

3 **Is correlation dimension a reliable proxy for the number**
4 **of dominant influencing variables for modeling risk of arsenic**
5 **contamination in groundwater?**

6 Jason Hill · Faisal Hossain · Bellie Sivakumar

7
8 © Springer-Verlag 2006

9 **Abstract** The correlation dimension (CD) of a time
10 series provides information on the number of dominant
11 variables present in the evolution of the underlying
12 system dynamics. In this study, we explore, using
13 logistic regression (LR), possible physical connections
14 between the CD and the mathematical modeling of
15 risk of arsenic contamination in groundwater. Our
16 database comprises a large-scale arsenic survey con-
17 ducted in Bangladesh. Following the recommendation
18 by Hossain and Sivakumar (Stoch Environ Res Risk
19 Assess 20(1–2):66–76, 2006a), who reported CD values
20 ranging from 8 to 11 for this database, 11 variables are
21 considered herein as indicators of the aquifer’s geo-
22 chemical regime with potential influence on the arsenic
23 concentration in groundwater. A total of 2,048 possible
24 combinations of influencing variables are considered as
25 candidate LR risk models to delineate the impact of
26 the number of variables on the prediction accuracy of
27 the model. We find that the uncertainty associated with
28 prediction of wells as safe and unsafe by LR risk model
29 declines systematically as the total number of influ-
30 encing variables increases from 7 to 11. The sensitivity

of the mean predictive performance also increases 31
noticeably for this range. The consistent reduction in 32
predictive uncertainty coupled with the increased sen- 33
sitivity of the mean predictive behavior within the 34
universal sample space exemplify the ability of CD to 35
function as a proxy for the number of dominant influ- 36
encing variables. Such a rapid proxy, based on non- 37
linear dynamic concepts, appears to have considerable 38
merit for application in current management strategies 39
on arsenic contamination in developing countries, 40
where both time and resources are very limited. 41

Keywords Nonlinear deterministic dynamics and 42
chaos · Correlation dimension · Arsenic contamination · 43
Logistic regression · Groundwater · Bangladesh 44

1 Introduction 45

Since the large-scale discovery of arsenic contamina- 46
tion in the alluvial Ganges aquifers of Bangladesh, 47
numerous studies have been conducted to better 48
understand the spatial variability of the contamination 49
scenario (e.g., Biswas et al. 1998; Burgess et al. 2000; 50
McArthur et al. 2001, 2004; Harvey et al. 2002; Muk- 51
herjee and Bhattacharya 2002; van Geen et al. 2003; 52
Yu et al. 2003; Ahmed et al. 2004; Hossain et al. 2006a, 53
b). Most of these studies have addressed the ‘spatial’ 54
pattern of arsenic using geo-statistical tools and the 55
classical notion of linear stochastic dynamics. For 56
example, in the first country-wide study toward spatial 57
(horizontal) characterization of the arsenic calamity, 58
conducted by the British Geological Survey (BGS) in 59
collaboration with the Department of Public Health 60
and Engineering (DPHE) of Bangladesh (hereafter 61

A1 J. Hill
A2 Department of Civil and Environmental Engineering,
A3 Tri-State University, 1 University Avenue, Angola
A4 IN 46703, USA

A5 F. Hossain (✉)
A6 Department of Civil and Environmental Engineering,
A7 Tennessee Technological University, Box 5015, Cookeville,
A8 TN 38505-0001, USA
A9 e-mail: fhossain@tntech.edu

A10 B. Sivakumar
A11 Griffith School of Engineering, Griffith University, Nathan,
A12 QLD 4111, Australia

62 called ‘BGS-DPHE’), an application of kriging (Journel and Huijbregts 1978) was reported to provide the
 63 ‘best’ estimate of the whole nation’s arsenic field at the
 64 regional scale with limited sampling information. The
 65 BGS-DPHE investigation involved the assumption
 66 that the arsenic concentration could be treated as a
 67 ‘regionalized’ linear stochastic random variable in
 68 space.
 69

70 It must be noted, however, that arsenic in ground-
 71 water is not a purely random occurrence and that
 72 (hidden) order and determinism may also exist, just as
 73 they do in any other natural or man-made phenome-
 74 non. Arguing that there existed profound geological
 75 and geochemical factors, with possible order, control-
 76 ling arsenic contamination dynamics (for details, see
 77 Hossain and Sivakumar 2006a; McArthur et al. 2004;
 78 Zheng et al. 2004), Hossain and Sivakumar (2006b)
 79 suggested that it was no longer defensible for the
 80 scientific community to continue to use purely geo-
 81 statistical (linear stochastic) approaches as stand-alone
 82 techniques for its spatial interpolation. Our under-
 83 standing of the role played by these physical factors in
 84 arsenic contamination of groundwater continues to be
 85 enhanced from recent studies by, for example, Zheng
 86 et al. (2004), Akai et al. (2004) and Ahmed et al.
 87 (2004). Traditional geostatistical tools are a ‘pattern-
 88 filling’ scheme based on the spatial correlation exhib-
 89 ited by two points in space separated by a lag h . This
 90 approach simplifies the spatial patterns manifested by
 91 the complex interactions between geology and time-
 92 sensitive fluid flow dynamics (Christakos and Li 1998).
 93 Concerns on the use of purely stochastic approaches
 94 and potential for alternative ones have been echoed by
 95 a few other studies as well (e.g., Faybishenko 2002;
 96 Sivakumar 2004a; Sivakumar et al. 2005).

