**RESEARCH ARTICLE**

**Key Points:**
- Transient data by itself can reveal the physical sources causing the transients
- Inverse methods can identify the sources of the transients
- Blind source separation based on nonnegative matrix factorization is applied

**Correspondence to:**
V. V Vesselinov,
vvv@lanl.gov

# Blind source separation for groundwater pressure analysis based on nonnegative matrix factorization

Boian S. Alexandrov[1] and Velimir V. Vesselinov[2]

[1]Theoretical Division, Physics and Chemistry of Materials Group, Los Alamos National Laboratory, Los Alamos, New Mexico, USA, [2]Earth and Environmental Sciences Division, Computational Earth Science Group, Los Alamos National Laboratory, Los Alamos, New Mexico, USA

**Abstract** The identification of the physical sources causing spatial and temporal fluctuations of aquifer water levels is a challenging, yet a very important hydrogeological task. The fluctuations can be caused by variations in natural and anthropogenic sources such as pumping, recharge, barometric pressures, etc. The source identification can be crucial for conceptualization of the hydrogeological conditions and characterization of aquifer properties. We propose a new computational framework for model-free inverse analysis of pressure transients based on Nonnegative Matrix Factorization (NMF) method for Blind Source Separation (BSS) coupled with $k$-means clustering algorithm, which we call NMF$k$. NMF$k$ is capable of identifying a set of unique sources from a set of experimentally measured mixed signals, without any information about the sources, their transients, and the physical mechanisms and properties controlling the signal propagation through the subsurface flow medium. Our analysis only requires information about pressure transients at a number of observation points, $m$, where $m \geq r$, and $r$ is the number of unknown unique sources causing the observed fluctuations. We apply this new analysis on a data set from the Los Alamos National Laboratory site. We demonstrate that the sources identified by NMF$k$ have real physical origins: barometric pressure and water-supply pumping effects. We also estimate the barometric pressure efficiency of the monitoring wells. The possible applications of the NMF$k$ algorithm are not limited to hydrogeology problems; NMF$k$ can be applied to any problem where temporal system behavior is observed at multiple locations and an unknown number of physical sources are causing these fluctuations.

## 1. Introduction

Most site hydrogeological studies include analysis of water-level transients; for example, pumping-test interpretation, evaluation of aquifer recharge, estimation of surface/subsurface water interactions, etc. However, identification of the physical sources causing spatial and temporal variation of aquifer water levels is challenging [cf. *Vasco et al.*, 2000]. Typically, there are multiple sources (forcings) that can cause the observed transients (signal observations) including municipal water-supply pumping, seasonal recharge fluctuations, barometric pressures propagating through vadose zone and boreholes, surface/subsurface water interactions, etc.

Source identification can be complicated because (1) some of these source signals may have similar temporal patterns, (2) some of the signals may interfere with each other, (3) signal propagation through the medium from the signal entry point to the signal observation point may be nonlinearly attenuated and subdued. Identification of factors that are (or are not) causing the observed transients can be crucial for conceptualization of the hydrogeological conditions and characterization of aquifer properties. If the original signals that cause the observed pressure transients at the observation wells can be successfully "unmixed" from the observations, decoupled physics models may then be applied to analyze the propagation of each signal independently; for example, to identify vadose zone properties controlling the pneumatic propagation of barometric pressures, or to estimate aquifer properties influencing hydrodynamics of water-pumping transients. The problem of identification, estimation, and removal of barometric pressure effects in observed groundwater levels is especially challenging [*Rasmussen and Crawford*, 1997; *Toll and Rasmussen*, 2007]. Even though, we typically know the barometric pressure fluctuations at the ground surface, we do not know the manifestation of the barometric pressure changes in the aquifer (at and below the water

table) after the pneumatic flow through the vadose zone. Barometric effects in the aquifer depend on the pneumatic (air) pressure fluctuations at the top of the regional aquifer and these are typically unknown.

Various methods can be applied for identification of sources causing pressure fluctuations. They can be divided into three broad categories. First, there are model-based methods where the physical factors are explicitly simulated, and their impact identified through a formal inverse analysis [cf. *Harp and Vesselinov*, 2011; *Halford et al.*, 2012]. These approaches may require development of complex physics models with multiple degrees of freedom where properties of the flow medium may also be unknown. The identification of medium properties is typically a part of the inverse analysis. As a result, the inverse problem may be ill-posed with many plausible solutions and computationally intensive to solve [cf. *Carrera et al.*, 2005]. The second category consists of a broad range of classical statistical methods (e.g., correlation analysis, regression analyses, principle component analysis, etc.) capable of analyzing hydrological data without (or with limited) prior information about the physical factors causing the transients. However, the hydrological records typically consist of mixtures of unknown individual signals that may be correlated which limits the applicability of these methods. The third category includes a broad range of Blind Source Separation (BSS) methods based on unsupervised (objective and adaptive) learning algorithms [cf. *Jutten and Herault*, 1991; *Zarzoso and Nandi*, 1999]. A classical BSS conundrum is the so-called "cocktail-party" problem [cf. *McDermott*, 2009]. Briefly, in the cocktail-party problem, several microphones are recording all the sounds in a ball-room (e.g., music, conversations, noise, etc.). Each of the microphones is recording a mixture of the available sounds, and each sound is recorded with different amplitude that depends on many factors including the distance from the sound source to the microphones. To be able to "unmix" and reconstruct the original sound sources (melodies, voices, etc.) from the records, a BSS algorithm is needed to obtain the best possible extraction of the original source signals from the mixtures. The BSS algorithms are capable of revealing hidden features and dependencies in large sets of observed data, and, based on these features, building a representation of the data that can contribute to understanding the physical mechanisms behind these data. However, the currently available BSS algorithms cannot identify the locations and the strengths (forcings) of the signal at the source locations. The BSS algorithms only identify the observed manifestation of the source forcings at the observation points (the "microphones"). Here, we focus on the methods in the third (BSS) category. We should also note that the obtained BSS results can be applied in pneumatic/hydraulic tomography studies [e.g., *Yeh and Liu*, 2000, *Vesselinov et al.*, 2001a, b, *Ni and Yeh*, 2008, *Cardiff et al.*, 2009, *Illman et al.*, 2009, *Berg and Illman*, 2011a, b] where the data about the identified source fssignals will be applied to calibrate a physical based models to characterize the vadose-zone/aquifer heterogeneity.

