

¹ Principal Component Geostatistical Approach for ² Large-Dimensional Inverse Problems

P. K. Kitanidis and J. Lee¹

P. K. Kitanidis, Department of Civil and Environmental Engineering, Stanford University, Stanford, CA, 94305, USA. (peterk@stanford.edu)

J. Lee, Department of Civil and Environmental Engineering, Stanford University, Stanford, 94305, USA. (jonghyun@stanford.edu)

¹Department of Civil and Environmental Engineering, Stanford University, Stanford, 94305, USA.

3 **Abstract.** The quasilinear geostatistical approach is for weakly nonlin-
4 ear underdetermined inverse problems, such as Hydraulic Tomography and
5 Electrical Resistivity Tomography. It provides best estimates as well as mea-
6 sures for uncertainty quantification. However, for its textbook implemen-
7 tation, the approach involves iterations, to reach an optimum, and requires
8 the determination of the Jacobian matrix, *i.e.*, the derivative of the obser-
9 vation function with respect to the unknown. Although there are elegant
10 methods for the determination of the Jacobian, the cost is high when the num-
11 ber of unknowns, m , and the number of observations, n , is high. It is also
12 wasteful to compute the Jacobian for points away from the optimum. Ir-
13 respective of the issue of computing derivatives, the computational cost of
14 implementing the method is generally of the order of m^2n , though there are
15 methods to reduce the computational cost. In this work, we present an im-
16 plementation that utilizes a matrix-free in terms of the Jacobian matrix Gauss-
17 Newton method and improves the scalability of the geostatistical inverse prob-
18 lem. For each iteration, it is required to perform K runs of the forward prob-
19 lem, where K is not just much smaller than m but can be smaller than n .
20 The computational and storage cost of implementation of the inverse pro-
21 cedure scales roughly linearly with m instead of m^2 as in the textbook ap-
22 proach. For problems of very large m , this implementation constitutes a dra-
23 matic reduction in computational cost compared to the textbook approach.
24 Results illustrate the validity of the approach and provide insight in the con-
25 ditions under which this method perform best.

1. Introduction

Hydrology and geophysics are among the many fields where one encounters large-dimensional parameter estimation or inversion problems. The relation between the unknown parameters and the observations is typically governed by partial differential equations. An important feature of such problems is that the information in the data does not suffice to constrain the answer to a unique solution that is not overly sensitive to the data. Examples include hydraulic tomography [Butler *et al.*, 1999; Yeh and Liu, 2000; Cardiff *et al.*, 2009, 2013] and electrical resistivity tomography [Slater *et al.*, 2000; Linde *et al.*, 2006; Pollock and Círpka, 2012], which lead to the formulation of such problems when a fine enough grid is used to discretize the unknown function. To summarize the main ideas consider the linear problem

$$\mathbf{y} = \mathbf{H}\mathbf{s} \tag{1}$$

where the n by m observation matrix \mathbf{H} has rank less than m in a numerical sense. Given \mathbf{y} , there are multiple solutions for \mathbf{s} or the solution may vary a lot in response to small changes in the data. The common approach to solving such problems is to introduce additional requirements, such as finding the flattest or smoothest solution consistent with the data, in a precisely specified sense.

This kind of problem can be addressed within a Bayesian framework [Gavalas *et al.*, 1976; Neuman, 1980; Kitanidis and Vomvoris, 1983; Carrera and Neuman, 1986; Woodbury and Rubin, 2000; Rubin *et al.*, 2010], where the multiplicity of solutions is explicitly recognized and each possible solution is assigned a probability based on how well it agrees

45 with the data (the likelihood function criterion) and how well it agrees with other infor-
46 mation (the prior probability criterion). A solution is ascribed a high probability if it
47 meets both criteria (logical conjunction) to a significant degree. There are many differ-
48 ent Bayesian approaches but the emphasis here is on the geostatistical approach (GA).
49 Within the context of Bayesian inference approaches, GA is an objective [*Berger, 2006,*
50 describes the basic ideas] and empirical Bayes [e.g., *Carlin and Louis, 2000; Kitanidis,*
51 2010] method. Objective is a technical term that refers to the need to rely as much on
52 data as the nature of the problem allows while empirical means that both the prior and
53 the likelihood functions are adjustable models informed by data.

54 Let us briefly review how the method works by considering just the best estimate. The
55 result is

$$\hat{\mathbf{s}} = \mathbf{\Lambda} \mathbf{y} \tag{2}$$

56 where $\mathbf{\Lambda}$ acts as a pseudo-inverse of \mathbf{H} . This $\mathbf{\Lambda}$ depends not just on \mathbf{H} but also on an
57 m by p drift matrix \mathbf{X} , an m by m covariance \mathbf{Q} of \mathbf{s} and an n by n covariance \mathbf{R} of
58 measurement error, as will be explained later.

59 1. The pseudo-inverse satisfies $\mathbf{\Lambda H X} = \mathbf{X}$, which means that $\mathbf{\Lambda}$ behaves like the true
60 inverse of \mathbf{H} w.r.t. \mathbf{X} . For example, in some applications (e.g., s is hydraulic head) the
61 datum is modeler dependent. By choosing $\mathbf{X} = \mathit{ones}(m, 1)$ we make sure that if the
62 datum is changed, the solution is not affected. In some other problems, we work with
63 $s = \log K$ (where K is hydraulic or electrical conductivity). By using $\mathbf{X} = \mathit{ones}(m, 1)$
64 we ascertain that the solution is scaled solely on the basis of data. The introduction of

65 the drift means that p characteristics of the solution depend only on the data and not on
66 the prior.

67 2. The selection of \mathbf{Q} describes the variability and continuity in the solution, in the
68 part not covered by the drift. A solution may be chosen to be smooth on the basis of what
69 is known about the unknown function. Or, based on the understanding that small scale
70 features cannot be estimated with certainty from data, one may choose to seek solutions
71 that are sufficiently smooth.

72 3. The selection of \mathbf{R} informs about the type and magnitude of observation error or
73 "noise" and allows to properly weight measurements. If \mathbf{R} is doubled, for example, the
74 data will receive less weight, the best solution will reproduce the data less closely, and
75 thus the effect of data noise on the result will be reduced.

76 4. In combination, \mathbf{Q} and \mathbf{R} have a critical role in error quantification. For example,
77 if both \mathbf{Q} and \mathbf{R} are multiplied by 4, the credible interval (Bayesian confidence interval)
78 doubles.

79 The geostatistical approach for nonlinear problems [*Kitanidis, 1995*] uses a Gauss-
80 Newton iterative method to obtain the "best estimate" and linearized uncertainty quantifi-
81 cation about the best estimate. However, the textbook implementation of this approach
82 is suited for problems where m and n are of modest size. By modest size we mean that
83 linear systems and matrix multiplications are computable within seconds with computa-
84 tional facilities available. At this point in time, this would mean $10^3 - 10^4$. However, we
85 are increasingly faced with problems where the number of unknowns is large, often larger
86 than 10^6 , whereas the number of observations is modest (say 10^3). For such problems,
87 each iteration requires about n runs of the forward problem (to construct the Jacobian

88 matrix via the adjoint-state method) plus operations of the order m^2n . Furthermore,
89 covariance matrices, of size m^2 , are very expensive to compute and store. Overall, the
90 cost of implementation increases roughly with m^2 , which means that the method becomes
91 impractical to use for a sufficiently large m . To address this problem, methods using FFT
92 or hierarchical matrices and fast multiple methods have been applied [Nowak *et al.*, 2003;
93 Saibaba *et al.*, 2012; Ambikasaran *et al.*, 2013] that speed up the cost of multiplication of
94 \mathbf{Q} with a vector from $O(m^2)$ to $O(m \log m)$, where the big O denotes order of. Such
95 methods reduce the overall cost quite significantly. Additionally, various other works
96 examine aspects of reducing the computational cost [e.g., Li and Cirpka, 2006; Liu and
97 Kitanidis, 2011; Saibaba and Kitanidis, 2012; Yadav and Michalak, 2013].

98 This work presents an implementation, termed Principal Component Geostatistical Ap-
99 proach (PCGA), that is mainly suited for cases with the following characteristics:

- 100 1. In terms of problem size, m is huge whereas n is of modest size.
- 101 2. The prior covariance corresponds to a reasonably smooth unknown function. More
102 technically the spectrum of the covariance matrix (*i.e.*, the eigenvalues) drops rapidly.
- 103 3. The information content of the observations is limited, due to observation error and
104 the ill-conditioning of the inverse problem. We will explain what this means in more
105 precise terms later.

106 Many problems encountered in practice, such as hydraulic tomography, fit this descrip-
107 tion. We will later discuss how one can relax 1. and 2. but requirement 3. is needed
108 for the method to perform most efficiently and to be competitive with other approaches.
109 The approach uses the idea expressed in Bui-Thanh *et al.* as follows: "... we exploit the
110 fact that the data are typically informative about low-dimensional manifolds of parameter

111 space to construct low rank approximations of the covariance matrix of the posterior pdf
112 via a matrix-free randomized method.” The method presented in this work also uses low-
113 rank matrix approximations and also avoids the computation and storage of the complete
114 Jacobian matrix but is different from the method in the aforementioned reference.

115 The main contribution of this paper that it presents a methodology for solving static
116 inverse problems, using the Geostatistical Approach, with the following characteristics:
117 (a) It is fully numerical and scales in computational cost and storage requirements roughly
118 linearly with m , the size of the vector of the unknowns. (b) It is a ”matrix-free” approach
119 in the sense that required derivatives are computed by performing matrix-vector multi-
120 plications, without computing the complete Jacobian matrix. Computing derivatives is
121 the bulk of computations in inverse methods. The method we present here allows signif-
122 icant savings, particularly at early iterations when computing a full Jacobian is usually
123 a waste of effort. (c) The method requires making calls to a forward model and does
124 not require an adjoint-state model. (d) The method has spectral convergence with the
125 number of components used, which, for many though not necessarily all problems, can be
126 much faster than the $1/\sqrt{K}$ convergence of ensemble or Monte-Carlo based methods.

