



1 **Improved logistic regression model based on a spatially weighted technique (ILRBSWT**
2 **v1.0) and its application to mineral prospectivity mapping**

3

4 Daojun Zhang^{1,2*}, Na Ren¹, Xianhui Hou¹

5 ¹College of Economics and Management, Northwest A&F University, Yangling 712100,
6 China

7 ²Center for Resource Economics and Environment Management, Northwest A&F University,
8 Yangling 712100, China

9 *Corresponding author: cugzdj@gmail.com (Zhang, D)

10 **Abstract:** Due to complexity, multiple minerogenic stages, and superposition during
11 geological processes, the spatial distributions of geological variables also exhibit specific
12 trends and non-stationarity. For example, geochemical elements exhibit obvious spatial
13 non-stationarity and trends because of the deposition of different types of coverage. Thus,
14 bias may clearly occur under these conditions when general regression models are applied to
15 mineral prospectivity mapping (MPM). In this study, we used a spatially weighted technique
16 to improve general logistic regression and developed an improved model called the improved
17 logistic regression model based on spatially weighted technique (ILRBSWT, version 1.0).
18 The capabilities and advantages of ILRBSWT are as follows: (1) ILRBSWT is essentially a
19 geographically weighted regression (GWR) model, and thus it has all its advantages when
20 dealing with spatial trends and non-stationarity; (2) the current software employed for GWR
21 mainly applies linear regression whereas ILRBSWT is based on logistic regression, which is
22 used more commonly in MPM because mineralization is a binary event; (3) a missing data
23 process method borrowed from weights of evidence is included to extend the adaptability
24 when dealing with multisource data; and (4) the differences of data quality or exploration
25 level can also be weighted in the new model as well as the geographical distance.



26 **Keywords:** anisotropy; geographical information system modeling; geographically weighted
27 logistic regression; mineral resource assessment; missing data; trend variable; weights of
28 evidence.

29

30 **1 Introduction**

31 The main distinguishing characteristic of spatial statistics compared with classical statistics is
32 that the former has a location attribute. Before the development of geographical information
33 systems, spatial statistical problems were often transformed into general statistical problems,
34 where the spatial coordinates were more like a sample ID because they only had an indexing
35 feature. However, even in non-spatial statistics, the reversal paradox or amalgamation paradox
36 (Pearson et al., 1899; Yule, 1903; Simpson, 1951), which is commonly called Simpson's
37 paradox (Blyth, 1972), has attracted much attention from statisticians and other researchers.
38 In spatial statistics, some spatial variables usually exhibit certain trends and non-stationarity.
39 Thus, it is possible for Simpson's paradox to occur when a global regression model is applied
40 and the existence of unknown important variables may make this condition even worse. The
41 influence of Simpson's paradox can be fatal. For example, due to the presence of cover and
42 other factors that occur after mineralization, the ore-forming elements in Area I are generally
43 much lower than those in Area II, but the actual probability of a mineral in Area I is higher
44 than that in Area II, and more deposits may be discovered in Area I (Agterberg, 1971). In this
45 case, a negative correlation will be obtained between the ore-forming elements and the
46 mineralization according to the classical regression model, whereas a high positive correlation
47 can be obtained in both areas if they are separated. Simpson's paradox is an extreme case of
48 the bias caused by using a global model and it is usually not so severe in practice. However,
49 this type of biased needs to be considered and we should take care when applying a classical
50 regression model to a spatial problem. Several solutions to this issue have been proposed



51 previously, which can be divided into three types.

52 (1) Locations are introduced as direct or indirect independent variables. Several studies
53 have employed spatial trend variables (Agterberg, 1964; Agterberg and Cabilio, 1969;
54 Agterberg, 1970; Agterberg and Kelly, 1971; Agterberg, 1971) to express linear or nonlinear
55 trends in space by adding coordinate variables or their functions in predictive models. In these
56 methods, the locations themselves are taken as independent variables as well as the normal
57 independent variables. For example, Reddy et al. (1991) performed logistic regression by
58 including trend variables for mapping the base-metal potential in the Snow Lake area,
59 Manitoba, Canada. In addition, Casetti (1972) developed a spatial expansion method where
60 the regression parameters are themselves functions of the x and y coordinates as well as their
61 combinations.

62 (2) Using local models to replace global models, i.e., geographically weighted models
63 (Fotheringham et al., 2002). Geographically weighted regression (GWR) is the most popular
64 model among the geographically weighted models. GWR was first developed at the end of the
65 20th century by Brunson et al. (1996) and Fotheringham et al. (1996, 1997, 2002) for
66 modeling spatially heterogeneous processes, and it has been used widely in the field of
67 geography.

68 (3) Reducing the trends in spatial variables. For example, Cheng developed a local
69 singularity analysis technique and spectrum-area (S-A) model based on fractal/multi-fractal
70 theory (Cheng, 1997; Cheng, 1999). These methods can remove spatial trends and prevent the
71 strong effects of the original high and low values of the variables on predictions, and thus they
72 are used widely to weaken the effect of spatial non-stationarity to some degree (e.g., Zuo et
73 al., 2016; Zhang et al., 2016; Xiao et al., 2017).

74 GWR can be readily visualized and understood, and it is particularly valid for dealing
75 with spatial non-stationarity, thus it has been used widely in geography and other areas that



76 require spatial data analysis. In general, GWR is a moving window-based model where
77 instead of establishing a unique and global model for prediction, it makes a prediction for
78 each current location using the surrounding samples, and a higher weight is given when the
79 sample is located closer. The theoretical foundation of GWR is based on Tobler's observation
80 that: "everything is related to everything else, but near things are more related than distant
81 things" (Tobler, 1970). In mineral prospectivity mapping (MPM), the dependent variables
82 are binary and logistic regression is used instead of linear regression, and it is necessary to
83 apply geographically weighted logistic regression (GWLR) instead. GWLR belongs to
84 geographically weighed generalized linear regression model (Fotheringham et al. 2002) and it
85 is included in the software module GWR 4.09 (Nakaya, 2016). However, GWLR can only
86 deal with the data in the form of a tabular dataset containing the fields of dependent and
87 independent variables, and the x-y coordinates. Therefore, the spatial layers must be
88 re-processed into two-dimensional tables and the resulting data needs to be transformed back
89 into a spatial form. Another problem with the application of GWR 4.09 for MPM is that it
90 cannot deal with missing data (Nakaya, 2016). Weights of evidence (WofE) is a widely used
91 model for MPM (Bonham-Carter et al., 1988, 1989; Agterberg, 1989; Agterberg et al., 1990),
92 which can avoid the effect of missing data. However, WofE was developed based on the
93 premise that an assumption of conditional independence is satisfied among the evidential
94 layers with respect to the target layer; otherwise, the posterior probabilities will be biased and
95 the number of estimated deposits will not be equal to the known deposits. Agterberg (2011)
96 combined WofE with logistic regression and proposed a new model that can obtain an
97 unbiased estimated of the number of deposits as well as avoiding the effect of missing data. In
98 the present study, this concept is employed to deal with missing data and we propose the
99 improved logistic regression model based on spatially weighted technique (ILRBSWT
100 v1.0) for MPM. The main features of ILRBSWT include the following: (1) a spatial



101 t -statistics method (Agterberg et al., 1993) is introduced to determine the best binary threshold
102 for independent variables, where binarization is performed based on a local window instead of
103 the global level, which can increase the effect of indicating the independent variables to the
104 target variable; and (2) a mask layer is included in the new model to deal with the data quality
105 and exploration level differences among samples.

106 The idea of this research is origin from the first author's doctoral thesis (Zhang, 2015)
107 in Chinese, which has been shown to have better efficiency for mapping intermediate and
108 felsic igneous rocks (Zhang et al., 2017). The contribution of this research is to elaborate
109 the principle of ILRBSWT, and provide a detailed algorithm for its design and
110 implementation process with the code and software module attached. In addition, the
111 processing of missing data is not covered by former researches. At last, the prediction of
112 Au ore deposits in western Meguma Terrain, Nova Scotia, Canada, is chosen as case study
113 to show the performance of ILRBSWT in MPM.

