



1 **An improved logistic probability prediction model for water shortage**
2 **risk in situations with insufficient data**

Longxia Qian¹, Ren Zhang^{1,2*}, Chengzu Bai¹, Yangjun Wang¹ and Hongrui Wang³

¹ Institute of Meteorology and Oceanography, National University of Defense Technology, Nanjing, China,
211101

3 ² Collaborative Innovation Center on Forecast Meteorological Disaster Warning and Assessment, Nanjing
4 University of Information Science & Technology, Nanjing, China, 210044

³ College of Water Sciences, Beijing Normal University, Key Laboratory for Water and Sediment Sciences,
Ministry of Education, Beijing, China, 100875

5 **Abstract.** In drought years, it is important to have an estimate or prediction of the
6 probability that a water shortage risk will occur to enable risk mitigation. This study
7 developed an improved logistic probability prediction model for water shortage risk in
8 situations when there is insufficient data. First, information flow was applied to select
9 water shortage risk factors. Then, the logistic regression model was used to describe
10 the relation between water shortage risk and its factors, and an alternative method of
11 parameter estimation (maximum entropy estimation) was proposed in situations
12 where insufficient data was available. Water shortage risk probabilities in Beijing
13 were predicted under different inflow scenarios by using the model. There were two
14 main findings of the study. (1) The water shortage risk probability was predicted to be
15 very high in 2020, although this was not the case in some high inflow conditions. (2)
16 After using the transferred and reclaimed water, the water shortage risk probability

* Correspondence to: Ren Zhang, Institute of Meteorology and Oceanography, National University of Defense
Technology, Nanjing, China, 211101
E-mail: zrpaper@163.com



17 declined under all inflow conditions (59.1% on average), but the water shortage risk
18 probability was still high in some low inflow conditions.

Keywords Information flow · Risk factors · Logistic regression model · Maximum
entropy estimation · Insufficient data

19

20 **1 Introduction**

21 Nowadays, water shortages have become a serious problem in many parts of the
22 world due to climate change, heightened demand of water and integrated urbanization,
23 and there is a negative impact on the security and sustainable development of water
24 resources (Giacomelli et al., 2008; Weng et al., 2015; Christodoulou 2011; Wang et al.
25 2012; Yang et al. 2015 Qian et al. 2014; Li et al. 2014). Risk is a measure of the
26 probability and severity of adverse effects (Haimes, 2009). It is important to have an
27 estimate or prediction of the probability that a water shortage risk will occur so that
28 effective measures for risk mitigation can be developed, particularly in the case of
29 precipitation deficits (drought).

30 Hashimoto et al. (1982) stated that risk can be described by the probability that a
31 system is in an unsatisfactory state. How to predict or estimate risk probability is still
32 an open issue with no definite solution. Mackenzie (2014) believed that an analyst
33 should first develop a probability distribution over the range of consequences that
34 fully describe the risk of an event. The simulation of probability distribution should be
35 based on a large number of data (Bedford and Cooke, 2001; Giannikopoulou et al.,
36 2015). Unfortunately, a full probabilistic assessment is generally not feasible, because



37 there is insufficient data to quantify the associated probabilities (Tidwell et al., 2005).
38 In some cases, frequency is often used as a substitute for probability in the risk
39 assessment of water resources (Hashimoto et al., 1982; Rajagopalan et al., 2009;
40 Sandoval-Solis et al., 2011), while in other cases, interval-valued probabilities and
41 fuzzy probabilities have been proposed to elaborate the concept of an imprecise
42 probability (Karimi and Hüllermeier, 2007). However, these approaches only consider
43 the probability of the hazard without consideration of the impact of risk factors. The
44 risk factors include characteristics of hazards and existing conditions of vulnerability
45 that could potentially harm exposed people, property, services and so on (UNISDR,
46 2009). There are many aspects of vulnerability arising from various physical, social,
47 economic, and environmental factors (Qian et al., 2016; Haimes, 2006; UNISDR,
48 2009). Therefore, it has been concluded that modeling risk probability requires a
49 consideration of vulnerability (Haimes, 2006). Although increasing attention has been
50 given to vulnerability assessment (Villagrán, 2006; Plummer, 2012), there have been
51 few studies of the relation between risk probability and water resources vulnerability.

52 A water shortage can either occurs or not occur, and therefore water shortage risk
53 is a binary categorical variable. According to statistical theory, a logistic regression
54 model is a nonlinear regression method of studying a binary categorical or
55 multi-categorical variable and its impact factors (Breslow, 1988). Therefore, a logistic
56 regression model can be used to describe the relation between water shortage risk and
57 its impact factors. However, the logistic regression model often requires a large
58 number of observed values of risk (i. e., samples that water shortage risk does or does



59 not occur) and risk factors for parameter estimation. The maximum likelihood
60 estimation is often used for parameter estimation; a large number of observed values
61 of risk and risk factors are required (Balakrishnan, 1992). However, the statistical data
62 about risk and its factors are insufficient in China. Therefore, the method of maximum
63 likelihood estimation is not applicable when the sample size is small. For this reason,
64 we proposed an improved logistic regression model for predicting water shortage risk
65 probability when data is insufficient (i.e. proposing an alternative method of
66 parameter estimation for a logistic regression model when data is insufficient).
67 Moreover, the backward mode is often applied for the selection of sensitive risk
68 factors, but it cannot unravel the cause-effect relation between the water shortage risk
69 and its factors.

70 The contributions of our paper are as follows. First, we used a logistic regression
71 model to predict water shortage risk probability. Then, we introduced an information
72 flow (Liang, 2014) for the selection of sensitive risk factors. Compared with the
73 backward mode, it was very easy to determine whether there was a cause and effect
74 between the water shortage risk and its factors. Finally, we proposed an alternative
75 method of parameter estimation (maximum entropy estimation) for a logistic
76 regression model in situations with a lack of data. The new method requires only a
77 few data, while maximum likelihood estimation requires a large amount of data.

78 The remainder of the paper is organized as follows. Section 2 presents the
79 principles and structure of the logistic probability prediction model for water shortage
80 risk. Section 3 presents the application of the model and the results of the research and



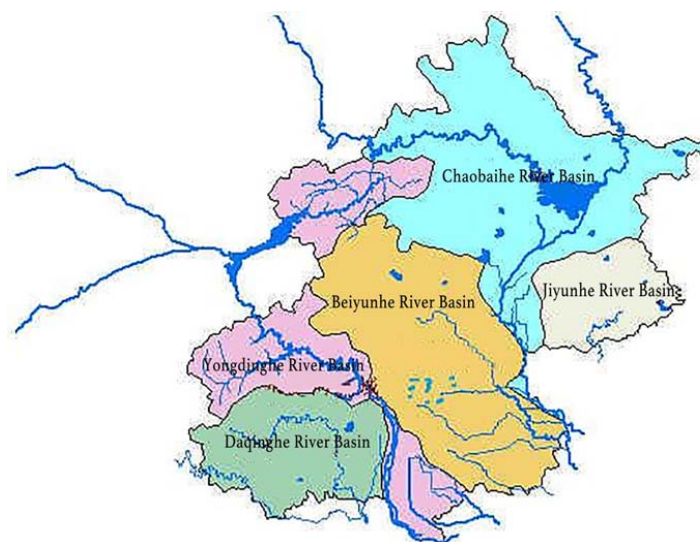
81 Section 4 presents some conclusions and proposes future work.