97 On the premise that the current ensemble of pro-
 98 posed ‘theories’ in scientific literature explaining
 99 arsenic mobility (e.g., Burgess et al. 2000; McArthur
 100 et al. 2001; Harvey et al. 2002; van Geen et al. 2003)
 101 can, in principle, be mathematically represented as the
 102 cumulative effect of a finite number of dominant pro-
 103 cesses comprising three or more partial differential
 104 equations, Hossain and Sivakumar (2006a) verified the
 105 existence of nonlinear deterministic and chaotic
 106 dynamic behavior in the spatial pattern of arsenic
 107 contamination in shallow wells (depth < 150 m) in
 108 Bangladesh. Employing the Grassberger–Procaccia
 109 correlation dimension (CD) algorithm (Grassberger
 110 and Procaccia 1983), their analysis revealed CD values
 111 (i.e., saturation of correlation exponents and a mani-
 112 festation of ‘determinism’) ranging anywhere from 8 to
 113 11 depending on the region and geology (see, for
 114 example, Fig. 1). Their findings suggested that the

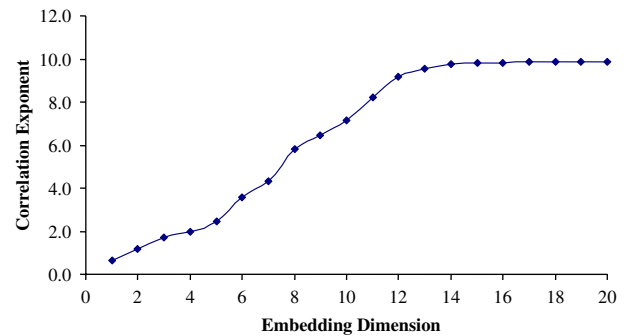


Fig. 1 Relationship between Correlation Exponent and Embedding Dimension for the whole Bangladesh based on BGS-DPHE (2001) arsenic data from shallow wells (after Hossain and Sivakumar 2006a)

115 arsenic contamination dynamics in space, from a cha-
 116 otic dynamic perspective, was a medium- to high-
 117 dimensional problem. While it is encouraging to note
 118 that the nonlinear CD analysis can reflect the influence
 119 of regional geology (and other factors) on arsenic
 120 contamination dynamics, the usefulness of the CD and
 121 other nonlinear deterministic dynamic techniques to
 122 understand the physics of the actual arsenic contami-
 123 nation phenomenon is far from clear, as explained
 124 next.

125 It is well known that the CD of (an attractor of) a
 126 time series generally provides information on the
 127 number of variables present in the evolution of the
 128 underlying system dynamics (e.g., Grassberger and
 129 Procaccia 1983; Hao 1984; Fraedrich 1986; Sivakumar
 130 2004b; Hossain and Sivakumar 2006a). However, cur-
 131 rent environmental literature is largely insufficient in
 132 the context of providing links between the CD and the
 133 actual physical mechanisms that take place in catch-
 134 ments/aquifers. While some studies have indeed con-
 135 ducted research in this direction, such have essentially
 136 been limited to the verification of the reliability of the
 137 CD estimate, and especially performed using nonlinear
 138 predictions of the respective time series. For example,
 139 Sivakumar et al. (2002c) investigated the reliability of
 140 the CD estimate of the monthly flow data observed at
 141 the Coaracy Nunes/Araguari River watershed in
 142 northern Brazil (see also Sivakumar et al. 2001a), using
 143 nonlinear local- (chaos theory-based) and global-
 144 (artificial neural networks-based) approximation tech-
 145 niques. The study, in fact, focused on the reliability of
 146 the CD in the context of short time series, since the
 147 data size requirement has been the primary subject of
 148 criticism on the reports of low-dimensional chaos in
 149 environmental time series (e.g., Ghilardi and Rosso
 150 1990; Schertzer et al. 2002; see also Sivakumar 2000,
 151 2005; Sivakumar et al. 2002a, for details). Similarly,
 152 nonlinear predictions of time series have served as the

153 basis, implicitly or explicitly, for verification of the CD
 154 estimate in other studies as well, albeit in different
 155 forms (e.g., Porporato and Ridolfi 1997; Lambrakis
 156 et al. 2000; Sivakumar et al. 2001b, 2002b).

157 With the encouraging results of their preliminary
 158 analysis (Hossain and Sivakumar 2006a) regarding the
 159 nonlinear deterministic nature of arsenic contamina-
 160 tion, Hossain and Sivakumar (2006b) subsequently
 161 discussed the potential role the nonlinear deterministic
 162 dynamic and related concepts can play in improving
 163 our understanding of arsenic contamination patterns in
 164 space. They especially highlighted their potential utili-
 165 ty in providing improved cost-effectiveness of envi-
 166 ronmental management in rural and resource-limited
 167 settings of developing countries, such as Bangladesh,
 168 Vietnam and India. In a related development, Serre
 169 et al. (2003) have reported that the spatial interpola-
 170 tion of arsenic contamination, if approached from the
 171 conventional paradigm of geostatistical mapping, can
 172 be challenging in Bangladesh as most of the variability
 173 in arsenic concentration occurs within short distances
 174 (2–5 km). Certainly acknowledging the fact that the
 175 traditional linear stochastic approaches had generally
 176 yielded fairly good and reliable results, Hossain and
 177 Sivakumar (2006b) also called for a much-needed
 178 change in the current state-of-the-art for spatial inter-
 179 polation of arsenic contamination, stating that: ‘While
 180 there is no structural, or even philosophical, flaw in
 181 using the conventional geo-statistical approach, there is
 182 indeed ample room to argue that the geo-statistical
 183 treatment of arsenic contamination in space as a
 184 regionalized random (or stochastic) variable may con-
 185 stitute only an incomplete analysis of its spatial vari-
 186 ability (even if system-dependent). Incompleteness can
 187 potentially arise from the fact that geo-statistics often
 188 fails to recognize the random looking but deterministic
 189 behavior that may be present due to self-similar (scale-
 190 invariant) factors in the continuum of the sub-surface.’

