



1 **Depth-to-Bedrock Map of China at a Spatial**
2 **Resolution of 100 Meters**

3 Fapeng Yan¹, Wei Shangguan^{2*}, Jing Zhang¹ and Bifeng Hu^{3,4,5}

4 ¹ College of Global Change and Earth System Science, Beijing Normal University, Beijing, China

5 ² Guangdong Province Key Laboratory for Climate Change and Natural Disaster Studies, School
6 of Atmospheric Sciences, Sun Yat-sen University, Guangzhou, China.

7 ³ Unité de Recherche en Science du Sol, INRA, Orléans 45075, France

8 ⁴ InfoSol, INRA, US 1106, Orléans F-4075, France

9 ⁵ Sciences de la Terre et de l'Univers, Orléans University, 45067 Orleans, France

10

11 *Correspondence to:* Wei Shangguan (shgwei@mail.sysu.edu.cn)



Abstract. Depth to bedrock serves as the lower boundary of soil, which influences or controls many of the Earth's physical and chemical processes. It plays important roles in geology, hydrology, land surface processes, civil engineering, and other related fields. This paper describes the materials and methods to produce a high-resolution (100 m) depth-to-bedrock map of China. Observations were interpreted from borehole log data (ca. 6,382 locations) sampled from the Chinese National Important Geological Borehole Database. To fill in large sampling gaps, additional pseudo-observations generated based on expert knowledge were added. Then, we overlaid the training points on a stack of 133 covariates including climatic images, DEM-derived parameters, land-cover and land-use maps, MODIS surface reflectance bands, vegetation index images, and the Harmonized World Soil Database. Spatial prediction models were developed using the random forests and gradient boosting tree, and ensemble prediction results were then obtained by these two independently fitted models. Finally, uncertainty estimation was generated by the quantile regression forest model. The 10-fold cross-validation showed that the ensemble models explain 57% of the variation in depth to bedrock. Based on comparison with depth-to-bedrock maps of China extracted from previous global predictions, our predictions showed higher accuracy. More observations, especially those in data-sparse areas, should be added to training data, and more covariates with high precision should be used to further improve the accuracy of spatial predictions. The resulting maps of this study are available on Figshare at the following DOI: <https://doi.org/10.6084/m9.figshare.7011524.v1>. And they are also available for download at <http://globalchange.bnu.edu.cn/>.

1 Introduction

Soil is the loose layer on the surface of the geosphere. It is the foundation of the whole terrestrial ecosystem (van Breemen and Buurman, 2002). The International Union of Soil Sciences (IUSS) divides the soil profile into six main genetic horizons: O (organic horizon), A (humus horizon), E (eluvial horizon), B (illuvial horizon), C (parent rock horizon), and R (hard rock). Of these, the bedrock (i.e., the R horizon) is the consolidated solid rock underlying unconsolidated surface materials, such as soil or other regolith (Jain, 2014). Depth to bedrock (DTB) is the depth to the R horizon, which is equivalent to the total thickness of the solum and weathered rocks; DTB controls or influences many physical and chemical processes of the Earth (Jain, 2014).



41 DTB information plays an important role in many fields of Earth system science. In geology,
42 DBT has been used for applications such as mineral exploration, earthquake modeling, and landslide
43 risk assessment (Schenk and Jackson, 2005; Fan et al., 2013). In land surface modeling, DTB is an
44 important input parameter that affects the energy, water, and carbon cycles. However, in most land
45 surface models, DTB has been set as a constant value because of a lack of data, which limits the
46 performance of land surface modeling (Gochis et al., 2010). DTB Information is also indispensable
47 to civil engineering in building homes, roads, railways, and bridges (Price, 2009). Furthermore,
48 DTB is of great importance to the study and applications of hydrology, ecology, agriculture, and
49 other relevant fields (Tromp-van Meerveld et al., 2007; Fu et al., 2011).

50 Although DTB is often considered equal to the thickness of the soil, there are great differences
51 between different measurement results. Soil thickness is mostly determined based on soil profiles
52 from soil surveys and borehole profiles from geological surveys. The observed depth of a soil profile
53 is generally less than 2 meters, and the thickness of the soil is therefore recorded as a value lower
54 than 2 meters (Shangguan et al., 2017). However, in reality, the DTB (the depth to the R horizon)
55 ranges from 0 meters to more than 1 kilometer, which is much greater than the average depth of soil
56 profiles. Limited by external factors such as equipment and technological constraints, traditional
57 soil surveys cannot reach bedrock in most cases. However, in contrast to traditional soil surveys,
58 geological borehole drillings usually reach depths of hundreds of meters or even deeper, and most
59 boreholes reach bedrock. Thus, borehole drilling logs are the most effective sources of DTB data.
60 Ground observations of DTB, which include soil profiles from soil surveying and borehole drilling
61 log data such as water well records and other measurements, have been widely used as training data
62 to produce spatial predictions of DTB (Tesfa et al., 2009; Shafuque et al., 2011; Miller and White,
63 1998; Hengl et al., 2014; Shangguan et al., 2017). Various mapping methods, which include
64 physically based models, interpolation from samples, and empirical-statistical models (Kuriakose
65 et al., 2009), have been employed for this purpose. Pelletier and Rasmussen (2009) proposed a
66 geomorphically based model that uses digital elevation model data to predict soil thicknesses based
67 on a hypothesis that there is a long-term balance between soil production and erosion. Karlsson et
68 al. (2013) developed a simplified regolith model modified from a trigonometric approach to estimate
69 regolith thickness based on slopes, outcrops, and distance to outcrops in eight directions, and
70 compared the results with those of linear regression and inverse distance weighting interpolation.



71 Shafique et al. (2011) proposed a multivariate linear model based on elevation, landform, and
72 distance to stream information to predict regolith thickness in a data-sparse environment. Hengl et
73 al. (2014) used zero-inflated models to predict global depth to bedrock based on a compilation of
74 major international soil profile databases. Dahlke et al. (2009) used a soil landscape model to predict
75 soil depth based on class means of merged spatial explanatory variables. Tesfa et al. (2009) applied
76 generalized additive and random forest models based on topographic and land-cover attributes to
77 predict soil depth at the watershed scale. Shangguan et al. (2017) predicted global depth to bedrock
78 using the random forest and gradient boosting tree models. Based on previous studies, machine
79 learning methods, especially random forest (RF) and gradient boosting tree (GBT) methods, showed
80 better performance than traditional interpolation methods under normal circumstances, and are
81 available in the “*randomForest*” (Breiman, 2001) and “*xgboost*” (Chen et al., 2016) packages in the
82 R software.

83 Although information about DTB is very important, to date, information about DTB in China is
84 very deficient, and there is no independent map of depth to bedrock in China. However, researchers
85 have advanced toward this target. Globally, there are several existing maps of DTB covering the
86 area of China (FAO, 1996; Hengl et al., 2014; Pelletier et al., 2016, Shangguan et al., 2017). The
87 earliest global distribution of DTB was produced by the FAO (Food and Agriculture Organization)
88 (1996); the depth was limited to the uppermost 2 meters and mapped using expert rules, and was
89 primarily based on soil unit classification, soil phase, and slope class. Hengl et al. (2014) developed
90 a global depth-to-bedrock map at 1-km resolution based on zero-inflated models using a compilation
91 of major international soil profile databases and 75 global environmental covariates representing
92 soil-forming factors. Pelletier et al. (2016) produced a global data set of the average thicknesses of
93 soil, intact regolith, and sedimentary deposits by representing uplands using soil data and lowlands
94 using water well data, with topographic, climatic, and geological data used as input. In China,
95 Shangguan et al. (2013) developed a comprehensive 30×30 arc-second resolution gridded data set
96 of soil characteristics that included soil depth derived from soil profiles and the Soil Map of China
97 (1:1,000,000), but the soil-depth data quality was relatively low because there were fewer
98 observations of deep soil. In addition, Shangguan et al. (2017) produced another global map of depth
99 to bedrock based on machine learning, using soil profile data, borehole data, and pseudo-
100 observations.



