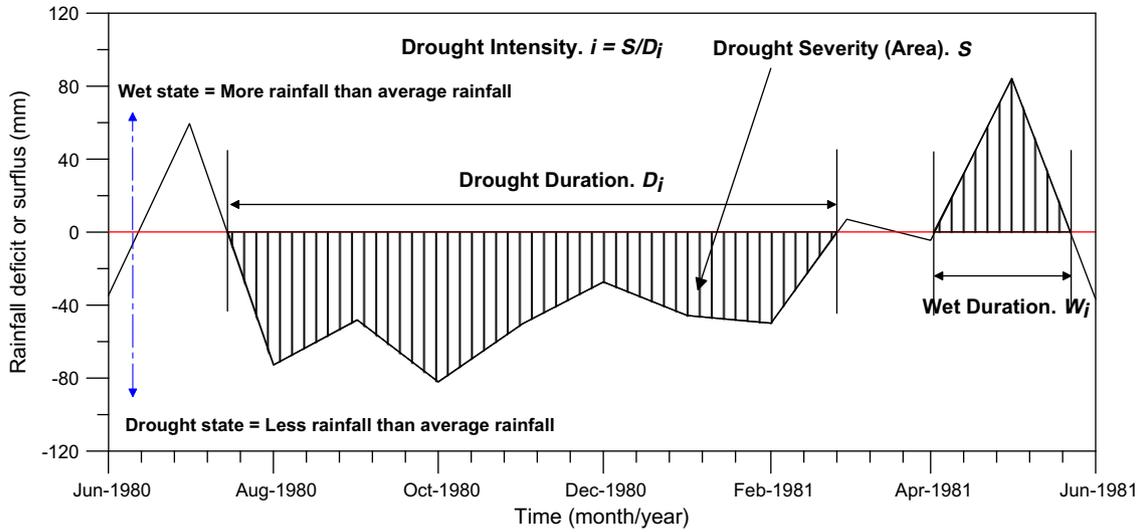


Fig. 1. Definition of drought properties.

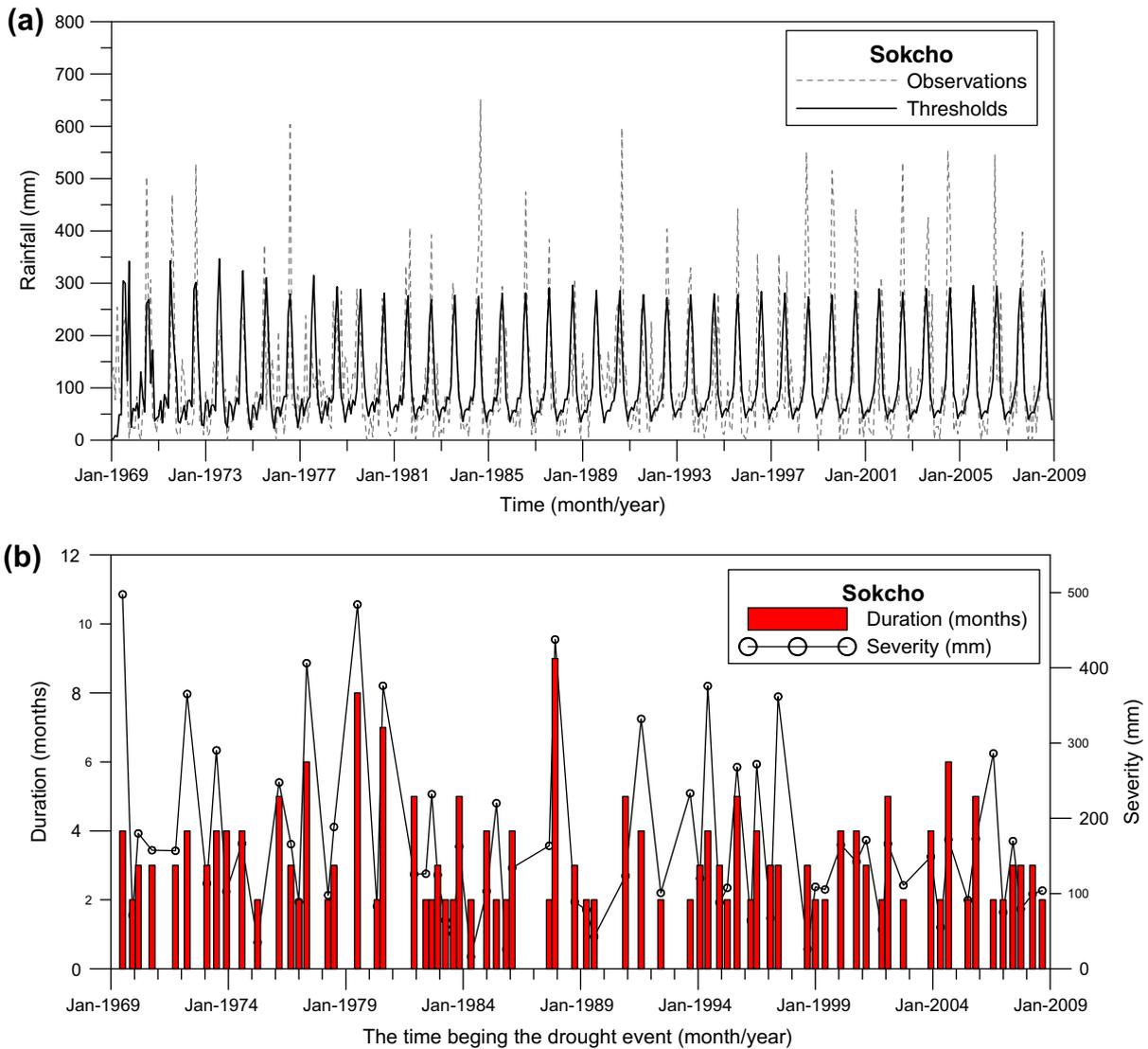


Fig. 2. Sokcho site, (a) time series of rainfall observations and monthly mean rainfall, (b) data series of drought properties.

