# Manifold Methods for Assimilating Geophysical and Meteorological Data in Earth System Models and their Components

Ammar Safaie<sup>1</sup>, Chinh Dang<sup>2</sup>, Han Qiu<sup>1</sup>, Hayder Radha<sup>2</sup>, Mantha S. Phanikumar<sup>1\*</sup>

#### Abstract

A novel manifold method of reconstructing dynamically evolving spatial fields is presented for assimilating data from sensor networks in integrated land surface – subsurface, oceanic / lake models. The method was developed based on the assumption that data can be mapped onto an underlying differential manifold. In this study, the proposed method was used to reconstruct meteorological forcing over Lake Michigan, the bathymetry of an inland lake (Gull Lake), and precipitation over the Grand River watershed in Michigan. In the first case study, hourly interpolated meteorological forcing data were used to run a three-dimensional hydrodynamic model of Lake Michigan to quantify the improvement that results from the use of the new interpolation method. In the second example, the bathymetry of Gull Lake was reconstructed from measured scatter point data using the manifold technique. A hydrodynamic model of Gull Lake was developed and further improved using improved bathymetry. In the last case study, daily participation data over a six-year period were interpolated over the Grand River watershed and used as input to an integrated, distributed hydrologic model. All three examples illustrate the superior performance of the manifold method over standard methods in terms of accuracy and computational efficiency. Our results also indicate that evaluating the relative performance of interpolation methods using the cross-validation method can lead to misleading conclusions about the relative performance of the methods.

#### 1 Introduction

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

Water continuously circulates between the Earth's surface and the atmosphere and moves through watersheds to become stored underground as groundwater or in surface reservoirs such as lakes. This circulation is a key aspect of Earth system models or components thereof and meteorological forcing plays an important role in modeling coupled physical - biogeochemical processes. Process-based modeling approaches describe land-lake-atmosphere interactions by explicitly considering the spatiotemporal variability of meteorological forcing fields (Xue et al., 2015). There are a variety of numerical ocean models, such as POM (Blumberg, and Mellor, 1987), ROMS (Haidvogel and Beckmann, 1999), and FVCOM (Chen et al., 2006), and their performance is highly dependent on how realistic the distributions of surface forcing fields (wind stress, heat flux, precipitation, evaporation) are. These forcing fields can be obtained from observational data, the output of a weather forecast model, or a combination thereof (Xue et al., 2015). In all cases it is crucial to use a suitable interpolation method in order to assimilate observations into the models, and estimate variables at unsampled locations and/or times. For instance, currents in large lakes such as the Laurentian Great Lakes are mostly controlled by wind. Therefore, by improving the representation of wind fields in models of lake circulation, we expect to describe coupled biophysical processes in lakes more accurately. For example, Safaie et al. (2016) demonstrated that improved representation of meteorological fields based on natural neighbor interpolation of weather station data produced superior results for currents and bacterial concentrations relative to similar results based on a nearest neighbor interpolation of the same data.

*In situ* observations generally have sparse and inhomogeneous distribution in space and time, and it is often infeasible to accurately reconstruct the true field from the data. However, more information about the structure of the field and its evolution, allows for better approximations (Barth et al., 2008). Various deterministic [e.g., nearest neighbor, natural neighbor, inverse distance weighting (IDW), spline, polynomial and geostatistical (e.g. kriging) interpolation methods have been developed to generate spatial fields. There have been numerous efforts to compare different spatial interpolation methods in order to identify the best method for a given model application. Many researches have used cross-validation for assessing the performance of the interpolation methods. In this method, a subset of the original dataset is withheld to be used later for validating the interpolated field constructed from the rest of the observational data. Mean error (ME), root mean square error (RMSE) and the coefficient of determination (R<sup>2</sup>) are commonly used to evaluate the performance of each interpolation method (Suparta and Rahman, 2016). However, every problem has a unique method of interpolation that works best for a given distribution of observations and the intended use of the interpolated data. Density of a sensor network, spatial variability of the variable of interest and its distribution, and observational errors, all influence the accuracy of the interpolated field (MacEachren and Davidson, 1987). For example, Luo et al. (2008) compared seven spatial interpolation techniques to identify which method produced the best estimation of the wind speed data recorded across England and Wales. Their study showed that kriging is the best method, and that the thin plate spline method had higher ME and RMSE values. However, in (Suparta and Rahman, 2016) the performance of the thin plate spline interpolation based on the RMSE and R<sup>2</sup> values was found to be better than kriging for less dense data points over the selected interpolation surface. Therefore, comparing interpolation methods using the cross-validation method without considering the data structure

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

and the purpose of interpolation is not guaranteed to produce the best representation of the underlying data.

Precipitation is another important component of the water cycle. Spatial distribution of precipitation influences the hydrologic response of a watershed (Daly et al., 2002). Basin responses of rainfall-runoff processes are closely related to the spatial variability of precipitation (Anquetin et al., 2010; Bell and Moore, 2000; Beven and Hornberger, 1982; Obled et al., 1994; Schuurmans and Bierkens, 2007; Syed et al., 2003; Tetzlaff and Uhlenbrook, 2005). Hydrologists usually use rain gauge data to obtain rainfall patterns over a watershed. Bell and Moore (2000) found that the responses of a distributed hydrologic model are sensitive to the locations of the rain gauges within the catchment and hence to the spatial variability of rainfall. Nicótina et al. (2008) performed numerical experiments to study the effects of different spatial resolutions of rainfall on various catchments. They found that the catchment response is sensitive to the spatial distribution of rainfall only when water residence time in the channels is comparable to the hillslope travel time; thus rainfall spatial heterogeneity likely plays a more important role in affecting the runoff response of large watersheds (typically larger than 103 km²) than smaller watersheds.

Various efforts have been devoted to improving the representation of spatially-distributed rainfall fields in hydrologic modeling. For example, Ly et al. (2011) compared seven interpolation methods for daily rainfall and found that geostatistical methods such as kriging and IDW algorithms significantly outperformed the Thiessen polygon method, which is also known as the nearest neighbor method. Sun et al. (2000) tested different rainfall estimation methods on a 1060 km<sup>2</sup> catchment in Australia to evaluate the flood modeling capabilities of a hydrologic model and found that blending radar and rain gauge data in a co-kriging framework provided

better performance compared with an approach based on kriging of the rain gauges data alone. Masih et al. (2011) used a semi-distributed model, Soil Water Assessment Tool (SWAT), to simulate the Karkheh River basin in Iran with two different precipitation interpolation methods. Their results, based on a comparison of simulated and observed discharges using the metrics of coefficient of determination and the Nash-Sutcliffe efficiency (NASH), showed that precipitation interpolation with inverse distance and elevation weighting technique produced better performance relative to the SWAT default method of nearest neighbor method.

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

Accurate representation of geophysical features such as topography and bathymetry is also important in earth system models and their components, and model performance depends on the interpolation method used to assign the topographic information over a numerical mesh in processed-based models. Yan et al. (2014) compared different interpolation methods, including IDW, global polynomial interpolation, local polynomial interpolation, radial basis functions, ordinary kriging (OK), simple kriging (SK), universal kriging (UK), and co-kriging (CK) to determine the water/land boundary point elevation based on in situ water level data from 14 control stations in Dongting Lake. They used a cross-validation method to select the optimal method, which was found to be the OK method. Merwade (2009) studied the effect of spatial trend on interpolation of river bathymetry, and compared the performance of different interpolation methods. The number of measurements and their spatial arrangement, as well as channel morphology and geology were found to influence the accuracy of the interpolation results (Merwade, 2009). Due to the effects of these and other factors on the performance of various methods, comparisons of different spatial interpolation methods could not point out the best universal interpolation method (Li and Heap, 2008; Šiljeg et al., 2015).

In this paper, we propose a novel manifold method to assimilate different types of spatiotemporal data in integrated earth system models based on the hypothesis that an environmental dataset (including independent variables such as longitude, latitude, and time, and the measured variables of interest) can be mapped onto an underlying differential manifold. A manifold (M) is an *n*-dimensional topological space such that each point of M and its neighborhood can be approximated by a small flat piece in the Euclidean space,  $\mathbb{R}^n$ . We can think of a manifold as a set of low-dimensional curves and surfaces within higher dimension Euclidean spaces (Victor and Pollack, 2010). Some typical examples of manifolds are smooth surfaces, such as a torus (Fig. 1a) or a sphere (Fig. 1b), where each point and its neighborhood can be approximated by a small flat linear-subspace within the three-dimensional Euclidean space. Another example of a manifold in a high dimensional space is a Calabi-Yau manifold which has found important applications in theoretical physics (e.g. superstring theory). Fig. 1c shows a two-dimensional cross-section of a six-dimensional Calabi-Yau manifold. Surfaces of all these three manifolds are not a Euclidean space. The laws of the Euclidean geometry, however, are valid locally.

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

Working directly in the high dimensional space generally involves dealing with complex algorithms. Modeling the high dimensional data using manifolds with fewer degrees of freedom has captured a great deal of attention recently (Zhang et al., 2016). The use of low-dimensional manifolds not only reduces computational load for further processing, but also helps visualize the entire dataset, which is an important initial step to make sense of the data before proceeding with more goal-directed modeling and analyses (Ma et al., 2011).