101 Among above-mentioned maps of DTB, most have relatively coarse resolutions (1 km or
102 coarser), except the map produced by Shangguan et al. (2017) (250 m resolution). In addition,
103 observations of DTB (FAO, 1996; Shangguan et al., 2013; Hengl et al., 2014) have been based
104 solely on soil data; thus, the predictions are often limited to soil surfaces with depths limited to
105 several meters. This depth is not consistent with the actual distribution of DTB. In addition, most
106 samples (Pelletier et al., 2016; Shangguan et al., 2017) were located in North America, whereas no
107 samples or only a small number of samples were located in China, which resulted in high uncertainty
108 for predictions in China. However, a large number of borehole logs produced by geologists in China
109 provide DTB information and are now available. In addition, several environmental covariates with
110 high resolution have been produced, which can be used to produce a high-resolution DTB map of
111 China. These data sources provide the cornerstone for producing a new map of DTB with higher
112 accuracy and resolution.

113 In this study, we aim to estimate DTB in China using machine learning methods. Observations
114 interpreted from geological borehole profiles and pseudo-observations of DTB are used as training
115 points. An extensive list of remote-sensing-based covariates, including DEM-derived parameters,
116 climatic images, MODIS products, land cover/land use, and the latest lithological/soil maps of
117 China are used as covariates. The objective of this paper is to (1) produce a DTB map of China at a
118 high spatial resolution of 100 meters; (2) compare and evaluate this map with observations and
119 existing DTB maps; and (3) estimate the uncertainty of the DTB map and discuss the outlook for
120 generating more accurate DTB maps in the future.

121 **2 Materials and methods**

122 **2.1 Borehole data**

123 A total of 6,382 borehole logs sampled from the Chinese National Important Geological Borehole
124 Database (NIGBD <http://zkinfo.cgsi.cn>) were used in our study. The NIGBD comprises about 80
125 million boreholes from across China (except Taiwan province). In every borehole log, geographic
126 coordinates and detailed lithological records are provided in the form of scanned images. Therefore,
127 the DTB of each borehole can be interpreted by finding the boundary between the regolith and fresh
128 bedrock.

129 **2.1.1 Observations sampled from the NIGBD**



130 The DTB of every borehole must be interpreted manually, and interpreting more than 80 million
131 boreholes logs therefore demands an immense amount of work and has high costs. However, many
132 boreholes that are located close to each other have similar DTB and environmental factors.
133 Therefore, we developed a sampling scheme to take a fraction of borehole drillings from the NIGBD
134 as the observation data sets in this study. Mapping methods, regardless of methods based on spatial
135 autocorrelation or soil environmental correlation, have requirements based on the number,
136 distribution, and typicality of the samples, which ensure global representation of the samples (Zhang
137 et al., 2012). To obtain representative samples from these boreholes, we used a sampling scheme
138 similar to stratified sampling to acquire our training points from the NIGBD.

139 The stratified sampling scheme includes designation of grid shape (such as a square grid,
140 triangular grid, or hexagonal grid) and grid size. A square grid is the easiest and most effective, and
141 is most widely used in sampling (Zhang et al., 2012). In general, smaller grid size leads to more
142 accurate predictions, but with greater sampling costs. Here, we used square grid sampling with a 0.2
143 $\times 0.2$ arc-degree grid, in consideration of the balance between representativeness and cost. Usually,
144 one observation or a number of observations are sampled at random locations from each grid.
145 However, the locations of boreholes in this study were determined in a previous geological survey.
146 Thus, we have taken one borehole randomly from each grid instead of one borehole from a random
147 location.

148 The depths of the boreholes range from 0 meters to more than 1 kilometer. Among these
149 boreholes, we were unable to determine the DTB from a few boreholes because of the limitations
150 of the records (see details in Sect. 2.1.2). This constraint resulted in vacancies of many grid cells
151 after the interpretation of all boreholes from the first sampling. To resolve this problem, we used an
152 additive sampling method; that is, additional samplings were taken multiple times until no new
153 observations could be added to the observation data sets. Thus, the latter samplings were aimed at
154 grids without DTB data based on the previous samplings. After a finite number of additive
155 samplings, the borehole logs of the NIGBD were considered efficiently used, and samples from all
156 the samplings were used in our study. The distribution of DTB observations interpreted from
157 boreholes is shown in Fig. 1.

158 **2.1.2 Interpretation of borehole records**

159 Interpreting DTB from borehole profiles sampled from the NIGBD was one of the crucial aspects



160 of this study. Borehole profiles, which were previously recorded by geologists, have longitudinal
161 verbal descriptions of soil layers and lithological layers with corresponding depths from the land
162 surface to the top and bottom of each layer. A typical simplified borehole profile diagram is shown
163 in Fig. 2.

164 Each borehole profile has several layers. Generally, the top layer of a borehole profile is pedolith,
165 where pedological processes have destroyed the original bedrock structure, principally through the
166 weathering of primary bedrock minerals and the formation and re-distribution of secondary
167 materials (National Committee on Soil and Terrain, 2009). Below is saprolite, referring to the zone
168 where the bedrock fabric is largely isovolumetrically weathered but primary bedrock structures are
169 still recognized. At the bottom is the unweathered bedrock. Because different boreholes were drilled
170 by different geological teams at different times, the details of stratification in the profiles often differ,
171 and the lithological description of each layer may be detailed or vague. These differences result in
172 inconsistencies or uncertainties in the borehole database, which were propagated into our DTB
173 observations.

174 To interpret the DTB from a borehole profile in the form of a scanned picture, we must manually
175 determine the boundary between the regolith and fresh bedrock based on lithological descriptions
176 and the dip angle of the borehole. The dip angles of a minority of boreholes whose dip angle were
177 not given were about 90° . Then, the DTB was calculated as the product of boundary depth and sine
178 of the dip angle. DTB can be interpreted from most sampled boreholes. However, some boreholes
179 are too shallow (several meters or less than 1 m) to reach the bedrock, and some have lithological
180 records that are unclear, which can make it is very difficult to determine the DTB (as described in
181 Sect. 2.1.1). Therefore, we used additive samplings. Because a number of boreholes went to depths
182 of more than 100 meters but still did not reach the bedrock, we could not obtain accurate DTB data
183 from these borehole profiles either. In this case, we regarded the depths of those boreholes as
184 approximations of the real DTB value. In addition, most research and applications focus on
185 relatively shallow depths.

186 2.2 Pseudo-observations

187 As shown in Fig. 1, DTB observations interpreted from borehole logs cover an extensive area across
188 China, except for the Qinghai-Tibet Plateau where boreholes are difficult to drill. Any purely data-
189 driven model fitted with large gaps in the covariate space is most likely to result in considerable



190 omissions, especially for areas that are often inaccessible or not of interest to soil surveys or
191 geological exploration. Therefore, we used pseudo-observations added to training data to fill such
192 gaps, which will avoid extrapolation for these areas (e.g., deserts and steep mountainous areas).
193 Deserts consist mainly of sand, and the DTB of such areas could be found in some publications.
194 Steep-slope areas without vegetation typically have very shallow or zero DTB; that is, rock outcrop.
195 Therefore, we used the following data sources to generate pseudo-observations to add to the training
196 points:

197 (1) The distribution map of deserts in China from the Data Center of Environmental and
198 Ecological Science in Western China (<http://zgsm.westgis.ac.cn>).