Accordingly, we focus on regionalization using a clustering algorithm in conjunction with predetermined rainfall properties as a starting point for partitioning weather stations of interests. Subsequently, regional drought frequency analyses are conducted for homogeneous regions. The appropriateness of utilizing a bivariate kernel estimator for drought severity and duration is also assessed as part of this study.

2. Drought properties

The Korean Meteorological Administration manages climatological data at over 70 weather stations throughout the Korean Peninsula. This study collected monthly precipitation data over 56 stations that cover more than 30 years (see Fig. 3). Drought thresholds were determined by the long-term mean of the monthly precipitation, as shown in Eq. (1). In this equation, the long-term mean of “year *n*” was calculated by using up to “year *n* – 1”. This method is identical to the Standardized Precipitation Index (SPI) proposed by Mckee et al. (1993).

$$t_{(n,i)} = \frac{X_{(s,i)} + X_{(s+1,i)} + \dots + X_{(n-1,i)}}{n - s} \tag{1}$$

where, $t_{(n,i)}$ is the threshold of a given *i*th month and *n*th year, and *s* is the first year of collected data.

The duration while successive rainfalls are less than the given threshold is regarded as drought duration, and the accumulated sum of rainfall deficit from the thresholds during drought duration is defined as drought severity. Fig. 1 illustrates a detailed graphical representation of defining drought properties, such as duration and severity. In this study, drought duration and severity were extracted from time series of monthly rainfall for each station. Typical examples are shown in Fig. 2. We calculated various basic descriptive statistics of duration and severity (e.g., mean, maximum value, minimum value, variance, and skewness), as well as the number of drought events, given the thresholds, during historical periods at each site (data not shown). Preliminary data analyses indicated that each site had a different period of recording, and that the number of drought events at most of the sites was too significantly small to reliably estimate the distribution function of drought properties. In addition, correlations between drought duration and severity are significant in more than 40 sites. These findings confirm that regional bivariate analysis is a more efficient method of analysis than point frequency analysis, given the limited data available on record in the area of study.

3. Methodology

Uncertainties in drought point frequency analysis are still problematic because drought extremes are rare by definition, and the data are often too short for reliable estimates of extremes. The main purpose of this study is to establish a way to overcome the limitations of drought point frequency analysis by increasing the data available for the study area. An essential way in augmenting

data is to group together weather stations that have similar drought properties. This approach helps to improve the accuracy of drought analysis through regional frequency analysis to reduce sampling errors in drought properties.

Drought events are defined based on different properties, such as duration and severity. Therefore, on-site univariate drought frequency analysis is unable to fully describe regional drought characteristics. In addition, the existence of significant correlations between drought properties is naturally addressed by regional bivariate frequency analysis.

The method of bivariate drought frequency analysis is divided into parametric methods and nonparametric methods. Parametric frequency analysis fits an assumed theoretical probability distribution to drought properties, while nonparametric frequency analysis does not require any assumptions with regard to probability distribution. For the sake of comparison, this study tried to apply the bivariate Gumbel model to drought properties at 56 stations, however, it is inappropriate for this study because the correlation coefficients between duration and severity were mostly higher than the upper limits for the application of parametric distributions. Alternatively, nonparametric approaches are considered in this study. Previous studies by Harrell and Davis (1982), Lall et al. (1993), and Moon and Lall (1994) concluded that nonparametric methods are more useful to analyze relatively small amounts of data in comparison to parametric methods. Accordingly, this study employed a bivariate kernel function (Kim et al., 2006), as shown in Eq. (2), for regional drought frequency analysis as follows:

$$f(x,y) = \frac{1}{nh_xh_y} \sum_{i=1}^n \left\{ K\left(\frac{x-x_i}{h_x}\right) K\left(\frac{y-y_i}{h_y}\right) \right\} \tag{2}$$

where *n* is the number of data pairs, and *h* is the bandwidth of variables, and *K* indicates the kernel function. Bandwidth selection is the most crucial issue in the nonparametric kernel density method. If a bandwidth is too small, it may lead to large variance and rough estimation. If a bandwidth is too large, there is too much bias, which may bring about the loss of information. This study used the optimal bandwidth selection method, as shown in Eq. (3).

$$h_{d,opt} = \left(\frac{4}{n(p+2)} \right)^{\frac{1}{(p+4)}} \sigma_d \tag{3}$$

Table 2
Heterogeneity measures.

Name	Duration			Severity		
	H(1)	H(2)	H(3)	H(1)	H(2)	H(3)
Cluster 1	-2.5340	-0.9459	-0.1525	-0.1385	-1.5881	-0.5105
Cluster 2	-2.4546	-0.3361	-0.1245	-0.1693	0.3435	-0.0972
Cluster 3	-1.6354	0.4268	-0.0921	-0.1842	-1.2308	-0.7460
Cluster 4	-3.1236	-1.6314	-1.3633	-0.2685	-2.3884	-1.4391
Cluster 5	-2.7503	0.1707	0.2020	-0.2211	-0.8825	-0.7905
Cluster 6	-3.4210	-2.8593	-1.5875	-0.2351	-1.9801	0.2872

Table 1
Basic statistics of regional droughts.