127 The next section reviews GA and discuss some important features regarding the stan-
128 dard or ”textbook” computational implementation of GA. Then we present the PCGA
129 algorithm, which is suited for problems for which the number of unknowns, m , is very
130 large, the number of observations, n , is moderate, and the evaluation of derivatives is too
131 expensive to be performed a very large number of times. Then, we present results that
132 illustrate the applicability of the method.

2. Overview of Geostatistical Approach

133 For the sake of completeness and for ease in referencing results, we will review the
134 quasilinear geostatistical approach (GA) [Kitanidis, 1995]. The observation equation,
135 which relates the vector of the unknowns \mathbf{s} to the vector of the data \mathbf{y} is

$$\mathbf{y} = \mathbf{h}(\mathbf{s}) + \mathbf{v} \quad (3)$$

136 where \mathbf{h} is the mapping from $R^{m \times 1}$ to $R^{n \times 1}$, \mathbf{v} is Gaussian with mean $\mathbf{0}$ and covariance
137 \mathbf{R} (often proportional to the identity matrix). The mapping must be sufficiently smooth
138 so that it can be approximated with a linear relation in the neighborhood where the
139 posterior probability of \mathbf{s} is concentrated. The prior probability of \mathbf{s} is Gaussian with
140 mean $\mathbf{X}\boldsymbol{\beta}$, where \mathbf{X} is an m by p known matrix and $\boldsymbol{\beta}$ a p by 1 vector of unknowns (*i.e.*,
141 to be determined by data), and generalized covariance matrix \mathbf{Q} .

142 The posterior pdf of \mathbf{s} and $\boldsymbol{\beta}$ are obtained through Bayes theorem and its cologarithm
143 (minus the logarithm), $-\ln p''(\mathbf{s}, \boldsymbol{\beta})$, is

$$\frac{1}{2}(\mathbf{y} - \mathbf{h}(\mathbf{s}))^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{h}(\mathbf{s})) + \frac{1}{2}(\mathbf{s} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{Q}^{-1}(\mathbf{s} - \mathbf{X}\boldsymbol{\beta}) \quad (4)$$

144 The maximum a posteriori or most likely values are obtained by minimizing this expression
145 with respect to \mathbf{s} and $\boldsymbol{\beta}$ vectors. This is a nonlinear optimization problem that is
146 commonly solved through iterative methods. For $n \ll m$, which is usually the case,
147 a convenient form of the Gauss-Newton method is in the form of the so-called cokriging
148 equations, to be described next.

149 Based on the initial guess (or most recent "good solution") \mathbf{s}_0 , we update to a new
 150 solution using a Newton-type iterative approach. First, define the n by m matrix \mathbf{H} as
 151 the Jacobian matrix of \mathbf{h} at \mathbf{s}_0 :

$$\mathbf{H} = \left. \frac{\partial \mathbf{h}}{\partial \mathbf{s}} \right|_{\mathbf{s}=\mathbf{s}_0} \quad (5)$$

152 Then, assuming that the actual $\hat{\mathbf{s}}$ is close to \mathbf{s}_0 , approximate

$$\mathbf{h}(\hat{\mathbf{s}}) = \mathbf{h}(\mathbf{s}_0) + \mathbf{H}(\hat{\mathbf{s}} - \mathbf{s}_0) \quad (6)$$

153 and, after some matrix manipulations, one obtains the updated solution to the minimiza-
 154 tion of (4)

$$\hat{\mathbf{s}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{QH}^T \boldsymbol{\xi} \quad (7)$$

155 where $\hat{\boldsymbol{\beta}}$ and the n by $\mathbf{1}$ vector $\boldsymbol{\xi}$ are found from the solution of a system of :

$$\begin{bmatrix} \mathbf{HQH}^T + \mathbf{R} & \mathbf{HX} \\ (\mathbf{HX})^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\xi} \\ \hat{\boldsymbol{\beta}} \end{bmatrix} = \begin{bmatrix} \mathbf{y} - \mathbf{h}(\mathbf{s}_0) + \mathbf{H}\mathbf{s}_0 \\ \mathbf{0} \end{bmatrix} \quad (8)$$

Note also that the objective function that is minimized can be written, using (7), as

$$\begin{aligned} J &= \frac{1}{2} (\mathbf{y} - \mathbf{h}(\mathbf{s}))^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{s})) + \frac{1}{2} (\mathbf{s} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{Q}^{-1} (\mathbf{s} - \mathbf{X}\boldsymbol{\beta}) \\ &= \frac{1}{2} (\mathbf{y} - \mathbf{h}(\mathbf{X}\boldsymbol{\beta} + \mathbf{QH}^T \boldsymbol{\xi}))^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{X}\boldsymbol{\beta} + \mathbf{QH}^T \boldsymbol{\xi})) \\ &\quad + \frac{1}{2} \boldsymbol{\xi}^T \mathbf{HQH}^T \boldsymbol{\xi} \end{aligned} \quad (9)$$

156 This expression can be used to gauge the progress of the minimization and to make sure
 157 that the new solution is not worse than the previous.

158 Once the optimum is achieved, we can proceed to uncertainty quantification. The
 159 linearized approach treats the problem as approximately linear in the proximity of the
 160 best estimate. This usually assumes that the posterior error is somehow small and the
 161 nonlinearity not too strong so that a linear approximation is good in the neighborhood
 162 where \mathbf{s} is most likely to be. Under these conditions the posterior is approximately
 163 Gaussian.

164 One approach to uncertainty quantification involves generation of conditional realiza-
 165 tions, *i.e.*, samples from the posterior distribution. They are computed as follows:

$$\mathbf{s}_i = \zeta_i + \mathbf{X}\beta + \mathbf{Q}\mathbf{H}^T \boldsymbol{\xi} \quad (10)$$

166 where ζ_i are unconditional realizations, *i.e.*, Gaussian with mean $\mathbf{0}$ and covariance \mathbf{Q} ,
 167 and $\hat{\beta}$ and $\boldsymbol{\xi}$ are found from:

$$\begin{bmatrix} \mathbf{H}\mathbf{Q}\mathbf{H}^T + \mathbf{R} & \mathbf{H}\mathbf{X} \\ (\mathbf{H}\mathbf{X})^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\xi} \\ \hat{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{y} + \mathbf{v}_i - \mathbf{h}(\mathbf{s}_0) + \mathbf{H}(\mathbf{s}_0 - \zeta_i) \\ \mathbf{0} \end{bmatrix} \quad (11)$$

168 where \mathbf{v}_i is generated Gaussian mean $\mathbf{0}$ and covariance matrix \mathbf{R} .

169 An alternative approach is to solve

$$\begin{bmatrix} \mathbf{H}\mathbf{Q}\mathbf{H}^T + \mathbf{R} & \mathbf{H}\mathbf{X} \\ (\mathbf{H}\mathbf{X})^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Lambda}^T \\ \mathbf{M} \end{bmatrix} = \begin{bmatrix} \mathbf{H}\mathbf{Q} \\ \mathbf{X}^T \end{bmatrix} \quad (12)$$

170 where $\boldsymbol{\Lambda}$ is an $m \times n$ matrix of coefficients and \mathbf{M} is a $p \times n$ matrix of multipliers. Then,
 171 the best estimate is:

$$\hat{\mathbf{s}} = \boldsymbol{\Lambda} (\mathbf{y} - \mathbf{h}(\mathbf{s}_0) + \mathbf{H}\mathbf{s}_0) \quad (13)$$

172 while the posterior covariance matrix, from the linearized approach, is

$$\mathbf{V} = \mathbf{Q} - \mathbf{F} \tag{14}$$

173 where

$$\mathbf{F} = \mathbf{X}\mathbf{M} + \mathbf{Q}\mathbf{H}^T\mathbf{\Lambda}^T \tag{15}$$

174 This result shows that the posterior covariance matrix is the prior covariance minus a
 175 low-rank correction. Mathematically, the rank of the correction \mathbf{F} cannot be higher than
 176 $n + p$, given the dimensions of the factors of the two terms that comprise the correction.

177 In the textbook approach, the Jacobian matrix must be re-calculated at each Gauss-
 178 Newton iteration using the adjoint-state method, which is the method of choice when
 179 $n \ll m$, requiring n solutions of the forward problem. Usually, this is the most com-
 180 putationally expensive part of the method. Furthermore, when m is large, the cost of
 181 dealing with covariance matrices is large and can even be prohibitive. The computation
 182 of the \mathbf{QH}^T is of the order of m^2n , which is very burdensome when m is large. There
 183 are methods to reduce the computational cost, associated with the multiplication of the
 184 covariance matrix with vectors [*Nowak et al.*, 2003; *Saibaba et al.*, 2012]. Lastly, the
 185 computation and storage of the posterior covariance matrix can be prohibitive.

186 The method described next is well suited for cases of huge m and moderate n and
 187 becomes even better when the effective rank of the correction in (14) is $\ll n$. (Or, to put
 188 it differently, it can be approximated by a matrix of rank $\ll n$.) The method is termed
 189 Principal Component Geostatistical Approach (PCGA) because it is an implementation
 190 of the Geostatistical Approach that utilizes principal components associated with the
 191 covariance matrix and the drift matrix. The method has the potential to reduce markedly

192 the number of runs of the forward problem and the dominant term in the cost of matrix
193 manipulations drops from $O(m^2)$ to order $O(m)$, which is a dramatic improvement for
194 large m .