199 (2) Steep, bare surface areas generated using a slope map of China and remote-sensing-based
200 data.

201 (3) Previously published detailed geological maps reporting DTB or bedrock outcrops.

202 We generated a certain number of points in random positions within deserts based on the
203 distribution map of China's deserts. The DTB values of these points were obtained from existing
204 material and previous studies of the sand thickness of the deserts. We must note that the number of
205 points was limited to less than 10% of the whole number of observations to prevent adding too many
206 soft observations, and we only used points whose values had high credibility. In addition, several
207 points located in high-slope areas ($> 60^\circ$) were added to the observations with DTB values that
208 varied between 0 and 0.1 m.

209 **2.3 Environmental covariates**

210 In our study, a total of 133 related environmental layers, which cover five types of factors (climate,
211 topography, living organisms, water dynamics, and parent material) and represent the factors of soil
212 formation according to Jenny (1994), were selected to generate a DTB map of China. These
213 predictors were generalized into seven predictive “*scorpan*” factors (McBratney et al., 2003). The
214 133 covariates classified as “*scorpan*” factors included:

215 (1) Harmonized soil database images: percent coverage of Andosols, Histosols, and dozens of
216 other soil types.

217 (2) Climatic images: images indicating the values of 8-day MODIS day-time and night-time
218 local standard time (LST), long-term and monthly precipitation data, etc.

219 (3) Land use and land cover images: including vegetation maps, land cover and land use



220 classifications, biomass and yield maps, etc.

221 (4) Relief data, mainly derived from digital elevation models: slope maps, the topographic
222 wetness index, the topographic openness index, physiographic landform units, elevation and
223 secondary terrain attributes, etc.

224 (5) Geological and parent material maps: geological ages based on surface geology.

225 The complete list of the 133 environmental covariates is given in Supplement File A.

226 **2.4 Spatial prediction model**

227 The framework of our research is shown in Fig. 3. This framework consists of four main processes:

- 228 1. Overlaying observations of DTB and covariates to generate a regression matrix for modeling;
- 229 2. Obtaining the best parameters for modeling using cross-validation;
- 230 3. Fitting the prediction models based on the whole regression matrix;
- 231 4. Applying spatial prediction models using covariates and comparing the prediction with
232 existing maps.

233 **2.4.1 Model fitting**

234 In this study, we overlaid observations of DTB and covariates under the same coordinate reference
235 to generate a matrix including DTB and covariate columns. The matrix was used as input data for
236 machine learning. Then, we separately used RF and GBT to fit the prediction models. Finally, the
237 spatial predictions were generated using an ensemble model based on the two models. RF and GBT
238 are decision-tree-based ensemble methods. The RF model uses fully grown decision trees and
239 reduces error by reducing variance (Breiman, 2001). The GBT model uses shallow trees and reduces
240 error mainly by reducing bias, and to some extent by reducing variance by aggregating the outputs
241 from many models (Chen and Guestrin, 2016). RF and GBT were implemented respectively in the
242 “*randomForest*” and “*xgboost*” packages in the R environment. Parallel computing was employed
243 to improve data processing efficiency.

244 **2.4.2 Model validation and evaluation**

245 Ten-fold cross-validation was used to evaluate prediction accuracy. Comparison with previously
246 existing DTB maps was then employed to evaluate our results.

247 In cross validation, samples were divided into a training set (5,740 samples) and validation set
248 (642 samples). The training set was used to fit models, and the validation set was used to validate
249 model performance. Some widely used indicators such as the coefficient of determination (R^2 or the



amount of variation explained by the model), mean error (ME), and root mean square error (RMSE) were used to evaluate model performance. Of these indicators, the coefficient of determination is calculated by:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (1)$$

where SSR is the regression sum of squares, SST is the total variation sum of squares, and SSE is the residual sum of squares, which is the difference of SST and SSR. The variable y_i is the measured target value, \hat{y}_i is the prediction of each point, \bar{y} is the average of the measurements, and n is number of validation points. The value of R^2 is usually between 0 and 1; a value close to 1 indicates a perfect model, and values around 0 indicate a failed model. The RMSE, which is also called standard error, is calculated by:

$$RMSE = \sqrt{MSE} = \sqrt{SSE / n}, \quad (2)$$

where MSE is the mean squared error. RMSE estimates the deviation between predictions and observed values. A smaller RMSE indicates a better prediction.

Different covariates have different importance to DTB. Covariates with no or weak relations with DTB may produce noise in fitted models. This noise results in higher error of predictions. Our results based on modeling with different covariates showed that the noise has a certain degree of influence on the accuracy of the models, especially for the gradient boosting tree model. Therefore, we removed some covariates with low importance based on the random forests model to reduce prediction errors. The covariates we ultimately used are marked in Supplement File A.

In addition, to verify whether our predictions are more accurate than existing DTB maps of China, we compared our predictions with existing DTB maps using the validation set.

2.4.3 Model prediction and uncertainty estimation

The final model was fitted based on all samples with parameters selected by cross-validation. The final spatial predictions were generated using an ensemble model based on random forests and the gradient boosting tree method, which can avoid the overshooting effect (Sollich and Krogh, 1996). To predict DTB in China at 100 m resolution, we used the available environmental covariates at 100 m resolution.



277 Because any model for digital soil mapping inevitably suffers from different sources of error, it
278 is important to quantify the uncertainty associated with the produced maps (Poggio et al., 2016).
279 Analyzing and evaluating help data users to understand its existence and also can help to improve
280 decision quality (Liang et al., 2018). In this study, we used quantile regression forests to estimate
281 the uncertainty of estimations. Quantile regression forests are a tree-based ensemble algorithm for
282 estimation of conditional quantiles. This method is particularly suitable for high-dimensional data.
283 Quantile regression forests were implemented via the R environment in the “*quantregForest*”
284 package (Meinshausen, 2014). To estimate the uncertainty of predictions at every location, we
285 generated the uncertainty map of predictions by:

$$286 \quad \text{uncertainty} = \frac{qp_{0.9} - qp_{0.1}}{qp_{0.5}}$$

287 where $qp_{0.9}$ is the 0.9 quantile prediction of DTB, $qp_{0.1}$ is the 0.1 quantile prediction of DTB, and
288 $qp_{0.5}$ is the 0.5 quantile prediction of DTB. The uncertainty map is the reference when using the
289 DTB map of China.

290 All code used to generate predictions is available from the Github channels
291 (<https://github.com/yanfp/DTB100China>).

292 **3 Results**

293 **3.1 Model input statistical summary**

294 A summary of the DTB statistics is provided in Table 1. The DTB ranged from 0 to 1,106.91 m,
295 with a mean DTB of 36.62 m and a median value of 8.24 m. Fig. 4(a) shows the histogram of DTB
296 within 100 m. The DTB after logarithmic transformation had a distribution similar to a normal
297 distribution but with many zero values (i.e., outcrops) (Fig. 4(b)).

298 **3.2 Model accuracy and variable importance**

299 As is shown in Table 2, the GBT model had good ability to estimate DTB and yielded relatively
300 higher R^2 (0.81) and lower RMSE than the RF model (Table 2) based on the training set.