191 In essence, Hossain and Sivakumar (2006b) argued
 192 for the need to couple/integrate the linear and nonlinear
 193 concepts/tools, whenever and wherever deemed neces-
 194 sary or appropriate [see also Sivakumar (2004b) for an
 195 example of possible integration of different concepts/
 196 methods for environmental modeling]. This, however, is
 197 easier said than done, since there is still some convincing
 198 needed, going by the criticisms, on the utility of the
 199 relatively new nonlinear deterministic dynamic con-
 200 cepts for arsenic contamination and other environmen-
 201 tal problems in the first place. Roughly speaking, the
 202 nonlinear analyses and results need to be verified using
 203 the conventional linear techniques, so as to first bring
 204 reconciliation between linear and nonlinear concepts
 205 and then to bridge the gap between them. With partic-

ular reference to the study by Hossain and Sivakumar 206
 (2006a), this should obviously start with the verification 207
 of the CD values obtained for the arsenic concentration 208
 data using any of the available linear tools. 209

210 In this spirit, we herein explore possible physical
 211 connections between the CD and the mathematical
 212 modeling of risk of arsenic contamination in ground-
 213 water by applying (the linear) logistic regression (LR)
 214 risk assessment technique. Using 11 potentially influ-
 215 encing variables that largely define the geochemical
 216 regime of aquifers and, hence, the variability of arsenic
 217 concentration, we attempt to provide a possible
 218 insightful evidence that the CD can be a proxy for the
 219 number of dominant influencing variables required in
 220 an LR risk model to optimally predict risk of arsenic
 221 contamination at non-sampled wells. To the best of our
 222 knowledge, such an insight, although preliminary,
 223 constitutes an important finding, with potential impli-
 224 cations on the reduction of uncertainty of risk maps
 225 produced from conventional (linear stochastic) para-
 226 digms. Even though we pursue this task primarily from
 227 a data-based perspective, a larger goal of our mission is
 228 to encourage greater interactions with the research
 229 community traditionally engaged in a more mechanis-
 230 tic understanding of arsenic contamination. We believe
 231 that such interactions can play a vital role in the inte-
 232 gration of non-linear deterministic dynamic concepts in
 233 future groundwater management protocols (discussed
 234 in detail later in the paper). In the sections that follow,
 235 we provide a systematic overview of our exploratory
 236 research to understand the value of CD in modeling
 237 risk of arsenic contamination.

2 Study region, data, and CD analysis 238

239 We choose to study arsenic contamination over the
 240 entire region of Bangladesh, as had been first surveyed
 241 by the BGS-DPHE (2001) study comprising 3,534
 242 wells. This is conducted in the manner similar to
 243 Hossain and Sivakumar (2006a) for estimating the CD
 244 values. The dataset is available (at the time of writing
 245 this manuscript) at [http://www.bgs.ac.uk/arsenic/ban-](http://www.bgs.ac.uk/arsenic/bangladesh/datadownload.htm)
 246 [gladesh/datadownload.htm](http://www.bgs.ac.uk/arsenic/bangladesh/datadownload.htm). Wells deeper than 150 m
 247 (and consistently below the safe limits) are excluded
 248 from the analysis, thus resulting in a set of 3,085
 249 shallow wells. While it is possible that such an exclu-
 250 sion of data based on depth may incur an added bias to
 251 our analyses on the application of CD, we believe, to
 252 the best of our knowledge, that the impact would be
 253 insignificant to alter the overall conclusions of our
 254 study, particularly when our goal is to demonstrate a
 255 proof-of-concept application of CD in deterministic



256 modeling. For details on the study region and data, the
257 reader is referred to the works of Hossain et al. (2006b)
258 and Hossain and Sivakumar (2006a).

259 The CD method employed by Hossain and Sivaku-
260 mar (2006a) used the correlation integral or function
261 (Grassberger and Procaccia 1983) for distinguishing
262 between chaotic and stochastic behaviors (more spe-
263 cifically, between low- and high-dimensional systems).
264 Although, traditional applications of the phase-space
265 reconstruction and the Grassberger–Procaccia algo-
266 rithms have been limited to data series in the contin-
267 uum of time (e.g., Takens 1981; Theiler 1987;
268 Rodriguez-Iturbe et al. 1989; Porporato and Ridolfi
269 1997; Sivakumar et al. 2001b, 2002c, 2005), Hossain and
270 Sivakumar (2006a) argued that there was no compelling
271 logic that disqualified its application to a data series in
272 space. Their CD analysis revealed positive evidence
273 regarding medium-to-high dimensional chaotic
274 dynamics in arsenic contamination in space, with a
275 country-wide dimension value ranging between 8 and
276 11. This subsequently led Hossain and Sivakumar
277 (2006a, b) to comment subjectively that the minimum
278 number of variables and hence the number of dominant
279 processes required to model the spatial variability of
280 arsenic contamination should also range from 8 to 11.

281 It is appropriate to mention, at this point, that
282 questions may be raised regarding the suitability of this
283 data set for CD analysis. Such questions may be related
284 to, among others, the data size (insufficient length) and
285 data quality (presence of noise), as these could
286 potentially influence the CD estimation (e.g., Neren-
287 berg and Essex 1990; Schreiber and Kantz 1996). These
288 issues, and also others, have been and continue to be
289 extensively discussed and debated in the literature,
290 including in the environmental sciences [e.g., Ghilardi
291 and Rosso 1990; Tsonis et al. 1994; Sivakumar et al.
292 1999, 2001b, 2002a, c; Sivakumar 2000, 2005; Schertzer
293 et al. 2002; see also Sivakumar (2004a) for a review].
294 Due to space limitations, and also to avoid unnecessary
295 deviation from the main focus of our study, we choose
296 not to discuss such issues, and consequently direct the
297 reader to the above studies and the numerous refer-
298 ences therein. We, however, would like to briefly
299 highlight a few points herein, in regards to the
300 reliability of the CD estimates for this data set reported
301 by Hossain and Sivakumar (2006a).

