1    **Predicting groundwater arsenic contamination in Southeast Asia**

2    **from surface parameters**

3

4    Lenny Winkel, Michael Berg*, Manouchehr Amini, Stephan J. Hug, C. Annette
5    Johnson
6
7    **Eawag, Swiss Federal Institute of Aquatic Science and Technology, 8600 Dübendorf, Switzerland.**
8
9    * **Corresponding author. e-mail: michael.berg@eawag.ch**
10    **Phone: +41-44-823 50 78; Fax: +41-44-823 50 28**

11

12

17

18    **Arsenic contamination of groundwater resources threatens the health of**

19    **millions of people worldwide, particularly in the densely populated river deltas**

20    **of Southeast Asia. Although many arsenic-affected areas have been identified in**

21    **recent years, a systematic evaluation of vulnerable areas remains to be carried**

22    **out. Here we present maps pinpointing areas at risk of groundwater arsenic**

23    **concentrations exceeding 10 µgL$^{-1}$. These maps were produced by combining**

24    **geological and surface soil parameters in a logistic regression model, calibrated**

25    **with 1756 aggregated and geo-referenced groundwater data points from the**

26    **Bengal, Red River and Mekong deltas. We show that Holocene deltaic and**

27    **organic-rich surface sediments are key indicators for arsenic risk areas and that**

28    **the combination of surface parameters is a successful approach to predict**

29    **groundwater arsenic contamination. Predictions are in good agreement with the**

30    **known spatial distribution of arsenic contamination and further indicate**

31    **elevated risks in Sumatra and Myanmar where no groundwater studies exist.**

32

33  More than 100 million people worldwide ingest excessive amounts of arsenic (As)

34  through drinking water contaminated from natural geogenic sources. Many Asian

35  countries, in particular, are known to be affected by high groundwater As

36  concentrations as a result of chemically reducing aquifer conditions: Bangladesh[1-10],

37  India[3,11,12], China[13,14], Nepal[15], Cambodia[16-18] and Vietnam[17,19,20]. However, since As

38  analysis is expensive and time-consuming, groundwater resources of many regions

39  still remain to be tested. Therefore, maps pinpointing areas vulnerable to As

40  contamination can guide households at risk of arsenic contamination, as well as

41  scientists and policy-makers to initiate early mitigation measures and protect the

42  populations from chronic As poisoning.

43

44  Though the exact chemical conditions and reactions leading to As mobilization are still

45  under debate, it is generally accepted that microbial and/or chemical reductive

46  dissolution of As-bearing iron minerals in the aquifer sediments[1,4,21] is the main cause

47  for the release of As. Reducing conditions are often associated with the presence of

48  natural (bio)degradable organic carbon embedded in sediments[9,11,22-25]. Other

49  identified key characteristics of contaminated areas are rapidly buried, young

50  (Holocene) sediments and low hydraulic gradients in flat and low-lying areas[8,18,26,27].

51  Ideally an As prediction model for groundwater should be based on parameters that

52  indicate the key characteristics mentioned above in three dimensions. However, in the

53  absence of a 3-dimensional spatially continuous database of aquifer conditions to

54  depth, globally and regionally available (two dimensional) surface parameters can be

55  used as indicators for As-enrichment in underlying aquifers[28,29].

56

57  In the past, several geostatistical interpolation methods (e.g. kriging) have been used

58  to predict elevated As in groundwater on a regional scale[30-32]. However, for predictions

59   of areas where no groundwater quality data exist, interpolation methods are not

60   applicable[32] and models based on logistic regression are more appropriate[33,34]. In an

61   expert-based statistical model to delineate areas at risk of groundwater As

62   contamination on a coarse global scale, we found that geological information was of

63   crucial importance[29]. Here we focus on an in-depth assessment of depositional

64   environments in Southeast Asia. We use a logistic regression approach based on

65   relationships between sedimentary information, soil maps and measured groundwater

66   As data of Bangladesh, Cambodia and Vietnam, to assess the relative importance of

67   the different surface proxies in these countries. We apply these relationships to set up

68   prediction maps of As contamination in Southeast Asia including Indonesia (Sumatra)

69   and other countries where groundwater quality data is scarce (Myanmar and Thailand).

70   Furthermore, we verify the predicted risk in South Sumatra, where the groundwater

71   has not previously been tested for the presence of As.

72

73   **ARSENIC PREDICTION MODEL**

74   Our model is based on three assumptions. First, sedimentary depositional

75   environments are characterized by a unique combination of chemical, physical, and

76   biological properties[35] and can serve as indicators (proxies) for chemical and physical

77   conditions of the aquifers beneath the surface. Second, soil properties are proxies for

78   present and past drainage conditions and they are also indicators of recent

79   depositional environments. Third, soil textures, for example clay and silt are proxies

80   for the chemical maturity of the sediments, where clay is more mature than silt. An

81   important factor in the development of soil textures is topography[36], which allows the

82   delineation of areas where the model is applicable.