114

115 **2 Models**

116 Linear regression is commonly used for exploring the relationship between a response
117 variable and one or more explanatory variables. However, in MPM and other fields, the
118 response variable is binary or dichotomous, so linear regression is not applicable and thus a
119 logistic model can be advantageous.

120 *2.1 Logistic Regression*

121 In MPM, the dependent variable(Y) is binary since Y can only take the value of 1 and 0,
122 which means the mineralization occurs or not. Suppose that π represents the estimation of Y ,
123 $0 \leq \pi \leq 1$, then a logit transformation of π can be made, i.e., $\text{logit}(\pi) = \ln(\pi/(1-\pi))$. Logistic
124 regression function can be obtained as following.

$$125 \text{Logit } \pi(X_1, X_2, \dots, X_p) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \quad (1)$$



126 where X_1, X_2, \dots, X_p , comprises a sample of p explanatory variables x_1, x_2, \dots, x_p , β_0 is the
 127 intercept, and $\beta_1, \beta_2, \dots, \beta_p$ are regression coefficients.

128 If there are n samples, we can obtain n linear equations with $p+1$ unknowns based on
 129 equation (1). Furthermore, if we suppose that the observed values for Y are Y_1, Y_2, \dots, Y_n , and
 130 these observations are independent of each other, then a likelihood function can be
 131 established:

$$132 \quad L(\beta) = \prod_{i=1}^n (\pi_i^{Y_i} (1 - \pi_i)^{1-Y_i}), \quad (2)$$

133 where $\pi_i = \pi(X_{i1}, X_{i2}, \dots, X_{ip}) = \frac{e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}}}$. The best estimate can be obtained if
 134 and only if equation (2) takes the maximum. Then the problem is converted into solving
 135 $\beta_1, \beta_2, \dots, \beta_p$. Equation (2) can be further transformed into the following log-likelihood
 136 function.

$$137 \quad \ln L(\beta) = \sum_{i=1}^n (Y_i \pi_i + (1 - Y_i)(1 - \pi_i)) \quad (3)$$

138 The solution can be obtained by taking the first partial derivative of β_i ($i = 0$ to p),
 139 which should be equal to 0.

$$140 \quad \begin{cases} f(\beta_0) = \sum_{i=0}^n (Y_i - \pi_i) X_{i0} = 0 \\ f(\beta_1) = \sum_{i=0}^n (Y_i - \pi_i) X_{i1} = 0 \\ \vdots \\ f(\beta_p) = \sum_{i=0}^n (Y_i - \pi_i) X_{ip} = 0 \end{cases} \quad (4)$$

141 where $X_{i0} = 1$, i takes the value from 1 to n , and equation (4) is obtained in the form of
 142 matrix operations.

$$143 \quad \mathbf{X}^T(\mathbf{Y} - \boldsymbol{\pi}) = \mathbf{0} \quad (5)$$

144 The Newton iterative method can be used to solve the nonlinear equations:

$$145 \quad \hat{\boldsymbol{\beta}}(t+1) = \hat{\boldsymbol{\beta}}(t) + \mathbf{H}^{-1} \mathbf{U}, \quad (6)$$

146 where $\mathbf{H} = \mathbf{X}^T \mathbf{V}(t) \mathbf{X}$, $\mathbf{U} = \mathbf{X}^T(\mathbf{Y} - \boldsymbol{\pi}(t))$, t represents the number of iterations, and $\mathbf{V}(t)$, \mathbf{X} ,
 147 \mathbf{Y} , $\boldsymbol{\pi}(t)$, and $\hat{\boldsymbol{\beta}}(t)$ are obtained as follows:



$$\begin{aligned}
 148 \quad \mathbf{V}(t) &= \begin{pmatrix} \pi_1(t)(1 - \pi_1(t)) & & & \\ & \pi_2(t)(1 - \pi_2(t)) & & \\ & & \ddots & \\ & & & \pi_n(t)(1 - \pi_n(t)) \end{pmatrix}, \\
 149 \quad \mathbf{X} &= \begin{pmatrix} X_{10} & X_{11} & \cdots & X_{1p} \\ X_{20} & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n0} & X_{n1} & \cdots & X_{np} \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \boldsymbol{\pi}(t) = \begin{pmatrix} \pi_1(t) \\ \pi_2(t) \\ \vdots \\ \pi_n(t) \end{pmatrix}, \text{ and } \hat{\boldsymbol{\beta}}(t) = \begin{pmatrix} \hat{\beta}_1(t) \\ \hat{\beta}_2(t) \\ \vdots \\ \hat{\beta}_n(t) \end{pmatrix}.
 \end{aligned}$$

150 Hosmer et al. (2013) provided more information about the derivation from equations (1) to
 151 (6).

152 2.2 Weighted Logistic Regression

153 In practice, vector data is popularly used, and sample size (area) has to be considered. In this
 154 condition, weighted logistic regression modeling should be used instead of general logistic
 155 regression. In addition, it is preferable to use a weighted logistic regression model when a
 156 logical regression should be performed for large sample data, since weighted logical
 157 regression can greatly reduce the size of the matrix and improve the computational efficiency
 158 (Agterberg, 1992). Assuming that there are four binary explanatory variable layers and the
 159 study area consists of 1000×1000 grid points, the matrix size for normal logic regression
 160 modeling would be $10^6 \times 10^6$; however, if weighted logistic regression is used, the matrix size
 161 would be 32×32 at most. That is because sample classification process is contained in
 162 weighted logistic regression, and all samples are classified into the classes which own the
 163 same values at dependent and each independent variables. The samples with the same
 164 dependent and independent variables form certain continuous and discontinuous patterns in
 165 space, which are called “unique condition” units. Each unique condition unit is then treated as
 166 a sample, and the area (grid number) for it is taken as weight in weighed logistic regression.
 167 Thus, in the case of weighted logical regression, equations (2) to (5) in section 2.1 need to be
 168 changed as following Equations (7) to (10) respectively.

169



$$170 \quad L_{new}(\beta) = \prod_{i=1}^n (\pi_i^{N_i Y_i} (1 - \pi_i)^{N_i(1-Y_i)}), \quad (7)$$

$$171 \quad \ln L_{new}(\beta) = \sum_{i=1}^n (N_i Y_i \pi_i + N_i(1 - Y_i)(1 - \pi_i)) \quad (8)$$

$$172 \quad \begin{cases} f_{new}(\beta_0) = \sum_{i=0}^n (Y_i - \pi_i) X_{i0} = 0 \\ f_{new}(\beta_1) = \sum_{i=0}^n (Y_i - \pi_i) X_{i1} = 0 \\ \vdots \\ f_{new}(\beta_p) = \sum_{i=0}^n (Y_i - \pi_i) X_{ip} = 0 \end{cases} \quad (9)$$

$$173 \quad \mathbf{X}^T \mathbf{W} (\mathbf{Y} - \boldsymbol{\pi}) = \mathbf{0} \quad (10)$$

174 where N_i is the weight for the i -th unique condition unit, i takes the value from 1 to n , and n
175 is the total number of grid points. And \mathbf{W} is a diagonal matrix which can be expressed as
176 following.

$$\mathbf{W} = \begin{pmatrix} N_1 & & & \\ & N_2 & & \\ & & \ddots & \\ & & & N_n \end{pmatrix}$$

177 Besides, new \mathbf{H} and \mathbf{U} should be used in equation (6) to perform Newton iterative under
178 weighted logistic regression, i.e., $\mathbf{H}_{new} = \mathbf{X}^T \mathbf{W} \mathbf{V}(t) \mathbf{X}$, $\mathbf{U}_{new} = \mathbf{X}^T \mathbf{W} (\mathbf{Y} - \boldsymbol{\pi}(t))$.

