

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

² Collaborative Innovation Center on Forecast Meteorological Disaster Warning and Assessment, Nanjing
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. The parameters of a logistic regression model are often estimated
58 by a maximum likelihood estimation; a large number of observed values of risk (i. e.,

59 samples that water shortage risk does or does not occur) and risk factors are required
60 for parameter estimation (Balakrishnan, 1992). However, the statistical data about risk
61 and its factors are insufficient in China. Therefore, the method of maximum
62 likelihood estimation is not applicable when the sample size is small.  For this reason,
63 we propose an alternative method of parameter estimation for a logistic regression
64 model when data is insufficient. Moreover, the backward mode is often applied for the
65 selection of sensitive factors, but the calculation is very complicated. 

66 The contributions of our paper are as follows. First, we used a logistic regression
67 model to explore the nonlinear relation between water shortage risk and its factors.
68 Then, we introduced an information flow (Liang, 2014) for the selection of significant
69 risk factors. Compared with the backward mode, it was very easy to determine
70 whether there was a cause and effect between the water shortage risk and its factors. 

71 Finally, we proposed an alternative method of parameter estimation (maximum
72 entropy estimation) for a logistic regression model in situations with a lack of data.
73 The new method requires only a few data, while maximum likelihood estimation
74 requires a large amount of data.

75 The remainder of the paper is organized as follows. Section 2 presents the
76 principles and structure of the logistic probability prediction model for water shortage
77 risk. Section 3 presents the application of the model and the results of the research and
78 Section 4 presents some conclusions and proposes future work.

79 **2 Materials and methods**

80 **2.1 Study area**

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

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



97

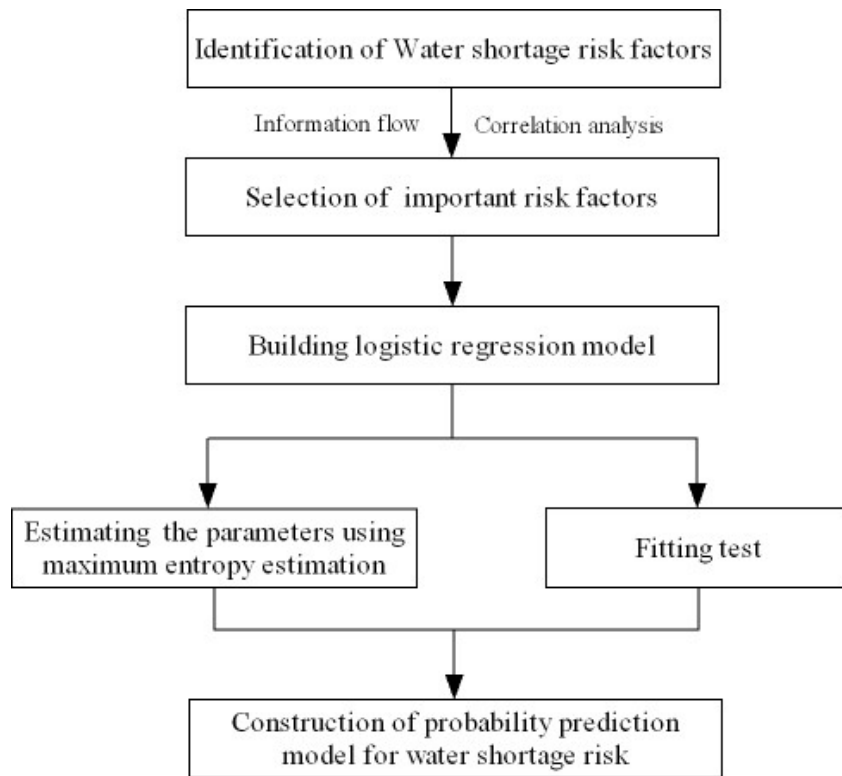
Figure 1. Distribution of water system of Beijing

98 **2.2 Data collection**

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

112 **2.3 Model development**

113 A flowchart showing the operation of the probability prediction model for water
114 shortage risk is given in Figure 2.