83

84   GIS-datasets were established from digital elevation data, countrywide geological

85   maps and global soil data (FAO), which were converted to a raster format using

86   ArcGIS (ver. 9.2). An overview of GIS data used in this study is provided in Table S1

87   (see Supplementary Information). Because each geological map applied a different

88   classification terminology, we created an uniformly classified geological map for all

89   regions (Figure 1). Although Bangladesh does geographically not belong to Southeast

90   Asia, it was included in the model because of the large number of available data

91   points[5]. Statistical relations between As concentrations and 30 parameters related to

92   soil properties, geology, climate, and hydrology (Table S2) were initially evaluated by

93   stepwise regression. The six parameters exhibiting a significance >95% and two

94   additional soil parameters were employed in the final model (see Methods and Table

95   1). Since young geological deposits and As groundwater contamination are rarely

96   observed in areas with steeper slopes, and groundwater As concentration data are

97   only available for regions with a flat topography (slope <0.1° ≈ 0.17%), areas with

98   slopes >0.1° were excluded from the model (see Supplementary Figure S1).

99

100  The prediction required three steps: i) Aggregation and binary-coding of measured As

101  concentrations to reduce spatial heterogeneities (dependent variables), ii) Logistic

102  regression to obtain weighting coefficients of independent variables (see Methods

103  section and corresponding results in Table 1), and, iii) Calculation of the probability of

104  As contamination based on the threshold value of 10 $\mu gL^{-1}$. The spatial datasets

105  considered as independent variables for the model are topographic data to delineate

106  the model area, sedimentary depositional environments as a proxy for aquifer

107  conditions, and soil variables as a proxy for drainage and chemical maturity of

108  sediments (see the Methods section).

109

## SURFACE PARAMETERS CONTRIBUTING TO THE MODEL

The weighting factors ($\lambda$) and significance of the eight independent variables retained in the final model are listed in Table 1. In general, variables describing sedimentary depositional environments have a larger contribution to the model than soil variables, presumably because geological variables come closest to describing As contamination in the aquifer itself. The presence of young deltaic deposits ($\lambda$=1.65) is a particularly significant indicator for As-contaminated aquifers. In Southeast Asia, delta initiation and progradation occurred simultaneously with the Holocene Climate Optimum[37]. Therefore, delta progradation resulted in the burial of organic matter at a high rate. The presence of relatively fresh organic carbon provides favourable conditions to establish reducing environments, which may lead to As enrichment in groundwater. Logistic regression confirms that organic-rich deposits ($\lambda$=1.11) play an important part in the model. Recent alluvial deposits ($\lambda$=0.55) are also indicative for elevated As concentrations in groundwaters (Table 1). Of the soil parameters, medium textured soils seem to be indicative of As-bearing aquifers evolved from rapid accumulation of young (Holocene) sediments. In contrast, the negative weighting coefficient ($\lambda$= -0.95) for floodplain deposits implies that these fine-grained deposits could overlie aquifers low in dissolved As (compare Figures 1 and 3). Fine-grained deposits (high clay content) thereby point to low-energy depositional environments with condensed sediments where Holocene aquifers are rare and where groundwater is likely drawn from older aquifers.

Pre-Holocene deposits, other Holocene deposits and tidal deposits (Figure 1) were found to be statistically insignificant (p >0.05) and they were excluded from the model in the first stepwise regression (see Table S2 and Table 1). Tidal sediments are

135  generally associated with aquifers abundant in sulfate that may be microbially reduced

136  to sulfide and re-precipitate As[38,39]. They might serve as proxies for low-risk areas, but

137  our As data in such aquifers were possibly too few to see such a relation.

138

139  **INFERRING DEPOSITIONAL ENVIRONMENTS AT DEPTH FROM THE SURFACE**

140  As mentioned in the introduction, our model is based on two- dimensional data (i.e.

141  surface maps). Nevertheless, geological information (sedimentary depositional

142  environments) inherently contains a three-dimensional component. The recent

143  sedimentary history of major Southeast Asian basins is characterized by delta

144  initiation, which occurred on a global scale at about 8500-6500 years BP and was

145  principally controlled by the deceleration of sea-level rise[40]. Delta progradation

146  resulted in the unconformable deposition of thick late Pleistocene-Holocene sediments

147  on Pleistocene and older sediments, e.g. as incised-valley-fill deposits[8,18]. As-

148  contaminated aquifers are mainly present in these Holocene aquifers, whereas deeper

149  Pliocene-Pleistocene aquifers are, to a large extent, free of As[7,8,26]. The boundary

150  between Pleistocene and Holocene sediments is not located at a constant depth[8] and

151  this can lead to misclassifications, since our model inherently assumes that the

152  underlying aquifer belongs to the same sedimentary depositional environment as on

153  the surface.

154

155  Two situations exist where the environment at the surface does not reflect the geology

156  at depth, i) when the As measurements used in the model were obtained from

157  tubewells tapping deep (Pleistocene) aquifers because shallow (Holocene) aquifers

158  are not present, or shallow groundwater is too saline for consumption (e.g. in coastal

159  regions), and ii) when Holocene sediments were deposited at sedimentation rates to

160  small to form a usable aquifer. In both situations Holocene depositional environments

161    are present at the surface and high-risk areas are indicated, although measured As

162    concentrations are low (false positive cases). Even though the probability maps would

163    be improved by including the thickness of Holocene sediments, the absence of

164    country-wide three-dimensional geological data rules this out. Furthermore, the

165    complexity of aquifer heterogeneity at a local scale makes it inevitable that

166    misclassifications occur.