The unmixing and reconstruction of the original signals in the BSS algorithms are usually based on some constrained and/or regularized optimization procedure minimizing an objective (cost) function together with a few imposed constraints, such as: maximum variability, statistical independence, non-negativity, smoothness, sparsity, simplicity, and others (further discussed below). The choice of the optimization constraints is usually based on a priori knowledge about the solved problem, and hence the constraints could be different for every particular case.

If the problem is solved in a temporally discretized framework, the goal of a BSS algorithm is to retrieve the original signals (sources), $\mathbf{S}$ ($\mathbf{S} \in \mathbf{M}_{p \times r}(\mathbb{R})$), that have produced observation records, $\mathbf{H}$ ($\mathbf{H} \in \mathbf{M}_{p \times m}(\mathbb{R})$), detected at a set of sensors (e.g., microphones), where $m$ is the number of the recording sensors, $r$ is the number of unknown signals, and $p$ is the number of discretized moments in time at which the signals are recorded at the sensors (the sources do not produce signals at discrete times; they are only recorded at discrete times). The time discretization does not need to be uniformly spaced; however, typically, the time steps are uniform and characterize the recording resolution of the sensors. We initially know only the matrix $\mathbf{H}$ containing the records of mixtures constructed from an unknown number of individual sources that are recorded by the sensors. Thus, in the simplest BSS problem, the recorded data, $\mathbf{H}$, is formed by a linear mixing of $r$ unknown original signals $\mathbf{S}$, blended by an unknown mixing matrix, $\mathbf{A}$ ($\mathbf{A} \in \mathbf{M}_{r \times m}(\mathbb{R})$), i.e.,

$$\mathbf{H}_{p \times m} = \mathbf{S}_{p \times r} \mathbf{A}_{r \times m} + \mathbf{E}_{p \times m}, \tag{1}$$

where $\mathbf{E}$ is a matrix ($\mathbf{E} \in \mathbf{M}_{p \times m}(\mathbb{R})$) describing possible noise or errors in each of the $m$ experimental records ($\mathbf{E}$ is also unknown). Since both factors $\mathbf{S}$ an $\mathbf{A}$ are unknown (we do not know even the exact size of these matrices, because we do not know how many original sources have been mixed), the main difficulty in

solving the BSS problem is that it is underdetermined (ill-posed). There are two widely used methods resolving the BBS underdetermination: Independent Component Analysis (ICA) [*Herault and Jutten*, 1986; *Amari et al.*, 1996], and Nonnegative Matrix Factorization (NMF) [*Paatero and Tapper*, 1994; *Lee and Seung*, 1997]. Although ICA and NMF approach the BBS underdetermination differently, in specific situations, they both successfully separate sensor-recorded data formed by mixtures of unknown signals with noise and/or measurement errors. When solving a real problem, the physical meaning and interpretation are often the key components to determine the concrete method for decomposition (unmixing) of the analyzed data. Below, we will briefly describe the basic principles of ICA and NMF.

ICA estimates the source, **S**, and mixing, **A**, matrices based on equation (1) by maximizing the statistical independence of the retrieved source signals in **S** (i.e., the matrix columns are expected to be independent). Typically, the source statistical independence is achieved by maximizing some high-order statistics for each source signal, such as the kurtosis or negentropy (negative entropy). The main idea behind ICA is that, while the probability distribution of a linear mixture of sources in **H** is expected to be close to a Gaussian (according to the Central Limit Theorem), the probability distribution of the original independent sources is expected to be non-Gaussian. As a result, ICA aims a maximization of the non-Gaussian characteristics in the estimated sources in **S** with the goal to find statistically independent non-Gaussian sources that reproduce the experimental data (equation (1)). ICA has applications in various unrelated fields, such as: neural computation, advanced statistics, text mining, telecommunications, signal processing, and many others [cf. *Hyvärinen and Oja*, 2000]. ICA was also utilized in a wide class of groundwater analyses [*Westra et al.*, 2007], and it was proposed as a tool for describing dependencies between water levels observed in surface waters and groundwater monitoring wells near radioactive storage facilities [*Nuzhny et al.*, 2008]. The ICA technique (combined with various optimization and simulation tools) was utilized for extracting and filtering of hydrological signals [*Frappart et al.*, 2010, 2011; *Forootan and Kusche*, 2012] from data detected by the Gravity Recovery and Climate Experiment (GRACE) satellites [*Tapley et al.*, 2004]. A significant benefit of ICA exploited in these studies was its ability to maximize the statistical independence of numerically generated univariate time-series, e.g., when these time series have been created specifically to mimic spatial and time dependencies of original stream flows or rainfalls [*Keylock*, 2012]. Although ICA was originally built to separate mixtures of stochastic signals, it was recently demonstrated that it can be also applied to separate deterministic trends mixed with stochastic signals as well [*Forootan and Kusche*, 2013].