The classic method of principal component analysis (PCA), including the more recent Robust PCA (Candès et al., 2011), is arguably the most popular framework for approximating a set of

high-dimensional data points by a low dimensional linear space. PCA does not work well when the underlying data structure is non-linear. Under such scenario, PCA methods require approximating the input data points using higher dimensional linear spaces to reduce the approximation error. Instead of assuming that the input data follow a linear structure, an alternative solution is learning the true underlying low-dimensional structure of the data. The problem of non-linear dimensionality reduction for a set of high dimensional data points is known as manifold learning. Examples of early works for non-linear dimensionality reduction include Isomap (Tenenbaum et al., 2000), local linear embedding (LLE) (Roweis, 2000), and Eigenmaps (Belkin and Niyogi, 2003), which have been used to learn the manifold structure of data. Since then, the manifold model has been exploited extensively in numerous applications such as face recognition, action classification, segmentation, image denoising, image/video super-resolution, and multi-scale image analysis (Carin et al., 2011; Dang et al., 2013, 2014; Dang and Radha, 2015).

Most of the above manifold learning methods have been inspired by linear techniques, mainly based on the assumption that non-linear manifolds can be approximated by locally linear parts (Mordohai and Medioni, 2010). Two pioneering works in this area are the Isomap approach (Tenenbaum et al., 2000) and the LLE algorithm (Roweis, 2000). The Isomap algorithm aims to preserve the geodesic distance among points from the input dataset. On the other hand, the LLE algorithm targets the local linear geometry of neighbors in a manifold. Numerous works on manifold learning have been further developed. A comprehensive review of prior works can be found in (van der Maaten et al., 2009).

Data assimilation methods seek to estimate both model parameters and model states and significant progress has been made in the development of joint state-parameter estimation

methods such as the ensemble Kalman filter (EnKF) and its variants (Moradkhani et al., 2005; Evensen, 2006; Pathiraja et al., 2016). The parameters estimated using these approaches are usually time-dependent although most dynamical models used in earth sciences use time-invariant parameters. Yang et al. (2007) describe an alternative optimization-assimilation approach for soil moisture in which they first estimate parameters within a long optimization window and then estimate model states within a short assimilation window. Due to the computational nature of the three-dimensional process-based models considered in the present study, we do not focus on the use of methods similar to the EnKF method. Instead we estimate optimal (relative to the observations) model parameters and states over the simulation period. Future papers will focus on the use assimilation methods such as the EnKF with the framework of manifold methods.

In this paper, the effectiveness of the presented manifold algorithm is evaluated through assimilation of geophysical and meteorological data in integrated land surface – subsurface and lake models, although the methods described are general and can be used in many other areas of computational geosciences. We first apply the proposed method to reconstruct wind fields (time-varying vector fields) over Lake Michigan. Since currents in Lake Michigan are primarily driven by wind, we expect to improve the simulation of hydrodynamic and biophysical variables of interest by improving the wind fields. Instead of relying on the cross-validation of interpolated wind data, however, we use a well-tested hydrodynamic model of Lake Michigan and compare current measurements with simulated currents to test the interpolation methods. The manifold methods are used to reconstruct meteorological data including wind fields, cloud cover, dew point, pressure, shortwave and longwave solar radiation, relative humidity, and air temperatures for improved simulation of circulation in Lake Michigan. Then the method is applied to

assimilate bathymetry data as a scalar field for use in a hydrodynamic model of Gull Lake - a relatively large (8 km² surface area) and deep (34 m maximum depth) clear water lake in the in Kalamazoo County in southwestern Michigan. In the third example, time-dependent fields of participation are interpolated over the Grand River watershed and used as input for an integrated, land surface – subsurface processes model (PAWS+CLM; Shen and Phanikumar, 2010; Shen et al., 2013). Grand River watershed is located in the middle of Michigan's Lower Peninsula and is the second largest watershed in Michigan. In this example, the manifold method is tested using stream discharge outputs of the PAWS+CLM model which has been tested in several catchments in the past (Niu et al., 2014; Niu and Phanikumar, 2015; Shen et al., 2014, 2013).

#### 2 Materials and Methods

#### 2.1 Manifold approach

Based on Einstein's theory of relativity, physical events are located on the continuum (manifold) of space-time. Therefore, station locations and times of observations form a space-time manifold viewed as a four-dimensional vector space. One way to handle spatiotemporal interpolation problems, inspired by this concept, is to integrate space and time simultaneously (Li and Revesz, 2004). An underlying assumption behind this approach is that time and space dimensions can be treated as equally important (Li et al., 2014). In order to add time as another dimension of space, time values are needed to be scaled for a spatiotemporal dataset by a scaling speed (Li et al., 2014; Schwab and Beletsky, 1998). For a point measurement, we can then define a four-vector  $P^{\mu} = (ct, \vec{x})$  where c is a time scale, t is the time coordinate and  $\vec{x}$  is a three-dimensional vector space. We assume that the set of high-dimensional data points P (and the estimated data points  $P_0$ ) belongs to a differential manifold M, which may be curved and have

a complicated topology, but the neighborhood of each point is approximately similar to a small piece of Euclidean space (resembles  $\mathbb{R}^D$ ). Since a traditional distance measure is built upon the geometry of Euclidean space, we adapt the calculation to a neighborhood or a small region of the assumed manifold.

An example of a one-dimensional curve in Fig. 2 illustrates the general idea of the manifold estimation approach. The set of points P in Fig. 2 includes sample data points where we have measured data as well as a point  $P_0$  where data are missing. For example, in the context of the wind field data, one full measurement (or data point) includes five components: time, longitude, latitude, wind speed, and wind direction. The partially missing data point may contain known components (time, longitude, latitude) and unknown or missing components (wind speed and wind direction).

Consider a smooth n-dimensional manifold M embedded in a D-dimensional Euclidean space. Suppose that it is desired to estimate the wind field for a data point  $P_0 \in \mathbb{R}^n$  from a set of training data points that belong to a manifold M. The space/time coordinates of the point (the independent variables) are known, however, the data (the dependent variable) are missing. We denote  $P_0 = \begin{bmatrix} P_0^\mu \\ P_0^\nu \end{bmatrix} \in \mathbb{R}^n$  as the data point using the superscript  $\mu$  to denote the independent variables and the superscript  $\nu$  to denote the dependent variable which is the missing component of interest here.  $P_0^\mu \in \mathbb{R}^{n_\mu}$  is the sub-vector of the known components, and  $P_0^\nu \in \mathbb{R}^{n_\nu}$   $(n_\mu + n_\nu = n)$  is the corresponding sub-vector (e.g., wind velocity vector) for the missing component where  $P_0^\nu = \vec{V} = (u, v)$  and u, v are the orthogonal components of the wind velocity

 $(\vec{V})$ . The training data points, for example  $P = \{P_1, P_2, ..., P_7\}$  in Fig. 2, also include the two components  $P_i = \begin{bmatrix} P_i^\mu \\ P_i^\nu \end{bmatrix} \in \mathbb{R}^n$ , but there is no missing component here since both dependent and independent variables are assumed to be known at the nearby stations. Given a point  $P_0^{\,\mu}\in\mathbb{R}^{n_\mu}$  , the algorithm locates a set of nearest points to  $P_0^{\mu}$  based on the distances  $d(P_i^{\mu}, P_0^{\mu})$  between pairs of points  $P_i^\mu$  and  $P_0^{\mu}$ . In order to determine local neighbors of  $P_0^{\mu}$ , we can calculate the distances between  $P_0^{\ \mu}$  and either all other points within a fixed radius  $\epsilon$ , or all of its k nearest neighbors (Tenenbaum et al., 2000). Then, a tangent space (linear subspace) of the manifold M at the point  $P_0$  is created from the set of nearest points (Fig. 3a), denoted by  $T_{P_0}(M) = \begin{bmatrix} T^{\mu} \\ T^{\nu} \end{bmatrix}$ where  $T^{\mu}, T^{\nu}$  denote the tangent spaces for the independent and dependent variables in the data at  $P_0^{\mu}$  and  $P_0^{\nu}$ . Finally, the point  $P_0 = \begin{bmatrix} P_0^{\mu} \\ P_0^{\nu} \end{bmatrix} \in \square^n$  will be located as the closest point that belongs to that tangent space.

To represent the closest distance between a point and a tangent space, we use the Euclidean distance of an orthogonal projection from that point to the tangent space. Since a tangent space is a linear space (or affine space in a more general case), one point can orthogonally project into that space. The question is how to define neighbors for each data point? The underlying idea is how to define similarity distance among the training data points, and then the overall similarity matrix. Several methods have been considered in the past, such as k-nearest neighbors (Press et al., 2007),  $\epsilon$ -ball method (Allard et al., 2012) or the use of sparse representation theory (Dang et al., 2014, 2013; Dang and Radha, 2015). To approximate the wind field, we do not focus on