179 2.3 Geographically Weighted Logistic Regression

180 GWLR is a local window-based model because logistic regression is established at each
181 current location in GWLR. The current location is changed using the moving window
182 technique with a loop program. If we suppose that \mathbf{u} represents the current location, which
183 can be uniquely determined by a pair of column and row numbers, \mathbf{x} denotes that p
184 explanatory variables x_1, x_2, \dots, x_p take values of X_1, X_2, \dots, X_p , respectively, and $\pi(\mathbf{x}, \mathbf{u})$
185 is the estimates of Y , i.e., the probability that Y takes a value of 1, then the following function
186 can be obtained.

$$187 \quad \text{Logit } \pi(\mathbf{x}, \mathbf{u}) = \beta_0(\mathbf{u}) + \beta_1(\mathbf{u})X_1 + \beta_2(\mathbf{u})X_2 + \dots + \beta_p(\mathbf{u})X_p, \quad (11)$$

188 where $\beta_0(\mathbf{u})$, $\beta_1(\mathbf{u})$, \dots , $\beta_p(\mathbf{u})$ denote that these parameters are obtained at the location of
189 \mathbf{u} . The predicted probability for the current location can be obtained under the condition that
190 the values of all the independent variables are known at the current location and all of the



191 parameters are also calculated based on the samples within the current local window.
192 According to equation (6) in section 2.1, the parameters for GWLR can be estimated with
193 equation (12):

$$194 \hat{\boldsymbol{\beta}}(\mathbf{u})_{t+1} = \hat{\boldsymbol{\beta}}(\mathbf{u})_t + (\mathbf{X}^T \mathbf{W}(\mathbf{u}) \mathbf{V}(t) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\mathbf{u}) (\mathbf{Y} - \boldsymbol{\pi}(t)), \quad (12)$$

195 where t represents the number of iterations; \mathbf{X} is a matrix comprising the values of all the
196 independent variable, and all of the elements in the first column are 1; $\mathbf{W}(\mathbf{u})$ is a diagonal
197 matrix where the diagonal elements are geographical weights, which can be calculated
198 according to distance, whereas the other elements are all 0; $\mathbf{V}(t)$ is also a diagonal matrix
199 and the diagonal element can be expressed as $\pi_i(t)(1 - \pi_i(t))$; and \mathbf{Y} is a column vector
200 representing the values taken by the dependent variable.

201 *2.4 Improved Logistic Regression Model based on Spatially Weighted Technique*

202 If a diagonal element in $\mathbf{W}(\mathbf{u})$ is only for one sample (grid point in raster data), section 2.3
203 can be seen as the improvement of section 2.1, i.e. samples are weighted according to its
204 location. If samples are reclassified firstly according to unique condition mentioned in section
205 2.2, and corresponding weights are then summarized according to each sample's geographical
206 weight, we can obtain an improved logistic regression model considering both sample sizes
207 and geographical distances. The new model can not only reflects the spatial distribution of
208 samples, but also reduce the matrix size, and it is to be discussed in following section.

209 In addition to geographic factors, the degree considered in the study can affect the
210 representativeness of a sample, e.g., differences in the level of exploration.

211 Suppose that there are n grid points in the current local window, S_i is the i -th grid, $W_i(g)$
212 is the geographical weight of S_i , and $W_i(d)$ represents the individual difference weight or
213 non-geographical weight (in some cases, there may be differences in quality or the
214 exploration level among samples, but $W_i(d)$ takes a value of 1 if there is no difference),
215 where i takes a value from 1 to n . Furthermore, if we suppose that there are N unique



216 conditions after overlaying all of the layers ($N \leq n$) and C_j denotes the j -th unique condition
217 unit, then we can obtain the final weight for each unique condition unit in the current local
218 window:

$$219 \quad W_j(t) = \sum_{i=1}^n [W_i(g) * W_i(d) * df_i], \quad (13)$$

220 where $\begin{cases} df_i = 1 & \text{if } S_i \in C_j \\ df_i = 0 & \text{if } S_i \notin C_j \end{cases}$, i takes a value from 1 to n , j takes a value from 1 to N , and

221 $W_j(t)$ represents the total weight (by combining both $W_i(g)$ and $W_i(d)$) for each unique
222 condition unit. We can use the final weight calculated in equation (13) to replace the original
223 weight in equation (12), which is one of the advantages of ILRBSWT.

224 *2.5 Missing data processing*

225 Missing data is a problem existing in all statistics-related research fields. In MPM, missing
226 data are also prevalent due to ground coverage, and limitations of exploration technique and
227 measurement accuracy. Agterberg and Bonham-Carter (1999) once compared following
228 commonly used missing data processing solutions: (1) removing variables containing missing
229 data, (2) deleting samples with missing data, (3) using 0 to replace the missing data, and (4)
230 replacing the missing data with the mean of the corresponding variable. From the point of
231 utilization efficiency of existing data, both (1) and (2) are clearly not good solutions since
232 more data will be lost. Solution (3) is superior to (4) for missing values due to the detection
233 limit of the measuring instrument; with respect to the missing data caused by the limitation of
234 geographical environment and the prospecting technique, solution (4) is obviously a better
235 choice. Missing data is mainly caused by the latter in MPM, and Agterberg (2011) pointed out
236 that missing data could be even better dealt with by performing WofE model. In WofE, the
237 evidential variable takes the value of positive weight (W^+) if it is favorable for the happening
238 of the target variable (e.g., mineralization); and the evidential variable takes the value of
239 negative weight (W^-) if it is unfavorable for the happening of the target variable; and



240 automatically the evidential variable takes the value of 0 if missing data happens.

$$241 \quad W^+ = \ln \frac{\frac{D_1}{D}}{\frac{A_1 - D_1}{A - D}} \quad (14)$$

$$242 \quad W^- = \ln \frac{\frac{D_2}{D}}{\frac{A_2 - D_2}{A - D}} \quad (15)$$

243 where A is an evidential layer, A1 means the area that A takes the value of 1, and A2 means
244 the area that A takes the value of 0; A3 means the area with missing data, and A1+A2 is
245 smaller than the total study area if missing data exists. D1, D2 and D3 are the area that the
246 target variable takes the value of 1 in A1, A2 and A3 respectively. In fact, A3 and D3 are not
247 used in equation (15) since no information is provided in area A3.

248 However, it is preferred to use 1 and 0 to represent the positive and negative condition of
249 the independent variable in logistic regression model. In this case, equation (16) can be used
250 to replace missing data in logistic regression modeling, which will cause an equivalent effect
251 just as missing data are processed in WofE.

$$252 \quad M = \frac{-W^-}{W^+ - W^-} = \frac{\ln \frac{D}{A-D} - \ln \frac{D_2}{A_2 - D_2}}{\ln \frac{D_1}{A_1 - D_1} - \ln \frac{D_2}{A_2 - D_2}} \quad (16)$$

253

254 **3 Design of the ILRBSWT Algorithm**

255 *3.1 Local Window Design*

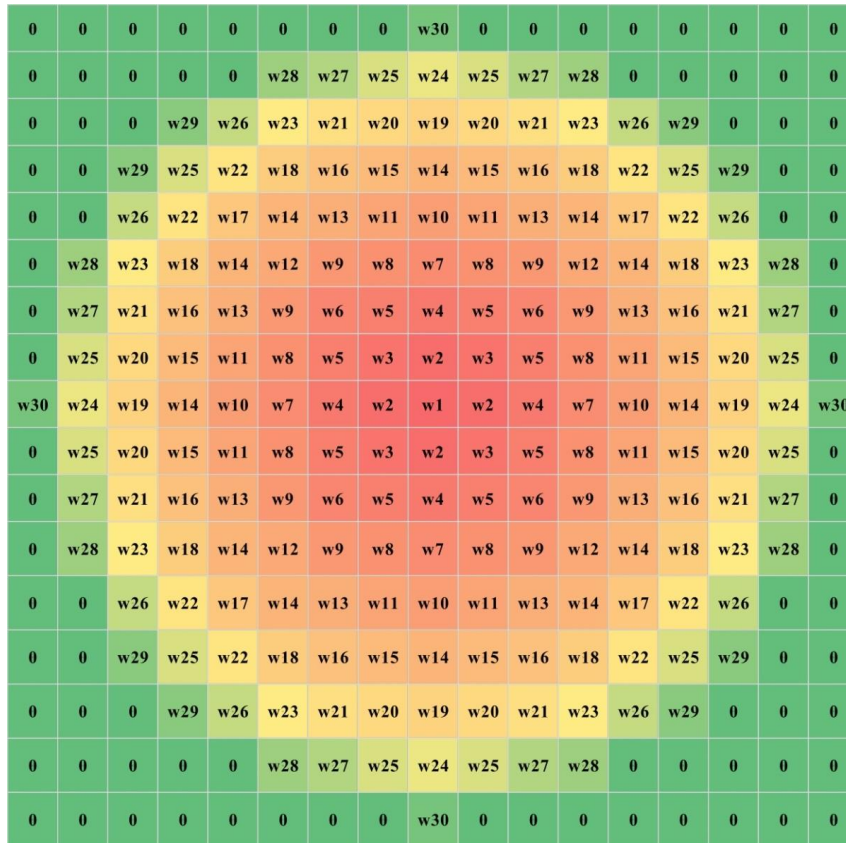
256 A raster data set is used for ILRBSWT modeling. With a regular grid, the distance between
257 any two grid points can be calculated easily and we can even obtain distance templates within
258 a certain window scope, which is highly efficient for data processing. The circle and ellipse
259 are used for isotropic and anisotropic local window designs, respectively.

