# Responses to Anonymous Referee #2:

**General comments**

The manuscript 'An improved logistic prediction model for water shortage risk in situations with insufficient data' Qian et al. is focused on the very important issue of water shortage prediction in a large metropolitan area as Beijing. The authors propose the approach of maximum entropy to estimate the parameters of a logistic model. Despite the indubitable interest of the issue, and the worth of the basic idea of the authors to refer to maximum entropy concept (to overcome the limitation of other approaches), the overall quality of the manuscript is rather poor. The presentation of the method is superficial, not well framed in the state of art of the existing risk assessment methods, and in particular for those that refer to the maximum entropy

approach. The formalization of the approach is rather rough and results are not described with sufficient detail to support the reliability of the proposed approach. Therefore, in my opinion the paper can not be accepted for publication in the present form. At the same time, I think that the idea is interesting and it is worth to be further elaborate on. Therefore I suggest to the authors to improve the quality of the manuscript and submit it again. My detailed suggestions can be found in the attached file.

Responses：Good suggestions. We have made a major revision to our paper. First, we have rewritten the introduction and cited more studies about parameter estimation of logistic regression model with a well expressed motivation of the background and objectives of the proposal. Second, we have added ten tests to evaluate the performance of maximum entropy estimation under different small sample size, compared with maximum likelihood estimation. The result shows that maximum entropy estimation performs much better than maximum likelihood estimation under small samples. More details are shown in the Section 3 of the revised manuscript and responses to your detailed suggestions. Third, we have revised the presentation of the method to make it well framed in the state of art of the existing risk assessment methods. Fourth, we have given more details about the approach and results to support the reliability of the proposed approach.

We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

**Detailed suggestions**

1. The introduction is almost confuse. The initial part from row 21 to 48 could be written in more concise way. A sentence that underlines the necessity to take into account the risk factors should be sufficient. The concepts expressed from row 49 to the end, should be better developed, justifying the different choices made, as for example:

Responses：Good suggestions. We have rewritten the introduction with a well expressed motivation of the background and objectives of the proposal. We have rewritten the initial part from row 21 to 48 in a more concise way. We have added some contents that underlines the necessity to take into account the risk factors. Moreover, we have cited more studies about the necessity to take into account the risk factors. The concepts expressed from row 49 to the end, have been better developed. More details are shown in the new introduction. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

a) What do you mean as water shortage?

Responses：Good suggestions. Water shortage refers to that the amount of water supply is less than that of water demand. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

b) Why do you use logistic regression? Was the only existing approach? Are there other approaches in literature to present the causal link between short water

probability and risk factors? In literature, have parameter estimation by a maximum likelihood estimation been never compared with other methods as maximum entropy estimation? Could your affirmation about the reduced number of data for parameter estimation by maximum entropy be supported by any tests? Since the main contribution of the paper is focused on parameter estimation by maximum entropy so it should be appropriate to cite more studies about this subject, concerning the use of such approach in the context of risk assessment.

Responses：Good suggestions. A water shortage can either occurs or not occur in a year, and it has two possible values: water shortage occurs or water shortage does not occur. Therefore, water shortage risk is a binary categorical variable. A logistic regression model, a method of probability prediction, is often applied to study the relation between a binary categorical variable and its factors (Scott and Wild, 1991). Therefore, this paper used a logistic regression model to present the causal link between water shortage risk and its factors. The merits of logistic regression model in risk predication applications have been discussed in many studies (Udevitz et al., 1987; Yerel and Anagun, 2010). The major purpose of this paper is to propose an alternative method of parameter estimation for logistic regression model in situations when insufficient data are available.

There are some approaches in literature to present the causal link between short water probability and risk factors such as discriminant analysis model and information diffusion model. However, these models can only obtain the level of risk rather than risk probability.

Maximum entropy estimation is proposed by us, so parameter estimation by a maximum likelihood estimation has been never compared with maximum entropy estimation in literature. In the Section 3 of the revised manuscript, we have added ten tests to evaluate the performance of maximum entropy estimation under different small sample size, compared with maximum likelihood estimation. The result shows that maximum entropy estimation performs much better than maximum likelihood estimation under small samples. Because maximum entropy estimation is proposed by us, there are no studies about this subject. Moreover, we cited more studies about maximum likelihood estimation and its shortcomings.