3. Important Preliminaries

195 This section describes some ideas and tools that will be used later.

3.1. Matrix-Free Approach

Instead of computing the complete Jacobian matrix \mathbf{H} and then computing the matrix-vector product $\mathbf{H}\mathbf{a}$, one can directly compute this matrix-vector product using a so-called matrix-free method. Consider the nonlinear vector function

$$\mathbf{y} = \mathbf{h}(\mathbf{s}) \tag{16}$$

where \mathbf{s} is m by 1 and \mathbf{y} is n by 1, with its Jacobian matrix at $\mathbf{s} = \mathbf{s}_0$ denoted as \mathbf{H} . Consider that we want to compute

$$\mathbf{H}\mathbf{a} \tag{17}$$

where \mathbf{a} is a vector with same dimensions as \mathbf{s} . Rather than follow the expensive process of first evaluating \mathbf{H} and then performing the matrix vector multiplication, we can take advantage of the following finite-difference approximation:

$$\mathbf{H}\mathbf{a} = \frac{\|\mathbf{a}\|}{\delta \|\mathbf{s}_0\|} \left[\mathbf{h} \left(\mathbf{s}_0 + \frac{\|\mathbf{s}_0\|}{\|\mathbf{a}\|} \delta \mathbf{a} \right) - \mathbf{h}(\mathbf{s}_0) \right] + O(\delta) \tag{18}$$

196 The ratio of norms (or vector lengths) $\|\mathbf{s}_0\| / \|\mathbf{a}\|$ is a normalization factor and δ is a
197 small dimensionless number. This approach with δ equal to about the square root of the
198 floating-point relative accuracy (for example, $\delta \approx 10^{-8}$ for "double precision", $\delta \approx 10^{-4}$ for
199 "single precision") is adequate, but one may want to experiment with choosing the δ that

200 is right for a specific application. The essential point is that each Jacobian times vector
201 computation involves one more run of the forward problem.

202 Note that this approach makes calls to a solver of the forward problem. If we needed
203 to compute $\mathbf{H}^T \mathbf{b}$, where \mathbf{b} is a vector with the same dimensions as \mathbf{y} , we could follow the
204 same procedure using a solver that solves the adjoint of the forward problem. The cost of
205 a forward run of the original problem and its adjoint are the same so that there is no real
206 benefit here in using the adjoint, unless the computations are arranged so that the number
207 of forward runs is reduced. In this work, computations have been arranged so that there
208 is no need to use an adjoint state solver, taking into account that most practitioners have
209 access to just forward-problem solvers that they can use in a "black box" fashion.

210 Note also that instead of using the finite-difference approximation one could solve a
211 linearized version of the forward problem. This approach would be more elegant and
212 would bypass the issue of choosing the right δ . However, again taking into account that
213 it is more practical to use the available forward model as a black box and not have
214 to write another PDE solver, which might be a nontrivial exercise, the finite-difference
215 approximation is recommended as the method of choice. After all, development of forward
216 problem solvers, which are typically PDE solvers, is quite advanced and one can find highly
217 optimized code, often approaching near linear scalability with problem size. The inverse
218 methodology presented herein has been structured to take advantage of such code.

3.2. Normalization of Prior Model

219 The drift matrix \mathbf{X} is m by p , where p is a small number. If $p = 1$, the normalized
220 version of \mathbf{X} is $\mathbf{U} = \mathbf{X}/\sqrt{m}$. For $p > 1$, one may use a singular value decomposition:

$$\mathbf{X} = \mathbf{U}\mathbf{S}_X\mathbf{V}_X^T \quad (19)$$

221 where \mathbf{U} is m by p , \mathbf{S}_X is diagonal p by p , and \mathbf{V}_X is p by p . The \mathbf{S}_X and \mathbf{V}_X can
 222 be discarded. Matrix \mathbf{U} is isomorphic to \mathbf{X} , *i.e.*, mathematically the same result in
 223 the geostatistical inversion is obtained whether one defines the drift in terms of \mathbf{X} or \mathbf{U} .
 224 However, since \mathbf{U} is orthonormal, meaning that $\mathbf{U}^T\mathbf{U}$ is identity matrix, it is preferable
 225 to use in computations. It is recommended to replace \mathbf{X} with \mathbf{U} when defining the model
 226 at the very beginning.

The detrending matrix \mathbf{P} is expressed

$$\mathbf{P} = \mathbf{I} - \mathbf{U}\mathbf{U}^T \quad (20)$$

This matrix is symmetric, has $m - p$ eigenvalues equal to 1 and p eigenvalues equal to
 0, is a projection to the space orthogonal to the drift matrix, $\mathbf{P}\mathbf{X} = \mathbf{0}$, and satisfies the
 idempotence property, meaning that it can be multiplied many times without changing
 the result beyond the initial multiplication, $\mathbf{P}\mathbf{P}\mathbf{a} = \mathbf{P}\mathbf{a}$. Furthermore the matrix vector
 product $\mathbf{P}\mathbf{a}$ can be computed with $O(m)$ operations using the relation

$$\mathbf{P}\mathbf{a} = \mathbf{a} - \mathbf{U}(\mathbf{U}^T\mathbf{a}) \quad (21)$$

The generalized covariance matrix \mathbf{Q} also has an issue of isomorphism. Whether
 one uses \mathbf{Q} or $\mathbf{Q} + \mathbf{X}\mathbf{b}\mathbf{b}^T\mathbf{X}$, where \mathbf{b} is an arbitrary p by 1 vector, the result of the
 geostatistical inversion method is mathematically identical. (The same is true of ordinary
 Kriging, an interpolation method.) This is a powerful, though perhaps not widely
 appreciated, feature of generalized covariance functions and matrices [*Matheron*, 1973;
Kitanidis, 1983, 1993]. The essence of isomorphism is that two covariance functions

or matrices may look different but act mathematically in exactly the same way, in the context of defining the prior with a specific drift matrix; and what matters is to identify the criterion for isomorphism and the invariants. The criterion is the following. Two covariance matrices \mathbf{Q}_1 and \mathbf{Q}_2 are isomorphic if

$$\mathbf{P}(\mathbf{Q}_1 - \mathbf{Q}_2)\mathbf{P} = \mathbf{0} \tag{22}$$

227 where \mathbf{P} is defined in (20). The definitions above lead to the following: (a) \mathbf{Q} and
 228 \mathbf{PQP} are isomorphic; (b) if \mathbf{Q}_1 and \mathbf{Q}_2 are isomorphic, then $\mathbf{PQ}_1\mathbf{P}$ is equal to $\mathbf{PQ}_2\mathbf{P}$
 229 and isomorphic to the original \mathbf{Q}_1 and \mathbf{Q}_2 ; and (c) thus, \mathbf{PQP} is the invariant.

One straightforward solution to removing the ambiguity is to replace \mathbf{Q} with \mathbf{PQP} . Note that whereas \mathbf{Q} may not be positive definite and may not look like a reasonable covariance matrix of a random vector, \mathbf{PQP} is nonnegative definite and has a solid probabilistic interpretation as the covariance of $\mathbf{P}\mathbf{x}$, where \mathbf{x} is a random vector. \mathbf{PQP} is "what matters" in the problems we are interested in. But one may avoid performing the multiplication, because it involves $O(m^2)$ operations, when applying the geostatistical approach. If, for example we multiply with a vector that is already detrended, by taking account of idempotence we have

$$(\mathbf{P}\mathbf{a})^T \mathbf{PQP} (\mathbf{P}\mathbf{a}) = (\mathbf{P}\mathbf{a})^T \mathbf{Q} (\mathbf{P}\mathbf{a}) \tag{23}$$

4. Low-Rank Approximation

230 A matrix of low rank can be factorized into the product of low-dimensional matrices.
 231 For example, a matrix \mathbf{Y} of size k by l and rank K can be factorized $\mathbf{Y} = \mathbf{A}\mathbf{B}^T$, where
 232 \mathbf{A} is k by K and \mathbf{B} is l by K . If $K \ll k, l$, the benefits are potentially tremendous in
 233 terms of storage and matrix-vector operations. Instead of storing $k \times l$ entries, one stores

234 $(k + l) K$. The matrix vector product $\mathbf{Y}\mathbf{y}$ instead of requiring $k \times l$ it takes only $(k + l) K$
 235 multiplications, since $\mathbf{Y}\mathbf{y} = \mathbf{A}(\mathbf{B}^T \mathbf{y})$. This crucial idea is the backbone of approaches
 236 in dealing with cases of large m .

237 In, for example, hierarchical matrix methods [Saibaba *et al.*, 2012; Ambikasaran *et al.*,
 238 2013, e.g., for applications in estimation, see], the structure of the dense prior covariance
 239 matrix is exploited to factorize it. Bui-Thanh *et al.* factorize a matrix appearing in the
 240 Hessian of the linearized problem to compute the low-order correction in the covariance.
 241 The method to be used here is closest works that employ orthogonal or Karhunen-Loeve
 242 decomposition of the covariance [Li and Cirpka, 2006; Marzouk and Najm, 2009]

243 In the approach developed herein, we start with the factorization of the prior covariance
 244 \mathbf{Q} with low-size matrices. There are many possible methods but we need a method with
 245 computational cost roughly linear in m and also one that does not require many repetitions
 246 of the procedure. The method we describe next is one of several possible computational
 247 approaches.

248 We start by generating K unconditional realizations; *i.e.*, ζ_k , where $k = 1 : K$, is m by
 249 1 generated with mean $\mathbf{0}$ and generalized covariance function \mathbf{Q} . For practical situations
 250 of generating realizations of random fields in 1, 2, or 3 dimensions using regular grids, one
 251 powerful approach is based on FFT or hierarchical matrices. The cost of generating a
 252 realization with such methods is $O(m \log m)$.

The simplest approximation to the covariance \mathbf{Q} is through

$$\mathbf{Q} \approx \frac{1}{K} \sum_{k=1}^K \zeta_k \zeta_k^T \quad (24)$$

This is factorization through the product of two rank K matrices

$$\mathbf{Q} \approx \mathbf{Z}\mathbf{Z}^T, \text{ where } \mathbf{Z}_{ij} = \frac{\zeta_{ij}}{\sqrt{K}}, i = 1, \dots, m, j = 1, \dots, K \quad (25)$$

253 This is the same factorization used in ensemble Kalman filtering [Evensen,
254 1994, 2003, 2006], which in some cases has been reported to work well even for K as
255 small as 50 or a 100, but in other cases requires large ensembles or additional enhance-
256 ments [Chatterjee *et al.*, 2012; Li *et al.*, 2014] Generally, the approximation error in (25)
257 is $O\left(\frac{1}{\sqrt{K}}\right)$. This means that to halve the error from the factorization one needs to
258 quadruple the number of realizations.