115

116

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


117

water shortage risk

118

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

119

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

120

2.3.1 Identification of water shortage risk factors

121

Water shortage risk factors include characteristics of hazards and existing conditions

122

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

123

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

124

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

125

2016). According to the study of Plummer et al. (2012), there are 50 different water

126

vulnerability assessment tools, and the water vulnerability indicators of these tools are

127

quite different. Therefore, a universal standard understanding of water resource

128

vulnerability indicators is difficult to develop. We established the indicators from

129 perspective of hydrological conditions, water resources, water supply and water use.
 130 The risk factors are: precipitation (P), water resources per capita (W_p), water
 131 consumption per GDP (W_c), satisfactory rate of water demand (S_r), and utilization
 132 rate of water resources (U_r), proportion of industrial water use (IW_p), proportion of
 133 agricultural water use (AW_p), proportion of domestic water use (DW_p) and the
 134 treatment rate of domestic sewage (DS_r). These indicators are defined as follows
 135 (Qian et al., 2014):

$$136 \quad W_p = \frac{W}{N} \quad (1)$$

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

$$138 \quad W_c = \frac{\text{the amount of water use}}{GDP} \quad (2)$$

$$139 \quad U_r = \frac{W_{ss} + W_{gs}}{W} = \frac{W_{as}}{W} \quad (3)$$

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

$$142 \quad DS_r = \frac{DS_t}{DS} \quad (4)$$

143 where DS_t is the amount of sewage treated and DS is the total amount of sewage
 144 discharged.

$$145 \quad S_r = \frac{W_{as}}{W_{td}} \quad (5)$$

146 where W_{as} is the water supply, and W_{td} is the water demand.

$$147 \quad IW_p = \frac{IW}{WU} \quad (6)$$

$$148 \quad AW_p = \frac{AW}{WU} \quad (7)$$

$$149 \quad DW_p = \frac{DW}{WU} \quad (8)$$

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

152 **2.3.2 Selection of important risk factors**

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


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

$$162 \quad T_{2 \rightarrow 1} = \frac{C_{11}C_{12}C_{2,d1} - C_{12}^2C_{1,d1}}{C_{11}^2C_{22} - C_{11}C_{12}^2} \quad (9)$$


163 where C_{ij} is the sample covariance between X_i and X_j , C_{i,d_j} is the covariance
 164 between X_i and \dot{X}_j , and \dot{X}_j is the difference approximation of $\frac{dX_j}{dt}$ using the Euler
 165 forward scheme.

$$166 \quad \dot{X}_{j,n} = \frac{X_{j,n+k} - X_{j,n}}{k\Delta t} \quad (10)$$

167 According to Liang (2014), with $k \geq 1$, for a general time series $k=1$ would be
 168 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
 169 cause X_1 , otherwise it is causal. A positive $T_{2 \rightarrow 1}$ means that X_2 functions to make X_1


170 more uncertain, while a negative value means that X_2 tends to stabilize X_1 . Liang
171 (2015) proposed a method of normalizing the causality between time series and the
172 range of value for $T_{2 \rightarrow 1}$ is 0 and 1. 

173 **2.3.3 Correlation analysis of selected risk factors**

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


179 **2.4 Risk probability prediction model using maximum entropy** 180 **estimation**

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

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

190 $p_i = p(y_i = 1 | x_{ij} (j=1,2,\dots,m))$ is the conditional probability when $y_i=1$ under

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

192
$$P_i = \frac{1}{1 + e^{-(\alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_m x_{im})}} \quad (11)$$
 

193 where $\alpha, \beta_1, \beta_2, \dots, \beta_m$ are the estimated parameters. The parameters are often

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

195 computing $\alpha, \beta_1, \beta_2, \dots, \beta_m$ is as follows:

196
$$\left\{ \begin{aligned} \frac{\partial L}{\partial \alpha} &= \sum_{i=1}^n \left[y_i - \frac{\exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)}{1 + \exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)} \right] = 0 \\ \frac{\partial L}{\partial \beta_j} &= \sum_{i=1}^n \left[y_i - \frac{\exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)}{1 + \exp\left(\alpha + \sum_{j=1}^m \beta_j x_{ij}\right)} \right] x_{ij} = 0 \quad j = 1, 2, \dots, m \end{aligned} \right. \quad (12)$$

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

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

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

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

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

202 parameter estimation is therefore required.

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

204 entropy principle. The new method is named after maximum entropy estimation. The


205 new method does not require the observed values of risk, and it requires only some

206 observed values of the factors. Its principle is as follows.