The details of the tests are as follows.

We used random number generator of logistic regression to generate a sequence with six parameters and its sample size is 1000. Ten simulations were performed to evaluate the performance of maximum entropy estimation under different small sample size. The small sample size are 100, 90, 80, 70, 60, 50, 40, 30, 20 and 10, respectively. Tables 1 and 2 shows a comparison of the values of the absolute percentage error (APE) generated from maximum likelihood estimation and maximum entropy estimation. The APE is calculated as follows.

$$APE = \frac{|p_i - P|}{P} \qquad (1)$$

where $p_i$ is the value of the parameter generated from different parameter methods under sample size of 100, 90, 80, 70, 60, 50, 40, 30, 20 and 10, and $P$ is the value of the parameter generated from maximum likelihood estimation under sample size of 1000.

Table 1. Comparison of the APE values generated from maximum likelihood estimation under different sample size

| APE | 100 | 90 | 80 | 70 | 60 | 50 | 40 | 30 | 20 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $\alpha$ | 0.8% | 21.0% | 198% | 132100% | 132100% | 131800% | 131800% | 131800% | 134000% | / |
| $\beta_1$ | 62.4% | 49.2% | 86.7% | 84900% | 85500% | 85500% | 85500% | 85500% | 85500% | / |
| $\beta_2$ | 1.3% | 50.1% | 83.3% | 260400% | 260800% | 260800% | 260800% | 260800% | 268600% | / |
| $\beta_3$ | 34.8% | 73.9% | 8.7% | 214900% | 215200% | 215200% | 215200% | 215200% | 220700% | / |
| $\beta_4$ | 18100% | 20100% | 16100% | 30306100% | 30332100% | 30332100% | 30332100% | 30332100% | 31026100% | / |
| $\beta_5$ | 88.4% | 27.3% | 335% | 144900% | 144800% | 144800% | 144800% | 144800% | 153200% | / |

Table 2. Comparison of the APE values generated from maximum entropy estimation under different sample size

| APE | 100 | 90 | 80 | 70 | 60 | 50 | 40 | 30 | 20 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $\alpha$ | 2.3% | 2.5% | 3.6% | 3.6% | 3.6% | 3.6% | 3.6% | 3.6% | 3.6% | 43.1% |
| $\beta_1$ | 0.6% | 0.1% | 0.2% | 0.2% | 0.2% | 0.2% | 0.2% | 0.2% | 0.2% | 21.2% |
| $\beta_2$ | 0.9% | 0.7% | 0.9% | 0.9% | 0.9% | 0.9% | 0.9% | 0.9% | 0.9% | 26.2% |
| $\beta_3$ | 10.4% | 14.3% | 11.6% | 11.6% | 11.6% | 11.6% | 11.6% | 11.6% | 11.6% | 11.6% |
| $\beta_4$ | 0% | 45.0% | 55.0% | 55.0% | 55.0% | 55.0% | 55.0% | 55.0% | 55.0% | 55.0% |
| $\beta_5$ | 2.3% | 0.7% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 18.8% |

Tables 1 and 2 show that maximum entropy estimation performs much better than maximum likelihood estimation. For example, maximum entropy estimation can obtain a satisfactory result when the sample size is greater than 20, while maximum likelihood estimation performs so badly. Moreover, maximum entropy estimation still provided an acceptable result with only 10 samples, while maximum likelihood

estimation was inapplicable with 10 samples.

We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

Reference：

Scott, A.J. and Wild, C.J.: Fitting logistic regression model in stratified case-control studies, Biometrics 47(2):497–510, 1991.

Udevitz, MS., Bloomfield, P., and Apperson, CS.: Prediction of occurrence of four Species of Mosquito Larvae with logistic regression on water-chemistry variables, Environmental Entomology, 16(1), 281–285, 1987.

Yerel, S., and Anagun, AS.: Assessment of water quality observation stations using cluster analysis and ordinal logistic regression technique, International Journal of Environment & Pollution, 42(4), 344–358, 2010.