301 The importance of covariates measured based on the residual sum of squares of the random
302 forests model is shown in Fig. 5. The four most important covariates for DTB in this study were the
303 topographic wetness index, physiographic landform units, the topographic openness index, and
304 slope. In contrast, the most important covariate for the DTB according to Pelletier et al. (2016) and
305 Shangguan et al. (2017) was precipitation. The relationships between DTB and four important



covariates are shown in Fig. 6. This figure shows that DTB had a positive correlation with the topographic wetness index. The topographic wetness index is a secondary terrain attribute related to the geomorphometry of the surface or landform classification. In addition, DTB showed a positive correlation with the topographic openness index and elevation, and a negative correlation with the slope. These relations are consistent with our knowledge about DTB.

3.3 Estimation accuracy

The cross-validation summary statistics of interpolation for models based on RF and GBT are shown in Table 3 and Fig. 7. These statistics show that RF produced more accurate estimations than GBT. Because the GBT model showed relatively higher R^2 and lower RMSE than the RF model based on the training set (Sect. 3.2), this result means that the GBT model had a large degree of overfitting. Our results showed significant overestimation in lower values of DTB, which is a common problem in regression, especially when the model is not able to explain > 50% of variability in the target variable (Shangguan et al., 2017).

3.4 Prediction results

Output estimations of DTB by the ensemble model based on RF and GBT at 100 meters resolution are shown in Fig. 8. Our estimated results reveal that the predicted mean DTB was 54.42 m. High values of DTB were mainly distributed in desert areas, the North China Plain (including areas in Hebei province, Henan province, and Jiangsu province) and the Northeast China Plain (including areas in Heilongjiang province, Jilin province, and Liaoning province). Relatively lower values of DTB were mainly located in hilly and mountainous areas, such as Sichuan province, Chongqing city, Guangxi province, and the mountainous areas of Northeast China. The spatial pattern of the DTB map of this study is similar to those of the maps produced by Pelletier et al. (2016) and Shangguan et al. (2017).

In addition, estimations of three percentiles (0.1 (Fig. 9(a)), 0.50 (Fig. 9(b), and 0.9 (Fig. 9(c)) were produced by the quantile regression forests model. The mean values of the estimated DTB for the three percentiles were 4.95 m, 31.22 m, and 99.56 m, respectively. The maps show that the spatial pattern of DTB predicted by the quantile regression forests model was similar to that of the ensemble model based on the random forest and gradient boosting tree methods.

The uncertainty map of the prediction of DTB is shown in Fig. 10. The uncertainty in the predictions in part depends on the density of sampling (Zhou et al., 2018). In our study it was low



336 in deserts, sandy areas, the North China Plain, and the Northeast China Plain, where the topography
337 is relatively simple and sampling was relatively dense. In the Tibetan Plateau and western Inner
338 Mongolia, where sampling was sparse and DTB is low, the uncertainty was high. The uncertainty
339 was also relatively high in the Yun-Gui Plateau where the topography is complex with widespread
340 karst landforms.

341 **3.5 Comparison with existing study results**

342 We compared our results with existing maps produced by Pelletier et al. (2016) and Shangguan et
343 al (2017). Our results show similar spatial patterns with these maps. Of course, DTB values in
344 deserts, sandy areas, and the North China Plain were relatively high, and values in hilly and
345 mountainous areas, such as Chongqing City and Yunnan province, were relatively low in the map
346 of this study and in maps from global predictions. The estimated mean DTB was 54.42 m in our
347 study, whereas the mean values predicted by Pelletier et al. (2016) (Fig. 11 (a)) and Shangguan et
348 al. (2017) (Fig. 11 (b)) were 11.81 m and 26.64 m. The correlation coefficient between DTB
349 observations and predictions in our study is 0.75, which is significantly higher than the estimation
350 results of Pelletier et al. (2016) and Shangguan et al. (2017) (Table 4). In addition, compared with
351 the prediction results of Pelletier et al. (2016) and Shangguan et al. (2017), our estimation results
352 had obviously lower RMSE (47.57) and ME (1.82).

353 In addition, our prediction result shows similar spatial patterns to the maps produced by Pelletier
354 et al. (2016) and Shangguan et al. (2017), but revealed more detailed information than previous
355 predictions. There are more jumping points in the map of Shangguan et al. (2017) than others, and
356 the map predicted by Pelletier et al. (2016) shows low continuity in space with high values and low
357 values in a wide range. From the comparison in a typical region in the North China Plain (Fig. 12),
358 our map revealed more spatial details, especially in high DTB areas, than the maps by Shangguan
359 et al. (2017) and Pelletier et al. (2016) (Fig. 12(a)). In contrast, the map estimated by Pelletier et al.
360 (2016) shows abrupt change between highland and lowland areas (Fig. 12(c)).

361 **4 Data availability**

362 The resulting maps are available on Figshare at the following DOI:
363 <https://doi.org/10.6084/m9.figshare.7011524.v1>. And they are also available for download at
364 <http://globalchange.bnu.edu.cn/>.



365 **5 Discussion**

366 **5.1 Success and limitations of the data set**

367 Our training observations were selected by using square grid sampling with a 0.2×0.2 arc-degree
368 grid. We sampled at least one observation within each grid cell. Under this condition, the training
369 data are most representative under the current sampling method, which will produce the most
370 accurate predictions. However, boreholes have uneven spatial distribution. Very few boreholes were
371 located in inhospitable areas such as deserts and mountainous areas (Fig. 13). In addition, we were
372 unable to interpret the DTB from some borehole profiles. These limitations resulted in vacancies of
373 observations in many grid cells. Lack of observations will increase the uncertainty of predictions in
374 these areas.

375 The reliability of training data and covariates together determines the accuracy of predictions.
376 Although observations in this study were less heavily distributed in western China, which may limit
377 the accuracy of our predictions, the number of observations in China is far greater than that in other
378 studies. In addition, the DTB values interpreted from borehole profiles were more accurate than
379 those from soil profiles. Therefore, the DTB maps produced from borehole profiles were also more
380 accurate than maps solely based on soil profiles, especially for deep-DTB areas. In addition, the
381 predictions show a higher correlation coefficient with observations than did previous DTB maps
382 based on the validation set. The amount of variation explained by models for the DTB is about 57%,
383 which means that more than half of the variation is explained. We produced the DTB maps of China
384 at a resolution of 100 m. Although only a few covariates had spatial resolution of 100 m because of
385 the lack of data, the spatial resolution of most covariates was about 1 kilometer. Thus, spatial
386 variation at 100-meter scale may not be fully explained. However, covariates with high correlation
387 with DTB, such as DEM-derived parameters and land cover, have high resolutions (Fig. 5 and 6).
388 More observations and more covariates with high precision should be used in the future to improve
389 prediction accuracy.

390 **5.2 Error from interpretation of borehole records**

391 As described above, the DTB observations were visually interpreted from every borehole profile.
392 Because different borehole profiles were mapped by different organizations, the basis of layer
393 stratification differed slightly for different profiles. This issue contributes to the disunity of DTB



394 observations. In addition, the level of detail for different borehole profile stratifications is discrepant
395 because of their original uses. Furthermore, lithological records of some borehole profiles that give
396 vague information about soil and lithology were not distinct enough for us to interpret the DTB
397 accurately. All these factors contributed to errors in our DTB observations.

398 **5.3 Models built from different topographic partitions**

399 The DTB was determined based on many covariates including factors of topography, climate,
400 geology, vegetation, age, and human activity. Soils at the surface of the Earth are formed under the
401 combined effects of those factors (Zhou et al., 2016). However, the mechanisms of soil formation
402 and the importance of each covariate still are not completely clear (Li et al., 2004). The most
403 important covariates related to the DTB may be different in different geographic partitions.
404 Therefore, a model based on observations over the whole area of China may not be able to capture
405 the major factors in some regional areas. Models built from regional partitions may produce more
406 accurate predictions than global models within the partitions. Pelletier et al. (2016) distinguished
407 global land surfaces into three landform components, upland hillslope, upland valley bottom, and
408 lowland, and used different models for each component to estimate the DTB. Peng et al. (2018)
409 divided training data into subsets according to the similarity of the predicted variables and attain the
410 independent prediction model, which improved the prediction accuracy. In the future, different
411 models should be built and spatial predictions should be applied separately in different topographic
412 partitions.