207 For an observation, we can define its entropy to evaluate its degree of uncertainty.

208 According to Jones and Jones (2000), the entropy of the i th observation of water
 209 shortage risk is

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

211 where C is a positive value and $p_i = p(y_i = 1 | x_{ij} (j = 1, 2, \dots, m))$ is the
 212 conditional probability when $y_i = 1$ under the conditions of
 213 $x_{ij} (i = 1, 2, \dots, n; j = 1, 2, \dots, m)$. According to the maximum entropy principle, if the
 214 values of $H(P_i)$ reaches a maximum, the optimal parameters are obtained (Jones
 215 and Jones, 2000). The reasons for obtaining a solution based on the maximum entropy
 216 principle are as follows. ① It conforms to the principle of entropy increase, which
 217 states that the entropy of an isolated system tends to reach a maximum. ② It accords
 218 with the principle that the solution should be in line with the sample/data and the least
 219 hypotheses must be constructed regarding the unknown parts when the data is
 220 insufficient. ③ It fits the maximum multiplicity principle. The multiplicity of a state
 221 refers to the number of possible ways in which a system can evolve to that state. The
 222 maximum multiplicity principle states that the greater the multiplicity of a state, the
 223 larger the possibility that a system is in this state. 

224 **2.4.1 Parameter estimation**

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

$$226 \quad \max H_i = -C \left\{ \frac{\left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right]} - \ln \left(1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right) \right\} \quad (14)$$

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

$$229 \quad \begin{cases} \frac{\partial H_i}{\partial \alpha} = \frac{\left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\} + \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \cdot \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right]}{\left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\}^2} - \frac{\exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)} = 0 \\ \frac{\partial H_i}{\partial \beta_j} = \frac{x_{ij} \cdot \left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\} + x_{ij} \cdot \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \cdot \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right]}{\left\{ 1 + \exp \left[- \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right) \right] \right\}^2} - \frac{x_{ij} \cdot \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)}{1 + \exp \left(\alpha + \sum_{j=1}^m \beta_j x_{ij} \right)} = 0 \end{cases} \quad (15)$$

230 The optimal estimation $\alpha, \beta_j (j = 1, 2, \dots, m)$ can be obtained by solving Eq.

231 (15). Numerical approaches are often used to obtain an approximate solution of Eq.

232 (15) rather than its exact solution. Therefore, we made use of the optimization

233 function of Matlab to estimate the parameters, i.e., the fminsearch function. If there

234 are n observations, there are $n H_i (i = 1, 2, \dots, n)$. It is impossible to find the parameters

235 that make all the $H_i (i = 1, 2, \dots, n)$ reach the maximum value. According to the

236 maximum entropy principle, the greater the entropy is, the larger the uncertainty of an

237 observation is. Therefore, the maximum value of the sequences $\{H_i, i = 1, 2, \dots, n\}$ was

238 taken as the objective function of the optimization model. 

239 **2.4.2 Goodness-of-fit test**

240 According to Brown (1982), a goodness-of-fit test should be made for evaluating the

241 fitting effect of the logistic regression model and its ability to identify water shortage
 242 risk. In this study, the Kolmogorov-Smirnov Test (K-S) test and Pearson χ^2 test are
 243 used.

244 **2.4.2.1 K-S test (t)**

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

248 Let $F_{n_1}(x)$ be the cumulative probability distribution of the samples that do not
 249 encounter a water shortage. $F_{n_2}(x)$ is the cumulative probability distribution of the
 250 samples that encounter a water shortage. A two independent samples test is then
 251 applied to compare whether the empirical distribution functions of two samples are
 252 the same. The test is as follows:

253
$$H_0 : F_{n_1}(x) = F_{n_2}(x) \quad H_1 : F_{n_1}(x) \neq F_{n_2}(x) \quad (16)$$

254 The value of K-S is:

255
$$K - S = \max |F_{n_1}(x) - F_{n_2}(x)| \quad (17)$$

256 When $N \rightarrow \infty$, the cumulative distribution curve and probability density curve of
 257 two samples can be obtained. The value of K-S is the maximum value of the
 258 cumulative distribution functions. When the value of K-S is greater than 0.35, the
 259 logistic regression model is applicable. The international classification standard of the
 260 logistic model is shown in Table 1 (Brown, 1982).