260 (1) Circular Local Window Design

261 If we suppose that W represents a local circular window where the minimum bounding
262 rectangle is R , then the geographical weights can be calculated only inside R . Obviously, the



263 grid points inside of R but outside of W should be weighted as 0, and the weights for grid
264 points inside W should be calculated according to the distances between themselves and the
265 current location. R should be a square so we can also assume that there are n columns and
266 rows in R , where n is an odd number. If we take east and south as the orientations of the x -axis
267 and y -axis, respectively, and the position of the northwest corner grid is defined as $(x = 1, y =$
268 $1)$, then a local rectangular coordinate system can be established and the position for the
269 current location grid can be expressed as $O(x = \frac{n+1}{2}, y = \frac{n+1}{2})$. The distance between any
270 grid inside W and the current location grid can be expressed as
271 $d_{o-ij} = \sqrt{\left(i - \frac{n+1}{2}\right)^2 + \left(j - \frac{n+1}{2}\right)^2}$, where i and j take values ranging from 1 to n . The
272 geographical weight is a function of distance, so it is convenient to calculate w_{ij} with
273 d_{o-ij} . Figure 1 shows the weight template for a circular local window with a half-window
274 size of nine grid points.



275

276

277

278

279

Fig. 1 Weight template for a circular local window with a half-window size of nine grid points, where w1 to w30 represent different weight classes that decrease with distance and 0 denotes that the grid is weighted as 0. Gradient colors ranging from red to green are used to distinguish the weight classes for grid points.

280

281

282

283

284

285

286

If we suppose that there are T_n columns and T_m rows in the study area, and *Current* (T_i, T_j) represents the current location, where T_i takes values from 1 to T_n and T_j takes values from 1 to T_m , then the current local window can be established by selecting the range of rows $T_i - \frac{n-1}{2}$ to $T_i + \frac{n-1}{2}$ and columns $T_j - \frac{n-1}{2}$ to $T_j + \frac{n-1}{2}$ based on the total research area. Next, we establish a local rectangular coordinate system according to the steps in the last paragraph, where the x and y coordinates for the northwest corner are defined as the coordinate origin by subtracting $T_i - \frac{n-1}{2}$ and $T_j - \frac{n-1}{2}$ from the x and y coordinates,



287 respectively, for all of the grid points in the range. The corresponding relationship can then be
288 established between the weight template and the current window. Global weights can also be
289 included via the matrix product between the global weight layer and local weight template
290 within the local window. In addition, special care should be taken when the weight template
291 covers some area outside the study area, e.g., $T_{.i} - \frac{n-1}{2} < 0$, $T_{.i} + \frac{n-1}{2} > T_{.n}$, $T_{.j} - \frac{n-1}{2} <$
292 0 , and $T_{.j} + \frac{n-1}{2} > T_{.m}$.

293 (2) Elliptic Local Window Design

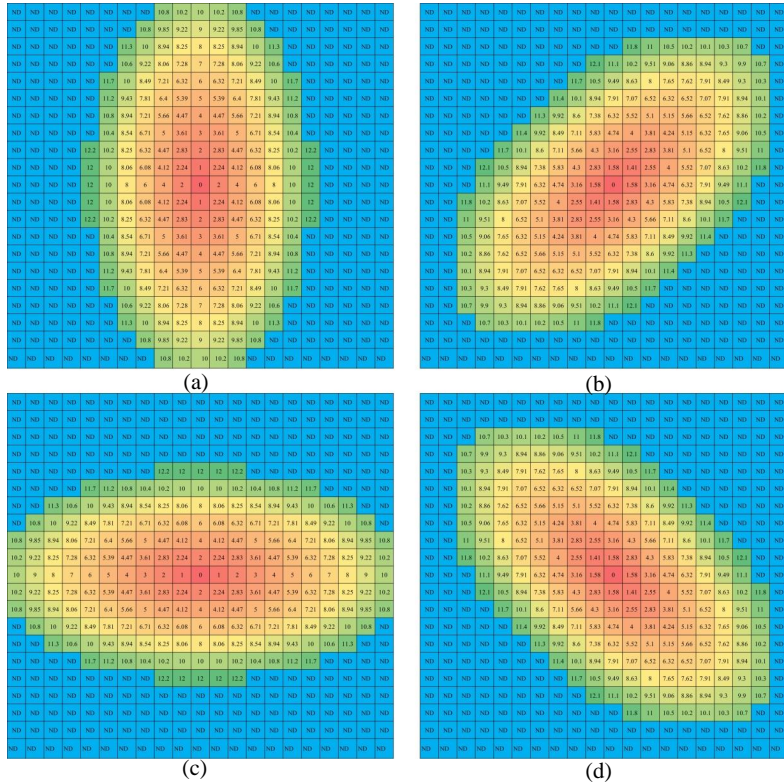
294 In most cases, the spatial weights change to variable degrees in different directions and
295 an elliptic local window may be better for describing the changes in the weights in space. In
296 order to simplify the calculation, we can convert the distances in different directions into
297 equivalent distances and an anisotropic problem then becomes an isotropic problem. For any
298 grid, the equivalent distance is the semi-major axis length of the ellipse that passes through
299 the grid and that is centered at the current location, where the parameters for the ellipse can be
300 determined using the kriging method.

301 We still use W to represent the local elliptic window and a , r , and θ are defined as the
302 semi-major axis, the ratio of the semi-minor axis relative to the semi-major axis, and the
303 azimuth of the semi-major axis, respectively. Then, W can be covered by a square R , where
304 the side length is $2a-1$ and the center is the same as W . There are $(2a-1) \times (2a-1)$ grid
305 points in R . We establish the rectangular coordinates as described above and we suppose that
306 the center of the top left grid in R is located at $(x=1, y=1)$, and thus the center of W should
307 be $O(x_0=a, y_0=a)$. According to the definition of the ellipse, two of the elliptical foci
308 are located at $F_1(x_1 = a + \sin(\theta)\sqrt{a^2 - (a*r)^2}, y_1 = a - \cos(\theta)\sqrt{a^2 - (a*r)^2})$ and
309 $F_2(x_2 = a - \sin(\theta)\sqrt{a^2 - (a*r)^2}, y_2 = a + \cos(\theta)\sqrt{a^2 - (a*r)^2})$. The summed
310 distances between a point and the two focus points can be expressed as



311 $l_{ij} = \sqrt{(i - x_1)^2 + (j - y_1)^2} + \sqrt{(i - x_2)^2 + (j - y_2)^2}$, where i and j take values from 1 to
312 $2a - 1$. According to the elliptical focus formula, we can decide whether a grid in R is located
313 in W . For any grid in R , if the sum of the distances between the two focal points and a grid
314 center is greater than $2a$, then the grid is located in W , vice versa. For the grid points outside
315 of W , the weight is assigned as 0, and the equivalent distances should be calculated for the
316 grid points within W . As mentioned above, the parameters for the ellipse can be determined
317 using the kriging method. In the ellipse W where the semi-major axis is a , we keep r and θ
318 as constants, so we can obtain countless ellipses centered at the center of W , and the
319 equivalent distance is the same on the same elliptical orbit. Thus, the equivalent distance
320 template can be obtained for the elliptic local window. Figure 2 shows the equivalent distance
321 templates under the conditions that $a = 11$ grid points, $r = 0.5$, and the azimuths for the
322 semi-major axis are 0° , 45° , 90° , and 135° , where the weight template can also be calculated
323 based on Fig. 2.