Name	Number of sites	Number of drought events	Duration				Severity				Optimal bandwidth
			Average (months)	Max (months)	Variance	Skewness	Average (mm)	Max (mm)	Variance	Skewness	
Cluster 1	6	583	3.40	12	3.56	1.79	188.5	989.8	26628.1	1.41	35.7
Cluster 2	8	650	3.34	13	3.13	1.90	163.5	720.8	18430.8	1.26	28.9
Cluster 3	6	527	3.38	14	2.87	1.80	127.7	621.7	9428.3	1.30	21.8
Cluster 4	11	738	3.23	11	2.71	1.70	139.1	736.1	12791.1	1.29	23.3
Cluster 5	9	644	3.27	11	2.87	1.82	145.6	581.4	11391.1	1.29	22.8
Cluster 6	13	1045	3.38	13	3.16	1.72	148.3	647.7	15997.8	1.29	23.9

where σ_d refers to the standard deviation in d dimensional distribution, and p is dimension. In this study, $p = 2$ because there are two variables, namely duration and severity.

The regional bivariate distributions can be constructed in this study using the bivariate Gaussian kernel function with the optimal bandwidth for homogeneous regions. The homogeneous regions are predetermined by the cluster analysis based on the attributes of drought across weather stations. Finally, regional bivariate drought frequency curves are derived for each cluster to estimate the severity of drought associated with drought duration.

The detailed procedure for drought frequency analysis using cluster analysis and bivariate probability distribution is illustrated in the next section.

4. Drought frequency analysis

4.1. Cluster analysis

K-means cluster analysis using drought properties was applied to partition regions into mutually exclusive clusters, to which each

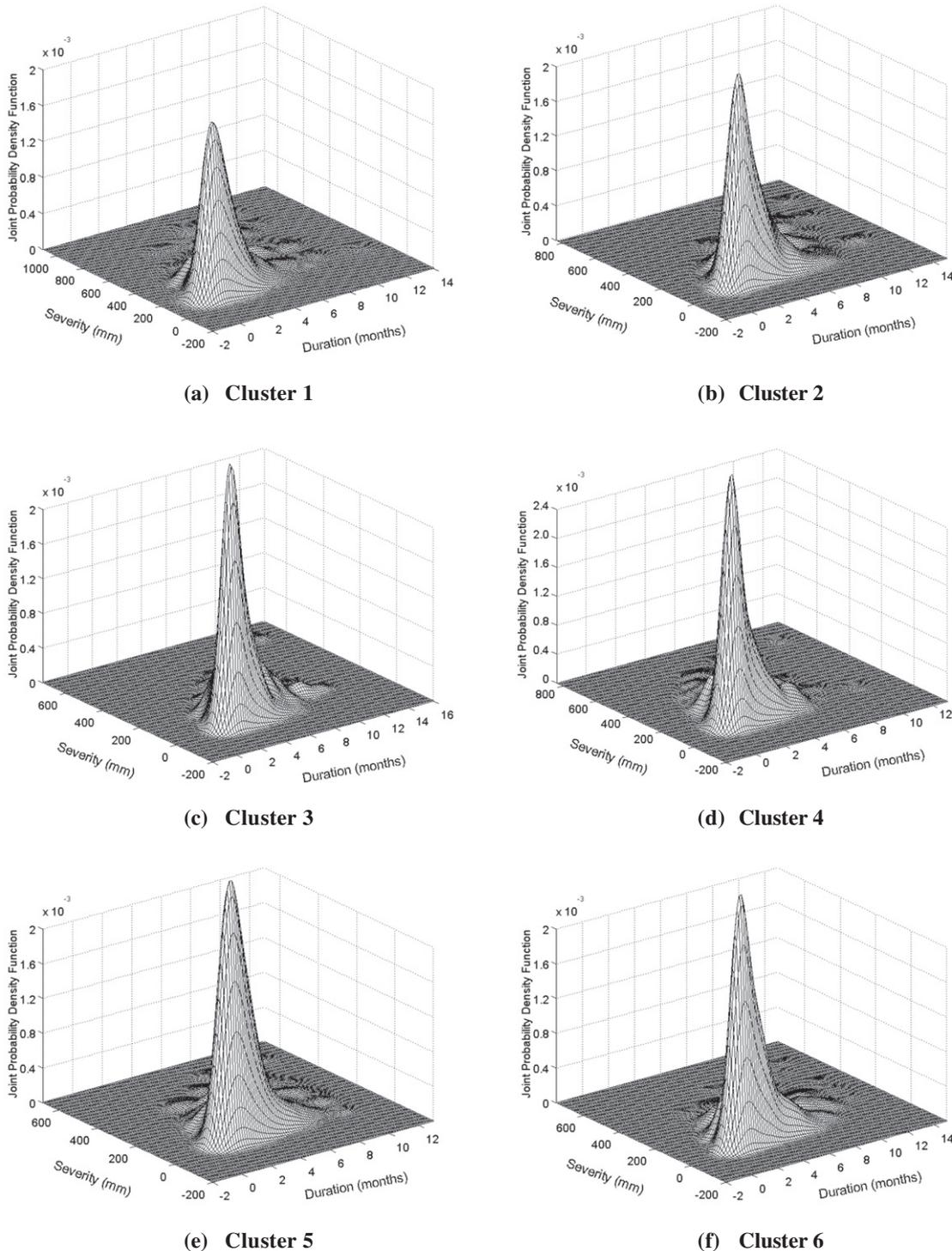


Fig. 4. Joint probability density function estimated by kernel estimator.

of the weather stations within the study area were assigned. The aim of the *K*-means algorithm is to partition points into groups, such that the sum of squares from points to their assigned cluster centers is minimized. The general procedure is to search for a partition with an optimal within-cluster sum of squares by moving points from one cluster to another (Hartigan, 1975). In this study,

six main clusters were distinguished through the regionalization procedure, as illustrated in Fig. 3. The descriptive statistics of drought properties associated with the clusters are presented in Table 1.