261 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

262

263 **2.4.2.2 Pearson χ^2 test**

264 The test is as follows:

$$265 \quad H_0: \text{the fitting is good} \quad H_1: \text{the fitting is bad} \quad (18)$$

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

$$267 \quad \chi^2 = \sum_{j=1}^l \frac{(O_j - E_j)^2}{E_j} \quad (19)$$

268 where $j = 1, 2, \dots, l$, l is the number of covariant types, O_j is the observed
 269 frequency of the j th covariant type, and E_j is the predicted frequency of
 270 the j th covariant type. The degree of freedom is the difference between the number of
 271 covariant types and parameters.

272 **3 Results and discussion**

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

276 **3.1 Construction of the Logistic probability prediction model**

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

282 Therefore, there are only 34-year data.

283 3.1.1 Determination of water resources vulnerability indicators

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

291 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

292

293 According to Liang (2014), a positive value of the information flow means that

294 the factor makes water shortage risk more uncertain, while a negative value means

295 that the indicator tends to stabilize water shortage risk. Therefore, all the factors tend

296 to make water shortage risk more uncertain. Furthermore, the impact of P , W_p ,

297 W_c are very significant.

298 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

299

300 A correlation analysis was performed on the remaining factors. The values of the
301 Pearson correlation coefficients are shown in Table 4.

302 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

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

306 **3.1.2 Construction of the logistic risk probability predication model**

307 The data for the risk and selected factors (W_c , W_p , U_r , P and DS_r) from 1979 to
308 2012 (Table 2) are used to construct the logistic risk predication probability model.

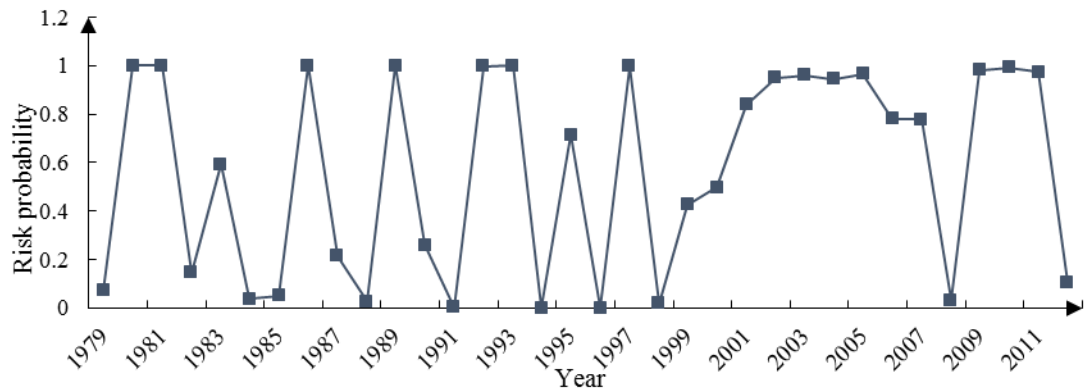
309 Because there is only 34 samples, it is impossible to estimate the parameters by the 

310 maximum likelihood estimation. Substituting the sequences of W_c , W_p , U_r , P and DS_r
311 from 1979 to 2012 (Table 2) into Eq. (14), the values of parameters obtained by
312 maximum entropy estimation can be obtained. The estimated values for
313 $\alpha, \beta_1, \beta_2, \dots, \beta_5$ are 61.6386, 0.004, -0.1262, -12.4077, -0.012 and -29.0963.

314 Therefore, the logistic regression model based on the maximum entropy
315 estimation is as follows:

$$316 \quad \text{Predicted probability} = \frac{1}{1 + e^{-(61.6386 + 0.004W_c - 0.1262W_p - 12.4077U_r - 0.012P - 29.0963DS_r)}} \quad (20)$$