- analyzing a predetermined set of tangent spaces as was done earlier (Dang et al., 2014, 2013), but instead create a tangent space for each input data point as creating a tangent space for a given input data leads to a better approximation of the manifold.
- The estimation of  $P_0^{\nu}$  is performed using the following steps:
- Given a set of neighboring points, estimate the tangent space T<sup>μ</sup> at the point of interest, P<sub>0</sub>:
   Details of the method for creating a tangent space from a set of data points are described in
   Appendix-A and in Dang et al. (2014). One simple method is to create a tangent space using
   singular value decomposition (SVD, Press et al., 2007). By way of an example in Fig. 2, a
   tangent space (the line b) is created for P<sub>0</sub><sup>μ</sup> from a set of its neighboring points (P<sub>5</sub><sup>μ</sup> and P<sub>6</sub><sup>μ</sup>).
   This tangent space at P<sub>0</sub><sup>μ</sup> ∈ M is denoted by T<sup>μ</sup>.
- 260 2. Find the orthogonal projection of  $P_0^{\mu}$  onto the tangent space:
- The closest point  $P' \in T^{\mu}$  to the given point  $P_0^{\mu}$  is located at the intersection of the line b and the line perpendicular to it which passes through the point  $P_0^{\mu}$ . P' which is a projection of  $P_0^{\mu}$  onto the subspace  $T^{\mu}$  can be represented as an approximation of point  $P_0^{\mu}$ . The orthogonal projection of vector point  $P_0$  in a high-dimensional space onto a low-dimensional vector subspace is given by:

$$266 \qquad \prod_{\mathbf{T}^{\mu}} (P_0^{\mu}) = A(A^T A)^{-1} A^T P_0^{\mu} = A A^+ P_0^{\mu}$$
 (1)

where  $A = T^{\mu} \in \mathbb{R}^{D \times n}$  is a full rank matrix with n < D containing the set of points on the tangent space of  $P_0^{\mu}$  and  $\prod_{T^{\mu}} (P_0^{\mu})$  denotes the projection of  $P_0^{\mu}$  onto the subspace  $T^{\mu}$ . This

projection is derived from the solution of the normal equation  $A^TAx = A^TP_0^\mu$  which is equivalent to the associated least squares solution of  $Ax = P_0^\mu$ . Due to the difficulty associated with inverting a general matrix that may be singular or non-square depending on the number of neighboring points selected in the manifold method, the problem (1) can be posed as a minimization problem in which the Moore-Penrose pseudoinverse  $A^+$  (Golub and Van Loan, 2013) of the original matrix A is used. The pseudoinverse  $A^+$  generalizes the concept of matrix inverse and arises in the minimum norm (that is, approximate as opposed to exact) or best-fit (in a least squares sense) solution to a system of linear equations. The problem:  $\min_x \|Ax - P_0^\mu\|_2$  has the solution  $x = A^+ P_0^\mu$ . The pseudoinverse can be computed using SVD as follows: if  $A = U\Sigma V^T$ , where U, V denote unitary matrices and  $\Sigma$  is a diagonal matrix containing the singular values of A, then  $A^+ = V\Sigma^+U^T$ . We used the function pinv to compute the pseudoinverse in MATLAB.

- 280 3. Find a linear representation coefficient vector  $\alpha$  of that projection onto the tangent space:
- This coefficient is calculated by solving the following equation:

$$283 \qquad \prod_{\mathsf{T}^{\mu}} (P_0^{\mu}) = A \cdot \alpha \tag{2}$$

- 284 4. Estimate the missing components of the point  $P_0$  ( $P_0^{\nu}$ ):
- The last step is finding a point on the subspace  $T^{\nu}$  that is closest (in norm) to the point  $P_0$ . In
- order to do that,  $T^{\nu}$  is projected using the projection coefficient  $\alpha$  computed in step 3:

$$P_0^{\nu} \cong T^{\nu}. \alpha \tag{3}$$

The result of this projection is the closest point to  $P_0^{\nu}$  that belongs to its subspace. In this algorithm, high-dimensional coordinates of selected neighborhoods on the manifold are

projected to a low-dimensional subspace. An alternative to this approach is to use kernel regression to assign a weight to each neighbor based on the distance from  $P_0^{\mu}$  (Fig. 3b). A weight for each selected neighborhood can be computed using the following Gaussian kernel function:

294 
$$W_i = e^{-\frac{\left(P_i^{\mu} - P_o^{\mu}\right)^2}{2\sigma^2}}, \ \sigma = \sqrt{\text{var}\left(d\left(P_i^{\mu}, P_o^{\mu}\right)\right)}$$
 (4)

Examples of manifolds representing geophysical (bathymetry) and meteorological (wind) data are shown in Figures 4(a) and (b). These figures support the assumption that the manifold can be considered as being linear locally, but with complicated topology overall.

#### 2.2 Test case: Analytical function

Before applying the manifold method to reconstruct complex geophysical and meteorological data, we first evaluate the effectiveness of the method in reproducing an analytical function, since errors can be computed relative to the known function values; therefore, the F7 function suggested by Lazzaro and Montefusco (2002) and Renka and Brown (1999) is used:

303  
304 
$$F7(x, y) = 2\cos(10x) \cdot \sin(10y) + \sin(10x \cdot y)$$
 (5)  
305

where the domain of F7 is restricted to  $0 \le x \le 1$  and  $0 \le y \le 1$  (Fig. 5a). Three sets of sparse random points from a normal distribution were generated in the domain with numbers of sampling points of 30, 60, and 90. The F7 function was sampled randomly as shown in Fig. 5b and the manifold method was tested by withholding one point at a time and estimating its associated value from the remaining points using the manifold method, in addition to other

methods such as the natural neighbor, nearest neighbor, and IDW interpolations. Since known components of the scatter points are located in the two-dimensional X-Y plane, at least two neighboring points are needed to form a tangent space for the manifold method. Therefore, for simplicity, only two nearest neighbors are used in both manifold and IDW interpolation methods.

#### 2.3 Assimilating Meteorological Data for Improved Lake Circulation Modeling:

#### Lake Michigan

The proposed method was first applied for the reconstruction of wind fields (time-varying vector fields) over Lake Michigan. Hourly wind speed and direction data during April-September 2008 were obtained from the National Data Buoy Center (NDBC) weather stations surrounding the lake (Fig. 6). The wind measurements were adjusted to a 10 m anemometer height using the profile methods described in Schwab (1987). Since the aerodynamic roughness over the lake is much lower compared to its counterpart over the land, an empirical overland-overlake adjustment was applied to the wind speeds recorded by overland stations (Schwab and Beletsky, 1998). The datasets of wind speed and direction were converted to two coordinates in the Cartesian coordinate system (x and y directions).

Instead of using the cross-validation method to evaluate the interpolated wind data, results from the hydrodynamic model of the lake were compared with current measurements to test the applied method. To this end, a well-tested three-dimensional hydrodynamic model of the lake (Safaie et al., 2016) was used. The model was based on the unstructured grid Finite Volume Community Ocean Model (FVCOM; Chen et al., 2006) which was successfully used in the past in ocean (Li et al., 2014), lake (Nguyen et al., 2014) and river (Anderson and Phanikumar, 2011)

modeling. Details of the unstructured mesh used in the hydrodynamic model are presented in Table 1.

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

Wind fields from April to September 2008 were reconstructed at the locations of nodes in the numerical mesh. Other hourly meteorological observations related to heat flux fields, including air temperature, cloud cover, dew point, long-wave solar radiation, short-wave solar radiation, and relative humidity, obtained from the National Climatic Data Center (NCDC) and NDBC stations, were interpolated over the computational grid using a smoothed natural neighbor method with a smoothing radius of 30 km. Air pressure was assumed to be constant (10<sup>5</sup> Pa) through the course of the study and a constant startup water temperature with a value of 2.5 °C was used in the model. The overlake dew points were estimated from overland observations using an empirical formula described in (Schwab and Beletsky, 1998). Air temperature and cloud cover were used to estimate long-wave solar radiation (Parkinson and Washington, 1979) and short-wave solar radiation was modeled using the clear-sky value and cloud cover (Nguyen et al., 2014). Six arc-second bathymetric data obtained from the NOAA National Geophysical Data Center (NGDC) combined with two-meter resolution LIDAR data along the Indiana coast from the National Oceanic and Atmospheric Administration (NOAA) were interpolated to the numerical mesh using the natural neighbor method (Safaie et al., 2015).

Three bottom-mounted, upward-looking Acoustic Doppler Current Profilers (ADCPs) were deployed at stations M, BB and S (Fig. 6) in southern Lake Michigan from early June to late August 2008 to measure nearshore currents for model testing (Thupaki et al., 2013; Safaie et al., 2016). The hydrodynamic model was run from April to August 2008 to have a two-month spin-up period. Evaluation of the manifold method was carried out by comparing the simulated currents with data collected by the ADCPs. Comparisons between simulated and observed

currents can be improved by identifying an optimal set of parameters in the manifold method. These parameters include: an optimum number of the nearest neighbors to create a tangent space, the time scale c, and parameters of the Gaussian kernel function. In addition, the method used for creating a tangent space from a set of data points (Appendix-A) can be changed to improve the agreement between simulated and observed currents. The manifold method for the reconstruction of wind fields was directly applied to reconstruct the other six scalar observations to calculate the heat flux fields. This time, however,  $P^{\nu}$  is a scalar, rather than a vector.

# 2.4 Assimilating Geophysical Data for Improved Lake Circulation Modeling: Gull

Lake

In the second example, the bathymetry of Gull Lake was reconstructed using a manifold method. The lake bathymetry data were collected using a SonTek RiverSurveyor M9 system. The M9 system has an Acoustic Doppler Profiler (ADP) with two sets of four profiling beams and one vertical acoustic beam (0.5-MHz echo-sounder) for river discharge measurements and bathymetric surveys. The system was equipped with differential GPS with sub-meter precision and mounted on a SonTeck hydroboard to avoid high pitch and roll angles. The vertical acoustic beam has a range of 0.2 m to 80 m with an accuracy of 1% and a resolution of 0.001 m. The bathymetry survey was performed in four days (June 9 – June 12, 2015) by collecting data along longitudinal and transverse transects of the lake with an approximate interval of 200 m between each transect pair and sampling interval of 0.2 m- 2 m along the transects depending on the boat speed (Fig. 7).

In order to assimilate the bathymetry of the lake, a three-dimensional hydrodynamic model based on FVCOM has been developed for the lake during the period of thermal stratification (June-August of 2014). The hydrodynamic equations were solved by the numerical model on an unstructured grid and details are given in Table 1.