302 1. We are convinced that the data size, with 3,085
303 points, is more than sufficient to obtain reliable CD
304 estimates of arsenic contamination in space. In this
305 regard, we are particularly comforted by past
306 studies that have reported reliable CD estimates
307 for much smaller data sizes, albeit in the contin-

uum of time (e.g., Sivakumar 2000, 2005; Sivaku- 308
mar et al. 2002a, c). 309

2. While we do admit that the arsenic concentration 310
data are likely contaminated with noise (e.g., 311
measurement errors), we do not believe that it 312
significantly influences our CD estimates [see, for 313
example, Sivakumar et al. (1999)]. Even if it were 314
to influence, the result would be only an overesti- 315
mation of CD, not underestimation. Therefore, the 316
interpretations and conclusions by Hossain and 317
Sivakumar (2006a) regarding medium-to-high 318
dimensional chaotic pattern would not only stand 319
the test but also be more solidified. 320
3. Another factor possibly leading to underestimation 321
of CD is the presence of a large number of zeros 322
(or any one particular value) in the data set (e.g., 323
Tsonis et al. 1994). Since there are no zeros (or 324
repetition of a particular value) in the arsenic data 325
set, this problem is also completely eliminated. 326

3 Logistic regression 327

The method of LR has been extensively used in epide- 328
miological studies, and more recently, has become 329
a common technique in environmental research 330
on modeling risk of groundwater contamination 331
(Twarakavi and Kaluarachchi 2006). Common regres- 332
sion techniques, such as the classical linear regression, 333
relate the response variables to the influencing variables. 334
LR relates the probability of a response variable to be 335
greater than a threshold value (i.e., a risk) to a set of 336
influencing variables (Afifi and Clark 1984; Helsel and 337
Hirsch 1992). In an LR risk model, regression is linear 338
between the natural logarithm of the odds ratio for the 339
probability of response to be less than the threshold 340
value and influencing variables. Equation 1 mathemat- 341
ically summarizes the LR model used in this study: 342

$$\ln[p/(1-p)] = \text{logit}(p) = \alpha + \beta\mathbf{x} \quad (1)$$

where p is the probability of response to be greater 344
than the safety threshold, α is a constant, β is a vector 345
of slope coefficients, and \mathbf{x} is a vector of influencing 346
variables. For more details on the use of LR for 347
modeling risk of arsenic contamination, the reader is 348
referred to Twarakavi and Kaluarachchi (2006). 349

4 The potential influencing variables 350

Table 1 shows the influencing variables considered 351
herein for defining the geochemical regime of aquifers. 352

353 These variables were sampled by BGS-DPHE (2001)
 354 in Bangladesh. The minimum and maximum values of
 355 these variables (Table 1) indicate the range of vari-
 356 ability across Bangladesh. The variables chosen are: (1)
 357 depth of wells (m), (2) P (Phosphorus) (mg/L), (3) Fe
 358 (Iron) (mg/L), (4) Ba (Barium) (mg/L), (5) Mg (Mag-
 359 nesium) (mg/L), (6) Ca (Calcium) (mg/L), (7) SO₄
 360 (Sulfate) (mg/L), (8) Mean annual precipitation (mm/
 361 day), (9) Si (Silicon) (mg/L), (10) Na (Sodium) (mg/L),
 362 and (11) Mn (Manganese) (mg/L). Although our
 363 choice of variables is primarily dictated by literature
 364 reports on the causes of arsenic mobility (e.g., Welch
 365 et al. 2000; Harvey et al. 2002; van Geen et al. 2003;
 366 McArthur et al. 2004; Zheng et al. 2004) and the
 367 availability of reliable data, we must also point out to
 368 the reader that the selection herein is governed purely
 369 from a data-based and qualitative paradigm. As indi-
 370 cated earlier, the larger goal of our study is to
 371 encourage greater interactions between the research
 372 communities on mechanistic modeling of arsenic con-
 373 tamination and non-linear dynamic analysis. We admit
 374 that such a data-based selection without a deeper
 375 physical regard for the pertinent mechanics and geo-
 376 chemistry of contamination (as appropriate for
 377 Bangladesh) may have potential limitations. However,
 378 we also believe that such potential limitations alone
 379 should not hamper our ability to investigate the use-
 380 fulness of the CD value, and particularly so when our
 381 intention is to primarily conduct a preliminary explo-
 382 ration. We believe that if there is a weakness in our
 383 choice of potential influencing variables, as may be
 384 revealed in our results, it only lends greater credibility
 385 to our mission in inviting the research community on
 386 arsenic contamination to interact more closely with the
 387 non-linear deterministic dynamic research community.

388 As a preliminary step, we first conduct the Spear-
 389 man's Rank Correlation Coefficient test for these

Table 1 The selected influencing variables for Logistic Regression Modeling

| Variable | Mean | Minimum | Maximum |
|---------------------------|--------|---------|----------|
| Well depth (m) | 60.550 | 0.600 | 362.000 |
| Ba (ppb) | 87.340 | 2.000 | 1360.000 |
| Ca (mg/L) | 51.590 | 0.100 | 366.000 |
| Fe (mg/L) | 3.353 | 0.005 | 61.000 |
| Mg (mg/L) | 20.750 | 0.040 | 305.000 |
| Mn (mg/L) | 0.555 | 0.001 | 9.980 |
| Na (mg/L) | 88.936 | 0.700 | 2700.000 |
| P (mg/L) | 0.765 | 0.100 | 18.900 |
| Si (mg/L) | 20.519 | 0.030 | 45.200 |
| SO ₄ (mg/L) | 5.917 | 0.200 | 753.000 |
| Annual precipitation (cm) | 86.001 | 25.350 | 596.140 |
| As (ppb) ¹ | 55.205 | 0.500 | 1660.000 |

¹ Arsenic (As) is the dependent variable in the LR risk model

390 selected variables to identify their non-linear depen-
 391 dence with arsenic concentration. Because all possible
 392 combinations of influencing variables are considered
 393 during LR modeling of contamination risk (discussed
 394 next), results from the Spearman's test are not used in
 395 the ranking of the variables according to the order of
 396 influence. The precipitation data are obtained from the
 397 Bangladesh Meteorological Department (BMD) and
 398 Bangladesh Water Development Board (BWDB). The
 399 data are derived from a network of 100 recording
 400 rainfall gauges that registered less than 5% missing
 401 data for the year 2000. The choice of precipitation as
 402 an influencing variable is governed by reports that
 403 groundwater pumping for irrigation and recharge could
 404 be one of the causes of arsenic mobility in the shallow
 405 geologic stratum (see Harvey et al. 2002). Because
 406 recharge data are not readily available for our study,
 407 we choose mean rainfall as a proxy indicator of
 408 recharge of aquifers. For consistency, we select pre-
 409 cipitation data pertaining to the year 2000 when the
 410 BGS-DPHE (2001) survey was completed. The mean
 411 annual rainfall value for each well is computed by the
 412 method of Thiessen Polygons using the ArcGISTM
 413 software (Ormsby et al. 2004).