2. The sentence should be better clarified. You give a definition of water shortage in paragraph 2.4. It would be more appropriate to provide such definition here.

Responses：Good suggestions. The sentence has been better clarified. We have given a definition of water shortage here. More details are shown in Lines 55-57. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

3. Lines 57-62, since this is the main motivation of the proposed study it should be discussed more deeply.

Responses：Good suggestions. The main motivation of this paper has been discussed more deeply. More details are shown in Lines 55-76. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

4. Line 64-65, This sentence is not clear. It could be omitted or the issue should be discuss more deeply.

Responses：Good suggestions. The sentence has been omitted, because it is needless. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

5. Does information flow work with small sample size?

Responses：Good suggestions. Information flow has no special requirement for sample size. You could refer to Liang 2014. Unraveling the cause-effect relation between time series, Physical Review E, 90, 052150-1–052150-11. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

6. The different steps of flowchart should be shortly described.

Responses：Good suggestions. The different steps of flowchart has be shortly described in the revised manuscript. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

7. The choice of risk factors should be justified and discussed.

Responses：Good suggestions. We have added more discussions and justifications about the choice of risk factors in the revised manuscript. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

8. What is the difference between $W$ and $W_{as}$? Here is not clear if other water supplies, beyond surface water supply and groundwater supply, are taking into account.

Responses：Good suggestions. $W$ is the total amount of water resources, and $W_{as}$ refers to the sum of amount of surface water supply and groundwater supply. $W_{as}$ may be less than $W$ or greater than $W$. As for Beijing, the available surface water and groundwater is unable to meet the needs of the city's economic and social development. Some measures, such as the use of transferred and reclaimed water have been put in place to mitigate the water shortage. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

9. Lines 140-141, this can be omitted.

Responses：Good suggestions. We have omitted the sentence "and $W$ is the total amount of water resources". We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

10. Do you normalize the information flow between 0 and 1? May you explain better what you have really done?

Responses：Good suggestions. We normalized the information flow between 0 and 1. We have given some explanations in the revised manuscript. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

11. Lines 174-178, this sentence could be put at the end of the paragraph 2.3.1.

Responses：Good suggestions. This sentence has been put at the end of the paragraph 2.3.1. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

12. Line 188, shortage risk

Responses：Good suggestions. Shortagerisk has been revised as "shortage risk". We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

13. Line 192, I am not sure that this definition is correct.

Responses：Good suggestions. This definition (Eq. (11)) is correct. You can refer to Houwelingen and Cessie 2010. Logistic Regression, a review, Statistica Neerlandica, vol. 42, 215-232. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

14. Since you don't use the maximum likelihood estimation, I don't think the relationship 12 is necessary.

Responses：Good suggestions. We have deleted relationship 12 in the revised manuscript. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

15. This motivations are almost generic and critical. You could refer to Sigh 1997. The use of entropy in hydrological and water resources, Hydrological processes vol. 11, 587-626, to explain the reasons to apply maximum entropy estimation.

Responses：Good suggestions. We have referred to Sigh 1997. The use of entropy in hydrological and water resources, Hydrological processes vol. 11, 587-626, to explain the reasons to apply the principle of entropy estimation. More details are shown in the revised manuscript. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

16. Would be more appropriate to refer to the methodology approach rather than the name of matlab function.

Responses：Good suggestions. We have revised the methodology approach as the nonlinear optimization method. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

17. I don't understand why the belonging of Fn1 and Fn2 at a common

cumulative probability distribution is a proof of a good fitting of logistic regression.

Responses: Good suggestions. A K-S test is often applied as a fitting test. $K - S = \mathbf{max} \left| F_{n1}(x) - F_{n2}(x) \right|$, and it is used to judge whether $F_{n1}(x)$ and $F_{n2}(x)$ are significantly different. Therefore, the value of K-S is a reflection of the ability of the logistic model to identify water shortage risk. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

18. Is $l$ equal to $m$?

Responses: Good suggestions. $l$ is the number of covariate pattern. $m$ is the number of parameters. $l$ is not equal to $m$. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

19. How have the observed frequency been estimated from a sample of 34 data in a 4-dimensional space of independent variables?