82 **2 Materials and methods**

83 **2.1 Study area**

84 Beijing, China's capital, is located in the northwest of the North China Plain, and
85 consists of five river systems from the east to the west (Figure 1). The average annual
86 precipitation is 585 mm. Precipitation in summer accounts for 70% of the total for the
87 whole year. Beijing, with a population of more than 20 million, is faced with a severe
88 shortage of water resources. The amount of self-generated water resources is only
89 $37.39 \times 10^8 \text{ m}^3$. The amount of water resources per capita is about 200 m^3 , which is
90 about one eighth of the value of water resources per capita for China and one thirtieth
91 of the global value of water resources per capita.

92 The available surface water and groundwater is unable to meet the needs of the
93 city's economic and social development. Some measures, such as the use of
94 transferred and reclaimed water have been put in place to mitigate the water shortage.
95 In 2014, through the South-to-North Water Diversion Project, water was channeled
96 from the Danjiangkou Reservoir in central China's Hebei province to Beijing.
97 Reclaimed water is also essential for Beijing and is mainly used for agricultural
98 irrigation and toilet flushing.



99

100

Figure 1. Distribution of river system of Beijing

101 2.2 Data collection

102 The data used in this paper were obtained from various sources. The inflow and
103 precipitation sequences from 1956 to 2012 were provided by Beijing Hydrological
104 Station. The water demand for 2020 was based on the Beijing City National
105 Comprehensive Plan for Water Resources (Beijing Municipal Development and
106 Reform Commission and Beijing Municipal Bureau of Water Affairs, 2009). The
107 water supply sequence for 2020 in the inflow conditions of 1956–2012 was computed
108 by an analysis of the balance between water supply and water demand. The
109 population size and gross domestic product (GDP) from 1979 to 2012 were taken
110 from the Statistical Yearbook 2014 of Beijing City (Statistical Bureau of Beijing City,
111 2014). The total amount of water resources from 1979 to 2012 were provided by
112 Beijing Hydrological Station. The water use statistics and data regarding the treatment
113 of domestic sewage from 1979 to 2012 were taken from the Statistical Yearbook 2014

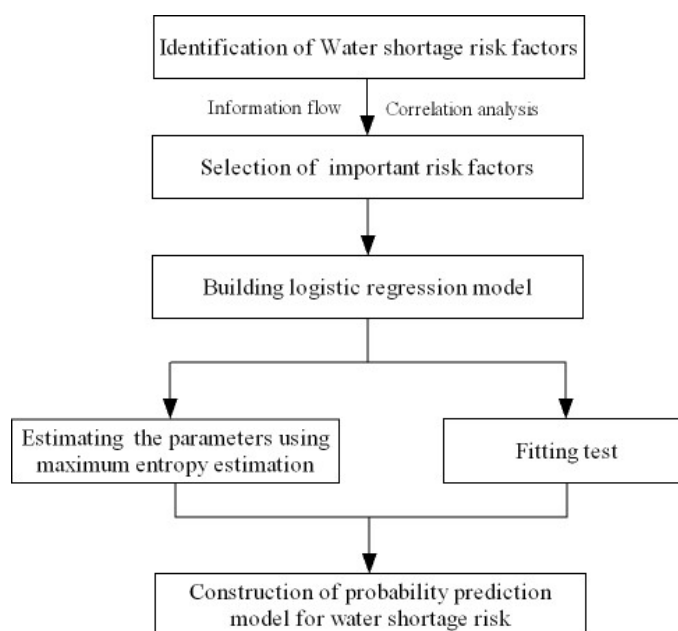


114 of Beijing City (Statistical Bureau of Beijing City, 2014).

115 **2.3 Model development**

116 A flowchart showing the operation of the probability prediction model for water

117 shortage risk is given in Figure 2.



118

119 Figure 2. Flowchart showing the operation of the improved probability prediction model for

120

water shortage risk

121 As can be seen from Figure 2 the model consists of a determination of water

122 shortage risk factors and the construction of a logistic probability prediction model.

123 **2.3.1 Identification of water shortage risk factors**

124 Water shortage risk factors include characteristics of hazards and existing conditions

125 of water resources vulnerability. Water resources vulnerability is referred to as the

126 manifestation of the inherent states (e.g., physical, social, and ecological) of the water

127 resources system that causes the system to be liable to a water shortage (Qian et al.,



128 2016). According to the study of Plummer et al. (2012), there are 50 different water
129 vulnerability assessment tools, and the water vulnerability indicators of these tools are
130 quite different. Therefore, a universal standard understanding of water resource
131 vulnerability indicators is difficult to develop. We established the indicators from
132 perspective of hydrological conditions, water resources, water supply and water use.
133 The risk factors are: precipitation (P), water resources per capita (W_p), water
134 consumption per GDP (W_c), satisfactory rate of water demand (S_r), and utilization
135 rate of water resources (U_r), proportion of industrial water use (IW_p), proportion of
136 agricultural water use (AW_p), proportion of domestic water use (DW_p) and the
137 treatment rate of domestic sewage (DS_r). These indicators are defined as follows
138 (Qian et al., 2014):

139
$$W_p = \frac{W}{N} \quad (1)$$

140 where W is the total amount of water resources, and N is the population size.

141
$$W_c = \frac{\text{the amount of water use}}{GDP} \quad (2)$$

142
$$U_r = \frac{W_{ss} + W_{gs}}{W} = \frac{W_{as}}{W} \quad (3)$$

143 where W_{ss} is the surface water supply, W_{gs} is the groundwater supply, and W is the
144 total amount of water resources.

145
$$DS_r = \frac{DS_t}{DS} \quad (4)$$

146 where DS_t is the amount of sewage treated and DS is the total amount of sewage
147 discharged.