5 Method of assessment

414 The dataset is divided randomly into two equal halves,
 415 with one half being employed for LR risk model cali-
 416 bration and the other half for validation. This random
 417 selection procedure is repeated 25 times within a
 418 Monte Carlo (MC) framework to assess the mean
 419 performance of the LR model. Using one-half of each
 420 randomly selected dataset, calibration of the LR model
 421 coefficients, α and β , is performed using ordinary least
 422 squares technique for a safety threshold of 50 ppb
 423 (Bangladesh limit). In the calibration phase, the ' p '
 424 values in Eq. 1 are assigned 0–1 binary values
 425 depending on the measured concentration of arsenic
 426 ($p = 1$ for exceeding the safety threshold; $p = 0$ for
 427 being below the threshold). During the validation
 428 phase, the LR model is assessed in terms of its ability
 429 to successfully predict contamination in 0–1 binary
 430 terms according to the safety threshold at non-sampled
 431 wells (i.e., over the other half of the dataset not used in
 432 calibration of the LR risk models). For this, we employ
 433 the notion of contamination risk associated with a pre-
 434 assigned probability (i.e., in this case, $p = 0.9$). For
 435 example, if the well is predicted by the LR risk model
 436 as unsafe with $p = 0.85$ for a given safety threshold,
 437 then that well would be flagged uncontaminated
 438 according to the high risk criterion of $p = 0.9$. The
 439



440 predictive power of the LR risk model for a given
 441 number of influencing variables is quantified by the
 442 probability of successful detection of a well's status as
 443 contaminated or uncontaminated at untested well
 444 locations. It should be noted that the pre-assignment of
 445 a probability value to denote risk category as high(low)
 446 is purely subjective and will linearly scale up(down)
 447 the predictive behavior of LR model without altering the
 448 response pattern to the number of influencing vari-
 449 ables. Hence, such a subjective assignment is consid-
 450 ered acceptable within the overall scheme of our study
 451 as the objective is to delineate the impact of the
 452 number of potentially influencing variables and not on
 453 the LR risk model performance per se.

454 The specific question we explore, using LR, in our
 455 study is: 'Is CD a reliable proxy for the number of
 456 dominant variables required to predict risk of arsenic
 457 contamination in groundwater?' We consider all possi-
 458 ble combinations of influencing variables from the
 459 total set of 11 as candidate LR models. This results in
 460 2,048 LR risk models being evaluated. Each evaluation
 461 is repeated 25 times within the MC framework and the
 462 mean and range of LR model prediction assessed. For
 463 a given number of influencing variables, the mean
 464 signified the most probable LR model performance
 465 while the range is an indicator of predictive uncertainty
 466 to expect. It is important to note that the predictive
 467 uncertainty (or range) has important implications for
 468 model complexity and parameter optimization. The
 469 wider the uncertainty, the more challenging naturally
 470 would be the optimization to converge to the best LR
 471 model configuration. We discuss this in more detail in
 472 the next section.

473 6 Results and discussion

474 Figure 2 shows the variation of probability of success-
 475 ful detection of wells, or the fraction of validation set
 476 wells correctly detected (as contaminated/uncontami-
 477 nated at the 50 ppb limit) as a function of the total
 478 number of influencing variables (Table 1) in the LR
 479 model. Basically, the terms 'contaminated/uncontami-
 480 nated' or 'unsafe/safe' refer to the wells with arsenic
 481 concentration exceeding/less than 50 ppb. The mean
 482 predictive ability (shown in red circles, Fig. 2) of the
 483 LR risk model, while remaining insensitive to number
 484 of influencing variables in the ranges of 1–7 variables,
 485 is found to noticeably increase in sensitivity when the
 486 number of variables is greater than 7. A systematic
 487 reduction in the predictive uncertainty is also observed
 488 as the number of variables is increased from 7 to 11
 489 (see Fig. 3). The probability of successful detection is

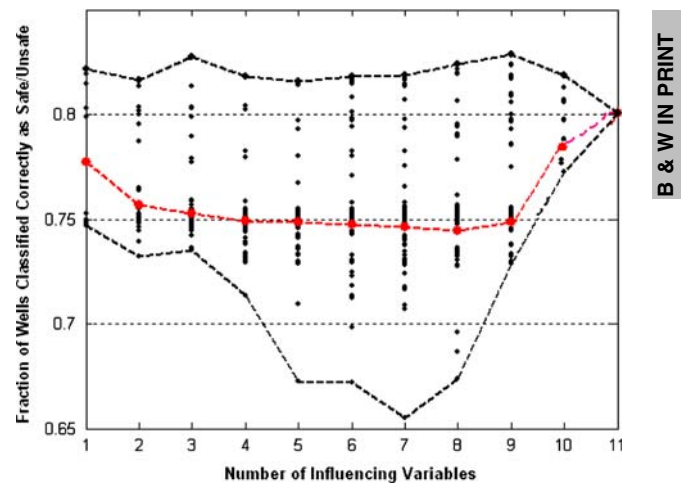


Fig. 2 Variation of fraction of wells correctly classified by LR model as safe/unsafe (i.e., probability of successful detection) with the number of influencing variables. The larger black circles with dashed line in the middle indicate mean values. The upper and lower dashed lines in black indicate the range of 25 Monte Carlo realizations for a given number of variables