An alternative idea is to factorize \mathbf{Q} in such a way that the factorization is "optimum"
for the specific number K , the rank of the approximation. From matrix theory [Golub
and Van Loan, 1989], it is known that the smallest achievable error when using a K -rank
approximation of a matrix \mathbf{Q} , where the spectral norm (aka 2-norm and denoted $\| \cdot \|$)
for matrices is used to quantify the error, is

$$\min_{\mathbf{Q}_K} \|\mathbf{Q} - \mathbf{Q}_K\| = \rho_{K+1} \quad (26)$$

259 where \mathbf{Q} is the full matrix, \mathbf{Q}_K is its rank K approximation, and ρ_{K+1} is the $(K + 1)$ -th
260 largest singular value of \mathbf{Q} . If the spectrum of \mathbf{Q} (*i.e.*, the set of singular values ρ_k ,
261 $k = 1 : m$) drops rapidly, it is possible to obtain a low-rank approximation that has a
262 small error. Most importantly, the optimum \mathbf{Q}_K is the one obtained from a singular
263 value decomposition of \mathbf{Q} where only the first K singular values are kept and the others
264 are set to zero.

265 To illustrate how the choice of model can make a difference, consider the following
266 example. Consider a one-dimensional domain of length $L = 1$ and points $m = 100$ on a
267 uniform grid and three models:

268 Model *I*: $p = 1, X_i = 1, Q_{ij}(x_i - x_j) = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$. This is the classic "nugget-
269 effect" model of geostatistics. No continuity is assumed and using this model in inverse
270 modeling can be interpreted as selecting as best estimate a small variance solution.

271 Model *II*: $p = 1, X_i = 1, Q_{ij} = -|x_i - x_j|$. This is the linear generalized covariance
272 function of geostatistics, corresponding to the commonly used linear semivariogram. The
273 unknown is assumed to belong to an ensemble of continuous but not differentiable func-
274 tions. The use of the model can be interpreted as selecting a best estimate that is a flat
275 solution.

276 Model *III*: $p = 2, X_{i1} = 1, X_{i2} = x_i, Q_{ij} = |x_i - x_j|^3$. This is the cubic generalized
277 covariance function of geostatistics. The unknown is assumed to belong to an ensemble
278 of continuous and differentiable functions and its use can be interpreted as selecting a
279 best estimate that is a smooth solution.

280 Note that none of these three models has a scale parameter and thus the choice of
281 domain length has no effect. The choice of m has a small effect but does not materially
282 affect the argument we will make. To account for the fact that the drift accounts for some
283 variability and the \mathbf{Q} is the variability around the drift, we examine the spectrum of \mathbf{PQP} ,
284 where \mathbf{P} is a symmetric projection matrix of rank $m - p$, as previously defined. Whenever
285 we need to compute a matrix vector product, \mathbf{Pa} , we notice that this is equivalent to
286 detrending \mathbf{a} ; *i.e.*, fitting to \mathbf{a} a trend $\mathbf{X}\boldsymbol{\beta}$, through least squares techniques [Golub, 1965]
287 and then subtracting it from \mathbf{a} to arrive at the desired result.

288 Because we are primarily interested in relative error in approximating \mathbf{Q} , we will plot
289 the normalized matrix spectrum, $\frac{\rho_k}{\rho_1}$, on figure 1.

290 The results clearly illustrate that for some covariance matrices, excellent accuracy in
291 approximating \mathbf{Q} can be obtained using approximations of low rank. Pretty much the
292 same behavior is for much larger m . In fact, for model *III*, even if m is 10^6 , a matrix
293 with rank as low as 20 suffices to achieve small relative error.

294 It is important to point out that these results only suggest the maximum number K
295 that is needed. The results will be better than what this analysis suggests, for a number
296 of reasons. Indeed, when K is large and approaches n , the results shown in figure 1
297 become less relevant. After all, matrices of interest like $\mathbf{H}\mathbf{Q}$ and $\mathbf{H}\mathbf{Q}\mathbf{H}^T$ are of rank
298 not larger than n , and sometimes considerably smaller than n because the low rank of
299 \mathbf{H} . Thus, the errors in \mathbf{Q} become less important. We will later explain how the action
300 of the forward operator and the presence of observation error make it unnecessary to use
301 a large K .

Once one decides what K value to use, one must proceed to compute the low-rank approximation. Among the many methods, each with advantages and disadvantages, we will focus on randomized methods [*Halko et al.*, 2011]. Here, we propose a low cost randomized method to construct an accurate symmetric approximation. The random driver is the m by K matrix \mathbf{Z} , equation (25), with the additional stipulation that each column (unconditional realization) has been detrended. We can perform a singular value decomposition (or a QR decomposition, it makes no difference in this application) of \mathbf{Z} ,

$$\mathbf{Z} = \mathbf{U}_z \mathbf{S}_z \mathbf{V}_z^T \tag{27}$$

Here, \mathbf{U}_z is m by K and satisfies $\mathbf{U}_z^T \mathbf{U}_z = \mathbf{I}_K$, \mathbf{S}_z is K by K diagonal with nonnegative elements, and \mathbf{V}_z is K by K unitary, $\mathbf{V}_z^T \mathbf{V}_z = \mathbf{V}_z \mathbf{V}_z^T = \mathbf{I}_K$. In what follows, only \mathbf{U}_z will be needed. Next compute

$$\mathbf{C} = \mathbf{U}_z^T (\mathbf{PQP}) \mathbf{U}_z \quad (28)$$

Since \mathbf{Z} has been detrended, so has \mathbf{U}_z and thus $\mathbf{P}\mathbf{U}_z = \mathbf{U}_z$, which means that we can simplify the notation

$$\mathbf{C} = \mathbf{U}_z^T \mathbf{Q} \mathbf{U}_z \quad (29)$$

Finally, the low-order approximation of the generalized covariance matrix, \mathbf{PQP} , is

$$\mathbf{PQP} \approx \mathbf{U}_z \mathbf{C} \mathbf{U}_z^T \quad (30)$$

302 which is an $(m \times k) (k \times k) (k \times m)$ factorization.

303 This is a randomized algorithm so one would like to evaluate the effect of randomness.
 304 For this purpose the process was repeated twenty times (*i.e.*, each time with a different
 305 seed number) for model *III* and $K = 30$ and the results were plotted on figure 2. By
 306 actual we mean the values computed with standard software without approximations
 307 about rank. The approximate ones are computed very efficiently from the small matrix
 308 \mathbf{C} . One can see from figure 2 that the error and the randomness are limited to the
 309 smaller computed singular values. This, together with other results, suggests that if
 310 one wants accuracy at the K -th singular value, one should use an approximation with a
 311 rank somewhat higher than that. Randomness, although small in the mean square sense,
 312 introduces unsightly small scale variability in the estimate and has a more pronounced
 313 material effect in computing variances.

314 This difficulty is fortunately easy to overcome. In this example, compute for 30 and
315 then keep only the first 15. All that is needed is to drop the last columns of \mathbf{A} and the
316 last columns and rows of \mathbf{C} .

317 In applications, one can evaluate retroactively how well one has done, without having
318 to compute the singular values of large \mathbf{Q} , by plotting all the singular values of the small
319 matrix \mathbf{C} , that should mirror the first K singular values of \mathbf{Q} . What one wants to see
320 is that they are dropping fast and the smallest is sufficiently small.

321 Summarizing the algorithm:

322 1. Generate K_s realizations, cost $O(K_s m \log m)$, create \mathbf{Z} matrix. Detrend, if not
323 already detrended, $O(K_s m)$ operations.

324 2. Do a singular value decomposition of the \mathbf{Z} matrix of rank K_s , produce m by K_s
325 matrix \mathbf{U}_z , with $O(K_s^2 m)$ operations.

326 3. Compute \mathbf{C} , with cost only $O(K_s m \log m)$ by using fast matrix vector multiplication
327 like *Fong and Darve* [2009] for 2-D.

328 4. Discard the "extra" columns and rows to reduce the size of \mathbf{A} to K columns and of
329 \mathbf{C} to K by K , like $K = K_s - 15$.

5. Principal Component Geostatistical Inversion

330 The static nonlinear inverse problem, Equation (3), can be solved through a variant to
331 the quasilinear approach, Equations (4)-(15), through the following steps.

332 Computation of \mathbf{HX}

For column \mathbf{X}_i , an $O(\delta)$ finite-difference approximation is

$$\mathbf{HX}_i = \frac{\|\mathbf{X}_i\|}{\delta \|\mathbf{s}_0\|} \left[h \left(\mathbf{s}_0 + \delta \frac{\|\mathbf{s}_0\|}{\|\mathbf{X}_i\|} \mathbf{X}_i \right) - h(\mathbf{s}_0) \right] \quad (31)$$

This can be computed from $O(\delta)$ approximation

$$\begin{aligned} &-\mathbf{h}(\mathbf{s}_0) + \mathbf{H}\mathbf{s}_0 \\ &= \frac{1}{\delta} [h((1 + \delta)\mathbf{s}_0) - (1 + \delta)h(\mathbf{s}_0)] \end{aligned} \quad (32)$$

334 **Computation of $\mathbf{H}\mathbf{Q}\mathbf{H}^T$ and $\mathbf{Q}\mathbf{H}^T$** 335 Consider the factorization of the covariance \mathbf{Q} is through

$$\mathbf{Q} = \mathbf{A}\mathbf{C}\mathbf{A}^T \quad (33)$$

336 which is the $(m \times k)(k \times k)(k \times m)$ factorization of equation (30) and we defined $\mathbf{A} =$
 337 \mathbf{U}_Z . The idea, of course, is that \mathbf{Q} is stored through an $m \times K$ matrix with orthonor-
 338 mal columns and a $K \times K$ symmetric and positive definite matrix, where K is much
 339 smaller than m . These matrices are used to perform operations involving \mathbf{Q} without ever
 340 computing \mathbf{Q} .