148
$$S_r = \frac{W_{as}}{W_{td}} \quad (5)$$



149 where W_{as} is the water supply, and W_{id} is the water demand.

$$150 \quad IW_p = \frac{IW}{WU} \quad (6)$$

$$151 \quad AW_p = \frac{AW}{WU} \quad (7)$$

$$152 \quad DW_p = \frac{DW}{WU} \quad (8)$$

153 where IW is the industrial water use, AW is the agricultural water use, DW is the
154 domestic water use and WU is total water use.

155 **2.3.2 Selection of important risk factors**

156 The purpose of this section was to select some important factors that have a
157 significant impact on water shortage risk. Liang (2014) reported that the cause and
158 effect between two time series can be measured by the time rate of information
159 flowing from one series to the other. Liang proposed a concise formula for causal
160 analysis. The causality is measured by information flow. Therefore, we can use the
161 information inflow to unravel the cause-effect relation between the risk factors and
162 water shortage risk.

163 According to Liang (2014), for series X_1 and X_2 , the rate of information flowing
164 (units: nats per unit time) from the latter to the former is

$$165 \quad T_{2 \rightarrow 1} = \frac{C_{11}C_{12}C_{2,d1} - C_{12}^2C_{1,d1}}{C_{11}^2C_{22} - C_{11}C_{12}^2} \quad (9)$$

166 where C_{ij} is the sample covariance between X_i and X_j , $C_{i,dj}$ is the covariance
167 between X_i and \mathbb{X}_j , and \mathbb{X}_j is the difference approximation of $\frac{dX_j}{dt}$ using the Euler
168 forward scheme.



169
$$X_{j,n}^{\&} = \frac{X_{j,n+k} - X_{j,n}}{k\Delta t} \quad (10)$$

170 According to Liang (2014), with $k \geq 1$, for a general time series $k = 1$ would be
171 suitable. If $T_{2 \rightarrow 1} = 0$ or the absolute value of $T_{2 \rightarrow 1}$ is less than 0.01, X_2 does not
172 cause X_1 , otherwise it is causal. A positive $T_{2 \rightarrow 1}$ means that X_2 functions to make X_1
173 more uncertain, while a negative value means that X_2 tends to stabilize X_1 . Liang
174 (2015) proposed a method of normalizing the causality between time series and the
175 range of value for $T_{2 \rightarrow 1}$ is 0 and 1.

176 ***2.3.3 Correlation analysis of selected risk factors***

177 In theory, a probability prediction model requires variables to be mutually
178 independent. Therefore, it is necessary to perform a correlation analysis. Because all
179 of the factors are continuous variables, Pearson correlation coefficients are often
180 applied. If the absolute correlation coefficient is greater than 0.5, there is a significant
181 correlation between two factors.

182 ***2.4 Risk probability prediction model using maximum entropy*** 183 ***estimation***

184 A logistic regression model is a nonlinear regression method of studying a binary
185 categorical or multi-categorical variable and its impact factors. Because a water
186 shortage either occurs or does not occur, water shortage risk belongs to a binary
187 categorical variable. Therefore, we can use a logistic regression model to simulate the
188 relation between water shortage risk and its factors. Suppose the risk factors
189 are $\{x_{ij} (i = 1, 2, L, n; j = 1, 2, L, m)\}$, where x_{ij} denotes the value of the j th factor in



190 the i th year. The risk sequence is $\{y_i (i = 1, 2, L, n)\}$,

191 where $y_i = \begin{cases} 0, & \text{water shortage risk does not occur} \\ 1, & \text{water shortage risk occurs} \end{cases}$, and is the observed value of the i th

192 year.

193 $p_i = p(y_i = 1 | x_{ij} (j = 1, 2, L, m))$ is the conditional probability when $y_i = 1$ under

194 the conditions of $x_{ij} (i = 1, 2, L, n; j = 1, 2, L, m)$. The logistic regression model is

$$195 \quad p_i = \frac{1}{1 + e^{-(\alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_m x_{im})}} \quad (11)$$

196 where $\alpha, \beta_1, \beta_2, L, \beta_m$ are the estimated parameters. The parameters are often

197 determined by a maximum likelihood estimation. The log likelihood equation of

198 computing $\alpha, \beta_1, \beta_2, L, \beta_m$ is as follows:

$$199 \quad \begin{cases} \frac{\partial L}{\partial \alpha} = \sum_{i=1}^n \left[y_i - \frac{\exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)}{1 + \exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)} \right] = 0 \\ \frac{\partial L}{\partial \beta_j} = \sum_{i=1}^n \left[y_i - \frac{\exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)}{1 + \exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)} \right] x_{ij} = 0 \quad j = 1, 2, L, m \end{cases} \quad (12)$$

200 According to Eq. (12), a large number of observed values of risk

201 ($y_i (i = 1, 2, L, n)$) and its factors are required for parameter estimation. Unfortunately,

202 the correlated samples between risk and its controlling factors are insufficient. It is

203 therefore far better to estimate the parameters. In this case, the maximum likelihood

204 estimation is not applicable for parameter estimation. An alternative approach for

205 parameter estimation is therefore required.

206 Thus, we proposed a new parameter estimation method based on the maximum



207 entropy principle. The new method is named after maximum entropy estimation. The
 208 new method does not require the observed values of risk, and it requires only some
 209 observed values of the factors. Its principle is as follows.

210 For an observation, we can define its entropy to evaluate its degree of uncertainty.
 211 According to Jones and Jones (2000), the entropy of the i th observation of water
 212 shortage risk is

$$\begin{aligned}
 H(p_i) &= -C [P_i \ln P_i + (1 - P_i) \ln (1 - P_i)] \\
 &= -C \left[P_i \ln \left(\frac{P_i}{1 - P_i} \right) + \ln (1 - P_i) \right] \quad (13) \\
 &= -C \left\{ \frac{\left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right]} - \ln \left(1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right) \right\}
 \end{aligned}$$

214 where C is a positive value and $p_i = p(y_i = 1 | x_{ij} (j = 1, 2, L, m))$ is the
 215 conditional probability when $y_i = 1$ under the conditions of
 216 $x_{ij} (i = 1, 2, L, n; j = 1, 2, L, m)$. According to the maximum entropy principle, if the
 217 values of $H(P_i)$ reaches a maximum, the optimal parameters are obtained (Jones
 218 and Jones, 2000). The reasons for obtaining a solution based on the maximum entropy
 219 principle are as follows. ① It conforms to the principle of entropy increase, which
 220 states that the entropy of an isolated system tends to reach a maximum. ② It accords
 221 with the principle that the solution should be in line with the sample/data and the least
 222 hypotheses must be constructed regarding the unknown parts when the data is
 223 insufficient. ③ It fits the maximum multiplicity principle. The multiplicity of a state
 224 refers to the number of possible ways in which a system can evolve to that state. The



225 maximum multiplicity principle states that the greater the multiplicity of a state, the
 226 larger the possibility that a system is in this state.