According to Hosking and Wallis (1997), heterogeneity measure (*H*) compares the dispersion of sample *L*-moments among sites

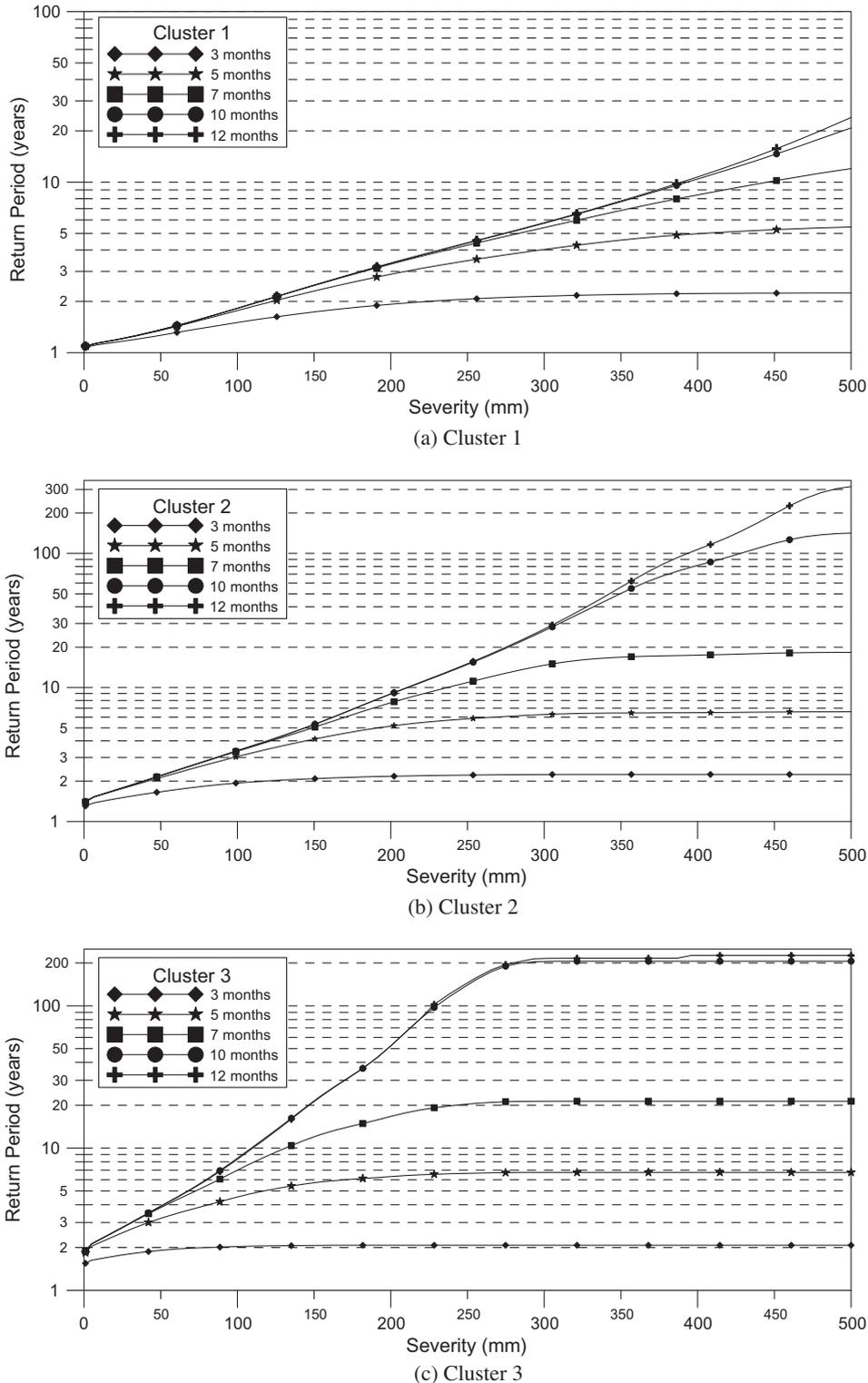


Fig. 5. Bivariate drought frequency curves for each cluster.