In contrast to ICA, NMF does not seek statistical independence or constrain any other statistical properties (i.e., NMF allows the estimated sources to be partially or entirely correlated); instead, NMF enforces a non-negativity constraint on the original sources in **S** and their mixing components in **A** (i.e. all the estimated matrix elements are greater than or equal to zero). These differences between NMF and ICA have important implications for the analyses presented in this paper, and, next, we will discuss them briefly.

First, in many situations, the specific physical properties of the analyzed systems unequivocally require a non-negativity constraint, e.g., if the data are composed of measurements of energies, masses, frequencies, densities, etc. The non-negativity means that the observed data have to be described only by additive signals that cannot cancel mutually. Such a representation of the signals from sensor-recorded data is called a "parts-based" representation [*Fischler and Elschlager*, 1973]. There are many situations where the parts-based representation is natural, for example, in the field of image recognition where the images are constructed by separate positive pixels that cannot cancel each other. The ICA application in parts-based cases is limited since the ICA idea of statistical independence leads to signals that can be subtracted in order to reproduce the observed data. As a result, ICA representation of parts-based systems is difficult to interpret [*Parra et al.*, 1999], and in such systems NMF usually gives better results [*Lee and Seung*, 1999].

Second, there are many cases where the ICA assumption that the original signals are statistically independent may contradict the physical conditions. Furthermore, it is apparent that if there are any dependencies (interactions) between the original signals (e.g., seasonality causing multiple separate signals to be temporally correlated), these signals could not be extracted by algorithms whose basis is to seek statistical independence [*Seung and Lee*, 2001].

Third, the ICA representation does not necessarily lead to sparsity (i.e., the extracted matrices **S** and **A** are typically full) because ICA represents the observed data by contributions from all the estimated signals. In contrast, because of the additivity and non-negativity requirements, NMF naturally leads to a sparseness in

both the signal, **S**, and mixing, **A**, matrices. The sparsity can be an important factor for finding patterns hidden in the observed data and estimated source signals [cf. *Cichocki et al.*, 2009]. Because of its sparsity, NMF representations are usually much simpler than ICA representations, and hence in many cases NMF solutions can be interpreted more easily.

In recent years, NMF have been successfully applied in various research fields: image processing, computer vision, medical imaging, spectra analyses, and many others [cf. *Cichocki et al.*, 2009; *Srivastava et al.*, 2008]. The simplicity of NMF interpretation makes its application very tempting to every problem where the statistical independence is not a requirement; NMF is the preferred method when partial signal correlations are expected or inevitable. However, NMF has a particular limitation; there is an ambiguity about how to define the number of the unknown signal sources, $r$. For each concrete problem, this uncertainty has to be resolved in a specific manner, and the solution is usually based either on subjective physical insights or objective criteria within the NMF algorithm.

Recently, a new NMF framework has been developed to address the estimation of the number of unknown sources (here we call it NMFc). In NMFc, a custom clustering algorithm was utilized to determine the number of unknown sources based only on the principle of parsimony. NMFc has been successfully applied for deciphering mutational signatures active in more than 7000 human cancer genomes [*Alexandrov et al.*, 2013a, 2013b].

In this paper, we develop an extended version of the NMFc framework, adapted for decomposing transient observations; we call this framework NMFk. The specifics of the NMFk framework are discussed in section 2. In section 3, the NMFk framework is applied to identify the unknown original sources causing the observed transients in the pressure fluctuations at monitoring wells. The analyzed pressure data set was collected at the Los Alamos National Laboratory site.

## 2. Methodology

### 2.1. Definitions

We consider an aquifer, subject to recording the changes in the hydraulic heads at $m$ monitoring (observation) wells. Let us assume that the aquifer is subjected to the influence of $r$ unknown physical sources causing pressure (water-level) fluctuations. These pressure fluctuations can be associated with different (independent or partially correlated; distributed or point) sources which cause propagation of pressure changes in the aquifer. The locations of the sources and their forcing transients are unknown. In our analyses, we do not make assumptions about which physical processes impact the signal propagation through the flow medium (the analysis is model-free, and there are no assumption about initial or boundary conditions either). Further, we make neither deterministic nor stochastic assumptions about the properties associated with these physical processes. The original source (forcing) signals are altered during their propagation through the flow media before reaching the observation wells. We focus only on the extraction (deconstruction) of the source signals as observed at the monitoring wells. The only assumption that we make is that source signals that are proportionally manifested at each observation well. The proportionality factor implicitly depends on the processes that occur in the flow medium and their properties. Therefore, NMFk allows for decoupling (deconstruction) of the source (forcing) signals as recorded at the observation points (wells) without models and model assumptions. The BSS cocktail-party problem analogy to the water-transients in aquifer is: the $m$ monitoring wells correspond to $m$ microphones, and the $r$ sources correspond to $r$ unknown sounds (barometric pressure, pumping, etc.) present in the ballroom (aquifer).