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

The meteorological observations for calculation of wind and heat flux fields were obtained from NCDC, Weather Underground (https://www.wunderground.com), and the Kellogg Biological Station Long-Term Ecological Research (KBS LTER, http://lter.kbs.msu.edu) stations, a total of 22 locations surrounding Gull Lake from May to August (Fig. 8). Instead of a constant air pressure, hourly air pressure data recorded by the KBS LTER stations were used to improve the performance of the model. This also helped in the calculation of water density in FVCOM based on a polynomial expression (Jackett and Mcdougall, 1995) that takes pressure into account. After applying overland-overlake adjustments, all observations were interpolated over the numerical mesh using a smoothed natural neighbor method with a smoothing radius of 15 km. This radius provided the best simulated results between the ranges of 0 to 30 km. Air temperatures were adjusted using the empirical formula of  $T_a = 0.4T_{la} + 0.6T_w$  (Schwab and Beletsky, 1998), where  $T_a$  is the adjusted air temperature over water,  $T_{la}$  is the air temperature reported by overland stations, and  $T_w$  is the surface water temperature. The surface water temperature was collected using an Onset HOBO Pro v2 sensor with an accuracy of 0.2°C. A linearly varying startup water temperature was used with a value of 12°C at the water surface and 4°C at the depth of 10 m. The hydrodynamic model was tested using observed current data measured using a Teledyne - RDI Sentinel-V ADCP (1000 kHz frequency with a bin size of 0.3 m) deployed in the nearshore waters of the lake in approximately 10 m of water (Fig. 7b). Finally, the bathymetry of the lake interpolated to grid nodes using the manifold method was assimilated into the model.

# 2.5 Assimilating Precipitation Data for an Integrated, Distributed Hydrologic

model

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

In the third example, we simulate the hydrology of a watershed in the Great Lakes region using the spatial distribution of six-year (2000-2005) daily participation data over the Grand River using an integrated, process-based hydrologic model, (GR) watershed in Michigan PAWS+CLM (Shen and Phanikumar, 2010; Shen et al., 2013; Niu et al., 2014). The model is able to simulate different hydrologic components and states including channel discharge, surface runoff, evapotranspiration, groundwater, soil moisture, soil temperature, and vegetation growth. PAWS+CLM uses a structured finite-volume grid to solve the governing partial differential equations for different hydrologic components. Governing equations and numerical details of PAWS have been described in (Shen and Phanikumar, 2010) and in Table 2 of (Niu et al., 2014). The Grand River (GR) watershed (Fig. 9) was selected as our study domain. GR watershed is located in the middle of Michigan's Lower Peninsula and it is the second largest watershed in Michigan. The watershed has a drainage area of 14, 431 km<sup>2</sup> and drains portions of 15 counties in Michigan. The GR stretches 420 kilometers to the outlet at Grand Haven on Lake Michigan and it is the longest river in Michigan. Shen et al. (2014) described the details of data input and integration algorithms of PAWS+CLM; thus we simply introduce the basic data input and processing information for our model in this section. We used a grid resolution of 1 km x 1 km for horizontal discretization which produced a 170×195 mesh for the GR watershed and 20 vertical layers to simulate the vadose zone dynamics by solving the Richards equation and 2 layers for the groundwater domain (unconfined and confined aquifers) (Table 1).

For topographic calculations (e.g. surface slope and overland flow), the 30 m resolution National Elevation Dataset (NED, http://ned.usgs.gov) from U.S. Geological Survey (USGS) was used as the Digital Elevation Model (DEM). For river network simulation, National Hydrography Dataset (NHD) from USGS was assimilated and reorganized as 'river segments' with a length of one kilometer. We used the 30 m resolution raster data provided by the Michigan Department of Natural Resources, i.e. the Integrated Forest Monitoring Assessment and Prescription (IFMAP) dataset as the land use and land cover (LULC) layout (MDNR, 2010). Soil information was obtained from Soil Survey Geographic (SSURGO) (Soil Survey Staff) database from U.S. Department of Agriculture. This information was processed by the pedotransfer functions provided in ROSETTA (Schaap et al., 2001) to provide soil properties of water retention and unsaturated conductivities. Climate driven data (e.g. precipitation, daily maximum temperature and minimum temperature, wind speed) are acquired as point input (Fig. 9) from National Climatic Data Center (NCDC, 2010) of the National Oceanic and Atmospheric Administration (NOAA). In this study, 14 rain gauges in the GR were selected to obtain the spatial distribution of rainfall over the watershed and for assimilation into the model. Previous applications of the PAWS+CLM model used the nearest neighbor method as the default for processing precipitation data. In our work, the manifold method was tested by evaluating the stream discharge outputs of PAWS+CLM against USGS data. The parameters of PAWS+CLM are listed in Table 1 of (Shen

et al., 2013), most of which are spatially distributed and have taken into account the spatial

heterogeneity of parameters such as the soil parameters and groundwater hydraulic

conductivities.

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

#### 3 Results and Discussion

#### 3.1 Analytical function

True values of the analytical function at each of the randomly selected sampling locations were compared with the estimated values obtained by the manifold method as well as other standard interpolation methods. The performance statistics for this example are provided in Table 2. For all methods, the approximation of the F7 function improved by increasing the number of sampling points. In this particular example, the results show that the manifold method produced better overall performance compared to the other three methods considered. However, the best method in this example might perform differently on another test function or for a different sampling point distribution; therefore, we examine the performance of the method for other datasets in the following sections.

#### 3.2 Lake Michigan

Due to the sparse distribution of weather stations around Lake Michigan, it was not clear a priori how many neighboring stations would provide an adequate representation of the data. Since choosing a relatively few (e.g., three) neighboring stations in this situation would involve using information from stations that are far apart as neighbors, we used kernel regression to assign weights to each station depending on the distance from the point of interest. For each node of the numerical grid of Lake Michigan, k number of nearest neighbors were selected and their assigned weights were projected to a low-dimensional subspace. The free parameters in the method are c (time scale),  $\sigma$  (the parameter used in kernel regression), and k. The standard deviation of weather station distances from the point of interest was used for the parameter  $\sigma$  in kernel regression. Performance of the manifold method as measured by a comparison of

simulated and observed currents in Lake Michigan is summarized in Table 3 relative to the other standard methods considered. We note that the manifold method based on kernel weighting considering all stations produced the best overall performance as measured by the root mean squared error (RMSE) between the observed and simulated currents. The performance of the method without kernel regression and with only three neighboring stations was comparable to the other methods but slightly inferior to the natural and nearest neighbor methods. Fig. 10 shows the RMSE and R<sup>2</sup> values for different numbers of nearest neighbors at different ADCP locations. Having all stations to create the tangent space for the manifold method resulted in a better representation of wind fields, and improved the results of the hydrodynamic model (Fig. 11 shows the comparison at station M). We can see that the manifold method performs better than the IDW method at two of the offshore stations (M and BB) but not at the nearshore location S. We believe that the reason for this has to do with the fact that in the nearshore region there are a number of additional processes (waves, wave-current interactions etc) which are not simulated in our model. Therefore model performance in that region cannot be directly attributed to the wind field. At the other two offshore stations M and BB, where the flow is predominantly wind-driven, an improvement in the simulated hydrodynamic fields can be seen.

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

Finally, cross-validation was used to compare the performance of the manifold method with other standard methods for the same Lake Michigan datasets. The performance metrics are summarized in Table 4. In this cross-validation method, one weather station was withheld to be used later for validating the manifold method, and all other stations surrounding the lake were used for the manifold training set. This process was repeated so that each weather station was given a chance to be part of this validation process. Based on these results the proposed manifold method with three nearest neighbors gave better results compared to other standard methods.

However, the performance of the hydrodynamic model based on these methods was relatively inferior compared to the performance of the model when the manifold method used all neighboring points. All different versions of the manifold methods had reasonable computational efficiency. The computational time for the O-kriging was high due to the time needed for finding the best variogram at each time step.

#### 3.3 Gull Lake

Wind and heat flux fields of the Gull Lake were interpolated over the numerical mesh of the lake using the natural neighbor method. Then, the bathymetry of the lake was interpolated over the mesh using the same natural neighbor method to develop the initial version of the lake hydrodynamic model. The raw bathymetry data, which has some regions of steep bathymetry change, created artificial currents in the model due to an error in the pressure gradient force introduced by the sigma-coordinate system of FVCOM (Mellor et al., 1998). Therefore, the interpolated bathymetry was smoothed with a radius of 100 m in order to reduce the errors. The results of the developed model using natural neighbor method and IDW with three nearest neighbors are presented in Fig. 12.

The model was used to assimilate the bathymetry of the lake based on the manifold method. The bathymetry data were reconstructed from the tangent space of the manifold with three nearest neighbors and smoothed with the same method described above. The hydrodynamic model was run with the reconstructed bathymetry. The final comparisons of the vertically-averaged velocity profiles at the ADCP location are presented in Fig. 12. The best value of  $\sigma$  used in kernel regression was equal to the standard deviation of distances of observational points where water depth values are available within a search radius of 50 m from the point of interest.

When the number of samples within this radius was smaller than 100,  $\sigma$  value was calculated based on locations of 100 nearest samples. This method is more accurate when enough samples are available around an estimated point, unless selecting 100 samples itself does a reasonable job. RMSE values (m/s) of eastward and northward velocities in Gull Lake for comparison of the manifold method with other standard methods used in limnology and oceanography are presented in Table 5. The statistics of cross-validation for all (=71) measured longitudinal and transverse transects are shown in Table 6. The cross-validation was performed by omitting one transect at each step and calculating the bathymetry for that transect from the rest of the observation data and repeating the process for all other transects.