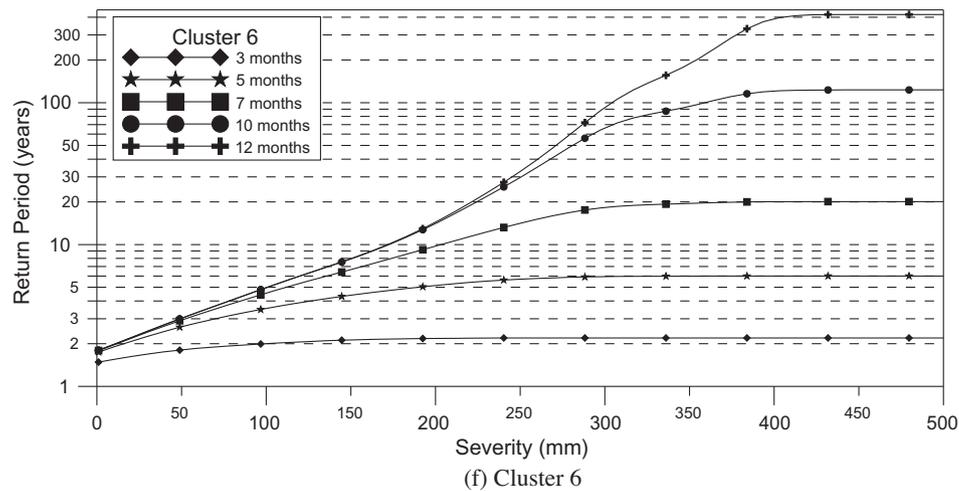
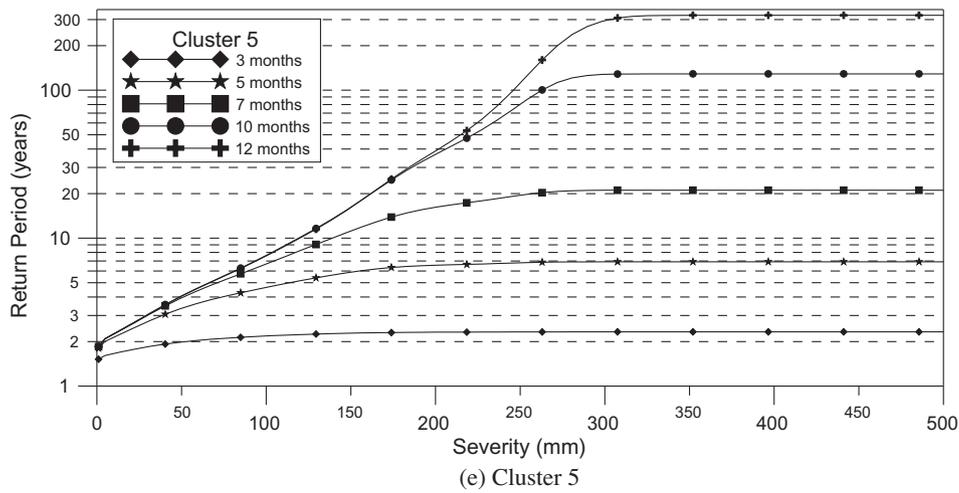
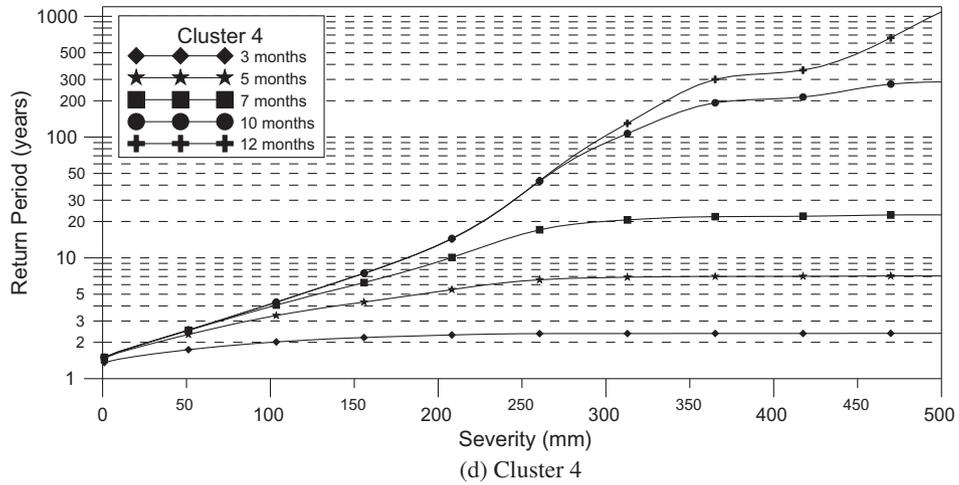


Fig. 5 (continued)

with the L -moment for the group of sites, so that is used to assess whether the regions might reasonably be treated as homogeneous region. A 4-parameter Kappa distribution based Monte Carlo simulation technique is employed to generate 500 homogenous regions with population parameters, which is equal to the regional average sample L -moment ratios. Finally the properties of the tar-

get region are compared to the simulated homogeneous region. The heterogeneity measure H statistic and V statistic for the sample and simulated regions can be defined as follow:

$$H = \frac{V - \mu_v}{\sigma_v} \tag{4}$$

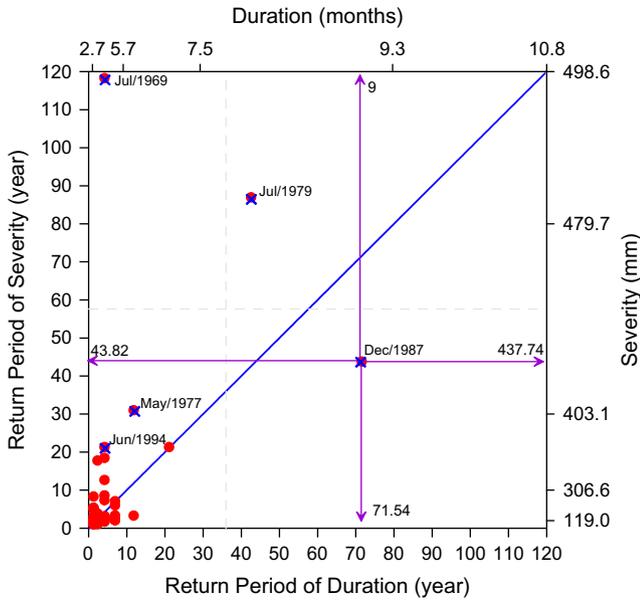


Fig. 6. Return periods corresponding to duration and severity at Sokcho.

$$V = \sqrt{\frac{\sum_{i=1}^N n_i (\tau^i - \tau^R)^2}{\sum_{i=1}^N n_i}} \quad (5)$$

where μ_V is the mean of simulated V values, σ_V is the standard deviation of simulated V values, n_i is record length at site i , τ^i is the sample L -moment at site i , and τ^R is the regional averaged sample L -moment. The L -coefficient of variation, the L -skewness, and the L -kurtosis are used as the L -moment in Eq. (5) for calculating $H(1)$, $H(2)$, and $H(3)$, respectively.