167

168    **PROBABILITY MAPS**

169    Predicted areas at risk of As contamination agree well with known spatial

170    contamination patterns, a finding which is supported by the results of the model

171    classification indicating the performance of model prediction (Figure 2) and the

172    Hosmer-Lemeshow goodness-of-fit statistics (Table S3). An absolute average

173    deviation of 7.3% is found between expected and modelled probabilities of As being

174    ≤10 or >10 $\mu gL^{-1}$. The model is further characterised by the receiver operating

175    characteristics curve area of ~0.7 (Figure S2), which is a good result considering that

176    neither depths of analyzed groundwater wells, nor aquifer hydrological data form part

177    of the model.

178

179    The probability maps for As concentrations exceeding 10 $\mu gL^{-1}$ are presented in

180    Figure 3 and supplementary Figure S3 (probability map of whole Southeast Asia). The

181    highest probabilities (0.7-0.8) for As contamination (>10 $\mu gL^{-1}$) are found in the south-

182    central part of Bangladesh, with a value of 0.5-0.6 in the north-eastern Sylhet basin

183    (Figure 3a). The probability of finding contaminated wells in the Red River delta

184    reaches a value of 0.7 (Figure 3e). The sedimentary depositional environments

185    present along the Mekong river differ from the deltaic environments of Bangladesh

186    and the Red River delta in that organic-rich deposits are found at close distance of the

187     modern Mekong and Bassac river courses, surrounded by extensive floodplain areas

188     free of As[18] (compare Figures 1 and 3). The probability map for the Mekong delta

189     shows values of up to 0.6 at close distance to the modern Mekong and Bassac river

190     courses and adjacent swampy marshes (Figure 3c). In addition, the large floodplain of

191     Lake Tonle Sap with organic-rich sediments is a risk area with probabilities ranging

192     between 0.4-0.6.

193

194     To compare our predicted low-risk (probability ≤0.4) and high-risk areas (probability

195     >0.4) (Figures 3b and 3d) with measured As contamination, a misclassification

196     analysis was performed based on the 1756 aggregated and binary-coded (≤10 or >10

197     $\mu gL^{-1}$) As measurement data for the three deltas. We report 71% (1023 aggregated

198     points), 59% (107 a.p.), and 75% (100 a.p.) of correctly classified cases for

199     Bangladesh, Red River and Mekong deltas respectively (average 70%). In our study,

200     the Red River and Mekong deltas have 6 and 9% false negative cases. In Bangladesh

201     false negative cases (13%) were found specifically in wells lying at a close distance to

202     rivers. The number of false negatives is outnumbered by false positives in all the three

203     deltas with 17, 35 and 16% for Bangladesh, the Red River and Mekong deltas,

204     respectively. Errors in the prediction can arise from the commonly reported well-to-well

205     variability where wells with low As levels are often present at close distance from wells

206     high in arsenic[5,6,18,19], as well as from uncertainties in measured As (estimated at 25%)

207     leading to misclassifications of concentrations being close to the threshold of 10 $\mu gL^{-1}$.

208     However, we interpret these misclassifications to be mainly an effect of modelling

209     three-dimensional processes based on two-dimensional data.

210

211 Apart from the three deltas discussed above, our Southeast Asia probability map

212 (supplementary Figure S3) also highlights risk areas that are largely unknown or

213 unreported, particularly in Sumatra, Myanmar and Thailand (Figure 3f).

214

215 **VERIFICATION OF PREDICTIONS FOR SUMATRA**

216 According to the modelled probability, an area of about 100,000 km$^2$ at Sumatra's east

217 coast is prone to high risk of As contamination (probability >0.4) (see Figure 4a and

218 supplementary Figure S3). To validate the Sumatra prediction map, 97 groundwater

219 samples were collected in 2007 in the Province of South Sumatra in the vicinity of

220 Palembang (see Methods section). This area was chosen because it is at the border

221 of a low- and high-risk area and not previously studied. Since there is no large

222 variation in geological and topographical features along Sumatra's east coast, the

223 study area is representative of the whole high risk area. Figure 4b shows As

224 concentrations (≤10 µgL$^{-1}$ and >10 µgL$^{-1}$) measured in the groundwater, imposed on

225 the binary probability map. The classification results for Sumatra (63% correctly

226 classified, 36% false positive and 1% false negative) are comparable to those of the

227 Bangladesh, Red River, and Mekong deltas. In total, 94% of the 50 tubewells located

228 in the low-risk area have As levels below 10 µgL$^{-1}$. Of the 12 contaminated wells (As

229 >10 µgL$^{-1}$) 75% are positioned in the high-risk area. However, in the high-risk area the

230 contaminated wells are clearly outnumbered by uncontaminated groundwater

231 measurements (9 contaminated wells vs. 38 uncontaminated wells), for reasons

232 explained below.