Based on the above description, we have (compare with the equation (1)),

$$h_i(t) = \mathcal{F}_i[s_1(t), s_2(t), s_3(t), \ldots, s_r(t)] + \epsilon_i(t), \forall i = 1, \ldots, m, \tag{2}$$

where $h_i(t)$ are the observed transient pressures at $i$-th well caused by $r$ sources with observation errors $\epsilon_i(t)$. The unknown sources are characterized by the transients (signals) $s_j(t), j = 1, \ldots, r$, they cause. These transients are propagated through the aquifer by some unknown mappings $\mathcal{F}_i[\cdot]$ (some of these mappings could be nonlinear) depending on aquifer properties and governing processes, and manifested in the observed transients $h_i(t), i = 1, \ldots, m$. The temporal characterization of $h_i(t), s_j(t)$ and $\epsilon_i(t)$ is commonly discretized based on the measurement frequencies. As a result, the observation period captured in equation (2) is represented by $p$ discrete measurements in time, and $h_i(t), s_j(t)$ and $\epsilon_i(t)$ are presented as rows

$h_{q,i}$, $s_{q,j}$, and $\epsilon_{q,i}$ ($q = 1, \ldots, p$) in the matrices $\mathbf{H}_{p \times m}$, $\mathbf{S}_{p \times r}$, and $\mathbf{E}_{p \times m}$, respectively. $\mathbf{H}_{p \times m}$ represents all the collected data, $\mathbf{S}_{p \times r}$ are the unknown source contribution at each discrete time moment ($s_{q,j} \in \mathbb{R}^+$), and $\mathbf{E}_{p \times m}$ are measurement errors or noise. Further, the simplest assumption is that the mapping $\mathcal{F}_i[\cdot]$ corresponds to a linear superposition, and hence it is assumed that the mapping can be represented with a multiplication factor, $a_{j,i}$, associated with propagation of each source $j$ to each monitoring well $i$. Thus, the general mapping, $\mathcal{F}_i[\cdot]$, is replaced by a mixing matrix $\mathbf{A}_{r \times m}$. The elements of the mixing matrix $a_{j,i}$ ($a_{j,i} \in \mathbb{R}^+$) do not depend on time but only on the governing processes, aquifer properties, and monitoring-well locations. Further, we assume that there are matrices $\tilde{\mathbf{S}}_{p \times r}$ and $\tilde{\mathbf{A}}_{r \times m}$, such as,

$$\mathbf{H}_{p \times m} \approx \tilde{\mathbf{S}}_{p \times r} \tilde{\mathbf{A}}_{r \times m}, \tag{3}$$

where the number of the sources, $r$, is unknown but we will assume that its value is less than the number of the monitoring wells, $m$, i.e., $r \leq m$. Here and henceforth $(\tilde{\cdot})$ is used to denote estimated quantities. The NMFk estimated noise will be equal to $\mathbf{E}_{p \times m} = \mathbf{H}_{p \times m} - \tilde{\mathbf{S}}_{p \times r} \tilde{\mathbf{A}}_{r \times m}$.

## 2.2. Source Identification

To be able to estimate the matrices $\tilde{\mathbf{S}}$ and $\tilde{\mathbf{A}}$ that reproduce the data $\mathbf{H}$ with some accuracy (equation (3)), we choose an objective function $\mathcal{O}$ that measures the discrepancy between the observation data, $\mathbf{H}$, and the NMF predictions, $\tilde{\mathbf{S}}_{p \times r} \tilde{\mathbf{A}}_{r \times m}$, based on Frobenius norm ($\| \cdot \|_F$) [cf. *Golub and Van Loan*, 1980],

$$\mathcal{O} = \frac{1}{2} \left( \| \mathbf{H} - \tilde{\mathbf{S}} * \tilde{\mathbf{A}} \|_F \right)^2 = \sum_{i=1}^{m} \sum_{q=1}^{p} \left( h_{q,i} - \sum_{j=1}^{r} s_{q,j} a_{j,i} \right)^2. \tag{4}$$

The minimization of the discrepancies between the observation data and the NMF predictions, $\mathbf{E} = \mathbf{H} - \tilde{\mathbf{S}}\tilde{\mathbf{A}}$, is equivalent to representing the discrepancies as a white noise, i.e., as independent Gaussian random variables. Indeed, minimizing $\mathcal{O}$ is the same as maximizing a likelihood function $\sim \exp\left( -\frac{\| \mathbf{H} - \tilde{\mathbf{S}} * \tilde{\mathbf{A}} \|_F^2}{2\sigma^2} \right)$ of normally distributed data points $\mathbf{E}$ with a standard deviation $\sigma$ [cf. *Cichocki et al.*, 2009].

To minimize our objective function $\mathcal{O}$, we apply the multiplicative update algorithm introduced by [*Lee and Seung*, 1999]. This algorithm consecutively updates the source components, $\tilde{\mathbf{S}}_{p \times r}$, while keeping $\tilde{\mathbf{A}}_{r \times m}$ fixed, and next updates the mixing components, $\tilde{\mathbf{A}}_{r \times m}$, while keeping the new $\tilde{\mathbf{S}}_{p \times r}$ fixed. When this algorithm is applied, if a matrix component becomes zero, its value remains at zero for all the successive iterations. To bypass this problem, the values of all the elements $s_{q,j}$ and $a_{j,i}$ are forced to be greater than a small positive constant, $\eta$ (e.g. $\eta = 10^{-16}$ is applied in the analyses presented here). Each step of the NMF solution is sought using the gradient descent approach in the following multiplicative update formulas, written as in Theorem 1 in [*Lee and Seung*, 1999],

$$a_{j,i}^* \leftarrow a_{j,i} \frac{\left[ \tilde{\mathbf{S}}^\mathsf{T} \mathbf{H} \right]_{j,i}}{\left[ \tilde{\mathbf{S}}^\mathsf{T} \tilde{\mathbf{S}} \tilde{\mathbf{A}} \right]_{j,i} + \eta}, \forall i = 1, \ldots, m, j = 1, \ldots, r, \tag{5}$$