The heterogeneity test is presented in Table 2. A region is regarded as acceptably homogeneous if $H < 1$, possibly heterogeneous if $1 < H < 2$, and definitely heterogeneous if $H > 2$. This study determined that the six clusters were possibly homogeneous regions, as presented in Table 2.

4.2. Nonparametric estimation of bivariate probability distribution

As mentioned in the previous section, the nonparametric kernel density estimator is better for estimating bivariate distributions in this study. This study aims to extend the bivariate kernel density estimator to incorporate regional drought features into frequency analysis. The proposed approach allows us to better understand complex features of drought severity and duration according to the clusters. The estimated joint probability density function of duration and severity is displayed in Fig. 4.

Bivariate drought return periods for each cluster were calculated by estimating cumulative distribution functions, as shown in Eq. (6):

$$T_{DS} = \frac{E(L)}{P(D \geq d, S \geq s)} = \frac{E(L)}{1 - F_D(d) - F_S(s) + F_{D,S}(d, s)} \quad (6)$$

where D and S denote duration and severity, respectively, T_{DS} refers to a bivariate return period for the case of $D \geq d$ and $S \geq s$, and $E(L)$ represents the interval of drought occurrences (Kim et al., 2006). The cumulative distribution functions were derived by integrating the probability density functions in Fig. 4. Fig. 5 illustrates bivariate drought return periods for each cluster. Probability distributions for each weather station were defined by regional kernel density functions, suggesting that underlying distributions within the clusters would be the same. In general, when one of two variables of interest is fixed, the joint cumulative distribution finally converges to the marginal distribution of the fixed variable as the values of the other variable become infinite. That is, $F_y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx$. Therefore, the flat curves in the frequency curves in Fig. 5 refer to the converged joint return periods which should be found in the range of marginal return periods (Kim et al., 2006). For example, for cluster 1 and duration of 3 months (closed diamond marks in Fig. 5a), the frequency curve illustrates that some the drought severity converges the marginal distribution around 400 mm so that the severity beyond 400 mm are insignificant for the duration of 3 months.

The return periods of drought duration were compared with return periods of drought severity by estimating univariate return periods. For example, the results at Sokcho, a site in cluster 5, showed that marginal return periods were inconsistent for either

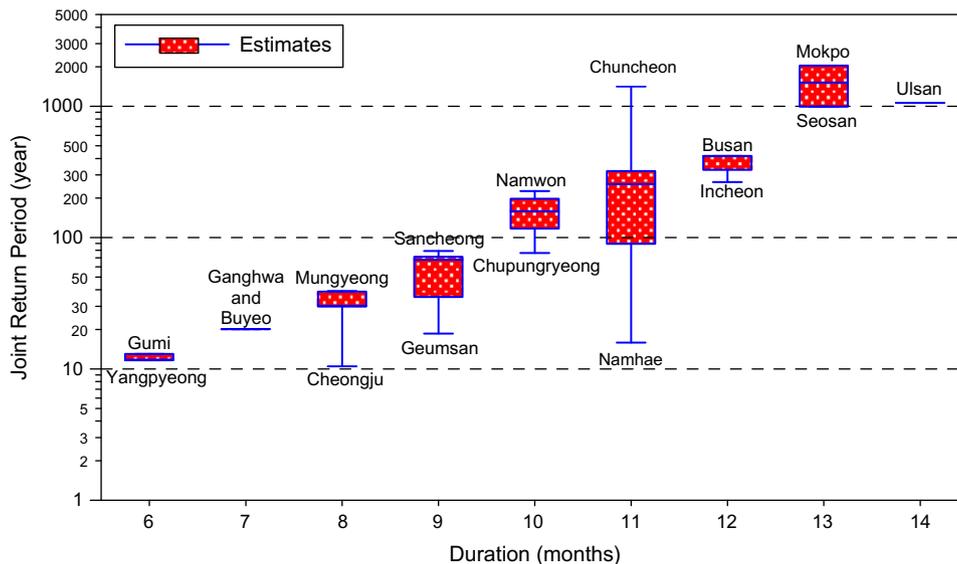


Fig. 7. Bivariate return periods for drought events given various durations.

duration or severity. For example, a drought on record began in December, 1982 and lasted for 9 months with a severity of 438 mm. The return periods, however, were inconsistent, showing a return period of 72 years by duration and a return period of 44 years by severity, as represented in Fig. 6. This inconsistency shows that bivariate drought frequency analysis is necessary for more consistent drought frequency analysis.