#### 3.4 Grand River watershed

The daily distributions of precipitation over the GR watershed for the period 2000 to 2005 were reconstructed from selected rain gauges using the manifold method. The precipitation fields over a period of six years were tested by comparing the simulated and observed stream discharges at selected USGS gaging stations within the watershed. In this example, kernel regression with a standard deviation of all rain gauge distances from the estimated points was used in the manifold method. The tangent space at each grid point of the numerical model was obtained from three-nearest neighbors around that point. The PAWS+CLM model was run with precipitation distributions built using the manifold, natural neighbor, nearest neighbor, and inverse distance methods. The final comparisons of simulated and observed stream flows for USGS gauges are presented in Figs. 13 and 14. Model performance metrics (NASH, Absolute Percent Bias (APB), and RMSE) were computed to evaluate the performance of the manifold method with other standard methods (Table 7). The manifold method provided superior results as can be seen from the improved representation of baseflow over the simulation period. This can be seen clearly from semi-log

plots of the stream hydrograph comparison. The cross-validation results for precipitation are presented in Table 8. In this method, data from one rain gauge was withheld and data from the rest of the gauges were used for manifold training. This process was repeated for all 14 rain gauges. Since rain gauges have a sparse distribution, the cross-validation method could not identify the best method for this case. The results of the integrated, distributed hydrologic model, however, clearly demonstrate the efficacy of the manifold method in the reconstruction of precipitation fields. In particular, the better simulated baseflow shows the strength of using precipitation data based on the manifold method to improve the simulation of heterogeneous partitioning of surface runoff and infiltration processes.

#### 4 Conclusions

We presented a novel manifold method of reconstructing spatio-temporal data for assimilating geophysical and meteorological data in integrated land surface subsurface, and oceanic/lake models. All three case studies illustrate the superior performance of the presented manifold algorithm over standard methods in terms of accuracy and computational efficiency. The hydrodynamic model of Lake Michigan based on the manifold method of reconstructing wind fields produced better performance relative to the other methods. The best results were obtained using kernel regression applied to all weather stations (neighbors). However, the cross-validation results show that the results of the three nearest neighbors were better than the other methods.

The Gull Lake model results indicated that the proposed method has the ability to reconstruct geophysical data at unsampled locations. The use of spatiotemporal precipitation fields constructed using the manifold method produced better stream discharge simulations compared to similar results from the nearest neighbor, natural neighbor and IDW interpolation methods in a

large watershed (> 10000 km²). Finally, all three examples show that evaluating the performance of interpolation methods using the cross-validation method without considering the data structure and the purpose of interpolation can lead to misleading conclusions about the relative performance of the methods considered. The comparisons presented here indicate that manifold methods show promise for modeling Earth system processes based on data from sensor networks. Our future work will combine manifold methods with approaches such as the EnKF method to further improve process-based modeling of land surface, subsurface and lake/ocean models.

Based on the results presented, we note that: (1) Details of the manifold method such as the tangent space estimation, the distance metric that defines spatiotemporal proximity and other details can be further improved to improve the performance of the manifold method; however, these topics are beyond the scope of the present paper. (2) We do not claim that the manifold method provides superior performance on all datasets and for all performance metrics but from the examples considered here it appears that the manifold method may offer an attractive method that is comparable or superior to other standard methods. More research is needed to understand the relative strengths and weaknesses of different manifold-based approaches compared to standard methods.

# Acknowledgments

This work was supported by a grant from the National Science Foundation, CyberSEES program (Award # 1331852). We thank Elena Litchman, Pam Woodruff, Jim Allen, Mike Gallagher, Andrew Fogiel, and Tuan D. Nguyen for their assistance with field data collection in Gull Lake.

- We gratefully acknowledge the use of data from the Kellogg Biological Station LTER which is
- 576 supported by NSF (DEB 1027253) and by Michigan State University AgBioResearch.

# **Appendix-A. Tangent Space Estimation**

- To understand the local geometry of the surface f(x) near a point  $x \in \mathbb{D}^n$ , we consider the first-
- order Taylor series expansion of the surface:

580 
$$f(\overline{x}) = f(x) + \frac{\partial f(x)}{\partial x} (\overline{x} - x) + O(\|\overline{x} - x\|^2) = f(x) + J_f(\overline{x} - x) + O(\|\overline{x} - x\|^2)$$
(A1)

- where  $J_f(x) \in \Box^{D \times n}$  is the Jacobian matrix of f at the point x. If the components of f(x) can
- be written as:  $f(x) = [f_1(x), f_2(x), f_3(x) \cdots f_D(x)]^T$  and  $x = [x_1, x_2, x_3 \cdots x_n]^T$ , then the Jacobian
- can be written as:

577

$$584 J_f(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \vdots & \vdots \\ \frac{\partial f_D}{\partial x_1} & \dots & \frac{\partial f_D}{\partial x_n} \end{bmatrix}$$
(A2)

- To understand the local shape of the surface in equation (A1), we seek to determine the space
- 586  $(\bar{x}-x)$ , such that as we move away from x, the value of the function doesn't change to within
- first order. This is equivalent to finding the space T such that:

588 
$$T = \left\{ (\overline{x} - x) \middle| J_f(x) (\overline{x} - x) = 0 \right\}$$
 (A3)

This space is the tangent space to the surface at point x and is the right null space of the Jacobian matrix  $J_f(x)$ . The space orthogonal to the tangent space is the row space of the Jacobian and orthogonal representations of these spaces can be obtained from SVD. The right null space of  $J_f$  is the columns of V corresponding to zero singular values. Therefore, the tangent space of the manifold M at y = f(x) is:

$$594 T(M) = \operatorname{span}(J_f(x)) (A4)$$

From a practical computation point of view, given a set of sample points  $y = \{y_1, y_2, y_3, \dots y_m\}$ , a simple method of constructing the tangent space is to approximate it as the line/surface obtained by joining the local neighboring points. The tangent space can also be directly estimated using SVD. If  $C^m$  denotes the local covariance matrix:

599 
$$C^m = \frac{1}{m} \sum_{i=1}^m y_i y_i^T = U \Sigma U^T$$
 (A5)

where  $U = [u_1, u_2, u_3, \dots u_D]$  and  $\Sigma = \text{diag}[\lambda_1, \lambda_2, \lambda_3 \dots \lambda_D]$  denote the eigenvector and eigenvalue matrices respectively, then the optimal (in a least-squares sense) *n*-dimensional linear subspace is the span of the *n*-largest eigenvectors in U:

603 
$$T(M) \cong \text{span}\{u_1, u_2, u_3, \dots u_n\}$$
 (A6)

#### 5. References

604

- Allard, W.K., Chen, G., Maggioni, M., 2012. Multi-scale geometric methods for data sets II:
- 606 Geometric Multi-Resolution Analysis. Appl. Comput. Harmon. Anal. 32, 435–462.
- 607 doi:10.1016/j.acha.2011.08.001

- Anderson, E.J., Phanikumar, M.S., 2011. Surface storage dynamics in large rivers: Comparing
- three-dimensional particle transport, one-dimensional fractional derivative, and multirate
- transient storage models. Water Resources Research 47, 1–15. doi:10.1029/2010WR010228
- Anguetin, S., Braud, I., Vannier, O., Viallet, P., Boudevillain, B., Creutin, J.-D., Manus, C.,
- 612 2010. Sensitivity of the hydrological response to the variability of rainfall fields and soils for the
- 613 Gard 2002 flash-flood event. J. Hydrol. 394, 134–147. doi:10.1016/j.jhydrol.2010.07.002
- Barth, A., Azcárate, A. A., Joassin, P., Jean-Marie, B., Troupin, C., 2008. Introduction to
- Optimal Interpolation and Variational Analysis. Presented at the SESAME Summer School,
- 616 SESAME Summer School, Varna, Bulgaria.
- Belkin, M., Niyogi, P., 2003. Laplacian eigenmaps for dimensionality. Speech Commun. 1, 349–
- 618 367.
- Bell, V.A., Moore, R.J., 2000. The sensitivity of catchment runoff models to rainfall data at
- different spatial scales. Hydrol. Earth Syst. Sci. 4, 653–667. doi:10.5194/hess-4-653-2000
- Beven, K.J., Hornberger, G.M., 1982. Assessing the Effect of Spatial Pattern of Precipitation in
- Modeling Stream Flow Hydrographs. J. Am. Water Resour. Assoc. 18, 823–829.
- 623 doi:10.1111/j.1752-1688.1982.tb00078.x
- Blumberg, A.F., Mellor, G.L., 1987. A description of a three-dimensional coastal ocean
- 625 circulation model. Am. Geophys. Union 1–16.
- 626 Candès, E.J., Li, X., Ma, Y., Wright, J., 2011. Robust principal component analysis? J. ACM 58,
- 627 1–37. doi:10.1145/1970392.1970395
- 628 Carin, L., Baraniuk, R.G., Cevher, V., Dunson, D., Jordan, M.I., Sapiro, G., Wakin, M.B., 2011.
- 629 Learning Low-Dimensional Signal Models. IEEE Signal Process. Mag. 28.
- 630 doi:10.1109/MSP.2010.939733
- 631 Chen, C., Beardsley, R., Cowles, G., 2006. An Unstructured Grid, Finite-Volume Coastal Ocean
- 632 Model (FVCOM) System. Oceanography 19, 78–89. doi:10.5670/oceanog.2006.92
- Daly, C., Gibson, W.P., Taylor, G.H., Johnson, G.L., Pasteris, P., 2002. A knowledge-based
- approach to the statistical mapping of climate. Clim. Res. 22, 99–113.
- Dang, C., Aghagolzadeh, M., Radha, H., 2014. Image Super-Resolution via Local Self-Learning
- Manifold Approximation. IEEE Signal Process. Lett. 21, 1245–1249.
- 637 doi:10.1109/LSP.2014.2332118
- Dang, C., Radha, H., 2015. Fast Image Super Resolution via Selective Manifold Learning of
- High Resolution Patches. Presented at the IEEE Proceedings of International Conference on
- Image Processing (ICIP15), IEEE Proceedings of International Conference on Image Processing
- 641 (ICIP15), Québec City, Canada.