We will employ this approximation to compute $\mathbf{H}\mathbf{Q}$ and $\mathbf{H}\mathbf{Q}\mathbf{H}^T$. Denote the columns of \mathbf{A} as \mathbf{a}_k , $k = 1 : K$,

$$\begin{aligned} \mathbf{H}\mathbf{Q} &= (\mathbf{H}\mathbf{A})\mathbf{C}\mathbf{A}^T \\ &= \left(\sum_{k=1}^K \mathbf{H}\mathbf{a}_k \right) \mathbf{C}\mathbf{A}^T \end{aligned} \quad (34)$$

Let us define

$$\mathbf{b}_k = \mathbf{H}\mathbf{a}_k \quad (35)$$

that can be computed as described previously. Form the n by K matrix \mathbf{B} with columns the computed \mathbf{b}_k . Then,

$$\mathbf{H}\mathbf{Q} = \mathbf{B}\mathbf{C}\mathbf{A}^T \quad (36)$$

$$\begin{aligned}
\mathbf{H}\mathbf{Q}\mathbf{H}^T &= (\mathbf{H}\mathbf{A})\mathbf{C}(\mathbf{H}\mathbf{A})^T \\
&= \mathbf{B}\mathbf{C}\mathbf{B}^T
\end{aligned} \tag{37}$$

342 One can proceed to solve system (12), obtain the estimate (13) and continue with
343 iterations. There are a number of possible refinements that reduce the computational
344 cost, such as one by taking advantage that some of the matrices we are working with
345 have rank K . For example, the matrix vector product $(\mathbf{H}\mathbf{Q}\mathbf{H}^T)\mathbf{a}$ can be computed
346 with $O(Kn)$ multiplications instead of $O(n^2)$. One can also solve a smaller system by
347 effectively compressing the data. Such refinements are important when n is large.

An important additional step, consistent with the objective of improving the scaling of
the computational problem with respect to m is the following. Form matrix

$$\mathbf{A}_p = [\mathbf{U}_z, \mathbf{U}]$$

which has $K + p$ orthonormal (i.e., $\mathbf{A}_p^T \mathbf{A}_p$ is identity matrix). Post-multiply $\mathbf{\Lambda}^T$ and
 \mathbf{M} with \mathbf{A}_p so that the system to solve is:

$$\begin{bmatrix} \mathbf{H}\mathbf{Q}\mathbf{H}^T + \mathbf{R} & \mathbf{H}\mathbf{X} \\ (\mathbf{H}\mathbf{X})^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} (\mathbf{\Lambda}^T \mathbf{A}_p) \\ (\mathbf{M} \mathbf{A}_p) \end{bmatrix} = \begin{bmatrix} \mathbf{H}\mathbf{Q}\mathbf{A}_p \\ \mathbf{X}^T \mathbf{A}_p \end{bmatrix} \tag{38}$$

348 By solving, we obtain the smaller matrices $(\mathbf{\Lambda}^T \mathbf{A}_p)$ and $(\mathbf{M} \mathbf{A}_p)$. Then approximate $\mathbf{\Lambda}^T$
349 and \mathbf{M} through K -rank approximations:

$$\mathbf{\Lambda}^T = (\mathbf{\Lambda}^T \mathbf{A}) \mathbf{A}_p^T \tag{39}$$

$$\mathbf{M} = (\mathbf{M} \mathbf{A}) \mathbf{A}_p^T \tag{40}$$

Thus, to compute the new \mathbf{s}

$$\hat{\mathbf{s}} = \mathbf{A}_p (\boldsymbol{\Lambda}^T \mathbf{A})^T (\mathbf{y} - \mathbf{h}(\tilde{\mathbf{s}}) + \tilde{\mathbf{H}}\tilde{\mathbf{s}}) \quad (41)$$

so that number of multiplications is $O(mK)$ rather than $O(mn)$, which is a slight improvement. The main benefit is in computing the covariance.

At the optimum, the posterior covariance correction satisfies

$$F \mathbf{A}_p = \mathbf{X} (\mathbf{M} \mathbf{A}_p) + \mathbf{Q} \mathbf{H}^T (\boldsymbol{\Lambda}^T \mathbf{A}_p) \quad (42)$$

Then, the correction to the covariance matrix, approximated through a rank $K + p$ matrix, is

$$\begin{aligned} F &\approx \mathbf{X} (\mathbf{M} \mathbf{A}_p) \mathbf{A}_p^T + \mathbf{Q} \mathbf{H}^T (\boldsymbol{\Lambda}^T \mathbf{A}_p) \mathbf{A}_p^T \\ &= \mathbf{X} (\mathbf{M} \mathbf{A}_p) \mathbf{A}_p^T + \mathbf{A} \mathbf{C} \mathbf{B}^T (\boldsymbol{\Lambda}^T \mathbf{A}_p) \mathbf{A}_p^T \end{aligned} \quad (43)$$

The sizes of the matrices involved are shown below

$$(m \times p) (p \times K) (K \times m), (m \times K) (K \times K) (K \times n) (n \times K) (K \times m) \quad (44)$$

while \mathbf{Q} is approximated through (33). The largest matrix to store is Km . Though the posterior covariance is not computed or stored, one can employ matrix-vector products to compute what is needed out of it. For example, say we want to compute the diagonal element $F(1,1)$ that is needed to find the variance of \mathbf{s}_1 and compute the Bayesian confidence intervals: One pre-multiplies by $[1 \ 0 \ 0 \ \dots \ 0 \ 0]$ and post-multiplies by the transpose of the same. Such operations can be arranged to be done quite efficiently.

Implementation

359 The approach requires a fast method to generate realizations with generalized covariance
360 matrix \mathbf{Q} . For practical situations of generating realizations of random fields in 1, 2, or
361 3 dimensions using regular grids, one approach is based on FFT [*Nowak et al.*, 2003].
362 For irregular grids, one can use methods from hierarchical matrices [*Saibaba et al.*, 2012;
363 *Ambikasaran et al.*, 2013]. It is also required to have an efficient algorithm for computing
364 unconditional covariance matrix times a vector. The same FFT and hierarchical methods
365 apply here as well. Note that the approach requires to use unconditional realizations and
366 the aforementioned products only at the start, before performing iterations.

367 If the standard Gauss-Newton method does not converge, one can introduce a between-
368 step search, either a line search or a trust region [for example, see *Zanini and Kitanidis*,
369 2009]. Assuming that the problem has been formulated correctly, there are really two
370 kinds of non-convergence issues. One is when the radius of convergence is large, compared
371 to the statistical error, but the starting solution is outside of it (the case of bad starting
372 point). In this case, one can start by using intermediate models, such as use a large
373 observation error (this makes the problem more linear-like and increases the radius of
374 convergence – linear problems have infinite radius of convergence) to get a very smooth
375 solution and gradually reduce the variance to the value that it should be. The other case is
376 that the radius of convergence is comparable to the statistical error of the solution (the case
377 of strong nonlinearity), and a plain Gauss-Newton simply does not consistently improve
378 the solution in successive steps even near where it should have converged. In this case,
379 it is essential to use a between-step procedure that guarantees that the objective function

380 can only improve. In the latter case, it is important to have a way to approximately
 381 compute the value of objective function. The objective function to be minimized is

$$\begin{aligned}
 J &= \frac{1}{2} (\mathbf{y} - \mathbf{h}(\mathbf{s}))^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{s})) + \frac{1}{2} (\mathbf{s} - \mathbf{X}\beta)^T \mathbf{Q}^{-1} (\mathbf{s} - \mathbf{X}\beta) \\
 &\simeq \frac{1}{2} (\mathbf{y} - \mathbf{h}(\mathbf{s}))^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{s})) + \frac{1}{2} \mathbf{s}^T \mathbf{A} \mathbf{C}^{-1} \mathbf{A}^T \mathbf{s}
 \end{aligned}
 \tag{45}$$

382 This expression is economical to compute, the second terms involving $O(Km)$ multipli-
 383 cations.

384 It should be noted that convergence is reached when the value of the J does not decrease
 385 much [e.g., *Kitanidis and Lane, 1985*]. For example, in the application shown in the
 386 next section, using a change of 0.01 would be enough for practical purposes.

6. Example

387 In this section, the methodology is tested on a toy problem. Application to a large-
 388 dimensional problem are reported elsewhere [*Lee and Kitanidis, 2014*].

389 Consider a problem of steady $1 - D$ flow, in a domain from 0 to 1, with variable and
 390 unknown conductivity, with constant recharge, and fixed-head conditions. The domain
 391 is discretized into 100 blocks, i.e., $m = 100$.

$$\frac{d}{dx} \left(K \frac{d\phi}{dx} \right) = -N \tag{46}$$

$$\phi(0) = \phi(1) = 1 \tag{47}$$

$$N = 10^{-5} \tag{48}$$

One may think of these as dimensionless quantities for a certain unit length and unit time. The model for the inversion is cubic generalized covariance function for the prior

and uncorrelated noise

$$Q_{ij} = 200 |x_i - x_j|^3, R_{ij} = (0.004)^2 \delta_{ij} \quad (49)$$

392 The drift matrix \mathbf{X} is m by 2; the first column of \mathbf{X} is ones and the second is the location
393 of blocks. The standard error in the 20 observations ($n = 20$) used, 0.004, is about 1%
394 of the difference in head between successive observation points.

395 For $K = 20$, the results for the textbook solution and the PCGA method are compared
396 against the true field. Between the two methods, the relative difference in norm between
397 the estimates is only 0.2% and in the computed covariance corrections 0.5%. Figure 3
398 shows the best estimates and the confidence intervals for the two methods.

399 When the procedure was repeated with fewer terms, like $K = 12$, PCGA gave subop-
400 timal results, but still surprisingly good since there is no way that all the information
401 in 20 measurements can be captured with just 12 components. The results are shown
402 in figure 4. The large-scale features are captured quite well, the method required fewer
403 iterations, and each iteration requires about 40% fewer calls to the forward problem. This
404 is an important advantage of the method. One can start with a small K performing the
405 first few iterations quickly and then to increase K when the largest features have been
406 identified and the solution is close to the optimum.