Responses: Good suggestions. The observed frequency refers to the number of the cases that risk occurs or does not occurs in each covariate pattern. We can count the observed frequency in each covariate pattern. Supposing the predicted frequency of the cases in the $jth$ covariant pattern is $E_j$. $E_j = n_j \times \hat{p}_j$, where $n_j$ is the number of cases in the $jth$ covariant pattern, and $\hat{p}_j$ is the predicted probability that cases in the $jth$ covariant pattern occur. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

20. You mean that $AW_p$ does not affect water shortage.

Responses: Good suggestions. Yes, we mean that $AW_p$ does not affect water shortage, because the value of the normalized information flow of $AW_p$ is only 0.003. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

21. Looking to the values of table 3, it is not clear why these factors are more significant than the others.

Responses: Good suggestions. We are sorry we made a mistake. The sentence "the impact of $P$, $W_p$ an $W_c$ are very significant" should be revised as "the impacts of $W_p$, $W_c$ and $U_r$ are very significant", because the values of normalized information flow of $W_p$, $W_c$ and $U_r$ are much larger than others. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

22. This session should be rewritten in more clear form. Results of Table 3 and table 4 discussed more deeply in order to provide a more clear explanation of the selected factors.

Responses: Good suggestions. We have rewritten this session in more clear form. We have discussed more deeply about the results of Tables 3 and 4. More details are shown in the revised manuscript. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

23. Lines 309-310, This sentence could be omitted.

Responses: Good suggestions. We have deleted it. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

24. To compare in the same figure the risk observed probability could be useful for an evaluation of the performance of the method, rather than the values of table 5.

Responses: Good suggestions. We have added the observed probability in Figure 3 for an evaluation of the performance of the method. The new figure 3 is as follows.
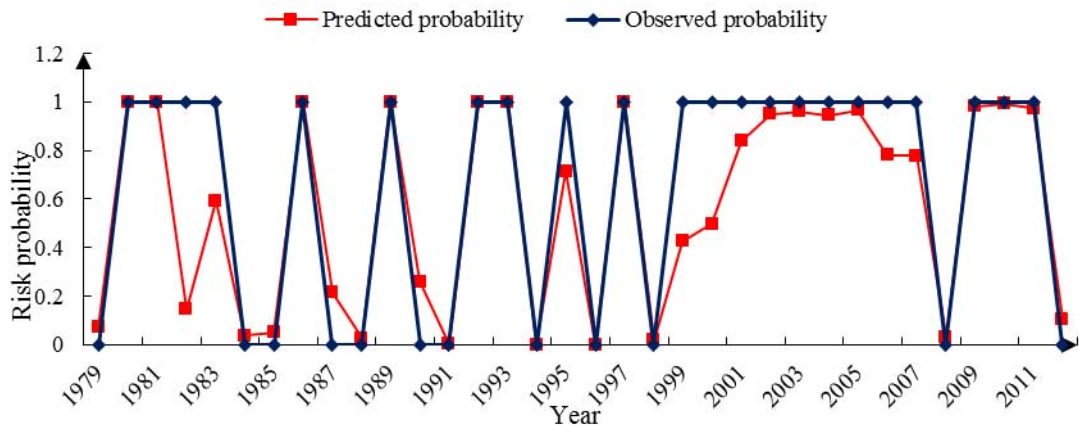


Figure 3. Comparison of the observed and predicted probability generated by the maximum entropy estimation from 1979 to 2012

According to Houwelingen and Cessie (2010) (Logistic Regression, a review, Statistica Neerlandica, vol. 42, 215-232.), table 5 is also a criteria for evaluation of the performance of logistic regression model, so we don't delete table 5.

We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

25. Lines 326-328, This sentence could be omitted.

Responses: Good suggestions. We have deleted it. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.

26. I don't think that the glossary is useful.

Responses: Good suggestions. We have deleted the glossary. We sincerely hope for your satisfaction with our revision. Thank you again for your kind suggestion.