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



320

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

322

to 2012

323

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

324

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

325

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

326

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

327

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

328

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

329

Table 5. The prediction accuracy using the maximum entropy estimation

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

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

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

341 ***3.2 Risk probability prediction in 2020 in Beijing***

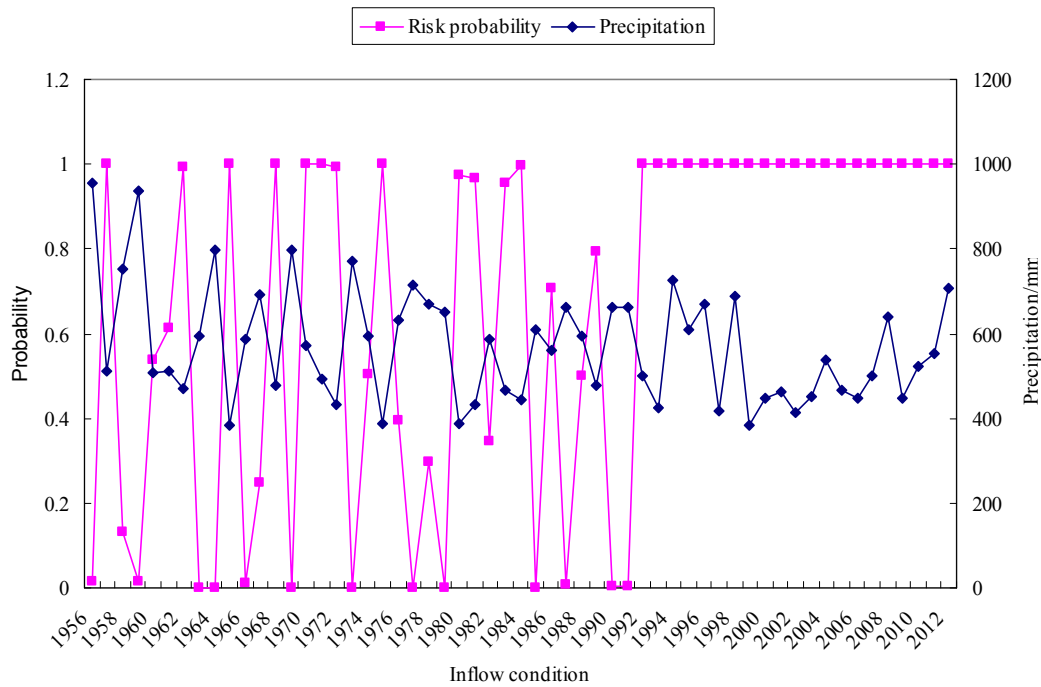
342 ***3.2.1 Risk probability prediction (without considering the use of*** 343 ***transferred and reclaimed water)***

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

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

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

367 In Figure 4, the horizontal axis represents the inflow conditions of 1956–2012.
368 Figure 4 shows that in 2020, the water shortage risk probability exceeded 0.95 under
369 33 different inflow conditions (accounting for 63.5% of all the inflow conditions) and
370 exceeded 0.5 under 38 different inflow conditions (accounting for 73.1% of all the
371 inflow conditions). In summary, there was a high probability of a water shortage risk
372 in 2020, although the probability was very low in some high precipitation periods.



373

374

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

375

3.2.2 Risk probability prediction after using transferred and reclaimed

376

water

377

According to Qian et al. (2016), 1.05 billion m³ of water will have been transferred

378

to Beijing in 2020 and the amount of reclaimed water used may reach 1 billion m³.

379

After using transferred and reclaimed water, the total amount of water resources

380

would increase, W_p and U_r would change and other indicators would remain

381

unchanged. Therefore, the sequences of W_p and U_r under the inflow scenarios of

382

1956–2012 had to be computed again. Substituting the sequences of

383

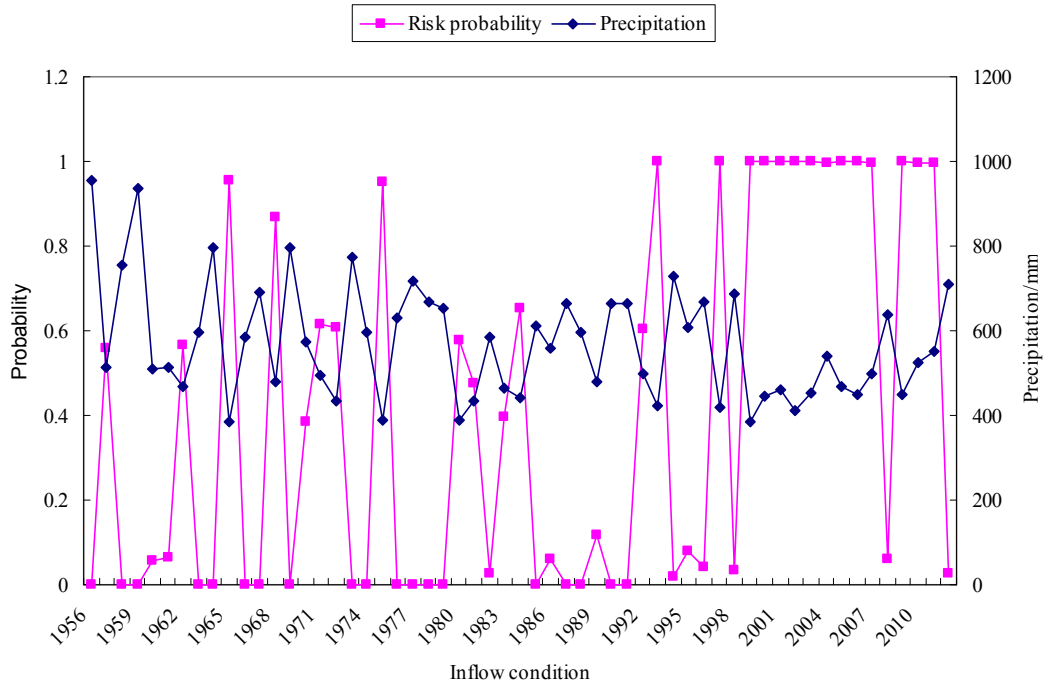
W_c, W_p, U_r, P and DS_r into Eq. (20), the water shortage risk probability in 2020