233

234 On average, both high and low risk areas in Sumatra are characterized by high DOC,

235 $NH_4^+$, $HCO_3^-$, $PO_4^{3-}$, and Fe(II) and low $SO_4^{2-}$ concentrations (Supplementary Table

236 S5). At least two-thirds of the sampled groundwaters are reducing in nature, especially

237 those located in the high-risk area. Since the chemical conditions of the aquifers would

238 permit the reductive dissolution of As, other explanations for the overall low As

239 concentrations must be considered.

240

241 To gain a greater insight into the characteristics of the aquifer, the local geology was

242 examined laterally and in depth. The studied area is geologically young (Tertiary-

243 Quaternary) (Figure 4c). The sediment contains deposits from peat swamp forests that

244 developed during the Holocene era (past 5,000 to 10,000 years) and unconformably

245 overlie older sediments[41]. The outline of the high-risk As-contamination zone mainly

246 follows the outline of these deposits. The peat deposits are usually found ranging from

247 4-8 meters although depths of up to 24 metres have been reported[42]. However, the

248 depths of sampled tubewells in Sumatra average 46 meters, which implies that most

249 of them tap groundwaters from aquifers below the Holocene peat deposits. This

250 shows that the prediction map is a useful tool for the identification of areas at risk of

251 As contamination, but that understanding the local geology as a function of depth is of

252 vital importance for specific areas.

253

254 **ARSENIC CONTAMINATION IN THAILAND AND MYANMAR**

255 The probability map in Figure 3f shows that the Chao Phraya basin in central Thailand

256 and the Irrawaddy delta in Myanmar have a risk of elevated As being present in

257 groundwater, whereas the Sittang basin (Myanmar) has a lower risk, and the Salween

258 basin (Myanmar) has virtually no risk at all. The problem of As contamination in the

259 Irrawaddy delta is partly known[43] although its spatial extent has not been investigated

260 to date, which is particularly worrying considering the size of the area at risk. In the

261 Chao Phraya basin in central Thailand, a groundwater survey was undertaken in 2001

262 using As field test kits to test wells with minimum depths of 80 m[44]. The range of

263    measured As concentrations in the 37 tested wells was <1 µgL⁻¹ to 100 µgL⁻¹ with an

264    average of 11 µgL⁻¹. These results correlate with our maps that predict a low to

265    moderate As contamination. Indeed, North-South geological profiles across the

266    basin[45] indicate maximal depths of 20 m for the pre-Holocene incisions implying that

267    sampled tubewells draw water from Pleistocene or older aquifers.

268

269    In contrast to the slow sedimentation rates of the Sittang and Chao Phraya deltas, the

270    Irrawaddy delta received massive amounts of sediments during the Cenozoic era[46].

271    The Irrawaddy River[47] still has a 10 times larger sediment load than the Chao Phraya

272    River[45]. In 2002 the Departments of Medical Research and Health (Lower Myanmar),

273    financially supported by UNICEF, conducted a groundwater sampling campaign in the

274    Irrawaddy division[43]. In total, 99 groundwater samples (90 shallow tubewells and 9

275    dug wells) were collected in 25 villages, and As was quatified by atomic absorption

276    spectroscopy. It was reported that 67% of the sampled wells had As levels >50 µgL⁻¹.

277    These results show that the risk of As contamination in the Irrawaddy delta, as

278    indicated by the probability map, should be taken seriously and that there is an urgent

279    need to test shallow tubewells in other townships or districts in the high risk area.

280

281    **ASSESSING ARSENIC ELSEWHERE**

282    In the study presented here, we identified regions in Southeast Asia, based on As

283    prediction maps, where tubewells should be tested for elevated As concentrations

284    (>10 µgL⁻¹). Although the use of such maps in other scientific investigations (e.g.

285    climate research) is a common practice, the prediction of As (and other groundwater

286    conditions) is still an emerging technique in the field of natural groundwater

287    contaminants. Our As model differs from earlier models in its ability to predict

288    contamination in areas of unknown groundwater quality and on a sub-continental

289     scale. The strength of the prediction lies rather in the combination of surface

290     parameters than in the individual parameters. As an example, the Sittang basin and

291     the vicinity of the lower Mekong river branches are both characterized by alluvial

292     deposits, but only the latter has a modelled high-risk because of the contribution of soil

293     properties. A limitation to be considered is that shallow, As-bearing groundwater can

294     only be expected where Holocene aquifers are present. Where this is not the case,

295     high risk areas may indicate the presence of reducing aquifers, but As concentrations

296     in groundwater could be lower than predicted.

297

298     Our approach provides a blueprint for further modelling and mapping of As-tainted

299     aquifers in and outside of Southeast Asia. The probability maps can further be

300     improved when data with higher spatial resolution or in three dimensions becomes

301     available, although we emphasize that it will not be possible to account for the local

302     heterogeneities of aquifers. The presented prediction maps are a valuable and

303     resource-saving tool that can serve both scientists and policy-makers to initiate early

304     mitigation measures in order to protect the people from As-related health problems as

305     well as to efficiently guide water resources management.