227 **2.4.1 Parameter estimation**

228 Based on the analysis above, an optimization model can be constructed as follows:

$$229 \quad \max H_i = -C \left\{ \frac{\left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right]} - \ln \left(1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right) \right\} \quad (14)$$

230 According to the extreme theory of multivariate function (Khuri 2003), we can
 231 obtain

$$232 \quad \begin{cases} \frac{\partial H_i}{\partial \alpha} = \frac{\left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\} + \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \cdot \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \cdot \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{\left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\}^2 \cdot 1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)} = 0 \\ \frac{\partial H_i}{\partial \beta_j} = \frac{x_{ij} \cdot \left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\} + x_{ij} \cdot \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \cdot \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \cdot x_{ij} \cdot \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{\left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\}^2 \cdot 1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)} = 0 \end{cases} \quad (15)$$

233 The optimal estimation $\alpha, \beta_j (j=1,2,L, m)$ can be obtained by solving Eq.
 234 (15). Numerical approaches are often used to obtain an approximate solution of Eq.
 235 (15) rather than its exact solution. Therefore, we made use of the optimization
 236 function of Matlab to estimate the parameters, i.e., the `fminsearch` function. If there
 237 are n observations, there are $n H_i (i=1,2,L, n)$. It is impossible to find the parameters
 238 that make all the $H_i (i=1,2,L, n)$ reach the maximum value. According to the
 239 maximum entropy principle, the greater the entropy is, the larger the uncertainty of an



240 observation is. Therefore, the maximum value of the sequences $\{H_i, i = 1, 2, L, n\}$ was
241 taken as the objective function of the optimization model.

242 **2.4.2 Goodness-of-fit test**

243 According to Brown (1982), a goodness-of-fit test should be made for evaluating the
244 fitting effect of the logistic regression model and its ability to identify water shortage
245 risk. In this study, the Kolmogorov-Smirnov Test (K-S) test and Pearson χ^2 test are
246 used.

247 **2.4.2.1 K-S test (t)**

248 A K-S test is often applied as a fitting test. It can be used to test the ability of the
249 model to identify water shortage risk. The value of K-S is between 0 and 1; the
250 greater the value is, the better the logistic model is. The idea is as follows.

251 Let $F_{n1}(x)$ be the cumulative probability distribution of the samples that do not
252 encounter a water shortage. $F_{n2}(x)$ is the cumulative probability distribution of the
253 samples that encounter a water shortage. A two independent samples test is then
254 applied to compare whether the empirical distribution functions of two samples are
255 the same. The test is as follows:

$$256 \quad H_0 : F_{n1}(x) = F_{n2}(x) \quad H_1 : F_{n1}(x) \neq F_{n2}(x) \quad (16)$$

257 The value of K-S is:

$$258 \quad K - S = \max |F_{n1}(x) - F_{n2}(x)| \quad (17)$$

259 When $N \rightarrow \infty$, the cumulative distribution curve and probability density curve of
260 two samples can be obtained. The value of K-S is the maximum value of the
261 cumulative distribution functions. When the value of K-S is greater than 0.35, the



262 logistic regression model is applicable. The international classification standard of the
 263 logistic model is shown in Table 1 (Brown, 1982).

264 Table 1. The international classification standard of the logistic model

K-S	The effect of the model
<0.2	Bad
0.2~0.4	General
0.4~0.5	Good
0.5~0.6	Better
0.6~0.75	Very good
0.75~1	Perfect

265

266 **2.4.2.2 Pearson χ^2 test**

267 The test is as follows:

268 H_0 : the fitting is good H_1 : the fitting is bad (18)

269 The expression of the χ^2 statistic is as follows.

270
$$\chi^2 = \sum_{j=1}^l \frac{(O_j - E_j)^2}{E_j} \quad (19)$$

271 where $j=1,2,L,l$, l is the number of covariant types, O_j is the observed
 272 frequency of the j th covariant type, and E_j is the predicted frequency of
 273 the j th covariant type. The degree of freedom is the difference between the number of
 274 covariant types and parameters.

275 **3 Results and discussion**



276 In this section, a logistic probability prediction model for water shortage risk is
 277 constructed and discussed, and the risk probability in 2020 in Beijing is predicted
 278 using the proposed model.

279 **3.1 Construction of the Logistic probability prediction model**

280 A sequence of risk factors were obtained for the period from 1979 to 2012, and were
 281 computed based on Eqs. (1)~(8). The risk sequence $\{y_i (i=1,2,L,34)\}$ from 1979
 282 to 2012 was obtained as follows. According to Qian and Zhang et al. (2016), a water
 283 supply is deemed inadequate if the supply is less than the demand, leading to a water
 284 shortage in the water supply system. $y_i = \begin{cases} 0, & \text{water shortage does not occur} \\ 1, & \text{water shortage occurs} \end{cases}$.

285 Therefore, there are only 34-year data.

286 **3.1.1 Determination of water resources vulnerability indicators**

287 Based on the risk factors sequences from 1979 into 2012 (Table 2) and the method of
 288 normalized information inflow (Liang, 2015), the values of normalized information
 289 flow from the factors to risk are shown in Table 3. According to the normalized
 290 information flow results (Table 3), the value of the normalized information flow
 291 from AW_p to water shortage risk is only 0.0031, and it is very little. It was concluded
 292 that the AW_p does not result in a water shortage risk. Therefore, AW_p was removed as
 293 risk factors.

294 Table 2. The values of the risk factors and risk from 1979 to 2012

Year	W_c (m ³ per CNY)	W_p (m ³ per capita)	U_r	P (mm)	DS_r (%)	AW_p	DW_p	IW_p	S_r	Risk
1979	0.36	426.15	1.12	652.00	10.20	0.56	0.10	0.33	0.71	0
1980	0.36	287.52	1.94	387.30	9.40	0.63	0.10	0.27	0.41	1