$$s_{q,j}^* \leftarrow s_{q,j} \frac{\left[ \mathbf{H} \tilde{\mathbf{A}}^{*\mathsf{T}} \right]_{q,j}}{\left[ \tilde{\mathbf{S}} \tilde{\mathbf{A}}^* \tilde{\mathbf{A}}^{*\mathsf{T}} \right]_{q,j} + \eta}, \forall j = 1, \ldots, r, q = 1, \ldots, p. \tag{6}$$

but with a small positive constant $\eta$ added to the denominator to avoid division by zero, and with an explicit marking of the updated matrices with "*." The solution is iterated until reaching the criteria for convergence, defined as a given number of iterations without a substantial change in the objective function $\mathcal{O}$, e.g. the change $\Delta \mathcal{O}$ is less than $10^{-10}$ over a predefined number of iterations (e.g., 10,000). The notation $[\mathbf{AB}]_{j,i}$ is equivalent to the $(j,i)$-th element of the matrix $\mathbf{C} = \mathbf{A} \times \mathbf{B}$. Matrix $\mathbf{A}^*$ with elements $a_{j,i}^*$, denotes the matrix $\mathbf{A}$ after the transformation (equation (5)). It has been proven that the Frobenius norm is: (i) nonincreasing under the update rules (5, 6), and ii) it is invariant under these updates if and only if $\tilde{\mathbf{S}}$ and $\tilde{\mathbf{A}}$ are at the norm stationary point [*Seung and Lee*, 2001]. Thus, a perfect reconstruction of $\mathbf{H}$ is necessarily a fixed point of the above update multiplicative rules (5, 6).

To estimate the number of the unknown sources, we apply the NMF algorithm described above to perform $m$ sets of independent analyses (recall $m$ is the number of the monitoring wells). Each of these $m$ sets of analyses is performed with a different number of predetermined sources, $r$; that is, we set the number of

unknown sources $r$, in equation (1), to range as $r=1,\ldots,m$. In addition, in each of these $m$ sets of analyses, we execute $n$ NMF runs with different random initial values (a white noise between 0 and 1) for $\tilde{\mathbf{S}}$ and $\tilde{\mathbf{A}}$. The number, $n$, of the NMF runs in each of the $m$ sets, is determined by the convergence of the final average solution (described below). As a result, the total number of the NMF runs becomes $m \times n$. We denote each set of $n$ solutions, corresponding to the same number, $r$, of predetermined (unknown) sources, as: $\mathcal{H}_r = ([\tilde{\mathbf{S}}_r^1 ; \tilde{\mathbf{A}}_r^1], [\tilde{\mathbf{S}}_r^2 ; \tilde{\mathbf{A}}_r^2], \ldots, [\tilde{\mathbf{S}}_r^n ; \tilde{\mathbf{A}}_r^n])$. In $\mathcal{H}_r$, each source solution-matrix $\tilde{\mathbf{S}}_r^i ; i=1,\ldots, n$ contains $r$ unique source signals represented by the columns of the matrix, $\tilde{\mathbf{S}}_r^i$. As a result, the total number of obtained possible solutions for the source signals becomes $n \times (1+2+\ldots+m) = \frac{n(m+1)m}{2}$.

### 2.3. Estimating the Optimal Number of the Unknown Source Signals

We apply an unsupervised clustering, based on a variation of the $k$-means clustering algorithm [*Hartigan and Wong*, 1979], to estimate the optimal number of the unknown sources, $r$. In the $k$-means clustering analysis, the number of the estimated clusters $k$ is predetermined. Here, we perform clustering analysis with a fixed number of clusters, $k$, to each of the $m$ sets of solutions, $\mathcal{H}_r$. Note that in each case, the number of the estimated clusters $k$ is equal to the number of the unknown sources $r$, where $r$ varies as: $r=1,\ldots, m$. Specifically, we cluster, based on their similarity, the columns of the matrices $\tilde{\mathbf{S}}_r^i ; i=1,\ldots, n$, in each set of solutions $\mathcal{H}_r$. Thus, each of the $n \times r$ columns in the estimated source matrices $\tilde{\mathbf{S}}$, within a given solution set $\mathcal{H}_r$, is assigned to exactly one of the $k$ clusters, derived via $k$-means clustering analysis. The similarity between any two $p$-dimensional source signals, $s_{q,j1}$ and $s_{q,j2}$, is measured by their cosine distance (also known as cosine similarity), $\rho(s_{q,j1}, s_{q,j2})$ [cf. *Pang-Ning et al.*, 2006]

$$\rho(s_{q,j1}, s_{q,j2}) = 1 - \frac{\sum_{q=1}^{p} s_{q,j1} s_{q,j2}}{\sqrt{\sum_{q=1}^{p} (s_{q,j1})^2} \sqrt{\sum_{q=1}^{p} (s_{q,j2})^2}}. \tag{7}$$