306

307     **METHODS**

308     Geological maps of Bangladesh, Cambodia, Thailand and Vietnam were available in

309     digital format (Supplementary Table S1). Maps of Myanmar and Sumatra were

310     digitized for this study. Geological maps of Malaysia and Laos were not available. The

311     sedimentary depositional environments employed in the model as independent

312     variables are deltaic deposits, organic-rich sediments (e.g. sediments deposited in

313     marshy environments), floodplains, alluvial deposits, tidal deposits, other Holocene

314     sediments and pre-Holocene sediments. Soil variables are percentages of silt, clay

315 and sand in both the topsoil (0-30 cm) and subsoil (30-100 cm), and coarse, medium

316 and fine soil textures.

317

318 We compiled >4600 data points of groundwater As-concentrations from Bangladesh

319 (median well depth 35 m), the Mekong delta (Cambodia and Southern Vietnam, m.w.d.

320 39 m) and the Red River delta (Northern Vietnam, m.w.d. 30 m). The data originates

321 from BGS and DPHE (Bangladesh, n=3534)[5], from Buschmann et al. (Mekong delta,

322 n=352)[18,20], and from Berg et al. (Red River delta, n=720)[17,25,48] and was used without

323 a restriction in well depths. To test the model, 97 tubewell were randomly sampled for

324 this study in Sumatra at about 5-10 km intervals at an average sampling density of 1

325 sample per 54 km$^2$. This study area is positioned at a latitude of 2°872S to 3°911S

326 and a longitude of 103°949E to 104°993E (sampling area 85 km by 65 km).

327 Procedures of sampling and analysis were carried out as described in Buschmann et

328 al.[18]. Concentrations of As and additional parameters measured in Sumatra

329 groundwater samples are provided in Table S5 (Supplementary Information).

330

331 Arsenic concentrations are point measurements within vertical depths of the wells,

332 while the other variables have coarser spatial resolution being generally greater than

333 30 arc seconds. Point data of measured As-concentrations were therefore aggregated

334 using the geometric mean to a resolution of one point per pixel with a size of 5 arc

335 minutes (~9.3 km at the equator), which is the pixel size of the global soil data (FAO).

336 Aggregation resulted in a decreased dataset of 1756 pixel-based data points

337 (Bangladesh 1443, Red River Delta 180, and Mekong delta 133 points). The

338 aggregated point-data were binary-coded using the WHO guideline value for As in

339 drinking water (10 µgL$^{-1}$) as a threshold. We acknowledge that several countries apply

340 a guideline of 50 µgL$^{-1}$, but adopting this threshold would result in significantly less

341    data points (992) for the calibration of the model. The binary variable of whether As

342    concentrations exceed the WHO threshold (1) or not (0) was used as the dependent

343    variable in this study.

344

345    Logistic regression was used to determine the weighting of the independent variables.

346    This is a common statistical method used in environmental research and allows for the

347    concurrent use of continuous and categorical variables[33,34]. The parameter P denotes

348    the probability of As concentrations exceeding the WHO threshold. Logistic regression

349    models log(odds), which is defined as the ratio of the probability that an event occurs

350    to the probability that it fails to occur log(P/(1−P)), as a linear combination of

351    independent variables[49].

352    $$\log(odds) = C + \sum_{i=1}^{n} \lambda_i X_i \tag{1}$$

353    where $C$ is the intercept of regression, $X_i$ are independent variables, and $\lambda_i$ are the

354    weighting coefficients that were obtained using the maximum likelihood procedure[49].

355    Exponential values of coefficients, Wald statistics, and p-values (see Table 1) indicate

356    the importance of each variable. Statistically insignificant independent variables were

357    excluded from the model during each of the subsequent regression steps (Table S2).

358    The threshold for maintaining a variable in the model was determined by the 95%

359    significance level ($p < 0.05$). The silt contents in the subsoil and medium textured soils

360    were kept in the model because of their good spatial match with known contaminated

361    areas, which is supported by the presence of silty sands at the surface of regions

362    exhibiting contaminated aquifers in West Bengal[50], and by reported elevated

363    groundwater As concentrations in aquifers capped with fine surface material (silt and

364    clay) in Bangladesh[10,27].

365

366  According to the calculated odds, the probability (P) of having an As concentration

367  above 10 µgL$^{-1}$ was calculated as follows:

368  $$P = \frac{\exp(C + \sum_{i=1}^{n} \lambda_i X_i)}{1 + \exp(C + \sum_{i=1}^{n} \lambda_i X_i)}$$  (2)

369  In addition, we tested how successful the model predicted the number of

370  contaminated cases for probability intervals of 0.1 (Figure 2 and Supplementary

371  Information).

372

373  Misclassification occurs when either a point with As concentration >10 µgL$^{-1}$ falls in an

374  uncontaminated area (false negative) or an As concentration <10 µgL$^{-1}$ falls in a

375  contaminated area (false positive). Based on the model classification results (Figure 2

376  and Table S4), a probability threshold of 0.4 resulted in a minimum misclassification

377  rate. The binary risk maps was hence categorized into low-risk areas (probability <0.4)

378  and high risk areas (probability >0.4).