1981	0.35	261.10	2.00	433.50	10.80	0.66	0.09	0.25	0.40	1
1982	0.30	391.44	1.29	585.10	10.90	0.61	0.10	0.29	0.62	1
1983	0.26	365.26	1.37	465.50	10.20	0.66	0.10	0.24	0.58	1
1984	0.18	407.36	1.02	442.10	10.00	0.55	0.10	0.36	0.79	0
1985	0.12	387.36	0.83	611.20	10.00	0.32	0.14	0.54	0.96	0
1986	0.13	262.94	1.35	560.30	8.90	0.53	0.20	0.27	0.59	1
1987	0.09	369.25	0.80	662.60	7.70	0.31	0.23	0.45	1.00	0
1988	0.10	369.27	1.08	594.70	7.40	0.52	0.15	0.33	0.74	0
1989	0.10	200.47	2.07	479.50	6.60	0.55	0.14	0.31	0.39	1
1990	0.08	330.20	1.15	662.40	7.30	0.53	0.17	0.30	0.70	0
1991	0.07	386.56	0.99	662.70	6.60	0.54	0.18	0.28	0.80	0
1992	0.07	203.63	2.07	500.00	1.20	0.43	0.24	0.33	0.39	1
1993	0.05	176.89	2.30	424.30	3.10	0.45	0.21	0.34	0.35	1
1994	0.04	403.73	1.01	727.70	9.60	0.46	0.23	0.32	0.79	0
1995	0.03	242.51	1.48	608.90	19.40	0.43	0.26	0.31	0.54	1
1996	0.02	364.22	0.87	669.40	21.20	0.47	0.23	0.29	0.92	0
1997	0.02	179.44	1.81	419.00	22.00	0.45	0.28	0.28	0.44	1
1998	0.02	302.67	1.07	687.40	22.50	0.43	0.30	0.27	0.75	0
1999	0.02	113.11	2.93	384.70	25.00	0.44	0.30	0.25	0.27	1
2000	0.01	123.64	2.40	446.60	39.40	0.41	0.33	0.26	0.33	1
2001	0.01	138.62	2.03	462.00	42.00	0.45	0.32	0.24	0.39	1
2002	0.01	113.13	2.15	413.00	45.00	0.45	0.34	0.22	0.37	1
2003	0.01	126.34	1.84	453.00	50.10	0.39	0.38	0.23	0.41	1
2004	0.01	143.36	1.52	539.00	53.90	0.39	0.39	0.22	0.50	1
2005	0.00	150.85	1.27	468.00	62.40	0.38	0.42	0.20	0.54	1
2006	0.00	154.97	1.14	448.00	73.80	0.37	0.45	0.18	0.57	1
2007	0.00	145.74	1.13	499.00	76.20	0.36	0.48	0.17	0.55	1
2008	0.00	201.77	0.74	638.00	78.90	0.34	0.51	0.15	0.78	0
2009	0.00	124.22	1.08	448.00	80.29	0.34	0.52	0.15	0.49	1
2010	0.00	117.64	0.99	524.00	81.00	0.32	0.42	0.14	0.52	1
2011	0.00	132.81	0.88	552.00	81.70	0.30	0.43	0.14	0.60	1
2012	0.00	190.89	0.58	708.00	83.00	0.26	0.45	0.14	0.88	0

295

296 According to Liang (2014), a positive value of the information flow means that
 297 the factor makes water shortage risk more uncertain, while a negative value means
 298 that the indicator tends to stabilize water shortage risk. Therefore, all the factors tend
 299 to make water shortage risk more uncertain. Furthermore, the impact of P , W_p ,



300 W_c are very significant.

301 Table 3. The values of information flow from the factors to water shortage risk

Factors	Information flow
W_c	0.3560
W_p	0.4823
U_r	0.3109
P	0.1575
DS_r	0.2413
IW_p	0.1320
AW_p	0.0031
S_r	0.1247
DW_p	0.1164

302

303 A correlation analysis was performed on the remaining factors. The values of the

304 Pearson correlation coefficients are shown in Table 4.

305 Table 4. Pearson correlation coefficients for the relations between various factors

Pearson correlation coefficients	W_c	W_p	U_r	P	DS_r	DW_p	IW_p	S_r
W_c	1	0.603	0.047	-0.066	-0.559	0.354	-0.780	0.047
W_p	0.603	1	-0.455	0.571	-0.682	0.654	-0.753	0.696
U_r	0.047	-0.455	1	-0.723	-0.268	0.026	-0.157	-0.869
P	0.066	0.571	-0.723	1	-0.100	0.219	-0.064	-0.820



DS_r	-0.559	-0.682	-0.268	-0.100	1	-0.802	0.902	-0.087
DW_p	0.354	0.654	0.026	0.219	-0.802	1	-0.715	0.354
IW_p	-0.780	-0.753	-0.157	-0.064	0.920	-0.715	1	-0.013
S_r	0.047	0.696	-0.869	0.820	-0.087	0.354	-0.013	1

306 Based on the results in Tables 3 and 4, AW_p , S_r , IW_p , and DW_p were
 307 removed as risk factors. Therefore, the selected factors for logistic regression model
 308 were W_c , W_p , U_r , P and DS_r .

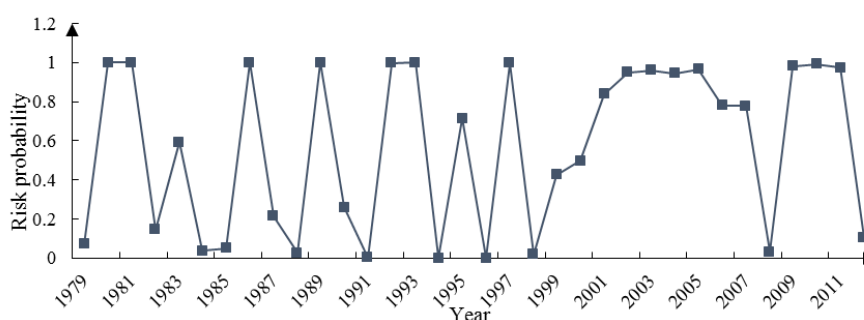
309 **3.1.2 Construction of the logistic risk probability predication model**

310 The data for the risk and selected factors (W_c , W_p , U_r , P and DS_r) from 1979 to
 311 2012 (Table 2) are used to construct the logistic risk predication probability model.
 312 Because there is only 34 samples, it is impossible to estimate the parameters by the
 313 maximum likelihood estimation. Substituting the sequences of W_c , W_p , U_r , P and DS_r
 314 from 1979 to 2012 (Table 2) into Eq. (14), the values of parameters obtained by
 315 maximum entropy estimation can be obtained. The estimated values for
 316 α , β_1 , β_2 , β_3 , β_4 , β_5 are 61.6386, 0.004, -0.1262, -12.4077, -0.012 and -29.0963.

317 Therefore, the logistic regression model based on the maximum entropy
 318 estimation is as follows:

$$319 \quad \text{Predicted probability} = \frac{1}{1 + e^{-(61.6386 + 0.004W_c - 0.1262W_p - 12.4077U_r - 0.012P - 29.0963DS_r)}} \quad (20)$$

320 Substituting the sequences of W_c , W_p , U_r , P and DS_r from 1979 to 2012 into Eq.
 321 (20), the predicted probability values of water shortage risk by the maximum entropy
 322 estimation is shown in Fig. 3.