490 shown for the mean of the 25 MC simulations on the y-
 491 axis of Fig. 2. Finally, we observe the best performance
 492 of the LR model when the number of influencing
 493 variables is 11. (Note that the lines all converge here to
 494 a point when the number of variables is 11 because the
 495 total number of possible LR model combinations is
 496 one. This observation should not be construed as an
 497 indication of no uncertainty for an LR model with 11
 498 variables, but rather as an indication of the last point of
 499 complex modeling within a set of 11 variables where
 500 only one possible model can be constructed). As
 501 evident from Figs. 2, 3, an a priori inclusion of CD
 502 value in assigning the minimum LR model complexity
 503 appears to guarantee global optimization of the model
 504 configuration with a considerably higher degree of

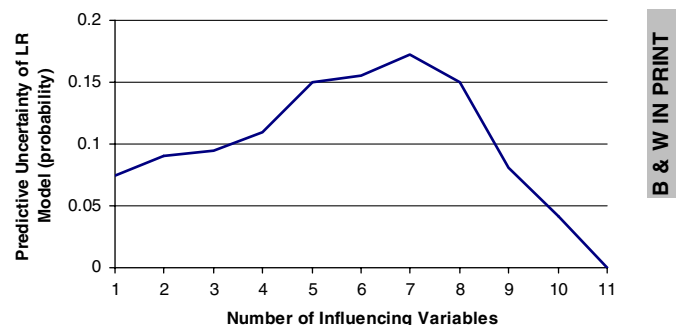


Fig. 3 Predictive uncertainty in terms of probability of successful detection (i.e., the range between upper and lower limits in Fig. 2) as a function of the number of influencing variables. (Note: the value when the number of influencing variables is 11 should be ignored.)

505 success. This empirical observation indicates consis- 554
 506 tency with the CD concept, according to which the 555
 507 inclusion of any additional variable deemed influential 556
 508 on the dynamics should yield either an improvement or 557
 509 simply no change (unless otherwise significantly influ- 558
 510 enced by noise) [see also, for example, Sivakumar et al. 559
 511 (2001b, 2002c)]. Overall, this preliminary finding seems 560
 512 to offer credence to the hypothesis that an acceptable 561
 513 number of variables to model the risk of arsenic con- 562
 514 tamination should range from 7 or 8 to 11 [The LR 563
 515 results also seem to strengthen our earlier point that 564
 516 the CD estimates reported by Hossain and Sivakumar 565
 517 (2006a) may only be an overestimation due to the 566
 518 presence of noise, if any, and not an underestimation]. 567

519 Currently, there are a number of maps available that 568
 520 characterize the probability of arsenic contamination 569
 521 in non-sampled regions based on kriging [see BGS- 570
 522 DPHE (2001) and McArthur et al. (2001), for exam- 571
 523 ple]. Preliminary findings of our study imply that an 572
 524 injection of the chaotic dynamic approach of LR 573
 525 modeling with variables equaling the CD could exped- 574
 526 ite refinement of the map toward reduction of 575
 527 uncertainty in risk of contamination at non-sampled 576
 528 locations than what would have otherwise been possi- 577
 529 ble by the kriging method alone. Although CD does 578
 530 not offer any physical insight on the variables that need 579
 531 to be chosen or the nature of their integration in risk 580
 532 assessment models, prior knowledge as a proxy for an 581
 533 acceptable number of variables required can be a 582
 534 valuable information that can potentially save consid- 583
 535 erable time during a rapid assessment of arsenic con- 584
 536 tamination for remediation management. 585

537 7 Conclusion 586

538 While applications of nonlinear dynamic concepts, 587
 539 such as the CD method, are gaining momentum in 588
 540 environmental sciences, their usefulness to understand 589
 541 the actual physical mechanisms occurring in our 590
 542 catchments and aquifers remains unclear. With the 591
 543 encouraging results reported recently by Hossain and 592
 544 Sivakumar (2006a) regarding the possible nonlinear 593
 545 deterministic nature of arsenic contamination phe- 594
 546 nomenon in Bangladesh (with CD values ranging from 595
 547 8 to 11), we herein have explored the possible physical 596
 548 connection between the CD and the mathematical 597
 549 modeling of risk of arsenic contamination in ground- 598
 550 water. We considered the LR model, with an aim to 599
 551 link the nonlinear CD technique with a linear analysis
 552 technique. Using 11 potential influencing variables that
 553 largely dictate the variability of arsenic concentration,

we observed that the CD may function as an accept- 554
 able proxy for the number of variables required in the 555
 LR model to accurately predict arsenic contamination 556
 at non-sampled wells. Given this preliminary finding, 557
 we believe it is time we considered more comprehen- 558
 sive investigations to assess the true merit of non-linear 559
 deterministic paradigms in conjunction with the more 560
 conventional linear stochastic methods, such as kriging, 561
 for reducing uncertainty of risk mapping for ground- 562
 water contamination in resource poor countries. 563

This study is not without its share of limitations. The 564
 two primary limitations that should be highlighted 565
 herein, so that findings from this study are not quoted 566
 out of context, are: (1) selection of potential influenc- 567
 ing variables from a purely data-based paradigm; and 568
 (2) maximum number of influencing variables being 569
 only 11 and barely exceeding the range of CD values. 570
 An earlier section (on 'The potential influencing vari- 571
 ables') in this paper has already discussed in detail the 572
 first limitation with a qualified disclaimer. On the sec- 573
 ond limitation, we unconditionally recognize that the 574
 value of CD could have been more convincingly 575
 demonstrated had more than 11 potential influencing 576
 variables been analyzed. However, inclusion of a 577
 higher number of variables is easier said than done, 578
 since there is paucity of quality-controlled data in a 579
 rural setting like Bangladesh. For example, an influ- 580
 encing variable such as soil cover is expected to influ- 581
 ence recharge and to ultimately affect the water table 582
 fluctuations, which may consequently be responsible 583
 for the mechanism that mobilizes arsenic (Twarakavi 584
 and Kaluarachchi 2006). However, such data are hard 585
 to obtain for the case of Bangladesh on a large scale. 586
 We believe that inclusion of a larger set of geochemical 587
 data is an important area of future study where we, as 588
 members of the non-linear deterministic community, 589
 should depend on effective feedback from the com- 590
 munity engaged in mechanistic understanding of ar- 591
 senic contamination in order to secure a more 592
 complete and appropriate dataset for CD integration. 593
 It must be noted, therefore, that more detailed studies 594
 are needed to verify the true limitations and strengths 595
 of the CD approach to designing LR models for rapid 596
 assessment of risk of arsenic contamination. Investi- 597
 gations in this direction are already underway, details 598
 of which will be reported elsewhere. 599