In our case, $0 \le \rho(s_{q,j1}, s_{q,j2}) \le 1$, because of the non-negativity of the sources signals (the components $s_{q,j}$ of $\tilde{\mathbf{S}}$ are positive). In each of the $m$ solution sets, $\mathcal{H}_r$, the $k$-means clustering analysis constructs exactly $k$ clusters ($k=1,\ldots, m$). Cluster centroids are calculated by averaging the solutions that belong to the same cluster. The $k$ cluster centroids, in each set $\mathcal{H}_r$, represent the average solution $\tilde{\mathbf{S}}_r^{\mathbf{a}}$ with $k$ sources ($k=r$), where the superscript $a$ is to emphasize that these are averaged matrices (solutions). The $k$-means clustering of the solutions also effectively cluster the set of the rows of the corresponding mixing matrices $\tilde{\mathbf{A}}_r^i ; i=1,\ldots, n$, in $\mathcal{H}_r$, and hence, we compute also the corresponding centroids forming the average mixing matrix $\tilde{\mathbf{A}}_r^a$. Each pair of matrices $\tilde{\mathbf{S}}_r^a$ and $\tilde{\mathbf{A}}_r^a$ provides the best estimate for representing the observed data $\mathbf{H}$ with the number of the unknown sources $r$, (where $r < m$ and $r = k$, in the $k$-means clustering). As a result, we obtain a series of $m$ average solutions $([\tilde{\mathbf{S}}_1^a ; \tilde{\mathbf{A}}_1^a], [\tilde{\mathbf{S}}_2^a ; \tilde{\mathbf{A}}_2^a], \ldots, [\tilde{\mathbf{S}}_m^a ; \tilde{\mathbf{A}}_m^a])$, each with a different number of unknown sources $r$ ($r=1,\ldots, m$).

The robustness of each of these $m$ solutions is evaluated by assessing the tightness of the corresponding clusters based on the Silhouette method [*Rousseeuw*, 1987]. For each cluster element (there are $n \times r$ cluster elements (columns) in each of the $m$ sets of solutions $\mathcal{H}_r, r=1,\ldots, m$), we compute the average dissimilarity (based on cosine distance; equation (7)) with all other elements within the same cluster ($a_d; d=1,\ldots, n \times r$) as well as the dissimilarity with all other elements within the other clusters ($b_d$). For each $k$-means ($k=r$) cluster solution, the Silhouette values ($c_d$) for each cluster element are computed as

$$c_d = \frac{b_d - a_d}{max[a_d, b_d]} \forall d=1,\ldots, n \times r, \tag{8}$$

where $-1 \le c_d \le 1$. If $c_d \approx 1$, the element is appropriately clustered. If $c_d \approx -1$, the element should belong to the closest neighboring cluster (the cluster for which $b_d$ is the lowest). If $c_d \approx 0$, the element is between two clusters. A Silhouette width is computed for each of the clusters by averaging all the Silhouette values ($c_d$) for all the cluster elements. Thus, the Silhouette width of a cluster demonstrates to what extent the elements within this cluster are similar to each other and dissimilar to the elements of the other clusters [cf. *Izenman*, 2008]. Further, we calculate the average Silhouette width of the $k$ clusters in each set of solutions $\mathcal{H}_r$. If this average Silhouette width is close to 1, all the clusters in the set of solutions $\mathcal{H}_r$ are tight and the NMF solutions are robust, i.e., the NMF algorithm is consistently extracting a similar set of $r$ solutions when applying different random initial values for the estimated matrices, $\tilde{\mathbf{S}}$ and $\tilde{\mathbf{A}}$. In contrast, if the average Silhouette width is much less than 1, the NMF solutions are not robust and we have a lack of reproducibility.

Hence, the average Silhouette width for each of the $m$ sets of $k$-means cluster solutions is a measure of the reproducibility and stability of the derived average solutions with a predetermined number of $r$ unknown sources.

In addition to the robustness, the average Frobenius reconstruction error (equation (4)) is used to evaluate the accuracy with which the derived average (cluster) solutions $[\tilde{\mathbf{S}}_{\mathbf{r}}^{\mathbf{a}}; \tilde{\mathbf{A}}_{\mathbf{r}}^{\mathbf{a}}]$ reproduce the observations $\mathbf{H}$. In general, the solution accuracy increases (while the solution robustness decreases) with the increase of the number of unknown sources. Hence, the average Silhouette width and the Frobenius reconstruction error for each of the $m$ $k$-means cluster solutions can be used to define the optimal number of sources, $\hat{r}$. Specifically, we select $\hat{r}$ to be equal to the minimum number of sources that accurately reconstruct the observations (i.e., the Frobenius reconstruction error is less than a given value), but the solutions are sufficiently stable (i.e., the Silhouette width is close to 1).

In general, the NMF (and any BSS model-free) solutions are not unique: any nonsingular matrix $C$ and its inverse, $C^{-1}$ can be used to transform the solutions $S$ and $A$ of the equation $\mathbf{H} \approx \tilde{S} * \tilde{A}$, if and only if $\tilde{S} * C > 0$ and $C^{-1} * \tilde{A} > 0$, because:

$$\mathbf{H} \approx \tilde{\mathbf{S}} * \tilde{\mathbf{A}} \equiv (\tilde{\mathbf{S}} * \mathbf{C}) * (\mathbf{C^{-1}} * \tilde{\mathbf{A}}), \tag{9}$$

In this case, the new (transformed) solutions will be: $\tilde{S} \rightarrow \tilde{S} * C$ and $\tilde{A} \rightarrow C^{-1} * \tilde{A}$. Control over this non-uniqueness of the NMF factorization can be obtained by enforcement of some constraints meaningful for the particular problem [cf. *Xu et al.*, 2003]. In our case, the only constraint is "non-negativity" of all the matrix components of $\mathbf{H}$, $\tilde{\mathbf{S}}$ and $\tilde{\mathbf{A}}$.