413 **6 Conclusions**

414 In this study, we demonstrated the use of an ensemble model to produce a DTB map of China at a
415 resolution of 100 meters using the most reliable ground observations of DTB interpreted from
416 borehole profiles. This study provides the final prediction map of DTB as well as an uncertainty
417 estimation map for China. The cross-validation showed that the R^2 of the ensemble model was 0.57,
418 and the comparison showed that our DTB map is more accurate than existing DTB maps. Even
419 though the shortage of data used in this study, including DTB observations and environmental
420 covariates, limited the precision of the DTB map at a scale of 100 meters, this data set provides
421 more accurate information for Earth system researches compared with previous maps of DTB.
422 Based on the spatial prediction framework, data processing, model fitting, and spatial prediction are



423 fully automated and can be updated easily. By adding more DTB observations and using more
424 accurate covariates, we will be able to produce more accurate DTB maps of China in the future.

425 **Author contributions.**

426 Wei Shangguan, Jing Zhang, and Fapeng Yan designed the experiment and control the planning.
427 Fapeng Yan collected and compiled DTB observation data, prepared a part of environmental
428 covariate data, built models and implemented the spatial prediction, and wrote the paper. Wei
429 Shangguan prepared the other part of environmental covariate data. Bifeng Hu contributed to the
430 process on the data compilation and data validation. Jing Zhang initiated and coordinated the work.
431 All authors contributed to the scientific discussion of the results, the editing, and revision of the
432 paper.

433 **Competing interests.**

434 The authors declare that they have no conflict of interest.

435 **Acknowledgments**

436 This work was supported by the Natural Science Foundation of China under grants 41575072 and
437 the National Key Research and Development Program of China under grants 2017YFA0604303.

438 **References**

- 439 Boer, M., Barrio, G. D., and Puigdefiibregas, J.: Mapping soil depth classes in dry Mediterranean areas
440 using terrain attributes derived from a digital elevation model, *Geoderma*, 72, 99-118, 1996.
- 441 Breemen, N. v., and Buurman, P.: Soil formation, Springer Science & Business Media, The Netherlands,
442 2002.
- 443 Breiman, L.: Random forests, *Machine learning*, 45, 5-32, 2001.
- 444 Chen, J., Li, M., and Wang, W.: Statistical Uncertainty Estimation Using Random Forests and Its
445 Application to Drought Forecast, *Mathematical Problems in Engineering*, 2012, 1-12,
446 10.1155/2012/915053, 2012.
- 447 Chen, T., and Guestrin, C.: Xgboost: A scalable tree boosting system, *Proceedings of the 22nd acm sigkdd*
448 *international conference on knowledge discovery and data mining*, 785-794, 10.1145/2939672.2939785,
449 2016.
- 450 Dahlke, H. E., Behrens, T., Seibert, J., and Andersson, L.: Test of statistical means for the extrapolation



451 of soil depth point information using overlays of spatial environmental data and bootstrapping techniques,
452 Hydrological Processes, 23, 3017-3029, 10.1002/hyp.7413, 2009a.

453 Hengl, T., de Jesus, J. M., MacMillan, R. A., Batjes, N. H., Heuvelink, G. B., Ribeiro, E., Samuel-Rosa,
454 A., Kempen, B., Leenaars, J. G., Walsh, M. G., and Gonzalez, M. R.: SoilGrids1km--global soil
455 information based on automated mapping, PLoS One, 9, e105992, 10.1371/journal.pone.0105992, 2014.

456 Hengl, T., Mendes de Jesus, J., Heuvelink, G. B., Ruiperez Gonzalez, M., Kilibarda, M., Blagotic, A.,
457 Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas, R.,
458 MacMillan, R. A., Batjes, N. H., Leenaars, J. G., Ribeiro, E., Wheeler, I., Mantel, S., and Kempen, B.:
459 SoilGrids250m: Global gridded soil information based on machine learning, PLoS One, 12, e0169748,
460 10.1371/journal.pone.0169748, 2017.

461 Jain, S.: Fundamentals of Physical Geology, Springer, New Delhi, 2014.

462 Jie, P., Asim, B., Qingsong, J., Ruiying, Z., Jie, H., Bifeng, H., Zhou, S., Estimating soil salinity from
463 remote sensing and terrain data in southern Xinjiang Province, China,
464 <https://doi.org/10.1016/j.geoderma.2018.08.006>, 2018.

465 Karlsson, C. S. J., Jamali, I. A., Earon, R., Olofsson, B., and Mörtberg, U.: Comparison of methods for
466 predicting regolith thickness in previously glaciated terrain, Stockholm, Sweden, Geoderma, 226-227,
467 116-129, 10.1016/j.geoderma.2014.03.003, 2014.

468 Kuriakose, S. L., Devkota, S., Rossiter, D. G., and Jetten, V. G.: Prediction of soil depth using
469 environmental variables in an anthropogenic landscape, a case study in the Western Ghats of Kerala,
470 India, Catena, 79, 27-38, 10.1016/j.catena.2009.05.005, 2009.

471 Li, T., Zhao, Y., Zhang, K., Zheng, Y., and Wang, Y.: Soil geography, China, 2004.

472 McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping, Geoderma, 117,
473 3-52, 10.1016/s0016-7061(03)00223-4, 2003.

474 Meinshausen, N.: Quantregforest: quantile regression forests, R package, 2014.

475 Miller, D. A., and White, R. A.: A Conterminous United States Multilayer Soil Characteristics Dataset
476 for Regional Climate and Hydrology Modeling, American Meteorological Society, 2, 1-26, 1998.

477 Osborne, J. W.: Improving your data transformations: Applying the Box-Cox transformation, Practical
478 Assessment, Research & Evaluation, 15, 2, 2010.

479 Pelletier, J. D., and Rasmussen, C.: Geomorphically based predictive mapping of soil thickness in upland
480 watersheds, Water Resources Research, 45, 10.1029/2008wr007319, 2009.