384

under the inflow scenarios of 1956–2012 (after using transferred and reclaimed

385

water) was predicted, and the results are shown in Figure 5.



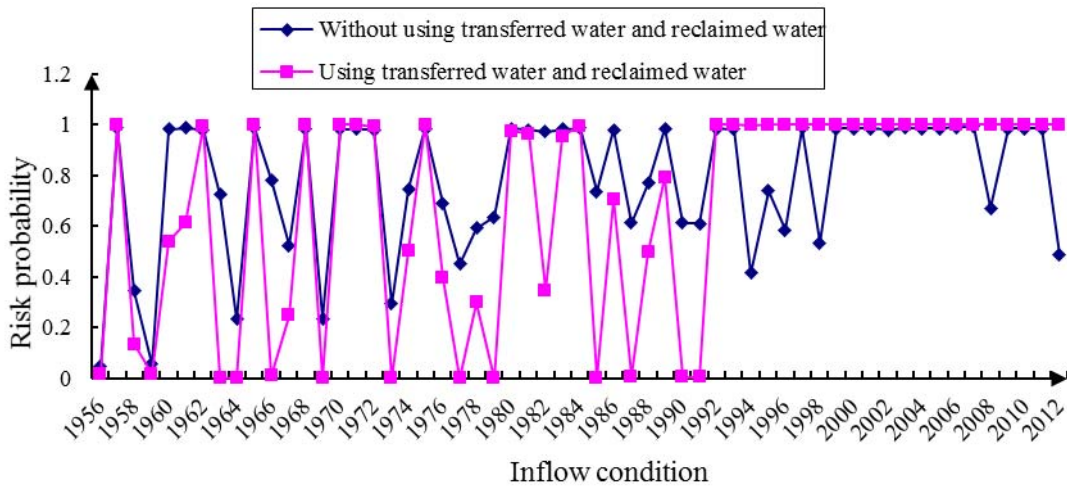
386

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

388

transferred and reclaimed water

389



390

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

392

water

393

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

394

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

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

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

408 **4 Conclusions**

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

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

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

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

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

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

433

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

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

437 **2. Maximum likelihood estimation.** It is a method of parameter estimation in
438 statistics.

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

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

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

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

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

452 8. **The third industry.** In China, the third industry is also known as the service
453 industry, and includes the traffic and transportation industry, communication industry,
454 and commercial industry.

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

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

457 2. **GDP** Gross domestic product

458 3. P . Precipitation.

459 4. W_p . Water resources per capita

460 5. W_c Water consumption per 10 thousand CNY GDP

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

469

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

472 **References**

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

483 Communications in Statistics, 11(10), 1087-1105, 1982.

484 Christodoulou, S.E.: Water resources conservancy and risk reduction under climatic instability,
485 Water Resources Management 25, 1059–1062, 2011.

486 Giacomelli, P., Rossetti, A., and Brambilla, M.: Adapting water allocation management to
487 drought scenarios, Nat. Hazards Earth Syst. Sci., 8, 293–302, 2008.

488 Giannikopoulou, A.S., Gad, F.K., and Kampragou, E.K.: Risk-based assessment of drought
489 mitigation options: the case of Syros Island, Greece, Water Resources Management, 31(2),
490 655–669, 2017.