- Dang, C.T., Aghagolzadeh, M., Moghadam, A.A., Radha, H., 2013. Single image super
- resolution via manifold linear approximation using sparse subspace clustering, in: Global
- 644 Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE. IEEE, Austin, Texas,
- 645 U.S.A., pp. 949–952.
- 646 Evensen, G., 2006, Data Assimilation: The Ensemble Kalman Filter, Springer, N. Y.
- Golub, G., Van Loan, C.F., 2013. Matrix Computations, 4th ed. Baltimore: Johns Hopkins, pp.
- 648 756.
- Haidvogel, D.B., Beckmann, A., 1999. Numerical Ocean Circulation Modeling. Imp. Coll. Press.
- Zhang, H., Mendoza-Sanchez, I., Miller, E.L., Abriola, L.M., 2016. Manifold Regression
- Framework for Characterizing Source Zone Architecture. IEEE Transactions on Geoscience and
- 652 Remote Sensing 54, 3–17. doi:10.1109/TGRS.2015.2448086
- Jackett, D.R., Mcdougall, T.J., 1995. Minimal Adjustment of Hydrographic Profiles to Achieve
- 654 Static Stability, J. Atmospheric Ocean, Technol. 12, 381–389. doi:10.1175/1520-
- 655 0426(1995)012<0381:MAOHPT>2.0.CO;2
- 656 Lazzaro, D., Montefusco, L.B., 2002. Radial basis functions for the multivariate interpolation of
- large scattered data sets. J. Comput. Appl. Math., Int. Congress on Computational and Applied
- 658 Mathematics 2000 140, 521–536. doi:10.1016/S0377-0427(01)00485-X
- 659 Li, J., Heap, A.D., 2008. A review of spatial interpolation methods for environmental scientists.
- 660 Geoscience Australia, Canberra.
- Li, L., Losser, T., Yorke, C., Piltner, R., 2014. Fast Inverse Distance Weighting-Based
- Spatiotemporal Interpolation: A Web-Based Application of Interpolating Daily Fine Particulate
- Matter PM2:5 in the Contiguous U.S. Using Parallel Programming and k-d Tree. Int. J. Environ.
- Res. Public. Health 11, 9101–9141. doi:10.3390/ijerph110909101
- 665 Li, L., Revesz, P., 2004. Interpolation methods for spatio-temporal geographic data. Comput.
- 666 Environ. Urban Syst. 28, 201–227. doi:10.1016/S0198-9715(03)00018-8
- 667 Li, R. et al., 2014. Observed wintertime tidal and subtidal currents over the continental shelf in
- the northern South China Sea, J. Geophys. Res. Oceans, 119(8), pp. 5289-5310, doi:10.1002/
- 669 2014JC009931.
- Luo, W., Taylor, M.C., Parker, S.R., 2008. A comparison of spatial interpolation methods to
- estimate continuous wind speed surfaces using irregularly distributed data from England and
- Wales. Int. J. Climatol. 28, 947–959. doi:10.1002/joc.1583
- Ly, S., Charles, C., Degré, A., 2011. Geostatistical interpolation of daily rainfall at catchment
- scale: the use of several variogram models in the Ourthe and Ambleve catchments, Belgium.
- 675 Hydrol. Earth Syst. Sci. 15, 2259–2274. doi:10.5194/hess-15-2259-2011

- MacEachren, A.M., Davidson, J.V., 1987. Sampling and isometric mapping of continuous
- 677 geographic surfaces. Am. Cartogr. 14, 299–320.
- Masih, I., Maskey, S., Uhlenbrook, S., Smakhtin, V., 2011. Assessing the Impact of Areal
- Precipitation Input on Streamflow Simulations Using the SWAT Model1. JAWRA J. Am. Water
- 680 Resour. Assoc. 47, 179–195. doi:10.1111/j.1752-1688.2010.00502.x
- Ma, Y., Niyogi, P., Sapiro, G., Vidal, R., 2011. Dimensionality reduction via subspace and
- 682 submanifold learning. IEEE Signal Process. Mag. 28, 14–126. doi:10.1109/MSP.2010.940005
- Mellor, G.L., Oey, L.Y., Ezer, T., 1998. Sigma coordinate pressure gradient errors and the
- seamount problem. Journal of Atmospheric and Oceanic Technology 15, 1122–1131.
- Merwade, V., 2009. Effect of spatial trends on interpolation of river bathymetry. J. Hydrol. 371,
- 686 169–181. doi:10.1016/j.jhydrol.2009.03.026
- Mordohai, P., Medioni, G., 2010. Dimensionality estimation, manifold learning and function
- approximation using tensor voting. J. Mach. Learn. Res. 11, 411–450.
- Moradkhani, H., Sorooshian, S., Gupta, H.V., Houser, P.R., 2005. Dual state-parameter
- 690 estimation of hydrological models using ensemble Kalman filter, Adv. Water Resour. 28, 135–
- 691 147.4
- Nguyen, T.D., Thupaki, P., Anderson, E.J., Phanikumar, M.S., 2014. Summer circulation and
- 693 exchange in the Saginaw Bay-Lake Huron system. J. Geophys. Res. Oceans 119, 2713–2734.
- 694 doi:10.1002/2014JC009828
- Nicótina, L., Alessi Celegon, E., Rinaldo, A., Marani, M., 2008. On the impact of rainfall
- patterns on the hydrologic response. Water Resour. Res. 44, 1–14. doi:10.1029/2007WR006654
- Niu, J., Phanikumar, M.S., 2015. Modeling watershed-scale solute transport using an integrated,
- 698 process-based hydrologic model with applications to bacterial fate and transport. J. Hydrol. 529,
- 699 35–48. doi:10.1016/j.jhydrol.2015.07.013
- Niu, J., Shen, C., Li, S.-G., Phanikumar, M.S., 2014. Quantifying storage changes in regional
- 701 Great Lakes watersheds using a coupled subsurface-land surface process model and GRACE,
- 702 MODIS products. Water Resour. Res. 50, 7359–7377. doi:10.1002/2014WR015589
- Obled, C., Wendling, J., Beven, K., 1994. The sensitivity of hydrological models to spatial
- rainfall patterns: an evaluation using observed data. J. Hydrol. 159, 305–333. doi:10.1016/0022-
- 705 1694(94)90263-1
- Oleson, K.W., Lawrence, D.M., Gordon, B., Flanner, M.G., Kluzek, E., Peter, J., Levis, S.,
- Swenson, S.C., Thornton, E., Feddema, J., others, 2010. Technical description of version 4.0 of
- the Community Land Model (CLM) (No. NCAR/TN-478+STR), NCAR Technical Note.
- National Center for Atmospheric Research, Boulder, Colorado.

- Parkinson, C.L., Washington, W.M., 1979. A large-scale numerical model of sea ice. J. Geophys.
- 711 Res. Oceans 84, 311–337. doi:10.1029/JC084iC01p00311
- Pathiraja, S., Marshall, L., Sharma, A., Moradkhani, H., 2016, Hydrologic modeling in dynamic
- 713 catchments: A data assimilation approach, Water Resour. Res. 52, 1-23.
- 714 doi:10.1002/2015WR01719
- Press, W.H., Flannery, B.P., Teukolsky, S.A., Vetterling, W.T., 2007. Numerical Recipes in
- 716 *C++: The Art of Scientific Computing*, third ed. Cambridge University Press, New York.
- Renka, R.J., Brown, R., 1999. Algorithm 792: Accuracy Test of ACM Algorithms for
- 718 Interpolation of Scattered Data in the Plane. ACM Trans Math Softw 25, 78–94.
- 719 doi:10.1145/305658.305745
- Roweis, S.T., 2000. Nonlinear dimensionality reduction by locally linear embedding. Science
- 721 290, 2323–2326. doi:10.1126/science.290.5500.2323
- Safaie, A., Wendzel, A., Ge, Z., Nevers, M.B., Whitman, R.L., Corsi, S.R., Phanikumar, M.S.,
- 723 2016. Comparative Evaluation of Statistical and Mechanistic Models of Escherichia coli at
- Beaches in Southern Lake Michigan. Environ. Sci. Technol., 50, 2442–2449.
- 725 doi:10.1021/acs.est.5b05378
- Schaap, M.G., Leij, F.J., van Genuchten, M.T., 2001. rosetta: a computer program for estimating
- soil hydraulic parameters with hierarchical pedotransfer functions. J. Hydrol. 251, 163–176.
- 728 doi:10.1016/S0022-1694(01)00466-8
- Schuurmans, J.M., Bierkens, M.F.P., 2007. Effect of spatial distribution of daily rainfall on
- 730 interior catchment response of a distributed hydrological model. Hydrol. Earth Syst. Sci.
- 731 Discuss. 11, 677–693.
- Schwab, D.J., 1987. Simulation and forecasting of Lake Erie storm surges. Mon. Weather Rev.
- 733 106, 1476–1487.
- Schwab, D.J., Beletsky, D., 1998. Lake Michigan Mass Balance Study: Hydrodynamic modeling
- project (No. NOAA Technical Memorandum ERL GLERL-108). Great Lakes Environmental
- 736 Research Laboratory, Ann Arbor, MI.
- Shen, C., Niu, J., Fang, K., 2014. Quantifying the effects of data integration algorithms on the
- outcomes of a subsurface—land surface processes model. Environ. Model. Softw. 59, 146–161.
- 739 doi:10.1016/j.envsoft.2014.05.006
- Shen, C., Niu, J., Phanikumar, M.S., 2013. Evaluating controls on coupled hydrologic and
- vegetation dynamics in a humid continental climate watershed using a subsurface-land surface
- 742 processes model. Water Resour. Res. 49, 2552–2572. doi:10.1002/wrcr.20189