407 A strategy to reduce computational cost is to start with a very small number, like $K = 2$,
408 obtain a solution, use it as starting value for $K = 4$, and so on every time doubling the
409 number of components till no improvement can be achieved.

410 Focusing on the case $K = 20$, we will consider the interaction between the prior co-
411 variance and the forward operator, expressed through the Jacobian matrix \mathbf{H} . We will
412 use the Jacobian matrix that corresponds to the best estimate. There are many ways

413 to examine interactions between \mathbf{H} and \mathbf{Q} . Here we will compare the eigenvalues of
414 $\mathbf{H}\mathbf{Q}\mathbf{H}^T$ with those of $\mathbf{R} = \sigma^2\mathbf{I}$. Note that $\mathbf{H}\mathbf{Q}\mathbf{H}^T$ quantifies variability in the obser-
415 vation attributable to variability in unknowns whereas \mathbf{R} variability due to observation
416 error ("noise"). When the first dominates, the measurements are informative about the
417 unknowns or, to put it another way, the measurements can change the estimates more.

418 In figure 5 we compare the eigenvalues of the three matrices. Eigenvalues of $\mathbf{H}\mathbf{Q}\mathbf{H}^T$
419 that are much smaller than the variance of the measurement error, σ^2 , have negligible
420 effect and they might as well be set equal to 0.

421 This simply is another way to understand why head observation are not informative
422 about small scale fluctuations in the conductivity and why even modest increases in the
423 observation error can affect resolution quite dramatically. In terms of the method pro-
424 posed here, the figure suggests the reason why, if the measurement error is significant, a
425 small number K of components is sufficient to give the solution. For example, if $\sigma = 0.04$,
426 one would expect that 12 components would be enough even though there are 20 obser-
427 vations. This is verified in figure 6, where one can see that the textbook solution and the
428 approximate with just 12 components are practically the same.

429 In an application the scaling of \mathbf{Q} versus \mathbf{R} is a critical part of the process of solving
430 the inverse problem. However, this issue is beyond the scope of this paper that focuses
431 on how to compute the solution.

7. Discussion

432 We have presented a method with computational cost that has near linear scaling with
433 the size of the vector of the unknowns. Thus, this method is a promising alternative for
434 problems with unknowns in the order of millions. The method has two phases.

435 In the first phase, the problem is parameterized in terms of orthogonal components.
436 The method proposed herein is purely algebraic and "black box", in the sense that there
437 is no need for analytical expansions. The user needs to provide the geostatistical model
438 in terms of the drift and covariance functions. Assuming that one has access to software
439 tools for fast covariance matrix vector multiplications (like FFT and hierarchical) that are
440 becoming increasingly easier to find, this phase involves computations with near linear
441 scaling. This is done once for every geostatistical model and involves no forward runs
442 and Gauss-Newton iterations.

443 The second phase is Gauss-Newton iterations. One advantage of the method is that
444 the number of runs of the forward model (typically a PDE solver) is reduced compared
445 to more standard approaches. Only the forward model used as a black box is required
446 and, in this version, there is no need for adjoint-state solvers. A second advantage is that
447 the overall computational cost of matrix operations (excluding the forward runs) scales
448 linearly rather than quadratically with the number of unknowns. The same holds true
449 for storage requirements, including the computation of the posterior covariance. For a
450 large problem, these savings can be by orders of magnitude.

451 One advantage of the approach is that the computations in the second phase, which
452 is usually the most demanding part computationally, are highly parallelizable. Like in
453 ensemble methods, one can run forward solvers in parallel.

454 The method is not meant as the solution to all types of inverse problems. It is best for
455 problems where we can focus on seeking relatively smooth solutions and the measurements
456 have limited information content, a situation that is often encountered.

457 This work focuses on the mathematics of the method. Only a simple and small problem
458 is presented here, to illustrate that the mathematical equations work as they are supposed
459 to and to give some insights into the problem. The scaling follows from a theoretical
460 analysis of the cost of each operation. Future work will present application to large
461 problems with analysis of actual cost and effectiveness.

462 **Acknowledgments.**

463 This material is based upon work supported by the Department of Energy under Award
464 Number DE-FE0009260. Additional support was received from the National Science Foun-
465 dation through its ReNUWit Engineering Research Center (www.renuwit.org; NSF EEC-
466 1028968).

References

- 467 Ambikasaran, S., Y. Li, E. Darve, and P. K. Kitanidis (2013), Large-scale stochastic linear
468 inversion using hierarchical matrices, *Computational Geosciences*, *published online*, doi:
469 10.1007/s10596-013-9364-0.
- 470 Berger, J. O. (2006), The case for objective Bayesian analysis, *Bayesian Analysis*, *1*(3),
471 385–402.
- 472 Bui-Thanh, T., C. Burstedde, O. Ghattas, J. Martin, G. Stadler, and L. C. Wilcox (),
473 Extreme-scale UQ for Bayesian inverse problems governed by PDEs, in *Proceedings of*
474 *the International Conference on High Performance Computing, Networking, Storage*
475 *and Analysis*, pp. 1–11, IEEE Computer Society Press.
- 476 Butler, J. J. J., C. D. McElwee, and G. C. Bohling (1999), Pumping tests in networks
477 of multilevel sampling wells; motivation and methodology, *Water Resources Research*,

- 479 Cardiff, M., W. Barrash, P. Kitanidis, B. Malama, A. Revil, S. Straface, and E. Rizzo
480 (2009), A potential-based inversion of unconfined steady-state hydraulic tomography,
481 *Ground Water*, 47(2), 259–270, doi:10.1111/j.1745-6584.2008.00541.x.
- 482 Cardiff, M., W. Barrash, and P. K. Kitanidis (2013), Hydraulic conductivity imaging from
483 3-d transient hydraulic tomography at several pumping/observation densities, *Water*
484 *Resour. Res.*, 49, 7311–7326, doi:10.1002/wrcr.20519.
- 485 Carlin, B. P., and T. A. Louis (2000), *Bayes and empirical Bayes methods for data anal-*
486 *ysis*, 2nd edition ed., Chapman and Hall CRC, Boca Raton.
- 487 Carrera, J., and S. P. Neuman (1986), Estimation of aquifer parameters under transient
488 and steady state conditions, 1. Maximum Likelihood method incorporating prior infor-
489 mation, *Water Resour. Res.*, 22(2), 199–210.
- 490 Chatterjee, A., A. M. Michalak, J. L. Anderson, K. L. Mueller, and V. Yadav (2012),
491 Toward reliable ensemble Kalman filter estimates of CO₂ fluxes, *Journal of Geophysical*
492 *Research: Atmospheres*, 117(D22), D22,306, doi:10.1029/2012jd018176.
- 493 Evensen, G. (1994), Sequential data assimilation with a non-linear quasi-geostrophic
494 model using Monte Carlo methods to forecast error statistics, *J. Geophys Res*, 99('C5'),
495 10.143–10.162, doi:10.1029/94JC00572.
- 496 Evensen, G. (2003), The ensemble Kalman filter: Theoretical formulation and practical
497 implementation, *Ocean Dynamics*, 53, 343–367, doi:10.1007/s10236-003-0036-9.
- 498 Evensen, G. (2006), *Data Assimilation: The Ensemble Kalman Filter*, Springer.
- 499 Fong, W., and E. Darve (2009), The black-box fast multipole method, *Journal of Com-*
500 *putational Physics*, 228(23), 8712–8725, doi:10.1016/j.jcp.2009.08.031.

501 Gavalas, G. R., P. C. Shah, and J. H. Seinfeld (1976), Reservoir history matching by
502 Bayesian estimation, *Soc. Petrol. Eng. J.*, 16, 337–350, box 3C.

503 Golub, G. H. (1965), Numerical methods for solving linear least squares problems, *Numer.*
504 *Math.*, 17, 206–216.

505 Golub, G. H., and C. F. Van Loan (1989), *Matrix Computations*, 462 pp., Johns Hopkins
506 Univ. Press, Baltimore.

507 Halko, N., P. G. Martinsson, and J. A. Tropp (2011), Finding structure with randomness:
508 Probabilistic algorithms for constructing approximate matrix decompositions, *SIAM*
509 *Review*, 53(2), 217–288.

510 Kitanidis, P. K. (1983), Statistical estimation of polynomial generalized covariance func-
511 tions and hydrologic applications, *Water Resources Research*, 19(4), 909–921, box 5B.

512 Kitanidis, P. K. (1993), Generalized covariance functions in estimation, *Mathematical*
513 *Geology*, 25(5), 525–540.

514 Kitanidis, P. K. (1995), Quasilinear geostatistical theory for inversing, *Water Resour.*
515 *Res.*, 31(10), 2411–2419.

516 Kitanidis, P. K. (2010), Bayesian and geostatistical approaches to inverse problems, in
517 *Large-scale inverse problems and quantification of uncertainty*, edited by L. Biegler,
518 G. Biros, O. Ghattas, M. Heinkenschloss, D. Keyes, B. Mallick, Y. Marzouk, L. Tenorio,
519 B. van Bloemen Waanders, and K. Willcox, pp. 71–85, Wiley.

520 Kitanidis, P. K., and R. W. Lane (1985), Maximum Likelihood parameter estimation of
521 hydrologic spatial processes by the Gauss-Newton method, *J. Hydrology*, 79, 53–71, box
522 4A.

523 Kitanidis, P. K., and E. G. Vomvoris (1983), A geostatistical approach to the inverse
524 problem in groundwater modeling (steady state) and one-dimensional simulations, *Wa-*
525 *ter Resour. Res.*, *19*(3), 677–690.

526 Lee, J., and P. K. Kitanidis (2014), Large-scale hydraulic tomography and joint inversion
527 of head and tracer data using the PCGA method, *Water Resour. Res.*, *submitted*.

528 Li, J. Y., S. Ambikasaran, E. Darve, and P. K. Kitanidis (2014), A Kalman filter powered
529 by H-matrices for quasi-continuous data assimilation problems, *Water Resour. Res.*, *in*
530 *review*.