5. Results and discussion

Frequency analysis combined with regionalization and bivariate distribution is known as regional bivariate frequency analysis in hydrological applications. The regional bivariate drought frequency analysis in this study allowed the quantification of various properties of extreme drought events for each cluster. Regional fre-

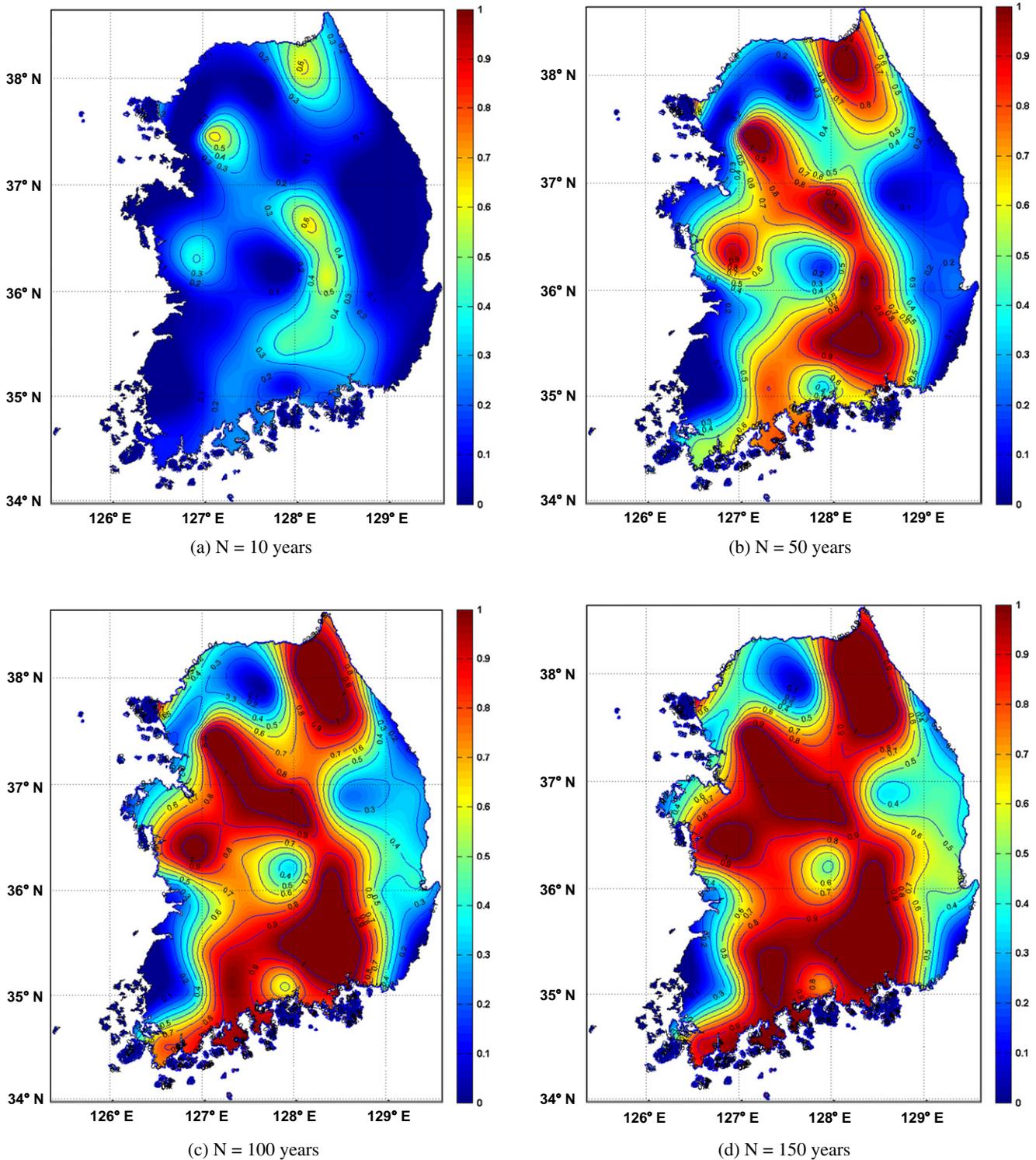


Fig. 8. Drought risk maps.

quency curves for each cluster were derived and used to determine return periods of severity given duration. The weather stations within the clusters were determined to be homogeneous regions, so that a representative probability density function was applied to all of the stations within each cluster. In addition, the bivariate distribution-based drought frequency analysis was able to estimate both univariate and bivariate return periods for the largest drought events to date.

These findings suggest that bivariate return periods of the largest drought events at each weather station are similar to univariate return periods of drought duration. The results also find that, depending on drought severity, bivariate return periods significantly differ among weather stations, as shown in Fig. 7, although the largest drought events of rainfall stations have the same return periods. For example, drought events at various weather stations (e.g., Chuncheon, Gangeung, Seoul, Suwon and Namhae) are on record as having lasted for 11 months, and among these, Namhae is likely to have the most severe drought risk because the lowest non-exceedance probability is assigned there. Regarding the univariate return period, the largest return period for drought duration at Namhae is 167 years, and the return period of drought severity is estimated at 11 years. In the case of the bivariate return period, the estimated return period is approximately 16 years.

Finally, relative differences in drought severity given duration can be effectively assessed, and drought risks can be estimated. The associated drought risks (R) can be calculated by Eq. (7) to impose the probability of extreme drought with T -year return period for N years (Chow et al., 1988).

$$R = 1 - \left(1 - \frac{1}{T}\right)^N \quad (7)$$

Fig. 8 illustrates spatial drought risks with return periods of maximum drought events observed at each site for 10, 50, 100, and 150 years of life span (N). The color red in Fig. 8 indicates a high drought risk compared to the color yellow. This study investigates drought risks via the largest drought events on record for 10 consecutive years. It appears that drought risks are much higher in some parts of the Han River basin, east coastal areas, and the Nakdong River basin.

Drought risks for the largest drought events over 50 and 100 consecutive years are relatively high in areas of 126.5–129° longitude. We see the highest drought risks in two regions in particular (36.5° to 37.5°N latitude, 127° to 128°W longitude and 35.2° to 35.7°N latitude, 128° to 128.7°W longitude). Accordingly, more concrete and comprehensive measures need to be established for proactive responses to future droughts, particularly in areas with higher drought risks.