References 600

- Afifi AA, Clark V (1984) Logistic regression in computer-aided 601
 multivariate analysis. Lifetime Learning Publications, Bel- 602
 mont 603

- 604 Ahmed KM, Bhattacharya P, Hasan MA, Akhter SH, Alam
605 SMM, Bhuyian MA, Imam MB, Khan AA, Sracek O (2004)
606 Arsenic enrichment in groundwater of the alluvial aquifers
607 in Bangladesh: an overview. *Appl Geochem* 19:181–200
608 Akai J, Izumi K, Fukuhara H, Masuda H, Nakano S, Yoshimura
609 T, Ohfuji H, Anawar MH, Akai K (2004) Mineralogical and
610 geomicrobiological investigations on groundwater arsenic
611 enrichment in Bangladesh. *Appl Geochem* 19:215–230
612 Biswas BK, Dhar RK, Samantha G, Mandal BK, Chakraborti D,
613 Faruk I, Islam KS, Chowdury M, Islam A, Roy S (1998)
614 Detailed study report of Samta, one of the arsenic-affected
615 villages of Jessore District, Bangladesh. *Curr Sci* 74:134–145
616 Burgess WG, Burren M, Perrin J, Ahmed KM (2000) Constraints
617 on sustainable development of arsenic-bearing aquifers in
618 southern Bangladesh. Part 1: A conceptual model of arsenic
619 in the aquifer. In: Hiscock, Rivett, Davison (eds) *Sustain-
620 able groundwater development*, vol 193. Geological Society
621 of London Special Publication, pp 145–163
622 Christakos G, Li X (1998) Bayesian maximum entropy analysis
623 and mapping: a farewell to kriging estimators? *Math Geol*
624 30(4):435–462
625 Faybishenko B (2002) Chaotic dynamics in flow through unsat-
626 urated fractured media. *Adv Water Resour* 25(7):793–816
627 Fraedrich K (1986) Estimating the dimensions of weather and
628 climate attractors. *J Atmos Sci* 43:419–432
629 Ghilardi P, Rosso R (1990) Comment on ‘Chaos in rainfall.’
630 *Water Resour Res* 26(8):1837–1839
631 Grassberger P, Procaccia I (1983) Measuring the strangeness of
632 strange attractors. *Physica D* 9:189–208
633 Hao B-L (1984) *Chaos*. World Scientific, Singapore
634 Harvey CF, Swartz CH, Badruzzaman ABM, Keon-Blute N, Yu
635 W, Ali MA, Jay J, Beckie R, Niedan V, Brabander D, Oates
636 PM, Ashfaq KN, Islam S, Hemond HF, Ahmed MF
637 (2002) Arsenic mobility and groundwater extraction in
638 Bangladesh. *Science* 298:1602–1606
639 Helsel DR, Hirsch RM (1992) *Statistical methods in water
640 resources*. Elsevier, New York
641 Hossain F, Sivakumar B (2006a) Spatial pattern of arsenic
642 contamination in shallow tubewells of Bangladesh: regional
643 geology and nonlinear dynamics. *Stoch Environ Res Risk
644 Assess* 20(1–2):66–76. DOI 10.1007/s00477-0055-0012-7
645 Hossain F, Sivakumar B (2006b) Spatial interpolation based on
646 complementary paradigms: a call for a change in attitude,
647 mathematical geology (in review; available online at [http://
648 www.iweb.tntech.edu/fhossain/papers/Math_GeolAs.pdf](http://www.iweb.tntech.edu/fhossain/papers/Math_GeolAs.pdf))
649 Hossain F, Bagtzoglou AC, Nahar N, Hossain MD (2006a)
650 Statistical characterization of arsenic contamination in
651 shallow tube wells of western Bangladesh. *Hydrol Process*
652 20(7):1497–1510. DOI 10.1002/hyp.5946
653 Hossain F, Hill J, Bagtzoglou AC (2006b) Geostatistically based
654 management of arsenic contaminated ground water in
655 shallow wells of Bangladesh. *Water Resour Manage* (in
656 press). DOI 10.1007/s11269-006-9079-2
657 Journel AG, Huijbregts CJ (1978) *Mining Geo-statistics*. Aca-
658 demic, San Diego
659 Lambrakis N, Andreou AS, Polydoropoulos P, Georgopoulos E,
660 Bountis T (2000) Nonlinear analysis and forecasting of a
661 brackish karstic spring. *Water Resour Res* 36(4):875–884
662 McArthur JM, Ravenscroft P, Safiullah S, Thirlwall MF (2001)
663 Arsenic in groundwater: testing pollution mechanisms for
664 sedimentary aquifers in Bangladesh. *Water Resour Res*
665 37(1):109–117
666 McArthur JM, Banerjee DM, Hudson-Edwards KA, Mishra R,
667 Purohit R, Ravenscroft P, Cronine A, Howarth RJ, Chat-
668 terjee A, Talukder T, Lowry D, Houghton S, Chadha DK
669 (2004) Natural organic matter in sedimentary basins and its
670 relation to arsenic in anoxic ground water: the example of
671 West Bengal and its worldwide implications. *Appl Geochem*
672 19:1255–1293
673 Mukherjee AB, Bhattacharya P (2002) Arsenic in groundwater
674 in the Bengal Delta plain: slow poisoning in Bangladesh.
675 *Environ Rev* 9:189–220
676 Nerenberg MAH, Essex C (1990) Correlation dimension and
677 systematic geometric effects. *Phys Rev A* 42(12):7065–7074
678 Ormsby T, Napoleon E, Burke R, Feaster L, Groessl C (2004)
679 Getting to know ArcGIS desktop, 2nd edn. ESRI Press,
680 Redlands (ISBN:1-58948-083-X)
681 Porporato A, Ridolfi L (1997) Nonlinear analysis of river flow
682 time sequences. *Water Resour Res* 33(6):1353–1367
683 Rodriguez-Iturbe I, De Power FB, Sharifi MB, Georgakakos KP
684 (1989) Chaos in rainfall. *Water Resour Res* 25(7):1667–1675
685 Schertzer D, Tchiguirinskaia I, Lovejoy S, Hubert P, Bendjoudi
686 H (2002) Which chaos in the rainfall-runoff process? A
687 discussion on ‘Evidence of chaos in the rainfall-runoff
688 process’ by Sivakumar et al. *Hydrol Sci J* 47(1):139–147
689 Schreiber T, Kantz H (1996) Observing and predicting chaotic
690 signals: is 2% noise too much? In: Krastov Yu A, Kadtko JB
691 (eds) *Predictability of complex dynamical systems*. Springer,
692 Berlin Heidelberg New York, pp 43–65
693 Serre ML, Kolovos A, Christakos G, Modis K (2003) An
694 application of the holistochastic human exposure method-
695 ology to naturally occurring arsenic in Bangladesh drinking
696 water. *Risk Anal* 23(3):515–528
697 Sivakumar B (2000) Chaos theory in hydrology: important issues
698 and interpretations. *J Hydrol* 227(1–4):1–20
699 Sivakumar B (2004a) Chaos theory in geophysics: past, present
700 and future. *Chaos, Solitons. Fractals* 19(2):441–462
701 Sivakumar B (2004b) Dominant processes concept in hydrology:
702 moving forward. *Hydrol Process* 18:2349–2353
703 Sivakumar B (2005) Correlation dimension estimation of hydro-
704 logic series and data size requirement: myth and reality.
705 *Hydrol Sci J* 50(4):591–604
706 Sivakumar B, Phoon KK, Liong SY, Liaw CY (1999) A
707 systematic approach to noise reduction in chaotic hydrolog-
708 ical time series. *J Hydrol* 219(3–4):103–135
709 Sivakumar B, Berndtsson R, Persson M (2001a) Monthly runoff
710 prediction using phase-space reconstruction. *Hydrol Sci J*
711 46(3):377–387
712 Sivakumar B, Sorooshian S, Gupta HV, Gao X (2001b) A
713 chaotic approach to rainfall disaggregation. *Water Resour
714 Res* 37(1):61–72
715 Sivakumar B, Berndtsson R, Olsson J, Jinno K (2002a) Reply to
716 ‘which chaos in the rainfall-runoff process?’ by Schertzer
717 et al. *Hydrol Sci J* 47(1):149–158
718 Sivakumar B, Jayawardena AW, Fernando TM GH (2002b)
719 River flow forecasting: use of phase-space reconstruction
720 and artificial neural networks approaches. *J Hydrol* 265(1–
721 4):225–245
722 Sivakumar B, Persson M, Berndtsson R, Uvo CB (2002c) Is
723 correlation dimension a reliable indicator of low-dimen-
724 sional chaos in short hydrological time series? *Water Resour
725 Res* 38(2). DOI 10.1029/2001WR000333
726 Sivakumar B, Harter T, Zhang H (2005) Solute transport in a
727 heterogeneous aquifer: a search for nonlinear deterministic
728 dynamics. *Nonlin Process Geophys* 12:211–218
729 Takens F (1981) Detecting strange attractors in turbulence. In:
730 Rand DA, Young LS (eds) *Dynamical systems and turbu-
731 lence*, lecture notes in mathematics, vol 898. Springer, Berlin
732 Heidelberg New York, pp 366–381
733 Theiler J (1987) Efficient algorithm for estimating the correla-
734 tion dimension from a set of discrete points. *Phys Rev A*
735 36(9):4456–4462