531 Li, W., and O. A. Cirpka (2006), Efficient geostatistical inverse methods for structured
532 and unstructured grids, *Water Resour. Res.*, *42*(W06402), doi:10.1029/2005WR004668.

533 Linde, N., A. Binley, A. Tryggvason, L. B. Pedersen, and A. Revil (2006), Improved hydro-
534 geophysical characterization using joint inversion of cross-hole electrical resistance and
535 ground-penetrating radar traveltime data, *Water Resources Research*, *42*(12), W12,404,
536 doi:10.1029/2006wr005131.

537 Liu, X., and P. K. Kitanidis (2011), Large-scale inverse modeling with an application in hy-
538 draulic tomography, *Water Resour. Res.*, *47*(2), W02,501, doi:10.1029/2010wr009144.

539 Marzouk, Y. M., and H. N. Najm (2009), Dimensionality reduction and polynomial
540 chaos acceleration of Bayesian inference in inverse problems, *Journal of Computational*
541 *Physics*, *228*(6), 1862–1902, doi:10.1016/j.jcp.2008.11.024.

542 Matheron, G. (1973), The intrinsic random functions and their applications, *Applied Prob-*
543 *ability*, *5*, 439–468, box 5B.

544 Neuman, S. P. (1980), A statistical approach to the inverse problem of aquifer hydrology,
545 3. Improved method and added perspectives, *Water Resour. Res.*, *16*(2), 331–346.

546 Nowak, W., S. Tenkleve, and O. Cirpka (2003), Efficient computation of linearized cross-
547 covariance and auto-covariance matrices of interdependent quantities, *Mathematical Ge-*
548 *ology*, *35*(1), 53–66.

549 Pollock, D., and O. A. Cirpka (2012), Fully coupled hydrogeophysical inversion of a lab-
550 oratory salt tracer experiment monitored by electrical resistivity tomography, *Water*
551 *Resour. Res.*, *48*(1), W01,505, doi:10.1029/2011wr010779.

552 Rubin, Y., X. Chen, H. Murakami, and M. Hahn (2010), A Bayesian approach for inverse
553 modeling, data assimilation, and conditional simulation of spatial random fields, *Water*
554 *Resour. Res.*, *46*(10), W10,523, doi:10.1029/2009wr008799.

555 Saibaba, A. K., and P. K. Kitanidis (2012), Efficient methods for large-scale linear inver-
556 sion using a geostatistical approach, *Water Resour. Res.*, *48*(5), W05,522.

557 Saibaba, A. K., S. Ambikasaran, J. Y. Li, P. K. Kitanidis, and E. F. Darve (2012),
558 Application of hierarchical matrices to linear inverse problems in geostatistics, *Oil Gas*
559 *Sci. Technol.*, *67*(5), 857–875, doi:10.2516/ogst/2012064.

560 Slater, L., A. M. Binley, W. Daily, and R. Johnson (2000), Cross-hole electrical imaging
561 of a controlled saline tracer injection, *Journal of Applied Geophysics*, *44*(2), 85–102,
562 doi:10.1016/S0926-9851(00)00002-1.

563 Woodbury, A. D., and Y. Rubin (2000), A full-Bayesian approach to parameter inference
564 from tracer travel time moments and investigation of scale effects at the Cape Cod
565 experimental site, *Water Resour. Res.*, *36*(1), 159–171.

566 Yadav, V., and A. M. Michalak (2013), Improving computational efficiency in large linear
567 inverse problems: an example from carbon dioxide flux estimation, *Geosci. Model Dev.*,
568 *6*(3), 583–590, doi:10.5194/gmd-6-583-2013.

569 Yeh, T.-C. J., and X. Liu (2000), Hydraulic tomography: Development of a new aquifer
570 test method, *Water Resour. Res.*, *36*(8), 2095–2105.

571 Zanini, A., and P. K. Kitanidis (2009), Geostatistical inversing for large-contrast transmis-
572 sivity fields, *Stochastic Environmental Research and Risk Assessment*, *23*(5), 565–577,
573 doi:10.1007/s00477-008-0241-7.

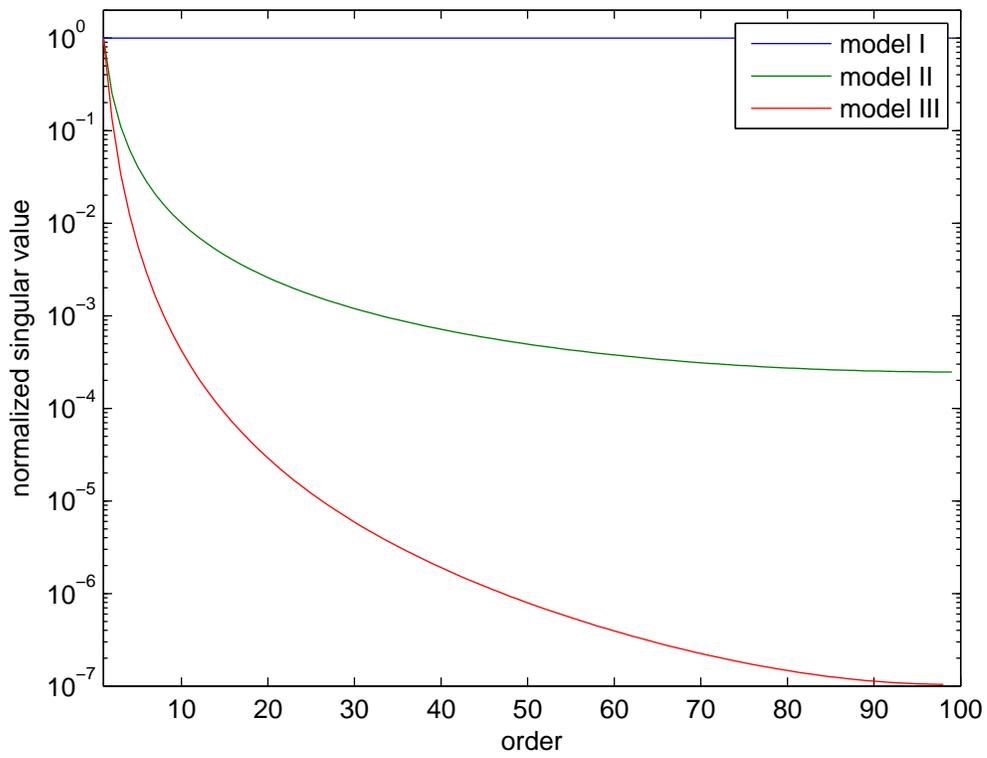


Figure 1. Normalized matrix spectrum (singular values) for three models.

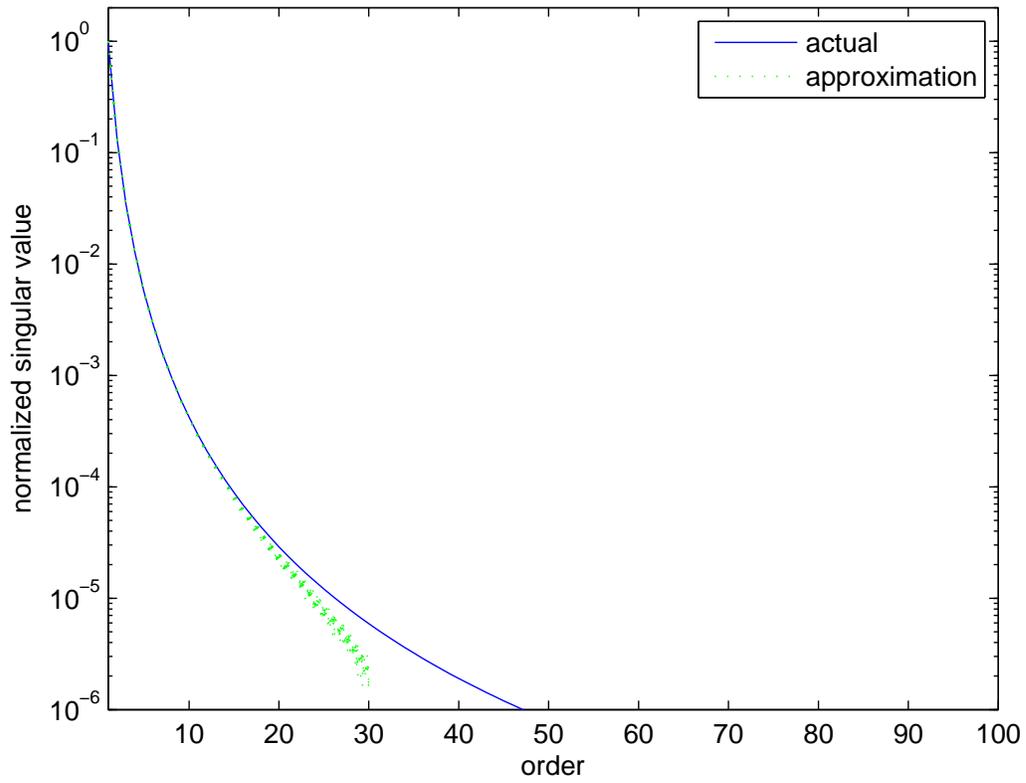


Figure 2. Results of twenty repetitions of the process for model III and $K = 30$.