6. Conclusions

This study developed a new approach for determining regional drought risks based on the cluster analysis and the bivariate distribution. To this end, we first employed a K -means clustering technique for regionalization of drought extremes. This method also demonstrated the application of bivariate kernel density estimation for an index of drought extremes. Regionalization was derived based on drought attributes, and the quality of regionalization was verified by measures of heterogeneity. Our proposed procedure identified six clusters, and showed relevant clustering performance given the heterogeneity criteria. Drought attributes in homogeneous regions were then incorporated into bivariate kernel density estimation.

The approach proposed in this study made it possible to overcome for sampling errors and to provide reliable and consistent

estimates of frequency curves for regions in which an individual site has limited data on record. Moreover, drought risks were further analyzed because our findings suggest that even though the largest drought events for weather stations have the same duration, bivariate return periods considering both duration and severity significantly differ among weather stations. With the proposed approach, spatio-temporal analyses for drought risks were successfully conducted and illustrated to forecast the potential for future drought risks over next some decades. The findings in the study confirm a possible use of bivariate kernel density function in regional frequency analysis to project future drought risks. This analysis can also be used to provide overall information about drought extremes at the national level for long-term risk management.

Acknowledgment

This work was supported by Grants from National Research Foundation of Korea (No. 2010-0016717) and the National Emergency Management Agency of Korea (NEMA-11-NH-40).

References

- Boo, K.-O., Kwon, W.-K., Oh, J.-H., Baek, H.-J., 2004. Response of global warming on regional climate change over Korea: an experiment with the MM5 model. *Geophys. Res. Lett.* 31, L21206.
- Choi, G., Kwon, W.-T., Boo, K.-O., Cha, Y.-M., 2008. Recent spatial and temporal changes in means and extreme events of temperature and precipitation across the Republic of Korea. *J. Korean Geograph. Soc.* 43 (5), 681–700.
- Chow, V.T., Maidment, D.R., Mays, L.W., 1988. *Applied Hydrology*. McGraw-Hill Book Company, New York.
- Clause, B., Pearson, C.P., 1995. Regional frequency analysis of annual maximum streamflow drought. *J. Hydrol.* 173, 111–130.
- Harrell, F.E., Davis, C.E., 1982. A new distribution free quantile estimator. *Biometrika* 69, 635–640.
- Hartigan, J.A., 1975. *Clustering Algorithms*. Wiley, New York.
- Hosking, J.R.M., Wallis, J.R., 1997. *Regional Frequency Analysis: An Approach Based on L-Moments*. Cambridge University Press, New York.
- Kao, S., Govindaraju, R.S., 2007. A bivariate frequency analysis of extreme rainfall with implications for design. *J. Geophys. Res.* 112, D13110. doi:10.1029/2007JD008522.
- Kim, D.H., Yoo, C., Kim, T.-W., 2011. Application of spatial EOF and multivariate time series model for evaluating agricultural drought vulnerability in Korea. *Adv. Water Resour.* 34 (3), 340–350.
- Kim, T.W., Valdés, J.B., Yoo, C., 2003. Nonparametric approach for estimating return periods of droughts in arid regions. *Journal of Hydrologic Engineering ASCE* 8 (5), 237–246.
- Kim, T.-W., Valdés, J.B., Yoo, C., 2006. Nonparametric approach for bivariate drought characterization using Palmer drought index. *Journal Hydrologic Engineering ASCE* 11 (2), 134–143.
- Lall, U., Moon, Y.I., Bosworth, K., 1993. Kernel flood frequency estimator: bandwidth selection and kernel choice. *Water Resour. Res.* 29 (4), 1003–1015.
- McKee, T.B., Doesken, N.J., Kleist, J., 1993. The relationship of drought frequency and duration to time scales. In: *Proceedings of the Eighth Conference on Applied Climatology*, January 17–22, Anaheim, California, pp. 179–184.
- Moon, Y.I., Lall, U., 1994. Kernel quantile function estimator for flood frequency analysis. *Water Resour. Res.* 30 (11), 3095–3103.
- Rossi, G., Benedini, M., Tsakiris, G., Giakoumakis, S., 1992. On regional drought estimation and analysis. *Water Resour. Manage* 6, 249–277.
- Salas, J.D., Fu, C., Cancelliere, A., Dustin, D., Bode, D., Pineda, A., Vincent, E., 2005. Characterizing the severity and risk of drought in the Poudre River, Colorado. *J. Water Resour. Plan. Manage.* ASCE 131 (5), 383–393.
- Yevjevich, V.M., 1967. An objective approach to definitions and investigations of continental hydrologic droughts. *Hydrologic Paper*. 23, Colorado State Univ., Fort Collins, CO.
- Yoo, C., Kim, S., 2004. EOF analysis of surface soil moisture field variability. *Adv. Water Resour.* 27, 831–842.
- Yue, S., 2000. The Gumbel mixed model applied to storm frequency analysis. *Water Resour. Manage.* 14, 377–389.
- Yue, S., Rasmussen, P., 2002. Bivariate frequency analysis: discussion of some useful concepts in hydrological application. *Hydrol. Process* 16, 2881–2898.
- Yue, S., Ouarda, T.B.M.J., Bobée, B., Legendre, P., Bruneau, P., 1999. The Gumbel mixed model for flood frequency analysis. *J. Hydrol.* 226, 88–100.
- Zhang, L., Singh, V.P., 2006. Bivariate flood frequency analysis using the copula method. *J. Hydrol. Eng.* ASCE 11 (2), 50–164.