379

380

381  **REFERENCES**

382  1.   Nickson, R. *et al.* Arsenic poisoning of Bangladesh groundwater. *Nature* **395**, 338
383       (1998).
384  2.   Smith, A. H., Lingas, E. O. & Rahman, M. Contamination of drinking water by
385       arsenic in Bangladesh: a public health emergency. *Bull. World Health Org.* **78**,
386       1093–1102 (2000).
387  3.   Chowdhury, U. K. *et al.* Groundwater arsenic contamination in Bangladesh and
388       West Bengal, India. *Environ. Health Perspect.* **108**, 393-397 (2000).
389  4.   McArthur, J. M., Ravenscroft, P., Safiulla, S. & Thirlwall, M. F. Arsenic in
390       groundwater: Testing pollution mechanisms for sedimentary aquifers in
391       Bangladesh. *Water Resour. Res.* **37**, 109-117 (2001).

392   5.   BGS and DPHE. *Arsenic contamination of groundwater in Bangladesh*. Eds.
393        Kinniburgh, D. G. & Smedley, P. L. (British Geological Survey, Keyworth, U.K.,
394        2001). www.bgs.ac.uk/arsenic/bangladesh.

395   6.   van Geen, A. *et al.* Spatial variability of arsenic in 6000 tube wells in a 25 km(2)
396        area of Bangladesh. *Water Resour. Res.* **39**(2003).

397   7.   Ahmed, K. M. *et al.* Arsenic enrichment in groundwater of the alluvial aquifers in
398        Bangladesh: an overview. *Appl. Geochem.* **19**, 181-200 (2004).

399   8.   Ravenscroft, P., Burgess, W. G., Ahmed, K. M., Burren, M. & Perrin, J. Arsenic in
400        groundwater of the Bengal Basin, Bangladesh: Distribution, field relations, and
401        hydrogeological setting. *Hydrogeol. J.* **13**, 727-751 (2005).

402   9.   Meharg, A. A. *et al.* Codeposition of organic carbon and arsenic in Bengal Delta
403        aquifers. *Environ. Sci. Technol.* **40**, 4928-4935 (2006).

404  10.  van Geen, A. *et al.* Flushing history as a hydrogeological control on the regional
405        distribution of arsenic in shallow groundwater of the Bengal Basin. *Environ. Sci.*
406        *Technol.* **42**, 2283-2288 (2008).

407  11.  McArthur, J. M. *et al.* Natural organic matter in sedimentary basins and its
408        relation to arsenic in anoxic ground water: the example of West Bengal and its
409        worldwide implications. *Appl. Geochem.* **19**, 1255-1293 (2004).

410  12.  Ahamed, S. *et al.* Arsenic groundwater contamination and its health effects in the
411        state of Uttar Pradesh (UP) in upper and middle Ganga plain, India: A severe
412        danger. *Sci. Total Environ.* **370**, 310-322 (2006).

413  13.  Smedley, P. L., Zhang, M., Zhang, G. & Luo, Z. Mobilisation of arsenic and other
414        trace elements in fluviolacustrine aquifers of the Huhhot Basin, Inner Mongolia.
415        *Appl. Geochem.* **18**, 1453-1477 (2003).

416  14.  Yu, G. Q., Sun, D. J. & Zheng, Y. Health effects of exposure to natural arsenic in
417        groundwater and coal in China: An overview of occurrence. *Environ. Health*
418        *Perspect.* **115**, 636-642 (2007).

419  15.  Shrestha, R. R. *et al.* Groundwater arsenic contamination, its health impact and
420        mitigation program in Nepal. *Environ. Sci. Health, Part A* **38**, 185-200 (2003).

421  16.  Polya, D. A. *et al.* Arsenic hazard in shallow Cambodian groundwaters. *Mineral.*
422        *Mag.* **69**, 807-823 (2005).

423  17.  Berg, M. *et al.* Magnitude of arsenic pollution in the Mekong and Red River
424        Deltas - Cambodia and Vietnam. *Sci. Total Environ.* **372**, 413-425 (2007).

425  18.  Buschmann, J., Berg, M., Stengel, C. & Sampson, M. L. Arsenic and Manganese
426      Contamination of Drinking Water Resources in Cambodia: Coincidence of Risk
427      Areas with Low Relief Topography. *Environ. Sci. Technol.* **41**, 2146–2152 (2007).

428  19.  Berg, M. *et al.* Arsenic contamination of groundwater and drinking water in
429      Vietnam: A human health threat. *Environ. Sci. Technol.* **35**, 2621-2626 (2001).

430  20.  Buschmann, J. *et al.* Contamination of drinking water resources in the Mekong
431      delta floodplains: Arsenic and other trace metals pose serious health risks to
432      population. *Environ. Int.* **34**, doi:10.1016/j.envint.2007.12.025 (2008).

433  21.  Ford, R. G., Fendorf, S. & Wilkin, R. T. Introduction: Controls on arsenic transport
434      in near-surface aquatic systems. *Chem. Geol.* **228**, 1-5 (2006).