324



325

326 **Fig. 2 Construction of the distance template based on an elliptic local window: $a = 11$ grid points,**

327 **$r = 0.5$, and the azimuths for the semi-major axis are 0° (a), 45° (b), 90° (c), and 135° (d).**

328 **3.2 Pseudocode for ILRBSWT**

329 The ILRBSWT method focuses mainly on two problems, i.e., spatial non-stationarity and
 330 missing data. We use the moving window technique to establish a local model, which can
 331 overcome the spatial non-stationarity better compared with the global model. The spatial
 332 t -value employed in the WoE method is used to binarize spatial variables based on the local
 333 window, which is quite different from binarization based on the global range, where the
 334 missing data can be handled well because positive and negative weights are used instead of
 335 the original “1” and “0” values, and the missing data can then be represented well as “0.”
 336 Both the isotropy and anisotropy window types are possible in our new proposed model. The
 337 geographical weights and the window size can be determined by the users themselves. If the



338 geographic weights are equal and there are no missing data, then ILRBSWT will yield the
339 same posterior probabilities as logistic regression; hence, the later can be treated as a special
340 case of the former. The core ILRBSWT algorithm is as follows.

341 Step 1. Establish a loop for all of the grid points in the study area according to both the
342 columns and rows. Determine a basic local window with a size of r_{\min} based on a variation
343 function or other method. In addition, the maximum local window with a size of r_{\max} is set,
344 with an interval of ΔR . If we suppose that a geographical weight model has already been
345 given in the form of a Gaussian curve determined by variations in the geostatistics, i.e.,
346 $W(g) = e^{-\lambda d^2}$, where d is the distance and λ is the attenuation coefficient, then we can
347 calculate the geographical weight for any grid in the current local window. The equivalent
348 radius should be used in the anisotropic situation. When other types of weights are considered,
349 e.g., the degree of exploration or research, it is also necessary to synthesize the geographical
350 weights and other weights (see equation 10).

351 Step 2. Establish a loop for all of the independent variables. In a circular (elliptical)
352 window with a radius (equivalent radius) of r_{\min} , apply the WofE (Agterberg, 1992) model
353 according to the grid weight determined in step 1, thereby obtaining a statistical table
354 containing the parameters of W_{ij}^+ , W_{ij}^- , and t_{ij} , where i is the i -th independent variable and
355 j denotes the j -th binarization.

356 Step 2.1. If a maximum t_{ij} exists and it is greater than or equal to the standard t -value
357 (e.g., 1.96), record the values of $W_{i-\max_t}^+$, $W_{i-\max_t}^-$, and $B_{i-\max_t}$, which denote the
358 positive weight, negative weight, and corresponding binarization, respectively, under the
359 condition where t takes the maximum value. Go to step 2 and apply the WofE model to the
360 other independent variables.

361 Step 2.2. If a maximum t_{ij} does not exist or it is smaller than the standard t -value, go to
362 step 3.



363 Step 3. In a circular (elliptical) window with a radius (equivalent radius) of r_{\max} , increase
364 the current local window based on r_{\min} according to the algorithm in step 1.

365 Step 3.1. If all of the independent variables have already been processed, go to step 4.

366 Step 3.2. If the size of the current local window exceeds the size of r_{\max} , then disregard
367 the current independent variable and go to step 2 to consider the remaining independent
368 variables.

369 Step 3.3. Apply the WofE model according to the grid weight determined in step 1 in the
370 current local window, which has increased. If a maximum t_{ij} exists and it is greater than or
371 equal to the standard t -value, record the values of $W_{i-\max_t}^+$, $W_{i-\max_t}^-$, $B_{i-\max_t}$, and r_{current} ,
372 which represents the radius (equivalent radius) for the current local window.

373 Step 3.4. If a maximum t_{ij} does not exist or it is smaller than the standard t -value, go to
374 step 3.

375 Step 4. Suppose that n_s independent variables are remaining.

376 Step 4.1. If $n_s \leq 1$, then calculate the mean value for the dependent variable in the
377 current local window with a radius size of r_{\max} and retain it as the posterior probability in the
378 current location. In addition, set the regression coefficients for all of the independent variables
379 as missing data. Go to step 6.

380 Step 4.2. If $n_s \geq 1$, then find the independent variable with the largest local window and
381 apply the WofE model to all the other independent variables, before recording the values of
382 $W_{i-\max_t}^+$, $W_{i-\max_t}^-$, and $B_{i-\max_t}$ for this time, and then go to step 5.

383 Step 5. Apply the logistic regression model based on geographic weights and for each
384 independent variable: (1) use $W_{i-\max_t}^+$ to replace all of the values that are less than or equal
385 to $B_{i-\max_t}$; (2) use $W_{i-\max_t}^-$ to replace all of the values that are greater than $B_{i-\max_t}$; and
386 (3) use 0 to replace no data (“-9999”). The posterior probability and regression coefficients
387 can then be obtained for all of the independent variables at the current location, and go to step



388 6.

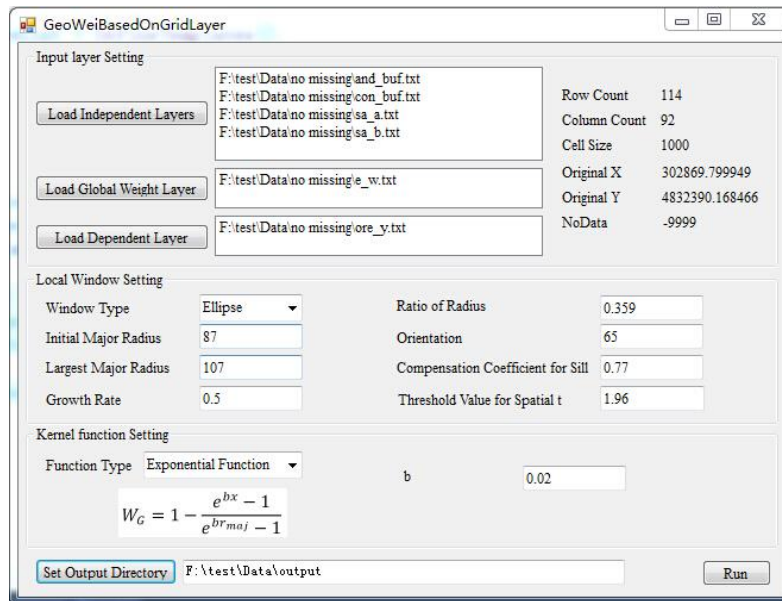
389 Step 6. Take the next grid as the current location and repeat steps 2–5.

390

391 **4 Interface Design**

392 In addition to the improved GWLR, we developed other modeling processes, where all of the
393 visualization and mapping procedures are performed using the ArcGIS 10.2 platform and
394 GeoDAS 4.0 software. The maps are stored in grid format, which are transformed into ASCII
395 files based on tools included in the Arc toolbox before the improved GLWR is performed.

396 As shown in Fig. 3, the main interface for the improved GLWR comprises four parts.
397 The upper left part is for the layer input settings, where independent variable layers,
398 dependent variable layers, and global weight layers should be assigned if they exist. Layer
399 information is shown at the upper right corner, including the row numbers, column numbers,
400 grid size, ordinate origin, and missing data. The local window can be defined in the middle.
401 Using the drop-down list, we can prepare a circle or ellipse to represent various isotropic and
402 anisotropic spatial conditions, respectively. The corresponding window parameters should be
403 set for each window type. For the ellipse, it is necessary to set parameters comprising the
404 initial length of the equivalent radius (initial major radius), the final length of the equivalent
405 radius (largest major radius), the increase in the length of the equivalent radius (growth rate),
406 the threshold of the spatial t -value used to determine the need to enlarge the window, the
407 length ratio of the major and minor axes, the orientation of the ellipse's major axis, and the
408 compensation coefficient for the sill. Next, it is necessary to define the attenuation function
409 and a variety of kernel functions, such as exponential model, logarithmic model, Gaussian
410 model, or spherical model, via the drop-down menu. More parameters can be set when a
411 certain model is selected. The output file settings are defined at the bottom and the execution
412 button is at the lower right corner.