- Shen, C., Phanikumar, M.S., 2010. A process-based, distributed hydrologic model based on a
- large-scale method for surface–subsurface coupling. Adv. Water Resour. 33, 1524–1541.
- 745 doi:10.1016/j.advwatres.2010.09.002
- 746 Šiljeg, A., Lozić, S., Šiljeg, S., 2015. A comparison of interpolation methods on the basis of data
- obtained from a bathymetric survey of Lake Vrana, Croatia. Hydrol. Earth Syst. Sci. 19, 3653–
- 748 3666. doi:10.5194/hess-19-3653-2015
- Sun, X., Mein, R.G., Keenan, T.D., Elliott, J.F., 2000. Flood estimation using radar and
- 750 raingauge data. J. Hydrol. 239, 4–18. doi:10.1016/S0022-1694(00)00350-4
- Suparta, W., Rahman, R., 2016. Spatial interpolation of GPS PWV and meteorological variables
- over the west coast of Peninsular Malaysia during 2013 Klang Valley Flash Flood. Atmospheric
- 753 Res. 168, 205–219. doi:10.1016/j.atmosres.2015.09.023
- 754 Syed, K.H., Goodrich, D.C., Myers, D.E., Sorooshian, S., 2003. Spatial characteristics of
- 755 thunderstorm rainfall fields and their relation to runoff. J. Hydrol. 271, 1–21. doi:10.1016/S0022-
- 756 1694(02)00311-6
- 757 Tenenbaum, J.B., De Silva, V., Langford, J.C., 2000. A global geometric framework for
- nonlinear dimensionality reduction. Science 290, 2319–2323.
- 759 Tetzlaff, D., Uhlenbrook, S., 2005. Significance of spatial variability in precipitation for process-
- oriented modelling: results from two nested catchments using radar and ground station data.
- 761 Hydrol. Earth Syst. Sci. 9, 29–41. doi:10.5194/hess-9-29-2005
- 762 Thupaki, P., M.S. Phanikumar and R.L. Whitman, 2013. Solute dispersion in the coastal
- boundary layer of southern Lake Michigan, J. Geophys. Res. Oceans, vol. 118, No. 3, pp. 1606-
- 764 1617, doi: 10.1002 / jgrc.20136 (2013)
- 765 Thupaki, P., Phanikumar, M.S., Schwab, D.J., Nevers, M.B., Whitman, R.L., 2013. Evaluating
- the role of sediment-bacteria interactions on *Escherichia coli* concentrations at beaches in
- southern Lake Michigan. J. Geophys. Res. Oceans 118, 7049–7065.
- 768 doi:10.1002/2013JC008919
- van der Maaten, L.J., Postma, E.O., van den Herik, H.J., 2009. Dimensionality reduction: A
- comparative review. J. Mach. Learn. Res. 10, 66–71.
- Victor, G., Pollack, A., 2010. Differential topology. American Mathematical Soc.
- Wendzel, A., 2014. Constraining mechanistic models of indicator bacteria at recreational
- beaches in Lake Michigan using easily-measurable environmental variables (M. Sc.
- 774 Dissertation). Michigan State University, East Lansing. MI.
- Xue, P., Schwab, D.J., Hu, S., 2015. An investigation of the thermal response to meteorological
- forcing in a hydrodynamic model of Lake Superior. J. Geophys. Res. Oceans 120, 5233–5253.
- 777 doi:10.1002/2015JC010740

- Yan, Y., Xiao, F., Du, Y., 2014. Construction of lake bathymetry from MODIS satellite data and
- 779 GIS from 2003 to 2011. Chin. J. Oceanol. Limnol. 32, 720–731. doi:10.1007/s00343-014-3185-4
- Yang, K., Watanabe, T., Koike, T., Li, X., Fujii, H., Tamagawa, K., Ma, Y., Ishikawa, H., 2007.
- Auto-calibration system developed to assimilate AMSR-E data into a land surface model for
- estimating soil moisture and the surface energy budget. J. Meteorol. Soc. Jpn. 85A, 229–242.

783

784

785

786

787

788

789

790

### **List of Figures**

- 791 Fig. 1. Some examples of manifolds (a) torus, (b) sphere, and (c) a two-dimensional cross-section of a
- 792 six-dimensional Calabi-Yau manifold.
- Fig. 2. Illustration of the proposed manifold approach for estimation of missing data at point  $P_0$ .
- 794 Fig. 3. A tangent space created from the set of nearest points using (a) coordinates of selected
- neighborhoods or (b) Kernel regression.
- 796 Fig. 4. Manifolds representing (a) bathymetry of Gull Lake and (b) wind components over Lake
- 797 Michigan in three dimensional space.
- 798 Fig. 5. (a) Analytical function used to test the manifold method for interpolation of scattered data.
- Random sampling was used to generate scatter points as shown in figures (b. 30 points), (c. 60 points) and
- 800 (d, 90 points) to reconstruct the function.

- Fig. 6. Locations of the ADCPs deployed during summer 2008 and weather stations surrounding Lake Michigan.
- Fig. 7(a). Bathymetry of Gull Lake. (b) Boat tracks generated during the sampling survey in Gull Lake.
- Fig. 8. Selected weather stations surrounding Gull Lake.
- Fig. 9. Map of Grand River watershed showing the locations of the USGS gauges, rain gauges and National Hydrography Dataset (NHD) streams.
- Fig. 10. Performance of the manifold method evaluated using observed and simulated currents at different stations in Lake Michigan. Different number of nearest neighbors were used to reconstruct the wind field using the manifold method with kernel regression.
- Fig. 11. Comparison of simulated (black lines) and observed (red lines) vertically averaged currents at the location M in Lake Michigan. (a) Alongshore velocity (b) Cross-shore velocity
- Fig. 12. Comparison of simulated (black lines) and observed (red lines) vertically averaged currents at the ADCP location in Gull Lake. (a) Eastward velocity and (b) Northward velocity
- Fig. 13. Comparison of simulated and observed stream flows for USGS gauge #04116000 in (a) a linear scale, and (b) logarithmic scale.
- Fig. 14. Comparison of simulated and observed stream flows for USGS gauge #04119000 in (a) a linear scale, and (b) logarithmic scale.

# List of Tables

818

819

820

821

822

823

824

#### **Table 1.** Properties of the numerical grids used for the hydrodynamic and hydrologic

Model	Grid Classification	Element shape	Grid Resolution		# Elements	#Vertical layers
FVCOM (Lake Michigan)	Unstructured	Triangle	4m -5 km	12,684	23,602	20
FVCOM (Gull Lake)	Unstructured	Triangle	8-100m	5,132	9,361	20
PAWS (Grand River)	Structured	Quadrilateral	1 km	33,150	32,786	22

**Table 2.** Cross-validation results for the analytical function based on different sampling points selected randomly.

Sample size	Method	$R^2$	RMSE	Fn	PBIAS	NASH	APB (%)
30	Manifold	0.667	0.778	0.710	14.420	0.416	0.652
	Natural neighbor	0.582	0.876	0.799	-8.866	0.259	0.713
	Nearest neighbor	0.619	0.882	0.804	-14.689	0.250	0.669
	IDW	0.577	0.870	0.793	35.847	0.270	0.727
60							
	Manifold	0.846	0.579	0.531	-33.924	0.703	0.472
	Natural neighbor	0.816	0.615	0.564	16.912	0.664	0.466
	Nearest neighbor	0.779	0.720	0.660	-34.524	0.540	0.512
	IDW	0.832	0.603	0.553	-44.292	0.677	0.469
90	Manifold	0.891	0.502	0.432	-14.103	0.791	0.400
	Natural neighbor	0.874	0.539	0.464	-10.666	0.759	0.344
	Nearest neighbor	0.867	0.571	0.491	-28.777	0.730	0.446
	IDW	0.859	0.567	0.487	-5.974	0.735	0.416

**Table 3.** RMSE values (m/s) of alongshore and cross-shore velocities for comparison of the manifold method with other standard methods used in limnology and oceanography

	Loaction: M		Locati	Location: BB		Location: S	
Method	RMSE u	RMSE v	RMSE u	RMSE v	RMSE u	RMSE v	
O-kriging	0.0385	0.0290	0.0590	0.0349	0.0540	0.0152	
Nearest Neighbor	0.0363	0.0286	0.0580	0.0348	0.0545	0.0152	
Natural Neighbor	0.0366	0.0275	0.0553	0.0334	0.0515	0.0158	
Manifold (3 NBR)	0.0383	0.0276	0.0594	0.0346	0.0568	0.0158	
Manifold+Kernel (3 NBR)	0.0371	0.0268	0.0576	0.0341	0.0559	0.0158	
Manifold+Kernel (all NBR)	0.0304	0.0265	0.0531	0.0312	0.0568	0.0154	
IDW (all NBR)	0.0328	0.0267	0.0535	0.0316	0.0498	0.0155	

**Table 4.** Cross-validation results for wind field over Lake Michigan.

Method	$R^2u$	$R^2v$	RMSEu	RMSEv	Computational time (s)
O-kriging	0.441	0.572	3.497	3.853	92463.8
Nearest Neighbor	0.666	0.743	2.792	3.044	18.6
Natural Neighbor	0.693	0.794	2.558	2.750	183.6
Manifold (3 NBR)	0.690	0.801	2.433	2.595	28.1
Manifold+Kernel (3 NBR)	0.710	0.806	2.392	2.566	55.1
Manifold+Kernel (all NBR)	0.547	0.681	2.884	3.129	77.3
IDW (3 NBR)	0.724	0.822	2.278	2.458	69.3

**Table 5.** RMSE values (m/s) of eastward and northward velocities in Gull Lake for comparison of the manifold method with other standard methods used in limnology and oceanography

Method	RMSE u	RMSE v
Natural Neighbor	0.0090	0.0205
Manifold+Kernel (3 NBR)	0.0098	0.0200
IDW (3 NBR)	0.0095	0.0204

 Table 6. Cross-validation results for Gull Lake bathymetry.