323

324 Figure 3. The predicted probability generated by the maximum entropy estimation from 1979

325

to 2012

326

If 0.5 is taken as threshold used to judge whether water shortage risk occurs, then

327

the prediction accuracy by using the maximum entropy estimation can be obtained,

328

and is shown in Tables 5. From Table 5, it can be seen that the average accuracy rate

329

using the maximum entropy estimation was very high (91.18%). The maximum

330

entropy estimation does not need observed values of risk ($y_i (i = 1, 2, L, n)$), whereas

331

the maximum likelihood estimation needs a large number of observed values of risk.

332

Table 5. The prediction accuracy using the maximum entropy estimation

	The prediction is that risk occurs	The prediction is that no risk occurs	Accuracy rate
Risk actually occurs	19	3	86.36%
Risk actually does not occur	0	12	100%
The average accuracy rate			91.18%



333 The K-S test and Pearson χ^2 test are performed and the results of the tests are
334 obtained. The value of K-S is 0.955 and according to Table 1, the logistic probability
335 prediction model was applicable. Moreover, the probability value was 0.000 (i.e., less
336 than 0.05), so the null hypothesis was rejected. Therefore, the ability of the logistic
337 regression model to predict water shortage is very strong.

338 Substituting the observed frequency and the predicted frequency into Eq. (19),
339 the value of the χ^2 statistics was 2.333 (the number of covariant type was 8). Because
340 the number of parameters was 6, there were 2 degrees of freedom. The $\chi_{0.1}^2(2)$ was
341 equal to 4.605 and was much greater than 2.333. Therefore, the null hypothesis was
342 accepted, i.e., the fitting of the model was very good. Based on the results of the K-S
343 test and Pearson χ^2 test, it was concluded that the model was applicable.

344 ***3.2 Risk probability prediction in 2020 in Beijing***

345 ***3.2.1 Risk probability prediction (without considering the use of*** 346 ***transferred and reclaimed water)***

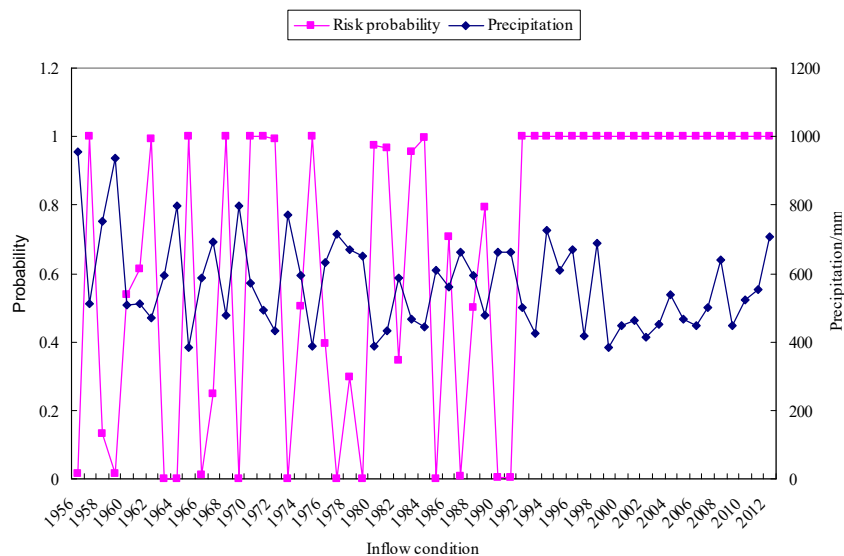
347 Because the inflow of 2020 is unknown, the inflow condition in 2020 was assumed to
348 be any annual inflow conditions from 1956 to 2012. In this section we predict the risk
349 probability of 2020 under different inflow conditions from 1956 to 2012. The
350 sequences for risk factors (W_c, W_p, U_r, P and DS_r) were obtained and computed as
351 follows. The precipitation in 2020 is assumed to be any annual precipitation from
352 1956 to 2012. First, an analysis of the balance between water supply and demand was
353 performed and the sequences of water supply and demand under the inflow scenarios
354 of 1956–2012 were obtained (Qian et al., 2016). The GDP of 2020 was the sum of the



355 gross agricultural product, gross industrial product, and gross product of the third
356 industry (details of the third industry are shown in Appendix A), using information
357 taken from the literature, and was estimated to be 4711.852 billion CNY (Qian et al.,
358 2016). N (the population size of 2020) was 24.43 million (Qian et al. 2016). The
359 total amount of water resources from 1956 to 2020 were considered to consist of
360 fifty-seven types of water resources in 2020. Substituting the total water resources
361 sequences and N of 2020 into Eq. (1), the sequence of W_p could be computed.
362 Substituting the water demand sequences and GDP of 2020 into Eq. (2), the sequence
363 of W_c could be computed. Substituting the sequence of the total water resources and
364 water supply for 2020 into Eq. (3), the sequence of U_r could be obtained. The DS_r of
365 2020 was about 90% (Beijing Municipal Development and Reform Commission and
366 Beijing Municipal Bureau of Water Affairs, 2009).

367 Substituting the sequences of W_c , W_p , U_r , P and DS_r into Eq. (20), the probability
368 that a water shortage risk will occur in 2020 under the inflow scenarios of 1956–2012
369 was predicted, and is shown in Figure 4.

370 In Figure 4, the horizontal axis represents the inflow conditions of 1956–2012.
371 Figure 4 shows that in 2020, the water shortage risk probability exceeded 0.95 under
372 33 different inflow conditions (accounting for 63.5% of all the inflow conditions) and
373 exceeded 0.5 under 38 different inflow conditions (accounting for 73.1% of all the
374 inflow conditions). In summary, there was a high probability of a water shortage risk
375 in 2020, although the probability was very low in some high precipitation periods.