- 481 Pelletier, J. D., Broxton, P. D., Hazenberg, P., Zeng, X., Troch, P. A., Niu, G.-Y., Williams, Z., Brunke,
482 M. A., and Gochis, D.: A gridded global data set of soil, intact regolith, and sedimentary deposit
483 thicknesses for regional and global land surface modeling, *Journal of Advances in Modeling Earth*
484 *Systems*, 8, 41-65, 10.1002/2015ms000526, 2016.
- 485 Piskin, K., and Bergstorm, R. E.: *Glacial Drift in Illinois Thickness and Character*, 1975.
- 486 Poggio, L., Gimona, A., Spezia, L., and Brewer, M. J.: Bayesian spatial modelling of soil properties and
487 their uncertainty: The example of soil organic matter in Scotland using R-INLA, *Geoderma*, 277, 69-82,
488 10.1016/j.geoderma.2016.04.026, 2016.
- 489 Richard, S. M., Shipman, T. C., Greene, L. C., and Harris, R. C.: *Estimated Depth to Bedrock in Arizona*,
490 *Arizona Geological Survey, Digital Geologic Map DGM-52*, layout scale, 1, 2007.
- 491 Roeset, J. M., Stokoe II, K. H., and Seng, C.-R.: *Determination of Depth to Bedrock from Falling Weight*
492 *Deflectometer Test Data*, *Transportation Research Record*, 1995.
- 493 Rue, H., and Martino, S.: Approximate Bayesian inference for latent Gaussian models by using integrated
494 nested Laplace approximations, *Royal Statistical Society*, 71, 319-392, 2009.
- 495 Shafique, M., der Meijde, M. v., and Rossiter, D. G.: Geophysical and remote sensing-based approach to
496 model regolith thickness in a data-sparse environment, *Catena*, 87, 11-19, 10.1016/j.catena.2011.04.004,
497 2011a.
- 498 Shanguan, W., Dai, Y., Liu, B., Zhu, A., Duan, Q., Wu, L., Ji, D., Ye, A., Yuan, H., Zhang, Q., Chen, D.,
499 Chen, M., Chu, J., Dou, Y., Guo, J., Li, H., Li, J., Liang, L., Liang, X., Liu, H., Liu, S., Miao, C., and
500 Zhang, Y.: A China data set of soil properties for land surface modeling, *Journal of Advances in Modeling*
501 *Earth Systems*, 5, 212-224, 10.1002/jame.20026, 2013.
- 502 Shanguan, W., Hengl, T., Jesus, J. S. M. d., and Dai, Y.: Mapping the global depth to bedrock for land
503 surface modeling, *Advances in Modeling Earth Systems*, 9, 65-88, 2016.
- 504 Sollich, P., and Krogh, A.: Learning with ensembles: How over-fitting can be useful, *Advances in Neural*
505 *Information Processing Systems*, 190-196, 1995.
- 506 Tesfa, T. K., Tarboton, D. G., Chandler, D. G., and McNamara, J. P.: Modeling soil depth from
507 topographic and land cover attributes, *Water Resources Research*, 45, 10.1029/2008wr007474, 2009a.
- 508 Wilford, J.: A weathering intensity index for the Australian continent using airborne gamma-ray
509 spectrometry and digital terrain analysis, *Geoderma*, 183, 124-142, 10.1016/j.geoderma.2010.12.022,
510 2012.



511 Y. Zhou, R. Webster, R.A. Viscarra Rossel, Z. Shi, and Chen, S.: Baseline map of soil organic carbon in
512 Tibet and its uncertainty in the 1980s, *Geoderma*, 334, 124-133, 2018.

513 Yin Z., Asim B., Zhiqiang M., Yanli L., Qiuxiao C., and Shi, Z.: Revealing the scale-specific controls of
514 soil organic matter at large scale in Northeast and North China Plain, *Geoderma*, 271, 71-79, 2016.

515 Zhang, S., Zhu, A., Liu, j., and Yang, L.: Summarization of Digital Soil Attribute Mapping Methods and
516 Sample Design Based on Samples, *Soils*, 44, 917-923, 2012.

517 Zongzheng, L., Songchao C., Yuanyuan Y., Ruiying Z., Zhou S., and Rossel, R. A. V.: Baseline map of
518 soil organic matter in China and its associated uncertainty, *Geoderma*, 335, 47-56, 2018.

519

520 Table 1: Summary statistics of depth to bedrock in meters

DTB	Number
=0	1026
0~2.00	585
2.00~10.00	1833
10.00~50.00	1768
50.00~100.00	427
100.00~300.00	630
>300.00	113

521

522 Table 2: Model fitting results for the depth to bedrock.

Model	Unit	R ²	RMSE	ME
Random forests	M	0.575	47.48	1.75
Gradient boosting tree	M	0.811	31.43	2.13

523

524 Table 3: Mapping performance for the depth to bedrock.

	Unit	R ²	RMSE	ME
Random forests	M	0.573	47.57	1.82
Gradient boosting tree	M	0.547	49.53	2.18
Ensemble	M	0.566	48.57	2.50



Table 4: Correlation index between observations and predictions

Study	Unit	R	RMSE	ME
This study	m	0.752	47.57	1.82
Pellertier et al. (2016)	m	0.486	81.98	36.52
Shangguan et al. (2017)	m	0.475	67.32	14.71

R denotes the correlation coefficient

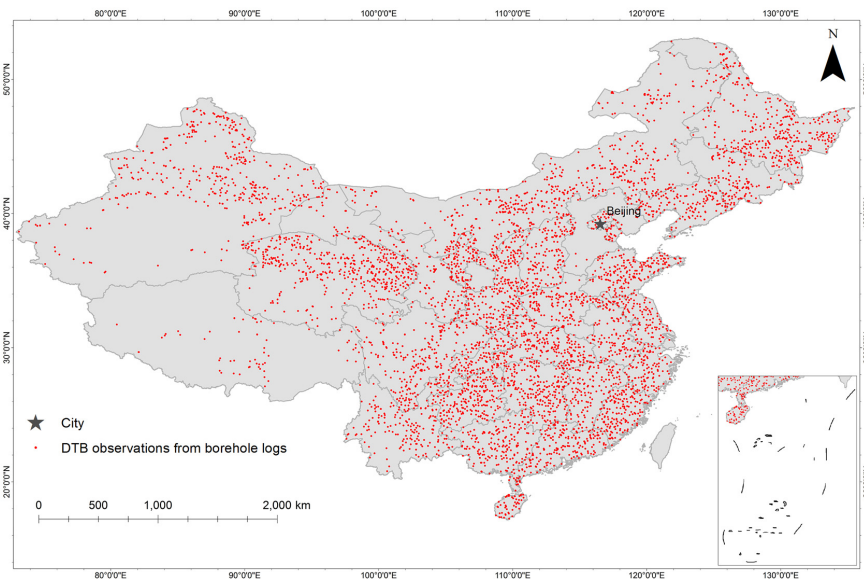


Figure 1: Distribution of DTB observations interpreted from boreholes.



Layer id	Depth(m)	Strip log	Lithology description
01	2.00		Quaternary surface soil , ...
02	10.00		Completely weathred basalt ...
03	20.00		Highly weathred basalt , ...
04	30.00		Moderately weathred basalt ... Slightly weathred basalt , ...
05	50.00		Fresh basalt , bedrock , ...

Figure 2: A typical borehole log sketch column. A borehole log describes the materials, color, and composition of each layer, and provides the depth, dip, and other relevant information. The original logs are in Chinese.

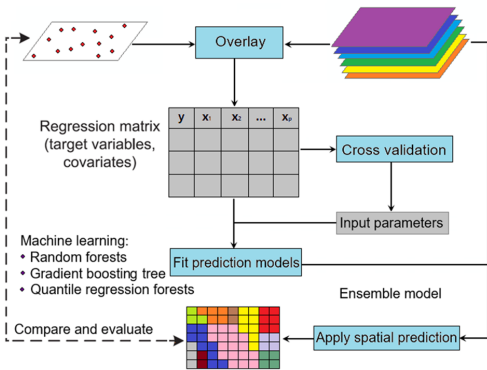
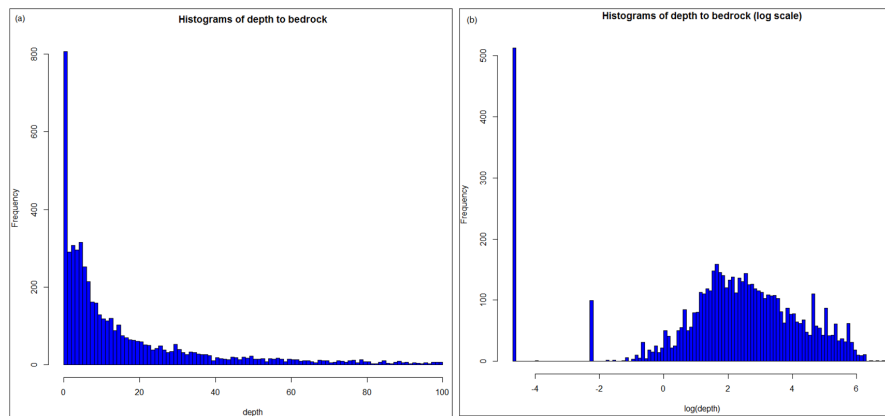


Figure 3: The spatial prediction framework used to fit models and apply spatial prediction of DTB in China at 100 m resolution.