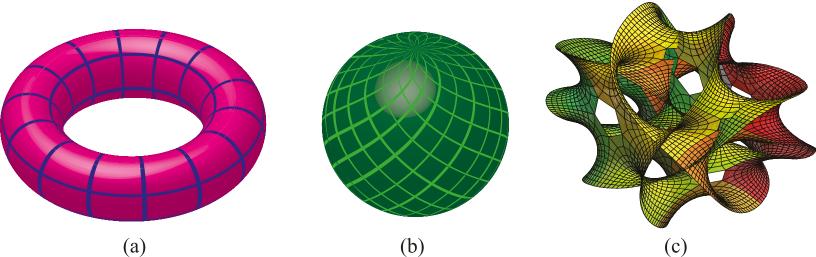
Method	$R^2$	RMSE (m)	Fn	NASH	PBIAS
Manifold	0.890	2.011	0.222	0.678	-14.016
Natural Neighbor	0.925	1.288	0.170	0.468	-13.301
Nearest Neighbor	0.888	2.039	0.230	0.670	-17.132
IDW (3 NBR)	0.839	3.282	0.6065	0.540	-15.918

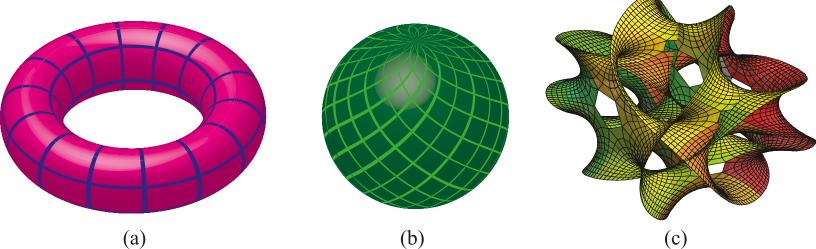
**Table 7.** Comparison of the manifold method with other standard methods for precipitation over the Grand River watershed.

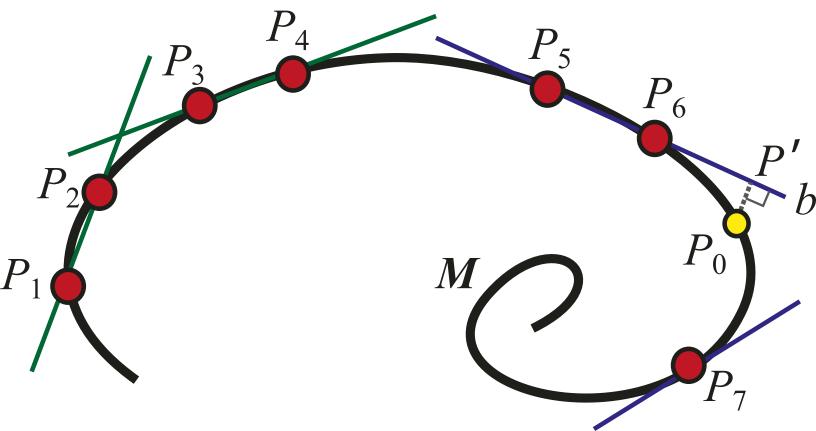
	USGS04116000			USGS04119000			
Method	NASH	RMSE	APB(%)		NASH	RMSE	APB(%)
Nearest Neighbor	0.56	37.13	36.47		0.58	58.15	30.54
Natural Neighbor	0.33	44.52	54.71		0.36	71.88	48.72
Inverse Distance	0.38	44.22	54.18		0.37	70.29	48.01
Manifold	0.59	35.56	33.76		0.63	54.57	27.87

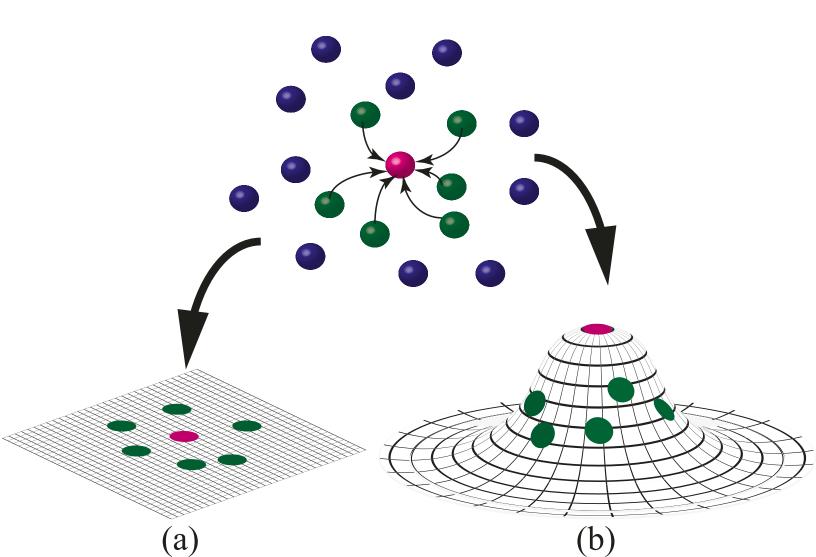
**Table 8.** Cross-validation results for precipitation over the Grand River watershed.

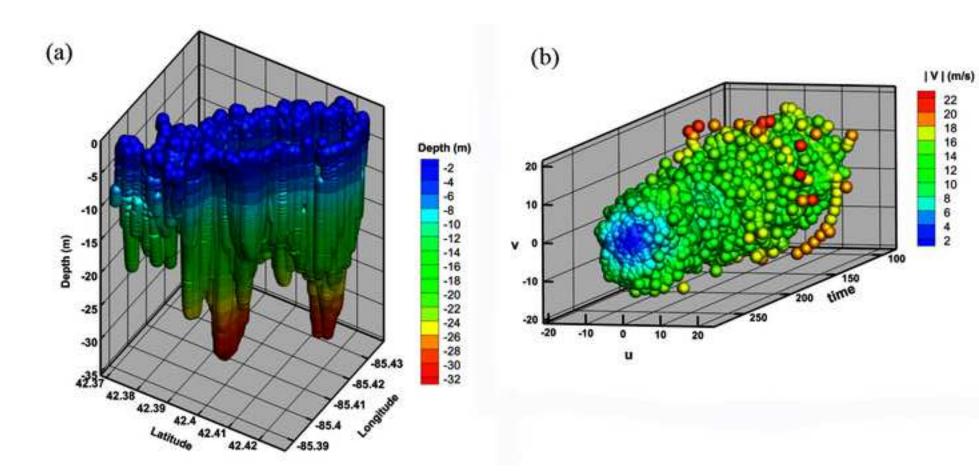
Method	$R^2$	RMSE (cms)	Fn	NASH	APB (%)
Manifold	0.543	7.933	0.910	0.624	72.4
Natural Neighbor	0.543	7.920	0.901	0.574	73.6
Nearest Neighbor	0.471	8.802	0.930	0.604	81.8
IDW (3 NBR)	0.567	7.542	0.907	0.619	70.0

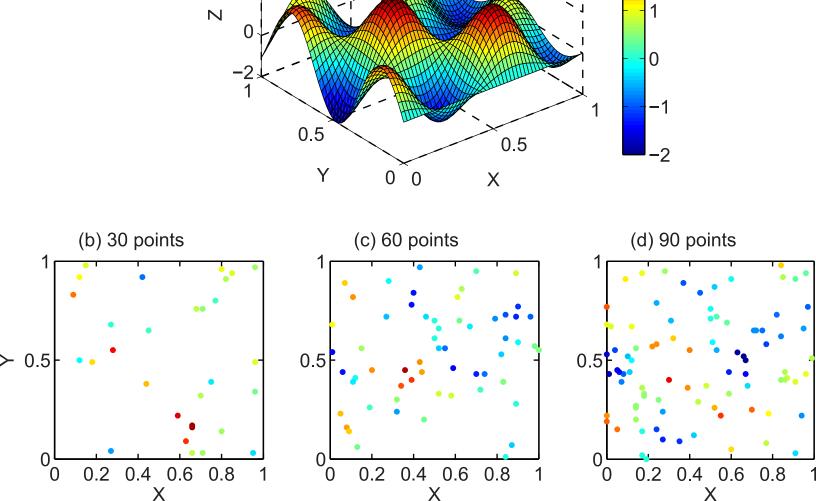












3

2

(a)

2