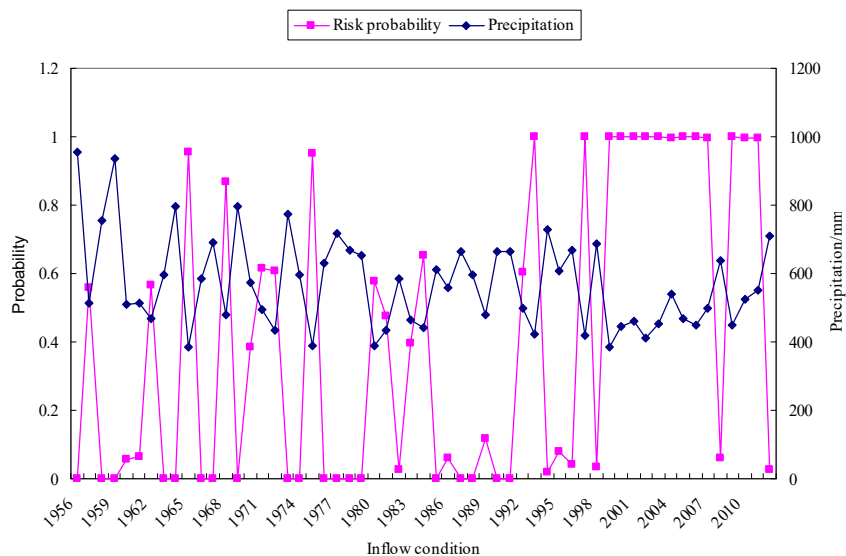
376

377

Figure 4. Risk probability under the inflow conditions of 1956–2012

378 **3.2.2 Risk probability prediction after using transferred and reclaimed** 379 **water**

380 According to Qian et al. (2016), 1.05 billion m³ of water will have been transferred
381 to Beijing in 2020 and the amount of reclaimed water used may reach 1 billion m³.
382 After using transferred and reclaimed water, the total amount of water resources
383 would increase, W_p and U_r would change and other indicators would remain
384 unchanged. Therefore, the sequences of W_p and U_r under the inflow scenarios of
385 1956–2012 had to be computed again. Substituting the sequences of
386 W_c, W_p, U_r, P and DS_r into Eq. (20), the water shortage risk probability in 2020
387 under the inflow scenarios of 1956–2012 (after using transferred and reclaimed
388 water) was predicted, and the results are shown in Figure 5.

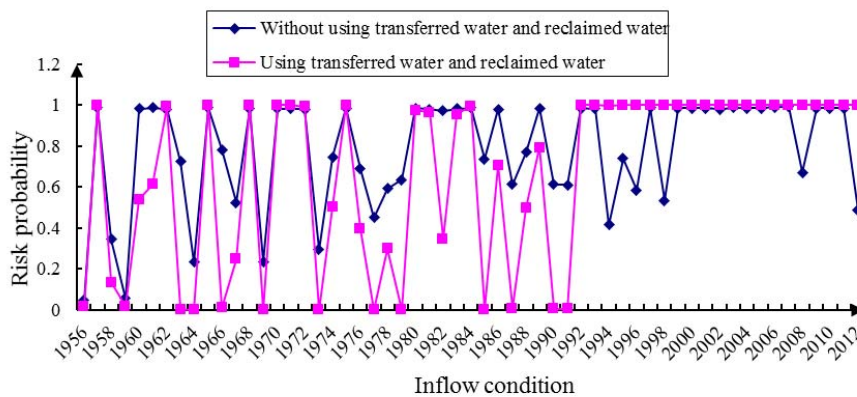


389

390 Figure 5. Values of risk probability under the inflow conditions of 1956-2012 after using

391 transferred and reclaimed water

392



393

394 Figure 6. Comparison of risk probability before and after using transferred and reclaimed

395 water

396 From Figures 5 and 6, it was concluded that the water shortage risk probability

397 would decline under all inflow conditions (59.1% on average). However, the water



398 shortage risk probability would still be high in some low inflow conditions. The risk
399 probability exceeded 0.5 under 24 different inflow conditions (accounting for 46.2%
400 of all inflow conditions). For example, the water shortage risk probability reached 1
401 under the inflow conditions of 1999–2008.

402 According to Qian et al. (2016), since 1999, Beijing has experienced drought in
403 ten consecutive years. This has had a strong effect on the water resources of Beijing,
404 including a significant reduction in surface water and severe over-exploitation of
405 groundwater. This means that a water shortage may occur in 2020 under the inflow
406 conditions of 1999–2008 although some measures have been taken. Moreover, water
407 resources vulnerability was still high in 2020 after using transferred and reclaimed
408 water (Qian et al., 2016). Therefore, we concluded that the water shortage risk
409 probability would still be high in 2020 after using transferred and reclaimed water,
410 especially in the case of precipitation deficits.

411 **4 Conclusions**

412 This study developed an improved logistic probability prediction model for water
413 shortage risk in situations when there is insufficient data. The model consists of the
414 following steps:

415 (1) Information flow was used to select some important factors that were likely
416 to have a significant impact on water shortage risk. This could determine the
417 cause-effect relation between the water shortage risk and its factors.

418 (2) The logistic regression model was applied to describe the nonlinear relation
419 between water shortage risk and its factors. A new parameter estimation method based



420 on the entropy principle, i.e. maximum entropy estimation, was proposed for
421 parameter estimation when insufficient data is available.

422 The results of the study were as follows. In 2020, the probability that a water
423 shortage risk will occur exceeded 0.95 under 33 different inflow conditions
424 (accounting for 63.5% of all inflow conditions) and exceeded 0.5 under 38 different
425 inflow conditions (accounting for 73.1% of all inflow conditions). After using the
426 transferred and reclaimed water, the water shortage risk probability declined under all
427 inflow conditions (by 59.1% on average), but the water shortage risk probability was
428 still high for some low inflow conditions. Risk probability exceeded 0.5 under 24
429 different inflow conditions (accounting for 46.2% of all inflow conditions).

430 However, some problems still exist with regard to the maximum entropy
431 estimation. Initial values of the parameters should be given for the optimization
432 function, but the optimization function belongs to local optimization, which was very
433 sensitive to the initial values. Therefore, we may obtain an unsatisfactory result if the
434 initial values are not correct. How best to search for a global optimum is an important
435 and difficult issue, and will be the focus of our further study.

436

437 **Appendix A. Glossary used in this paper**

438 1. **Logistic regression model.** It is nonlinear regression method of studying binary
439 categorical or multi-categorical variable and its impact factors.

440 2. **Maximum likelihood estimation.** It is a method of parameter estimation in
441 statistics.



442 3. **Maximum entropy estimation.** We propose a new parameter estimation method
443 for a logistic regression model when insufficient data is available. We called this new
444 method maximum entropy estimation.

445 4. **Backward.** It is a method of selecting the variables for a logistic regression model.
446 The methods of selecting the variables for a logistic regression model include enter,
447 forward and backward.

448 5. **Information flow.** Information flow, proposed and named by Liang (2014), is a
449 method for unraveling the cause-effect relation between time series.

450 6. **The extreme theory of multivariate function.** This is a theory used for
451 calculating extreme values in advanced mathematics.

452 7. **Two independent samples test.** This is one type of Kolmogorov-Smirnov (K-S)
453 test. The K-S test includes a one-sample K-S test, two independent sample test, and a
454 test for several independent samples.

455 8. **The third industry.** In China, the third industry is also known as the service
456 industry, and includes the traffic and transportation industry, communication industry,
457 and commercial industry.

458 **Appendix B. Abbreviations used in this paper**

459 1. **PLA** People's Liberation Army of China.