The NMF$k$ analysis also allows for the estimation of the standard Signal-to-Noise-Ratio (SNR) criteria for each of the reconstructed mixtures by using the formula (see, e.g., *Cichocki et al.* [2009]):

$$\mathrm{SNR}_i = \frac{\|s_{q,j} a_{j,i}\|_F^2}{\|h_{q,i} - s_{q,j} a_{j,i}\|_F^2}, i = 1, \ldots, m, \tag{10}$$

where $\| \cdot \|_F$ is the Frobenius norm. The SNR criteria represent the ratio between the signal and the noise in the analyzed data for each monitoring point. If the extracted signal and the noise are of the same magnitude, the SNR criteria will be equal or close to 1. If SNR criteria are substantially greater than 1 than the estimated signals, are substantially larger than the estimated noise in the data. Theoretically, the NMF$k$ will fail to extract clear signals if there are substantial measurement errors (random noise) in the analyzed data (here, substantial noise will be the case when the noise magnitude is greater or equal to the magnitude of the observed signals). However, if in the considered system, there is a correlated (systematic) noise (i.e., a noise that causes similar temporal pattern at each observation point) with a specific temporal pattern, then the NMF$k$ algorithm, by definition, will treat and eventually extract such a noise as a separate signal that participates in the forming of the NMF$k$-identified signal mixtures [cf., *Cichocki et al.*, 2009)

The $k$-means algorithm is applied to estimate the optimal number of sources characterizing the data. Here, the cluster centroids representing the average of the solutions within the cluster are applied to obtain the final solution. However, it is also possible to analyze all the solutions within a cluster to explore probabilistically the uncertainty as a result of multiple acceptable solutions within the cluster as well. In this case, the solution likelihood can be estimated based on the Frobenius norm.

The NMF$k$ algorithm is coded in MatLab based on the NMFc code utilized by [*Alexandrov et al.*, 2013b]. In the next section (3), we apply the algorithm to interpret the water-level transients observed in regional aquifer monitoring wells, and identify unknown source (forcing) signals causing these fluctuations. The only information provided to the algorithm are the water-level records. We expect that there are sufficient differences in the observed water-level transients to allow unique identification of multiple source (forcing) signals.

## 3. Site Data

The analyzed pressure data are collected at four monitoring wells within the regional aquifer beneath the LANL site. The aquifer is a complex stratified hydrogeologic structure [*Vesselinov*, 2004] which includes

unconfined zones (under phreatic conditions near the regional water table) and confined zones (the deeper regional aquifer zones). The regional aquifer is also potentially in hydraulic connection with surface waters of Rio Grande which flows East of the LANL site. The vadose zone and aquifer are composed of basin-fill sediments (alluvium, volcanic ash beds, tuffs) and fractured volcanic rocks (dacites and basalts). The aquifer study area where the monitoring and the water-supply wells are screened is predominantly within sedimentary units (sands and gravels) *Broxton and Vaniman* [2005]. The total thickness of the regional aquifer is unknown but it exceeds 1000m. The vadose zone between the ground surface and the regional water table has a thickness of ∼300m, and it is also a complex hydrogeologic structure including perched water-saturation zones and areas of focused recharge. Barometric pressures propagate through the vadose zone and impact aquifer pressures.

Due to concerns related to the migration of potential LANL-derived contaminants in the subsurface, an extensive monitoring network is established in the regional aquifer beneath LANL. The network includes more than 100 regional monitoring wells with about 350 monitoring screens [*Koch and Schmeer*, 2008]. Pressure fluctuations at each screen are automatically monitored using pressure transducers. The aquifer beneath LANL is an important source of water for LANL and neighboring municipalities. There are six water-supply wells in close vicinity to the study area, and 19 more water-supply wells are located nearby. The ultimate goal is to incorporate all the water-level data in the conceptualization, development, calibration, and analysis of the regional aquifer model.

Here, we analyze a subset of the data from monitoring wells, limiting our analysis to a small subarea of the LANL site. The pressure records considered here are collected from four monitoring wells labeled as R-11, R-13, R-15, and R-28. We denote the pressure records of R-11, R-13, R-15, and R-28 as $\mathbf{h}_1$, $\mathbf{h}_2$, $\mathbf{h}_3$, and $\mathbf{h}_4$, respectively. Figure 1 displays a map of the monitoring-well locations. The map also includes the six water-supply wells located in the vicinity; these wells are actively used for municipal water-supply pumping. The surface water in Rio Grande is also potentially hydraulically connected with the aquifer; Rio Grande flows about 8 km to the East from the monitoring wells. Table 1 lists well location coordinates and screen depths below the water table. Figure 2 presents the pressure records observed at the monitoring wells. The data record is selected to be about a 1.5 years to cover two summer periods when most of the water-level decline occur (Figure 2).

Previous analysis by *Harp and Vesselinov* [2011] demonstrated that the observed pressure transients in the LANL aquifer are influenced by (a) municipal water-supply pumping, (b) barometric pressure effects, and (c) long-term water-level declines with unknown origin (e.g., long-term climate changes, changes in regional discharge elevation, or regional aquifer overexploitation). The prior work was focused on the data from three of the wells (R-11, R-15, and R-28) and used a longer pressure record (2005–2009), and the analysis was based on an inverse model that predicts pressure changes caused by the known transients (daily) in the municipal water-supply [*Harp and Vesselinov*, 2011]. The pumping influences were modeled using a relatively simple analytical model based on the Theis equation [*Theis*, 1935] where the aquifer parameters were estimated in the inverse process. The long-term water-level declines were represented by a linear model with unknown slope (also estimated in the inverse process). The barometric pressure effects were removed assuming constant barometric efficiency equal to 100% for each of the three wells. The barometric efficiency of ∼100% was estimated based on analyses of short-term (several weeks) pressure observations during installation and development of the monitoring wells.