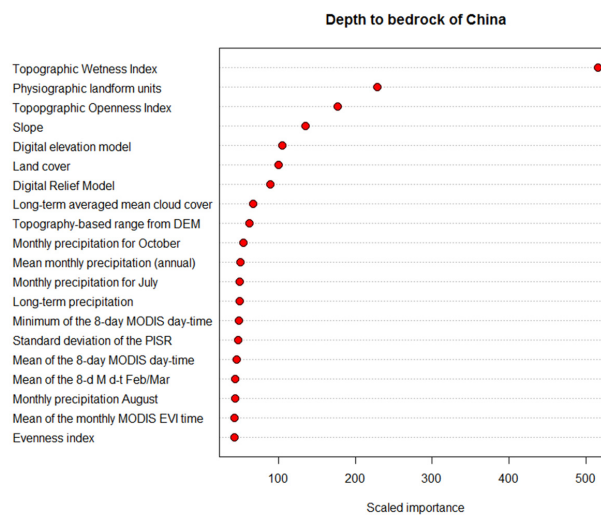
542



543

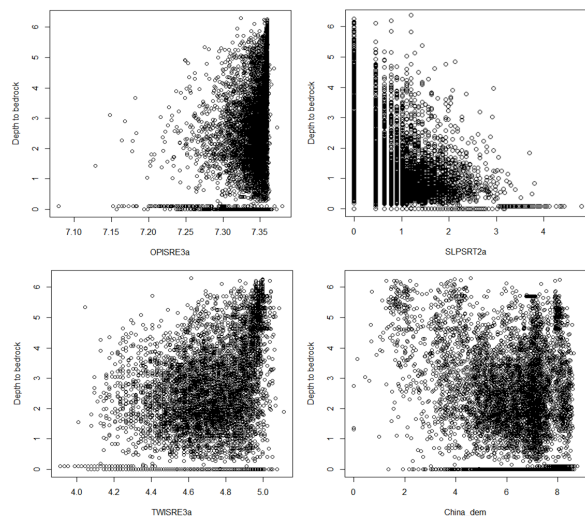
544 Figure 4: Histogram of depth to bedrock (a) and (b) after logarithmic transformation (values large
 545 than 100 m are not shown).

546

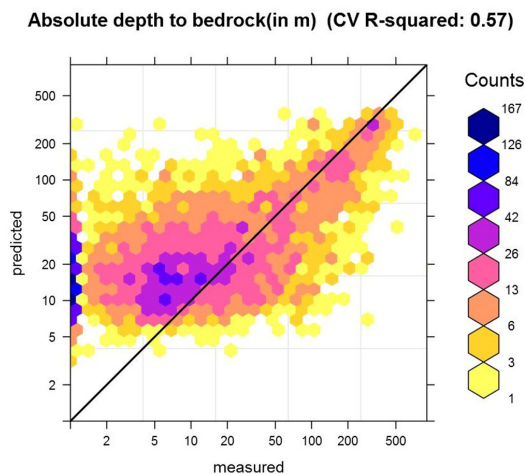


547

548 Figure 5: Importance of covariates for the depth to bedrock based on the random forest model.



549
 550 Figure 6: Relationships for target variables and the most important covariates (logarithmic scale).
 551 (*TWISRE3a* is the SAGA Topographic Wetness Index; *SLPSRT2a* is a slope map in percent;
 552 *OPISRE3a* is the SAGA Topographic Openness Index; *China_dem* is a digital elevation model of
 553 China.)
 554



555
 556 Figure 7: Plot showing cross-validation results for depth to bedrock on a logarithmic scale; R^2
 557 is calculated using Eq. (1).

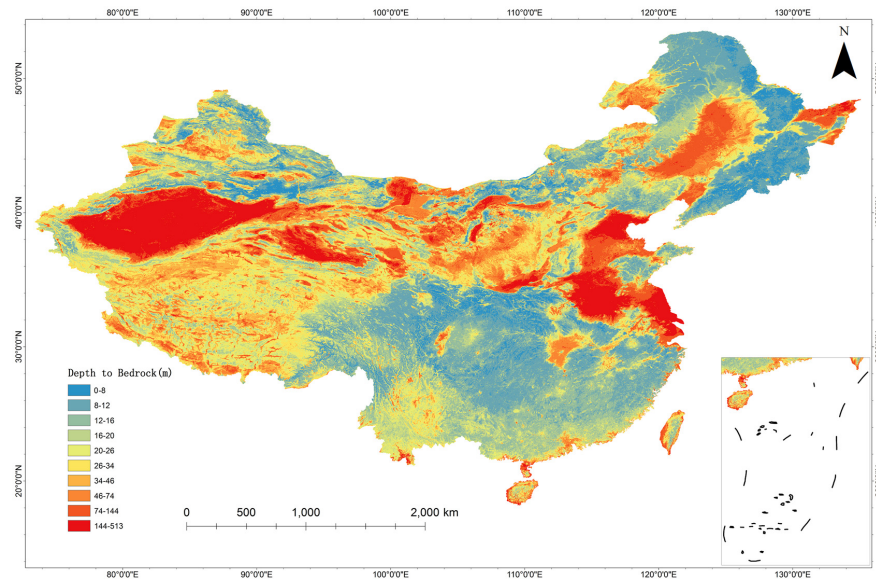


Figure 8: Final prediction of the depth to bedrock based on the ensemble model

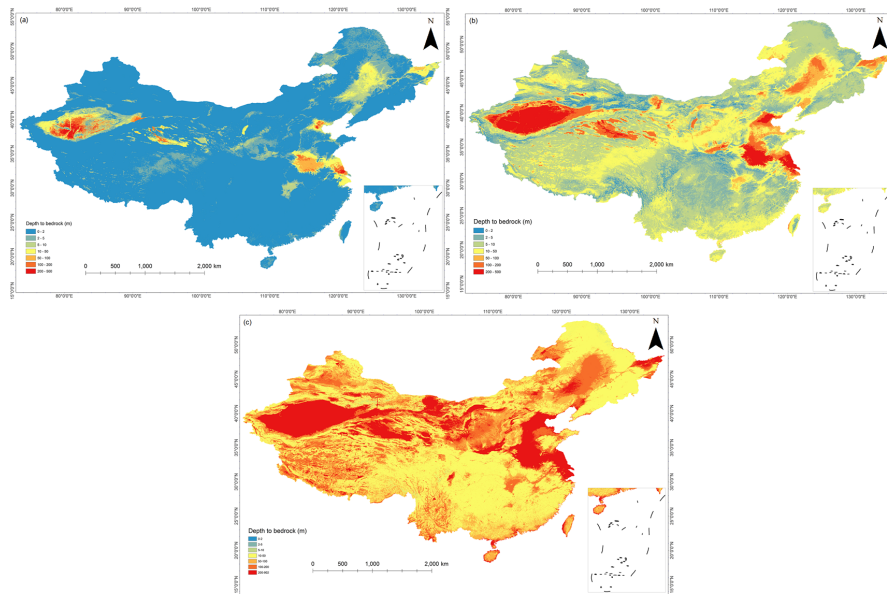


Figure 9: Depth to bedrock maps produced by the quantile regression forests model at the percentiles of 0.1 (a), 0.50 (b), and 0.9 (c).