413

414

Fig. 3 User interface design.

415

416 **5 Real Data Testing**

417 *5.1 Data source and preprocessing*

418 The test data used in this study were obtained from the case study reported by Cheng (2008).

419 The study area ($\approx 7780 \text{ km}^2$) was located in western Meguma Terrain, Nova Scotia, Canada.

420 Four independent variables were used in the WofE model for gold mineral potential mapping

421 by Cheng (2008), i.e., buffer of anticline axes, buffer for the contact of Goldenville–Halifax

422 Formation, and background and anomaly separated with the S-A filtering method based on

423 the loadings of the ore elements of the first component. More information about the data set

424 can be found in Cheng (2008).

425 Four independent variables mentioned above were also used for ILRBSWT modeling in

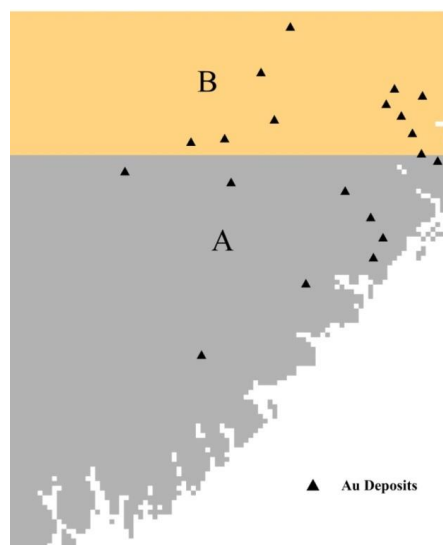
426 this study. In order to demonstrate the advantages of the new method when processing

427 missing data, we designed a situation where the geochemical data were missing for the

428 northern part of the study area, as shown in Fig. 4. In that case, grids in region A own values



429 at all of the four independent variables; however, grids in region B only own values at two
430 independent variables, and they have no values in the two geochemical variables.

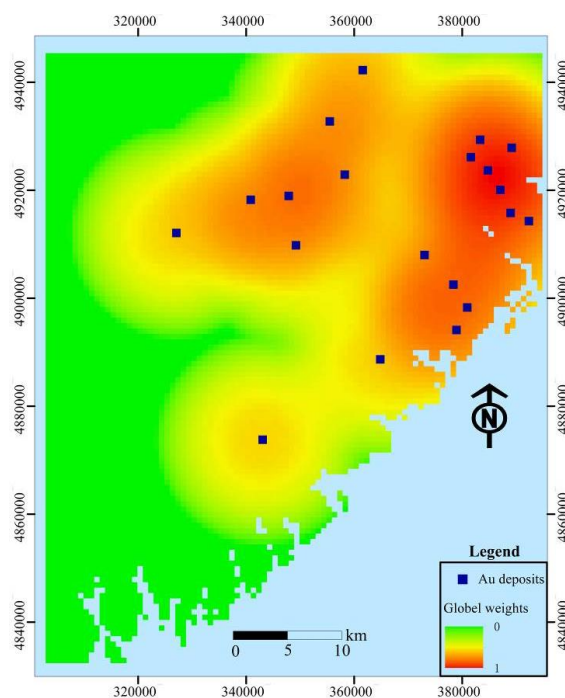


431

432 **Fig. 4 Study area (A and B) and the scope with missing geochemical data (B).**

433 *5.2 Mapping weights for the exploration level*

434 These types of weights can be determined based on prior knowledge according to differences
435 in the exploration data, e.g., different scales may exist throughout the whole study area. They
436 can also be obtained quantitatively. The density of known deposits is a good index for the
437 exploration level, where the degree of research is higher when more deposits are discovered.
438 The exploration level weights for the mapped study area obtained using the kernel density
439 tool provided by the ArcToolbox in ArcGIS 10.2 are shown in Fig. 5.



440

441

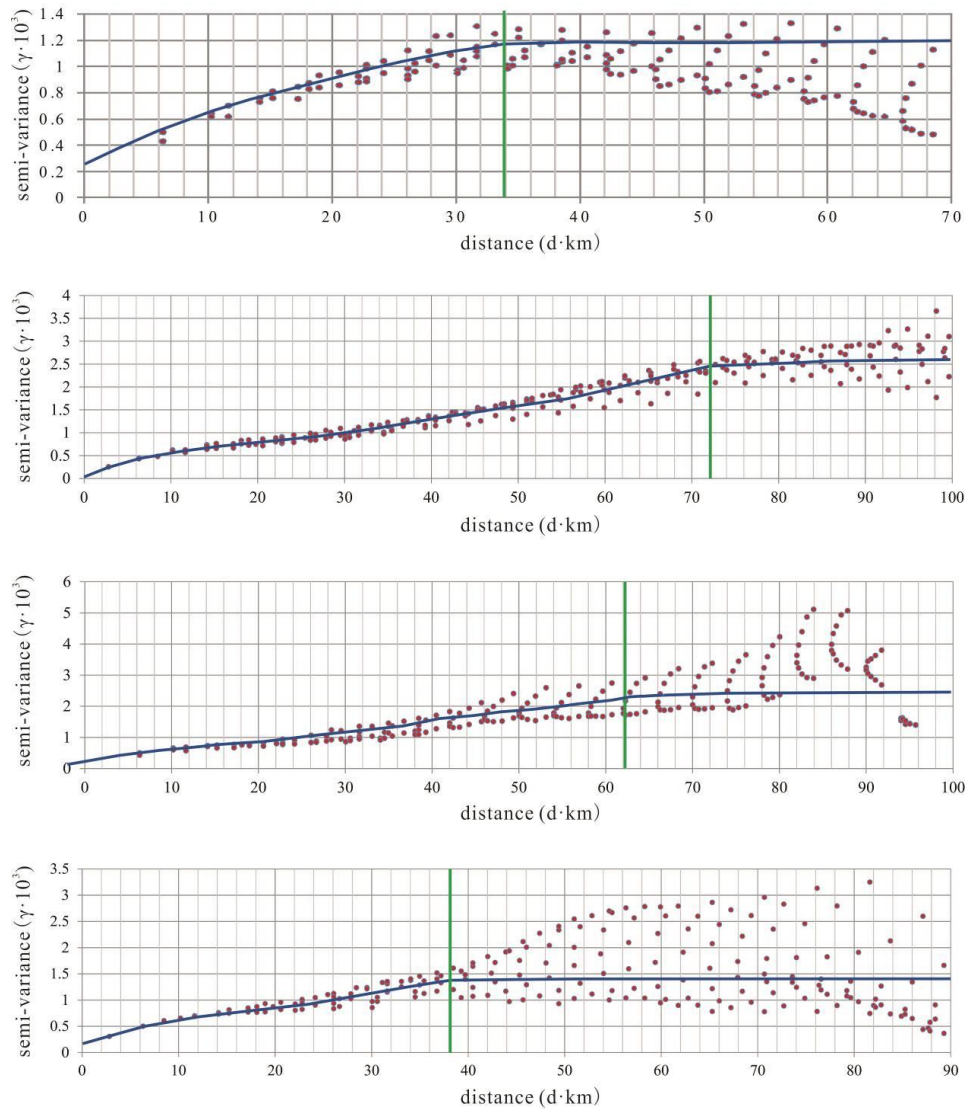
Fig. 5 Exploration level weights.

442 *5.3 Assignment of local window parameters and geographical weights*

443 Empirical and quantitative methods can be used to determine the local window parameters
 444 and the attenuation function for geographical weights. The variation function in geostatistics
 445 is an effective method for describing the structures and trends of spatial variables, so it was
 446 used in this study. In order to calculate the variation function for a dependent variable, it is
 447 necessary to first map the posterior probability using the global logistic regression method,
 448 before establishing the variation function to determine the local window type and parameters.
 449 Variation functions are established in four directions in order to detect anisotropic changes in
 450 space. If there are no significant differences among the various directions, a circular local
 451 window can be used for ILRBSWT, as shown in Fig. 1; otherwise, an elliptic local window
 452 should be used, as shown in Fig. 2. The specific parameters for the local window in the study
 453 area were obtained as shown in Fig. 6, and the final local window and geographical weight



454 attenuation were determined as indicated in Fig. 7 (a) and 7(b), respectively.