435  22.  Harvey, C. F. *et al.* Arsenic mobility and groundwater extraction in Bangladesh.
436      *Science* **298**, 1602-1606 (2002).

437  23.  Islam, F. S. *et al.* Role of metal-reducing bacteria in arsenic release from Bengal
438      delta sediments. *Nature* **430**, 68-71 (2004).

439  24.  Rowland, H. A. L. *et al.* The control of organic matter on microbially mediated
440      iron reduction and arsenic release in shallow alluvial aquifers, Cambodia.
441      *Geobiology* **5**, 281-292 (2007).

442  25.  Berg, M. *et al.* Hydrological and sedimentary controls leading to arsenic
443      contamination of groundwater in the Hanoi area, Vietnam: The impact of iron-
444      arsenic ratios, peat, river bank deposits, and excessive groundwater abstraction.
445      *Chem. Geol.* **249**, 91-112 (2008).

446  26.  Smedley, P. L. & Kinniburgh, D. G. A review of the source, behaviour and
447      distribution of arsenic in natural waters. *Appl. Geochem.* **17**, 517-568 (2002).

448  27.  Stute, M. *et al.* Hydrological control of As concentrations in Bangladesh
449      groundwater. *Water Resour. Res.* **43**(2007).

450  28.  Twarakavi, N. K. C. & Kaluarachchi, J. J. Arsenic in the shallow ground waters of
451      conterminous United States: Assessment, health risks, and costs for MCL
452      compliance. *J. Am. Water Resour. Assoc.* **42**, 275-294 (2006).

453  29.  Amini, M. *et al.* Statistical modeling of global geogenic arsenic contamination in
454      groundwater. *Environ. Sci. Technol.* **42**, 3669-3675 (2008).

455  30.  Goovaerts, P. *et al.* Geostatistical modeling of the spatial variability of arsenic in
456      groundwater of southeast Michigan. *Water Resour. Res.* **41**(2005).

457    31.  Lee, J. J., Jang, C. S., Wang, S. W. & Liu, C. W. Evaluation of potential health
458          risk of arsenic-affected groundwater using indicator kriging and dose response
459          model. *Sci. Total Environ.* **384**, 151-162 (2007).

460    32.  Hossain, F., Hill, J. & Bagtzoglou, A. C. Geostatistically based management of
461          arsenic contaminated ground water in shallow wells of Bangladesh. *Water*
462          *Resour. Manag.* **21**, 1245-1261 (2007).

463    33.  Twarakavi, N. K. C. & Kaluarachchi, J. J. Aquifer vulnerability assessment to
464          heavy metals using ordinal logistic regression. *Ground Water* **43**, 200-214 (2005).

465    34.  Ayotte, J. D. *et al.* Modeling the probability of arsenic in groundwater in New
466          England as a tool for exposure assessment. *Environ. Sci. Technol.* **40**, 3578-
467          3585 (2006).

468    35.  Reading, H. G. (ed.) *Sedimentary Environments: Processes, Facies and*
469          *Stratigraphy*, thrid edition (Blackwell Science, Oxford, 1996).

470    36.  Tan, K. H. *Environmental Soil Science*, second edition (M. Dekker, New York,
471          2000).

472    37.  Hori, K. *et al.* Delta initiation and Holocene sea-level change: example from the
473          Song Hong (Red River) delta, Vietnam. *Sediment. Geol.* **164**, 237-249 (2004).

474    38.  Kirk, M. F. *et al.* Bacterial sulfate reduction limits natural arsenic contamination in
475          groundwater. *Geology* **32**, 953-956 (2004).

476    39.  Lowers, H. A. *et al.* Arsenic incorporation into authigenic pyrite, bengal basin
477          sediment, Bangladesh. *Geochim. Cosmochim. Acta* **71**, 2699-2717 (2007).

478    40.  Stanley, D. J. & Warne, A. G. Worldwide Initiation of Holocene Marine Deltas by
479          Deceleration of Sea-Level Rise. *Science* **265**, 228-231 (1994).

480    41.  Wosten, J. H. M. *et al.* Interrelationships between hydrology and ecology in fire
481          degraded tropical peat swamp forests. *Int. J. Water Resour. Dev.* **22**, 157-174
482          (2006).

483    42.  Giesen, W. *Causes of Peatswamp Forest Degradation in Berbak National Park*
484          *and Recommendations for Restoration, Water for Food & Ecosystems*
485          *Programme*. (Arcadis Euroconsult, Arnhem, Holland, 2004).

486    43.  Tun, K. M. A. *Report on the Assessment of Arsenic Content in Groundwater and*
487          *the Prevalence of Arsenicosis in Thabaung and Kyonpyaw Townships,*
488          *Ayeyarwaddy Division*. (Department of Medical Research, Yangoon, Myanmar,
489          2002).

490    44.  Kohnhorst, A. Arsenic in groundwater in selected countries in south and
491          southeast Asia: A review. *J. Trop. Med. Parasitol.* **28**, 73-82 (2005).