| | | | |
|-----|--|---|-----|
| 736 | Tsonis AA, Triantafyllou GN, Elsner JB, Holdzkom JJ II, | Welch AH, Westjohn DB, Helsel DR, Wanty RB (2000) Arsenic | 749 |
| 737 | Kirwan AD Jr (1994) An investigation of the ability of | in ground water of the United States: occurrence and | 750 |
| 738 | nonlinear methods to infer dynamics from observables. Bull | geochemistry. <i>Ground Water</i> 38(4):589–604 | 751 |
| 739 | <i>Am Meteorol Soc</i> 75:1623–1633 | Yu WH, Harvey CM, Harvey CF (2003) Arsenic groundwater in | 752 |
| 740 | Twarakavi NKC, Kaluarachchi JJ (2006) Arsenic in ground | Bangladesh: a geo-statistical and epidemiological frame- | 753 |
| 741 | waters of conterminous United States: assessment, health | work for evaluating health effects and potential remedies. | 754 |
| 742 | risk, and costs for MCL compliance. <i>J Am Water Resour</i> | <i>Water Resour Res</i> 39(6):1146. DOI 10.1029/2002WR001327 | 755 |
| 743 | <i>Assoc</i> 42(2):275–294 | Zheng Y, Stute M, van Geen A, Gavrieli I, Dhar R, Simpson HJ, | 756 |
| 744 | van Geen A, Zheng Y, Vesteege R, Stute M, Horneman A, Dhar | Schlosser P, Ahmed KM (2004) Redox control of arsenic | 757 |
| 745 | R, Steckler M, Gelman A, Ahsan H, Graziano JH, Hussain | mobilization in Bangladesh groundwater. <i>Appl Geochem</i> | 758 |
| 746 | I, Ahmed KM (2003) Spatial variability of arsenic in 6000 | 19:201–214 | 759 |
| 747 | tube wells in a 25 km ² area of Bangladesh. <i>Water Resour</i> | | 760 |
| 748 | <i>Res</i> 39(5):1140. DOI 10.1029/2002/WR001617 | | |

UNCORRECTED PROOF