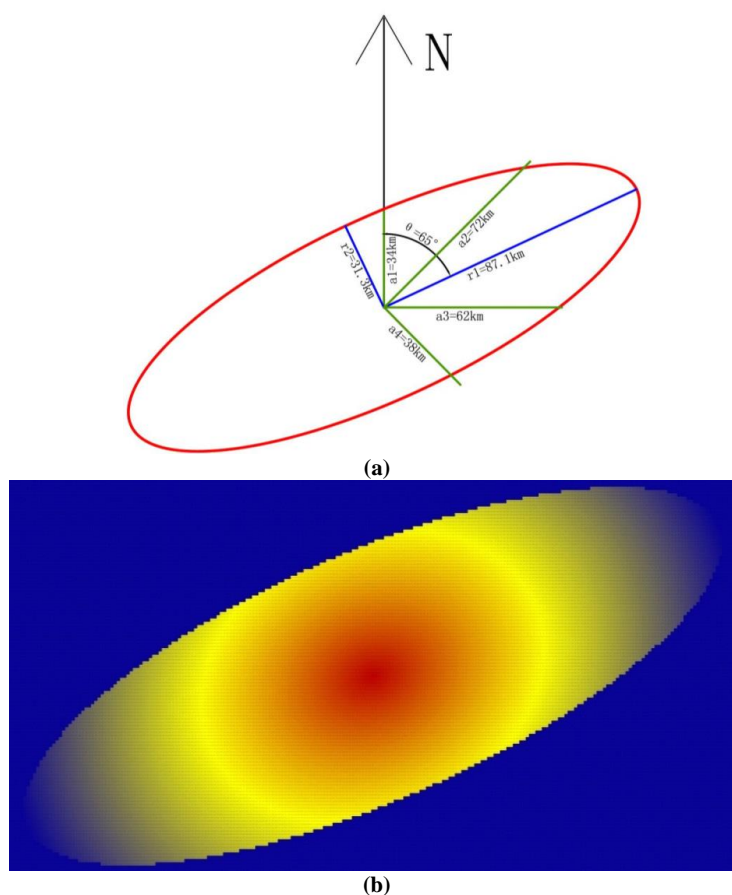
455

456 **Fig. 6** Experimental variogram fitting in different directions, where the green lines denote the

457 variable ranges determined for azimuths of (a) 0°, (b) 45°, (c) 90°, and (d) 135°.

458

459

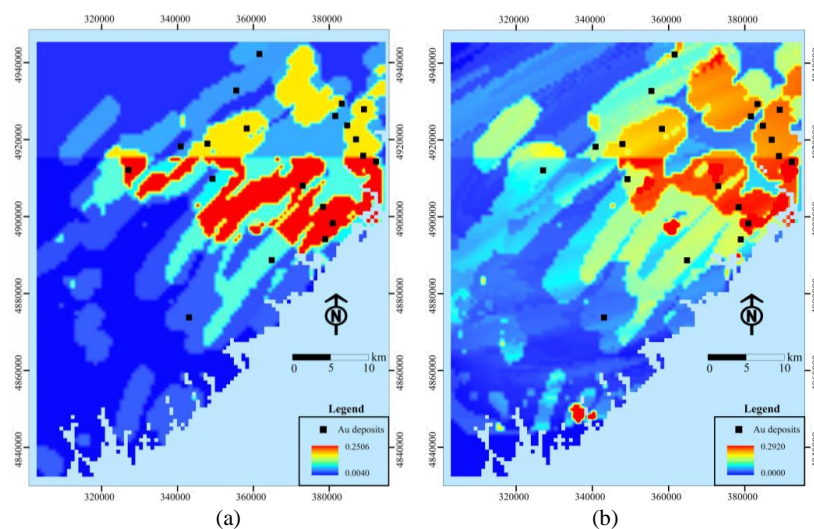


460

461 **Fig. 7 Nested spherical model for different directions. The green lines in (a) correspond to those**
462 **in Fig. 6, and (b) shows the geographical weight template determined based on (a).**

463 5.4 Data integration

464 Using the algorithm described in section 3.2, ILRBSWT was performed for the study area
465 according to the settings in Fig. 3. The estimated probability map obtained for intermediate
466 and felsic igneous rocks by ILRBSWT is shown in Fig. 8 (b), while Fig. 8 (a) presents the
467 results obtained by logistic regression. It can be seen from Fig. 8 that ILRBSWT can better
468 weak the effect of missing data than logistic regression, since the Au deposits in the north part
469 of the study area (where missing data exist) are well felled into the region with relatively
470 higher posterior probability in Fig. 8 (b) than in Fig. 8 (a).

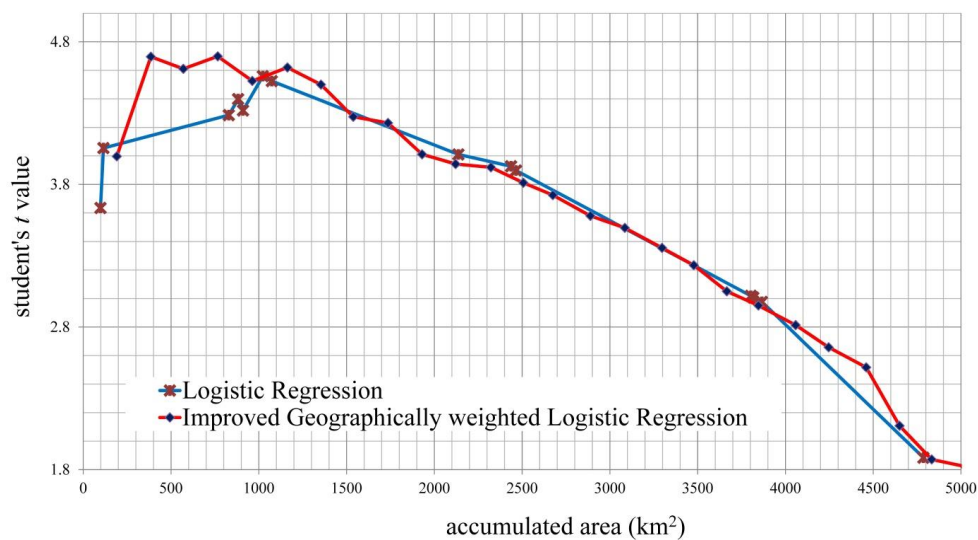


471

472 **Fig. 8 Posterior probability maps obtained for an Au deposit by (a) logistic regression and (b)**
473 **ILRBSWT.**

474 5.5 Comparison of the mapping results

475 In order to evaluate the predictive capacity of the newly developed method and the traditional
476 method, the posterior probability maps obtained by logistic regression and ILRBSWT shown
477 in Fig. 8(b) and 8(a), respectively, were divided into 20 classes by the quantile method and the
478 t -values were then calculated using WofE modeling (Fig. 9). Clearly, ILRBSWT performed
479 better because higher t -values were obtained, especially when a smaller area was delineated as
480 the target area, which is much more realistic. In the northern part of the study area, the known
481 deposits fitted better to the high posterior probability area shown in Fig. 8(b) than that in Fig.
482 8(a), which indicates that ILRBSWT can deal with missing data better than logistic
483 regression.



484

485 **Fig. 9 Student's t -values calculated for the spatial correlation between the known Au deposit**
 486 **layer and the predicted posterior probability layers obtained by logistic regression and ILRBSWT at**
 487 **different threshold levels.**

488

489 6 Conclusions

490 In this study, we developed an improved GWLR model ILRBSWT based on logistic
 491 regression, WofE, and the current GWR model. Furthermore, a software module was
 492 developed for ILRBSWT and a case study demonstrated its capacities and advantages.

493 Following objectives were achieved:

494 (1) A moving window technique is employed for spatial variable-parameter logistic
 495 regression, which can overcome or weaken the effect of spatial non-stationarity in MPM and
 496 improve the accuracy of mineral prediction.

497 (2) The variogram model in geostatistics is used to determine the spatial anisotropic
 498 parameters and geographical weight attenuation model, which makes the local window
 499 parameter design more objective and tenable.

500 (3) The spatial t -statistics method based on WofE is introduced to perform



501 binarization/discretization for the independent variables in each local window, and the new
502 model can better handle missing data.