492    45.  Tanabe, S. *et al.* Stratigraphy and Holocene evolution of the mud-dominated
493          Chao Phraya delta, Thailand. *Quat. Sci. Rev.* **22**, 789-807 (2003).

494    46.  Metivier, F., Gaudemer, Y., Tapponnier, P. & Klein, M. Mass accumulation rates
495          in Asia during the Cenozoic. *Geophys. J. Int.* **137**, 280-318 (1999).

496    47.  Robinson, R. A. J. *et al.* The Irrawaddy River sediment flux to the Indian Ocean:
497          The original nineteenth-century data revisited. *J. Geol.* **115**, 629-640 (2007).

498    48.  Berg, M. *et al.* Arsenic removal from groundwater by household sand filters:
499          Comparative field study, model calculations, and health benefits. *Environ. Sci.*
500          *Technol.* **40**, 5567-5573 (2006).

501    49.  Kleinbaum, D. G. & Klein, M. *Logistic Regression: A Self-Learning Text*, second
502          edition (Springer-Verlag, New York, 2002).

503    50.  Pal, T., Mukherjee, P. K., Sengupta, S., Bhattacharyya, A. K. & Shome, S.
504          Arsenic pollution in groundwater of West Bengal, India - An insight into the
505          problem by subsurface sediment analysis. *Gondwana Res.* **5**, 501-512 (2002).

506

507

508

509    Correspondence and requests for materials should be addressed to Michael Berg.

510

517

518

519    **Supplementary Information** accompanies this paper

520

521    **Competing financial interests statement**

522    The authors declare that they have no competing financial interests.

523 **Figure captions**

524 **Figure 1. Uniformly classified geological map of Southeast Asia.** It indicates

525 seven different sedimentary depositional environments in the mapped countries of

526 Bangladesh, Myanmar, Thailand, Cambodia, Vietnam, and Sumatra (Indonesia).

527

528 **Figure 2. Model classification results.** The graph shows the sensitivity (true

529 positives) and specificity (true negatives) of the model for different probability cutoff

530 values. The full classification table is provided in the Supplementary Information

531 (Table S4). A probability threshold of 0.4 was applied to delineate low- and high-risk

532 areas in the binary risk maps shown in Figures 3b, 3d, and 4b (see Methods).

533

534 **Figure 3. Modelled probability of As concentrations exceeding 10 µgL$^{-1}$ under**

535 **reducing aquifer conditions. a**, Continuous probability map of Bangladesh. **b**, Binary

536 map of Bangladesh (probability threshold 0.4) indicating high- and low-risk areas

537 overlain by aggregated As concentrations. Areas where groundwater is mainly drawn

538 from Pleistocene aquifers are sketched in brown. **c**, Continuous probability map of the

539 Mekong delta (Cambodia and Vietnam). **d**, Binary risk map of the Mekong delta

540 overlain by aggregated As concentrations. **e**, Continuous probability map of the Red

541 River delta (Vietnam). **f**, Continuous probability map of the Irrawaddy delta (Myanmar)

542 and Chao Phraya basin (Thailand).

543

544 **Figure 4. Maps of the model verification study area in Southeast Sumatra. a**,

545 Map of whole Sumatra depicting the probability of groundwater As concentrations

546 exceeding 10 µg L$^{-1}$. The corresponding colour code is given in Figure 3. **b**, Binary risk

547 map (probability cuttoff 0.4) and As concentrations measured for this study in the

548 vicinity of Palembang (South Sumatra). Swampy areas (high-risk area) are scarcely

549    populated and the number of sampled tubewells in this area is therefore limited. **c**,

550    Geological map (source: see Supplementary Table S1).

551

552 **Table 1. Results of logistic regression analysis.** Weighting coefficients of the

553 independent variables (λ) were used to calculate probabilities of As contamination.

554 Wald values (%) indicate the relative importance of the variables, and p-values the

555 absolute significance, where a value <0.05 indicates a significance of at least 95%.

556 Variables that were not statistically significant (based on a 95% level) were excluded

557 from the model (Table S2), with exception of medium textured soils and silt in subsoil

558 (see Methods). Excluded variables are: tidal deposits, other Holocene deposits, pre-

559 Holocene deposits (see Figure 1), coarse textured soils, soil sand and clay contents,

560 and climate.

561

| Variables | | λ | Wald | p-value |
|---|---|---|---|---|
| Sedimentary depositional environments | Deltaic deposits | 1.65 | 71.55 | <0.001 |
| | Organic-rich deposits | 1.11 | 72.20 | <0.001 |
| | Alluvial deposits | 0.55 | 12.51 | <0.001 |
| | Floodplain deposits | – 0.95 | 2.32 | 0.010 |
| Soil variables | Medium textured soils | 0.60 | 7.80 | 0.128 |
| | Fine textured soils | 0.24 | 1.00 | 0.005 |
| | Silt in subsoil | 0.10 | 6.62 | 0.317 |
| | Silt in topsoil | – 0.09 | 6.28 | 0.012 |
| – | Regression constant | – 1.54 | 20.67 | <0.001 |

562

563