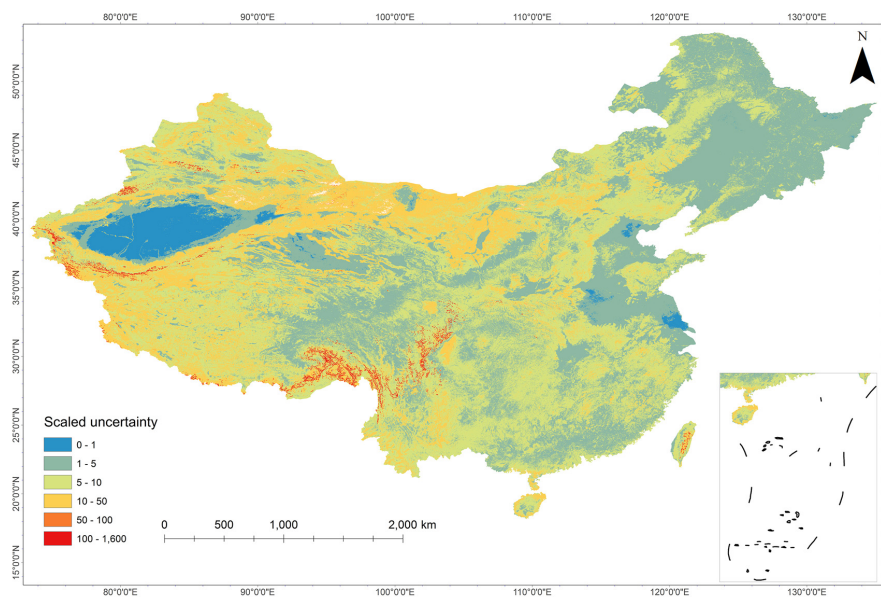


Figure 10: Uncertainty map of prediction of the depth to bedrock

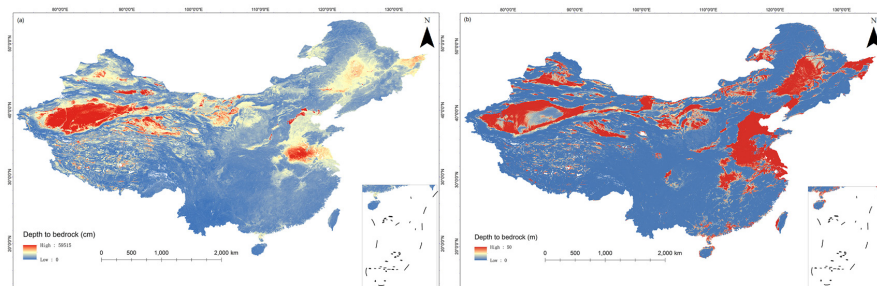
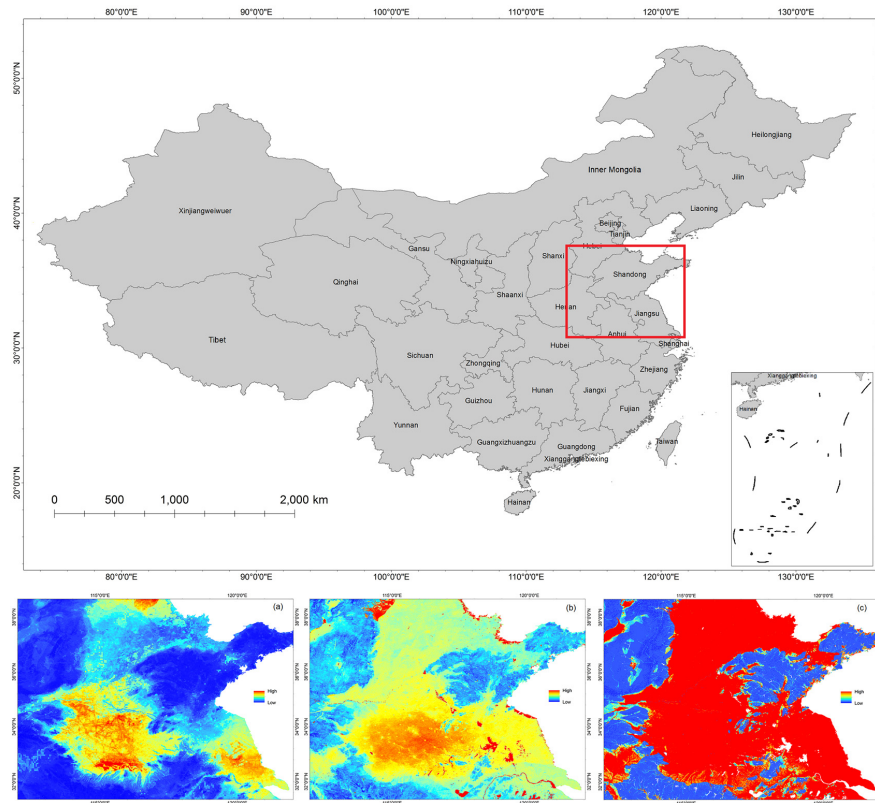


Figure 11: Extracted maps from global predictions of (a) Shangguan et al. (2017) and (b) Pelletier et al. (2016)



571
 572 Figure 12: Regional maps of (a) this study, (b) Shangguan et al. (2017), and (c) Pelletier et al.
 573 (2016).

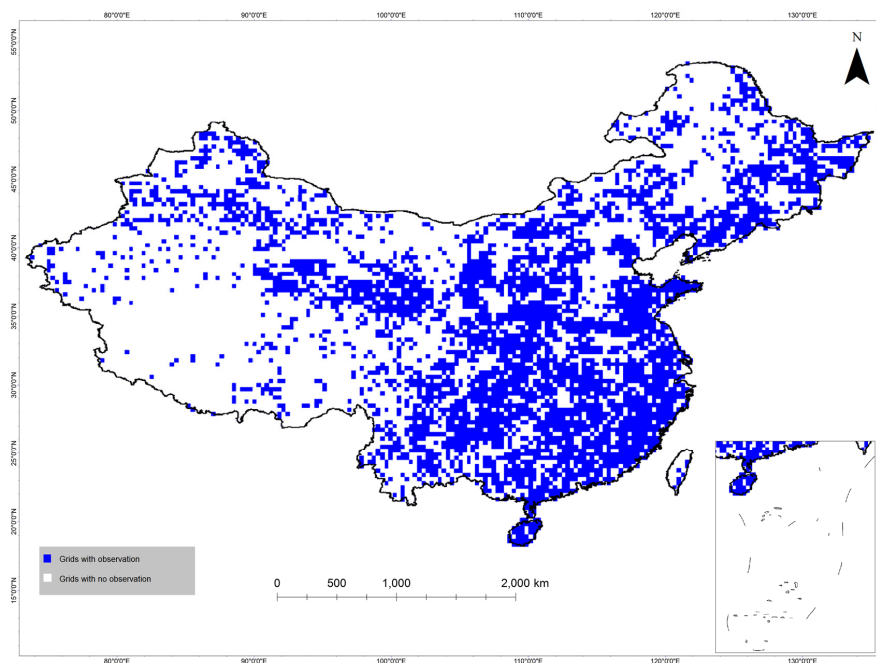


Figure 13: The distribution of 0.2×0.2 arc-degree grid with observation (blue color).