503 (4) The global weight layer in ILRBSWT can reflect differences in the data quality or
504 exploration level well.

505

506 ***Code availability***

507 The software tool ILRBSWT v1.0 in this research is developed by using C#, and the main
508 codes and key functions are prepared in file “Codes & Key Functions”. The executable
509 program files are placed in the folder “Executable Programs for ILRBSWT”. Please find them
510 in gmd-2017-278-supplement.zip.

511

512 ***Data availability***

513 The data used in this research is sourced from the demo data of GeoDAS software
514 (<http://www.yorku.ca/yul/gazette/past/archive/2002/030602/current.htm>), and this data is also
515 used by Cheng (2008). All spatial layers used in this work is included in the folder “Original
516 Data” in the format of ASCII file, which can be also found in gmd-2017-278-supplement.zip.

517 **Acknowledgments**

518 This study benefited from joint financial support by the Programs of National Natural Science
519 Foundation of China (Nos. 41602336 and 71503200), China Postdoctoral Science Foundation
520 (Nos. 2016M592840 and 2017T100773), Shaanxi Provincial Natural Science Foundation (No.
521 2017JQ7010), and the Fundamental Research Funds for the Central Universities (No.
522 2017RWYB08). The first author thanks former supervisor Drs. Qiuming Cheng and Frits
523 Agterberg for discussions about spatial weights and for providing constructive suggestions.

524

525

526 **References**

- 527 Agterberg, F.P., & Cabilio, P., 1969. Two-stage least-squares model for the relationship between mappable geological
528 variables. *Journal of the International Association for Mathematical Geology*, 1(2), 137-153.
- 529 Agterberg, F.P., & Kelly, A.M., 1971. Geomathematical methods for use in prospecting. *Canadian Mining Journal*, 92(5),
530 61-72.
- 531 Agterberg, F.P., 1964. Methods of trend surface analysis. *Colorado School Mines Quart*, 59(4), 111-130.
- 532 Agterberg, F.P., 1970. Multivariate prediction equations in geology. *Journal of the International Association for*
533 *Mathematical Geology*, 1970 (02), 319-324.
- 534 Agterberg, F.P., 1971. A probability index for detecting favourable geological environments. *Canadian Institute of Mining*
535 *and Metallurgy*, 10, 82-91.
- 536 Agterberg, F.P., 1989. Computer Programs for Mineral Exploration. *Science*, 245, 76 – 81.
- 537 Agterberg, F.P., 1992. Combining indicator patterns in weights of evidence modeling for resource evaluation. *Nonrenewal*
538 *Resources*, 1(1), 35–50.
- 539 Agterberg, F.P., 2011. A Modified WoE Method for Regional Mineral Resource Estimation. *Natural Resources Research*,
540 20(2), 95-101.
- 541 Agterberg, F.P., Bonham-Carter, G.F., & Wright, D.F., 1990. Statistical Pattern Integration for Mineral Exploration. in Gaál,
542 G., Merriam, D. F., eds. *Computer Applications in Resource Estimation Prediction and Assessment of Metals and*
543 *Petroleum*. New York: Pergamon Press: 1-12.
- 544 Agterberg, F.P., Bonham-Carter, G.F., Cheng, Q., & Wright, D.F., 1993. Weights of evidence modeling and weighted
545 logistic regression for mineral potential mapping. *Computers in geology*, 25, 13-32.
- 546 Blyth, C.R., 1972. On Simpson's paradox and the sure-thing principle. *Journal of the American Statistical*
547 *Association*, 67(338), 364-366.
- 548 Bonham-Carter, G.F., Agterberg, F.P., & Wright, D.F., 1988. Integration of Geological Datasets for Gold Exploration in
549 Nova Scotia. *Photogrammetric Engineering & Remote Sensing*, 54(11), 1585-1592.
- 550 Bonham-Carter, G.F., Agterberg, F.P., & Wright, D.F., 1989. Weights of Evidence Modelling: A New Approach to Mapping
551 Mineral Potential. In Agterberg F P and Bonham-Carter G F, eds. *Statistical Applications in the Earth Sciences*, 171-183.
- 552 Brunson, C., Fotheringham, A.S., & Charlton, M.E., 1996. Geographically weighted regression: a method for exploring
553 spatial nonstationarity. *Geographical analysis*, 28(4), 281-298.
- 554 Casetti, E., 1972. Generating models by the expansion method: applications to geographic research. *Geographical Analysis*, 4,
555 81-91.
- 556 Cheng, Q., 1997. Fractal/multifractal modeling and spatial analysis, keynote lecture in proceedings of the international
557 mathematical geology association conference, 1, 57-72.
- 558 Cheng, Q., 1999. Multifractality and spatial statistics. *Computers & Geosciences*, 25, 949–961.
- 559 Cheng, Q., 2008. Non-Linear Theory and Power-Law Models for Information Integration and Mineral Resources
560 Quantitative Assessments. *Mathematical Geosciences*, 40(5), 503-532.
- 561 Fotheringham, A.S., Brunson, C., & Charlton, M.E., 1996. The geography of parameter space: an investigation of spatial
562 non-stationarity. *International Journal of Geographical Information Systems*, 10, 605-627.
- 563 Fotheringham, A.S., Brunson, C., & Charlton, M.E., 2002. *Geographically Weighted Regression: the analysis of spatially*
564 *varying relationships*, Chichester: Wiley.
- 565 Fotheringham, A.S., Charlton, M.E., & Brunson, C., 1997. Two techniques for exploring nonstationarity in geographical
566 data. *Geographical Systems*, 4, 59-82.
- 567 Hosmer, D.W., Lemeshow, S, & Sturdivant, R.X., 2013. *Applied logistic regression*, 3rd edn. Wiley, New York
- 568 Nakaya, T., 2016. GWR4.09 user manual. WWW Document. Available online:
569 https://raw.githubusercontent.com/gwrtools/gwr4/master/GWR4manual_409.pdf (accessed on 16 February 2017).



- 570 Pearson, K., Lee, A., & Bramley-Moore, L., 1899. Mathematical contributions to the theory of evolution. VI. Genetic
571 (reproductive) selection: Inheritance of fertility in man, and of fecundity in thoroughbred racehorses. *Philosophical*
572 *Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 192,
573 257-330.
- 574 Reddy, R.K.T., Agterberg, F.P., & Bonham-Carter, G.F., 1991. Application of GIS-based logistic models to base-metal
575 potential mapping in Snow Lake area, Manitoba. *Proceedings of the Canadian Conference on GIS*, 18-22.
- 576 Simpson, E.H., 1951. The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society. Series B*
577 (Methodological), 238-241.
- 578 Tobler, W.R., 1970. A computer movie simulating urban growth in the Detroit region. *Economic Geography*, 46(2), 234-24.
- 579 Xiao, F., Chen, J., Hou, W., Wang, Z., Zhou, Y., & Erten, O., 2017. A spatially weighted singularity mapping method
580 applied to identify epithermal Ag and Pb-Zn polymetallic mineralization associated geochemical anomaly in Northwest
581 Zhejiang, China. *Journal of Geochemical Exploration*.
- 582 Yule, G.U., 1903. Notes on the theory of association of attributes in statistics. *Biometrika*, 2(2), 121-134.
- 583 Zhang, D., 2015. Spatially Weighted technology for Logistic regression and its Application in Mineral Prospectivity
584 Mapping (Dissertation). China University of Geosciences, Wuhan (in Chinese with English abstract).
- 585 Zhang, D., Cheng, Q., & Agterberg, F.P., 2017. Application of spatially weighted technology for mapping intermediate and
586 felsic igneous rocks in fujian province, china. *Journal of Geochemical Exploration*, 178, 55-66.
- 587 Zhang, D., Cheng, Q., Agterberg, F.P., & Chen, Z., 2016. An improved solution of local window parameters setting for local
588 singularity analysis based on excel vba batch processing technology. *Computers & Geosciences*, 88(C), 54-66.
- 589 Zuo, R., Carranza, E.J.M., & Wang, J., 2016. Spatial analysis and visualization of exploration geochemical
590 data. *Earth-Science Reviews*, 158, 9-18.