460 2. **GDP** Gross domestic product

461 3. **P** . Precipitation.

462 4. **W_p** . Water resources per capita

463 5. **W_c** Water consumption per 10 thousand CNY GDP



- 464 6. S_r Satisfactory rate of water demand
465 7. U_r Utilization rate of water resources
466 8. IW_p Proportion of industrial water use
467 9. AW_p Proportion of agriculture water use
468 10. DW_p Proportion of domestic water use
469 11. DS_r Treatment rate of domestic sewage
470 12. **CNY.** The Chinese Yuan
471 13. **K-S test.** Kolmogorov-Smirnov Test

472

473 **Acknowledgments** The study was supported by National Natural Science Foundation
474 of China (Grant Nos. 51609254, 51279006 and 51479003).

475 **References**

- 476 Balakrishna, N.: Approximate MLES for the location scale parameters of the half-Logistic
477 distribution with type right-censoring, IEEE Transactions on Reliability, (40), 140–145, 1991.
478 Bedford, T., and Cooke, R.M.: Probability Risk Analysis: Foundations and Methods, Cambridge
479 University Press, Cambridge, 2001.
480 Beijing Municipal Development and Reform Commission and Beijing Municipal Bureau of Water
481 Affairs: Beijing City comprehensive planning of water resources, China Water Power Press,
482 Beijing, 2009. (in Chinese)
483 Breslow, N.E. and Zaho, L.P.: Logistic regression for stratified case-control studies, Biometrics
484 (44):891–899, 1988.
485 Brown, C.C.: On a goodness-of-fit test for the logistic model based on score statistics,



- 486 Communications in Statistics, 11(10), 1087-1105, 1982.
- 487 Christodoulou, S.E.: Water resources conservancy and risk reduction under climatic instability,
488 Water Resources Management 25, 1059–1062, 2011.
- 489 Giacomelli, P., Rossetti, A., and Brambilla, M.: Adapting water allocation management to
490 drought scenarios, Nat. Hazards Earth Syst. Sci., 8, 293–302, 2008.
- 491 Giannikopoulou, A.S., Gad, F.K., and Kampragou, E.K.: Risk-based assessment of drought
492 mitigation options: the case of Syros Island, Greece, Water Resources Management, 31(2),
493 655–669, 2017.
- 494 Haimes, Y.Y.: On the definition of vulnerability in measuring risks to infrastructures, Risk
495 Analysis, 26(2), 293–296, 2006.
- 496 Haimes, Y.Y.: On the complex definition of risk: a systems-based approach, Risk Analysis, 29(12),
497 1647–1654, 2009
- 498 Hashimoto, T., Stedinger, J.R., and Loucks, D.P.: Reliability, resiliency and vulnerability criteria
499 for water resources system performance evaluation, Water Resources Research, 18(1), 14–20,
500 1982.
- 501 Jones, G.A., and Jones, J.M.: Information and coding theory, Springer-Verlag London Ltd, London,
502 2000.
- 503 Karimi, I., and Hüllermeier, E.: Risk assessment system of natural hazards: a new approach based
504 on fuzzy probability, Fuzzy Sets and Systems, 158, 987–999, 2007.
- 505 Khuri, A.I.: Advanced calculus with applications in statistics, John Wiley & Sons, Inc., Hoboken,
506 New Jersey, 2003.
- 507 Li, F.W., Qiao, J.L., Zhao, Y., and Zhang, W.: Risk assessment of groundwater and its application



- 508 part II: using a groundwater risk maps to determine control levels of the groundwater, *Water*
509 *Resources Management*, 28(13), 4875–4893, 2014
- 510 Liang, X.S.: Unraveling the cause-effect relation between time series, *Physical Review E*, 90,
511 052150-1–052150-11, 2014.
- 512 Liang, X.S.: Normalizing the causality between time series, *Physical Review E*, 92, 022126, 2015
- 513 Mackenzie, A.C.: Summarizing risk using risk measures and risk indices, *Risk Analysis*, 34(12),
514 2143–2162, 2014.
- 515 Plummer, R., Loë de Rob, and Armitage, D.: A systematic review of water vulnerability
516 assessment tools, *Water Resources Management*, 26, 4327–4346, 2012.
- 517 Qian, L.X., Wang, H.R., and Zhang, K.N.: Evaluation criteria and model for risk between water
518 supply and water Demand and its application in Beijing, *Water Resources Management*, 28,
519 4433–4447, 2014.
- 520 Qian, L.X., Zhang, R., Hong, M., Wang, H.R., Yang, L.Z.: A new multiple integral model for
521 water shortage risk assessment and its application in Beijing, China, *Natural Hazards*, 80(1),
522 43–67, 2016.
- 523 Wang, C.H., and Blackmore, J.M.: (2012) Supply–demand risk and resilience assessment for
524 household rainwater harvesting in Melbourne, Australia, *Water Resources Management*
525 26(15), 4381–4396, 2012.
- 526 Rajagopalan, B., Nowak, K., Prairie, J., Hoerling, M., Harding, B., Barsugli, J., Ray, A., and Udall,
527 B.: Water supply risk on the Colorado River: Can management mitigate? *Water resources*
528 *research*, 45, W08201, 2009.
- 529 Sandoval-Solis, S., McKinney, D.C., and Loucks, D.P.: Sustainability index for water resources



- 530 planning and management, *Journal of Water Resources Planning and Management*, 137(5),
531 381–390, 2011.
- 532 Statistical Bureau of Beijing City: *Statistical Yearbook 2014 of Beijing City*. China Statistics Press,
533 Beijing, 2014. (in Chinese)
- 534 Tidwell, V.C., Cooper, J.A., and Silva, C.J.: Threat assessment of water supply systems using
535 markov latent effects modeling, *Journal of Water Resources Planning and Management*,
536 131(3), 218–227, 2005.
- 537 UNISDR: *Terminology on disaster risk reduction*, United Nations, Geneva, 2009.
- 538 Villagrán, De León J: *Vulnerability: A conceptual and methodological review*, SOURCE
539 No.4.UNU-EHS, Bonn, 2006.
- 540 Weng, B.S., Yan, D.H., Wang, H., Liu, J.H., Yang, Z.Y., Qin, T.L., and Yin, J.: Drought assessment
541 in the Dongliao River basin: traditional approaches vs. generalized drought assessment index
542 based on water resources systems. *Nat. Hazards Earth Syst. Sci.*, 15, 1889–1906, 2015.
- 543 Yang, C.C., Yeh, C.H., and Ho, C.C.: Systematic quantitative risk analysis of water shortage
544 mitigation projects considering climate change, *Water Resources Management*, 29(4),
545 1067–1081, 2015.