491 Haimes, Y.Y.: On the definition of vulnerability in measuring risks to infrastructures, Risk
492 Analysis, 26(2), 293–296, 2006.

493 Haimes, Y.Y.: On the complex definition of risk: a systems-based approach, Risk Analysis, 29(12),
494 1647–1654, 2009

495 Hashimoto, T., Stedinger, J.R., and Loucks, D.P.: Reliability, resiliency and vulnerability criteria
496 for water resources system performance evaluation, Water Resources Research, 18(1), 14–20,
497 1982.

498 Jones, G.A., and Jones, J.M.: Information and coding theory, Springer-Verlag London Ltd, London,
499 2000.

500 Karimi, I., and Hüllermeier, E.: Risk assessment system of natural hazards: a new approach based
501 on fuzzy probability, Fuzzy Sets and Systems, 158, 987–999, 2007.

502 Khuri, A.I.: Advanced calculus with applications in statistics, John Wiley & Sons, Inc., Hoboken,
503 New Jersey, 2003.

504 Li, F.W., Qiao, J.L., Zhao, Y., and Zhang, W.: Risk assessment of groundwater and its application

505 part II: using a groundwater risk maps to determine control levels of the groundwater, *Water*
506 *Resources Management*, 28(13), 4875–4893, 2014

507 Liang, X.S.: Unraveling the cause-effect relation between time series, *Physical Review E*, 90,
508 052150-1–052150-11, 2014.

509 Liang, X.S.: Normalizing the causality between time series, *Physical Review E*, 92, 022126, 2015

510 Mackenzie, A.C.: Summarizing risk using risk measures and risk indices, *Risk Analysis*, 34(12),
511 2143–2162, 2014.

512 Plummer, R., Loë de Rob, and Armitage, D.: A systematic review of water vulnerability
513 assessment tools, *Water Resources Management*, 26, 4327–4346, 2012.

514 Qian, L.X., Wang, H.R., and Zhang, K.N.: Evaluation criteria and model for risk between water
515 supply and water Demand and its application in Beijing, *Water Resources Management*, 28,
516 4433–4447, 2014.

517 Qian, L.X., Zhang, R., Hong, M., Wang, H.R., Yang, L.Z.: A new multiple integral model for
518 water shortage risk assessment and its application in Beijing, China, *Natural Hazards*, 80(1),
519 43–67, 2016.

520 Wang, C.H., and Blackmore, J.M.: (2012) Supply–demand risk and resilience assessment for
521 household rainwater harvesting in Melbourne, Australia, *Water Resources Management*
522 26(15), 4381–4396, 2012.

523 Rajagopalan, B., Nowak, K., Prairie, J., Hoerling, M., Harding, B., Barsugli, J., Ray, A., and Udall,
524 B.: Water supply risk on the Colorado River: Can management mitigate? *Water resources*
525 *research*, 45, W08201, 2009.

526 Sandoval-Solis, S., McKinney, D.C., and Loucks, D.P.: Sustainability index for water resources

527 planning and management, *Journal of Water Resources Planning and Management*, 137(5),
528 381–390, 2011.

529 Statistical Bureau of Beijing City: *Statistical Yearbook 2014 of Beijing City*. China Statistics Press,
530 Beijing, 2014. (in Chinese)

531 Tidwell, V.C., Cooper, J.A., and Silva, C.J.: Threat assessment of water supply systems using
532 markov latent effects modeling, *Journal of Water Resources Planning and Management*,
533 131(3), 218–227, 2005.

534 UNISDR: *Terminology on disaster risk reduction*, United Nations, Geneva, 2009.

535 Villagrán, De León J: *Vulnerability: A conceptual and methodological review*, SOURCE
536 No.4.UNU-EHS, Bonn, 2006.

537 Weng, B.S., Yan, D.H., Wang, H., Liu, J.H., Yang, Z.Y., Qin, T.L., and Yin, J.: Drought assessment
538 in the Dongliao River basin: traditional approaches vs. generalized drought assessment index
539 based on water resources systems. *Nat. Hazards Earth Syst. Sci.*, 15, 1889–1906, 2015.

540 Yang, C.C., Yeh, C.H., and Ho, C.C.: Systematic quantitative risk analysis of water shortage
541 mitigation projects considering climate change, *Water Resources Management*, 29(4),
542 1067–1081, 2015.