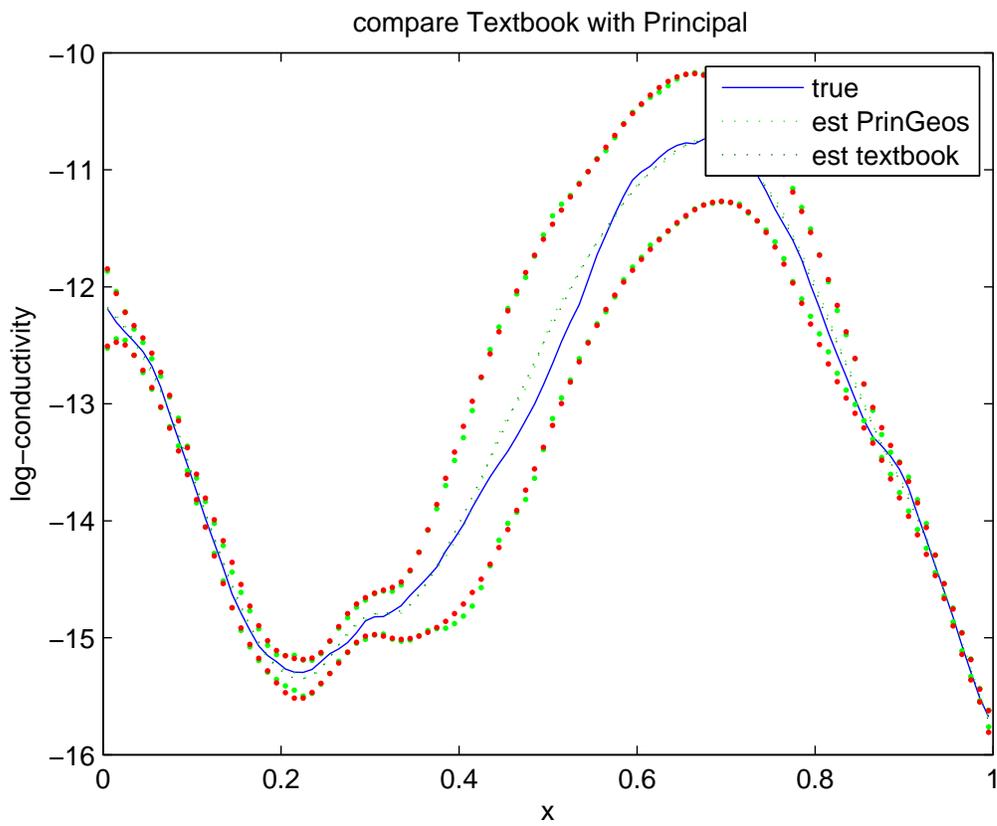


Figure 3. With $K=20$ components: The true, best estimate, and confidence intervals. Red is for textbook and green for the new method.

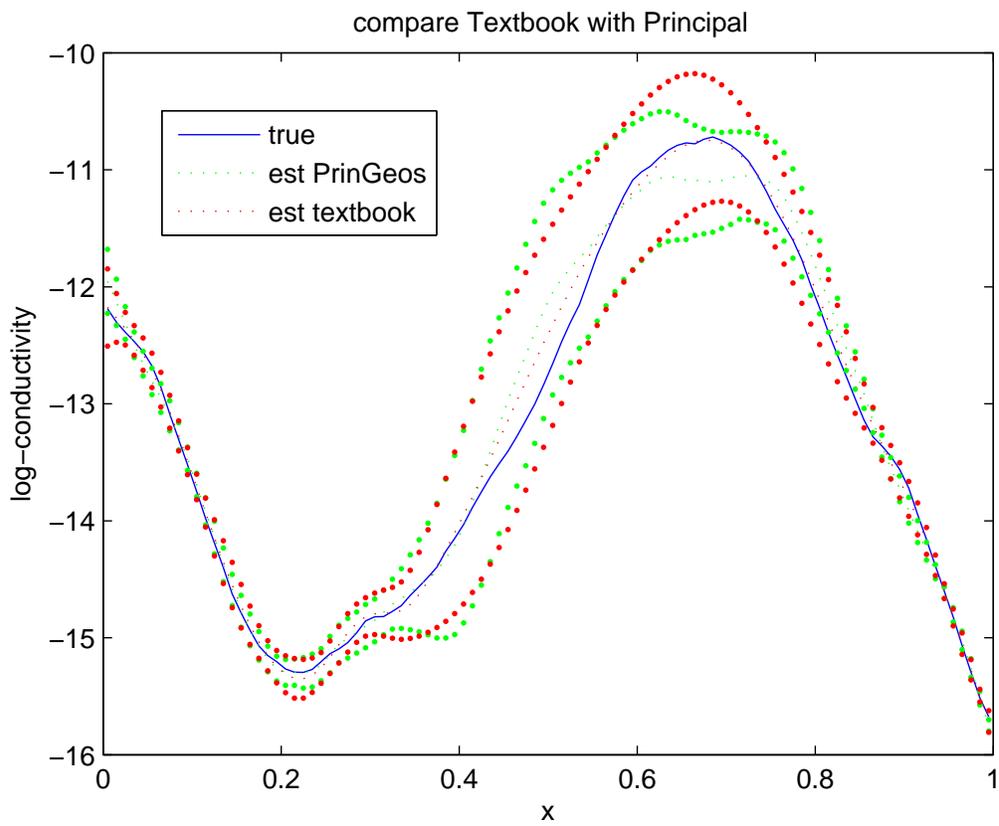


Figure 4. With $K=12$ components: The true, best estimate, and confidence intervals. Red is for textbook and green for the new method.

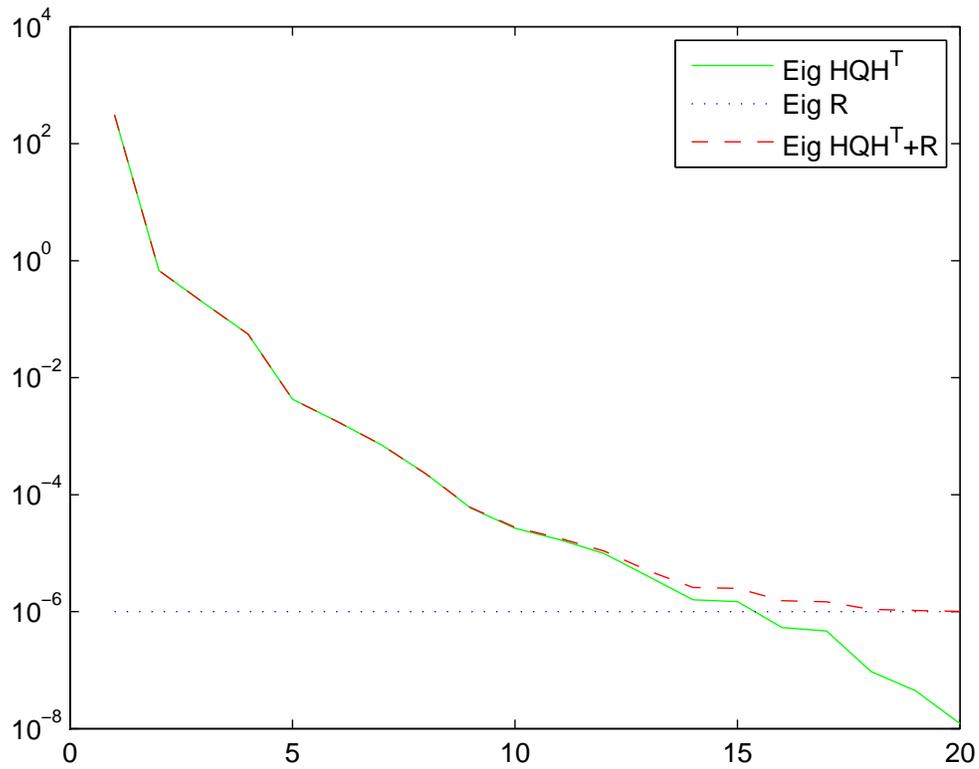


Figure 5. Comparison of eigenvalues of HQH^T and $HQH + R$.

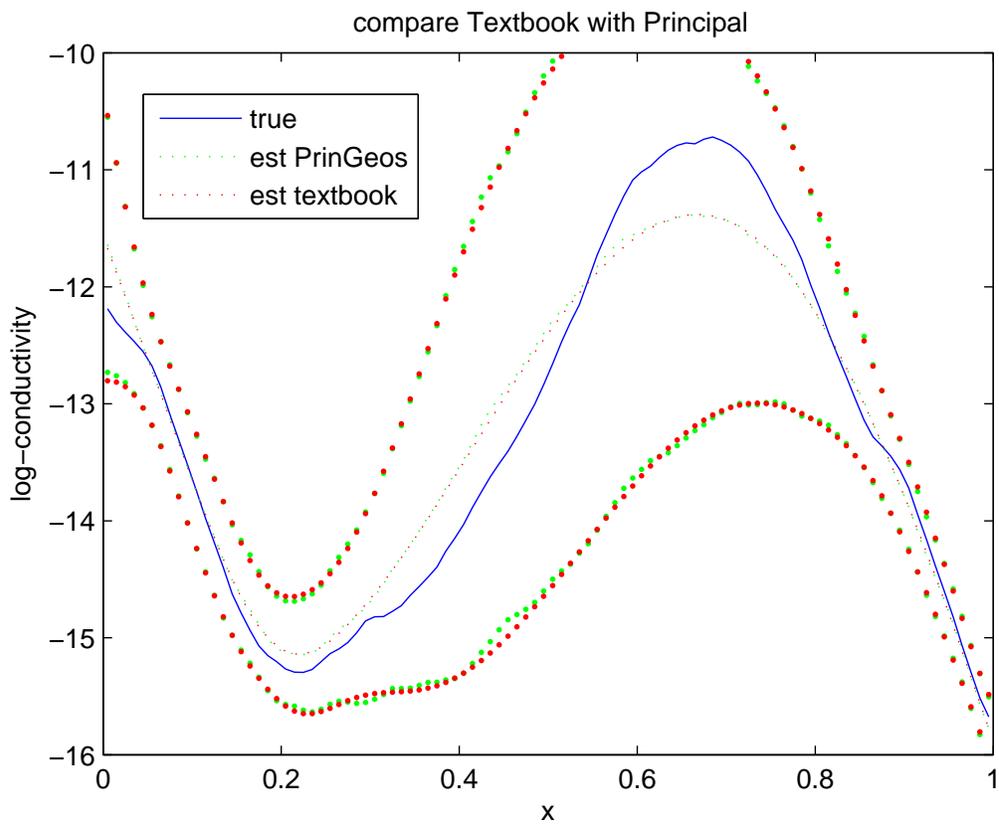


Figure 6. High-noise case with $K=12$ components: The true, best estimate, and confidence intervals. Red is for textbook and green for the new method.