The pressure data applied in the analyses presented here overlap with the data used by *Harp and Vesselinov* [2011]. This allow us to compare the results obtained with two distinct inverse methods relying on different assumptions. In contrast with *Harp and Vesselinov* [2011], here we use shorter pressure records, consider one additional well (R-13), and do not take into account the existing information about the transients in the municipal water-supply pumping and the barometric pressure. The goal here is to estimate pumping and barometric effects in the pressure data using the pressure data by itself only. No data about the municipal water-supply pumping rates or barometric pressure fluctuations are applied in the analysis. The shorter observation period is selected to test applicability of the developed NMF*k* methodology to identify source signals based on limited data records.

The analyzed data represent water levels observed between 6 May 2005 and 27 August 2006, at four monitoring wells ($m = 4$). For each day of this period, observations have been taken hourly at noncoincidental time points. We averaged all the data to obtain daily records with $p = 481$ data points.
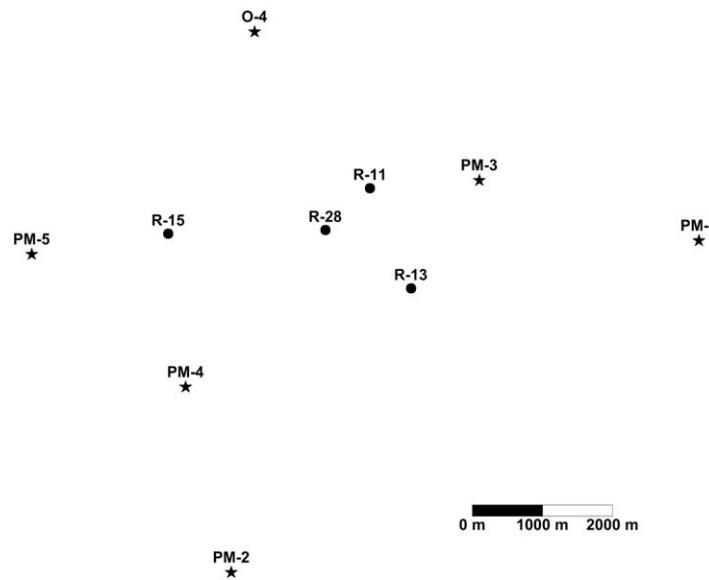
**Figure 1.** Monitoring (dots) and water-supply well (stars) locations.

Further, the daily-averaged water-level transients, $h'_{q,i}$, have been transformed for each well $i$ ($i = 1, \ldots, m$) as:

$$h_{q,i} = \frac{\Delta h'_{q,i}}{\sum_{w=1}^{p} \Delta h'_{w,i}}, q = 1, \ldots, p,$$

(11)

where $\Delta h'_{q,i} = h'_{q,i} - h'_{min,i}$, and $h'_{min,i}$ is the minimum level at each well $i$ during the observation period ($h'_{min,i} = \min (h'_{q,i}), q = 1, \ldots, p$) and represents a zero datum for each well. Note that this transformation imposes "non-negativity" and normalizes the data so that sum of daily records in the column vectors of the matrix **H** is equal to 1. In this way, we define the data matrix **H**, and this is the only input information provided to the NMF$k$ algorithm to identify unknown source signals in the data.

## 4. Results

Evaluation of the number of the unknown source signals causing the observed pressure transients in the Los Alamos aquifer was performed using the NMF$k$ algorithm. After all the NMF$k$ iterations, we obtain four solutions for $r = 1, 2, 3$, and 4 unknown, unique, source signals. The final number of NMF$k$ runs, $n$, for each of the four sets of solutions is equal to 200. For each of the solutions in the four solution sets, the number of NMF$k$ iteration steps is $10^6$. Figure 3 shows the average Silhouette widths (left vertical axis) and the average Frobenius reconstruction norm (right vertical axis) of the four solutions. For the case of one source, the average Silhouette width is equal to 1, and the solution is robust by definition. The Silhouette width for two-source solution is slightly below 1 which still suggests robust estimates. The addition of a third and forth unknown source, the average Silhouette width drops substantially (Figure 3). This defines low reproducibility of these NMF$k$ solutions, and the solutions can be considered unstable.

In contrast, average Frobenius reconstruction norm in Figure 3 is relatively high for the case of only one unknown source. The norm is relatively low and similar for two, three, and four sources. Based on the balance between the average Silhouette widths and the average Frobenius reconstruction norm in Figure 3, the NMF$k$ solution with two sources is selected to be the best because it provides a sufficiently accurate and robust approximation of the observed data.

To illustrate the tightness of the solutions in the two-source case, we represent schematically in Figure 4 the two final clusters representing the selected NMF$k$ solutions. The figure also shows the cluster maximum radii and the relative distance between the clusters; the separations between solutions and clusters in Figure 4 are based on the cosine distance (equation (7)).

The two source signals that were obtained are shown Figure 5. The two signals are characterized by similar large-scale (seasonal) fluctuations but very different small-scale (daily) transient profiles. Both signals have lows in approximate ranges of 80–150 and 380–450 days (the summer months) and highs between 250 and 350 days which are the late-winter/early spring months. Each of the source signals contribute differently to the four pressure records observed at the monitoring wells R-11, R-13, R-15, and R-28. The NMF$k$ reconstructions of the four observed pressure records (black dots) with two source signals are presented in Figure 6. Note that the analyses provide almost perfect reconstruction of the data (the residuals are normally distributed and uncorrelated). The correlation coefficients, $\rho_i$, between the NMF$k$-reconstructed mixtures and the

**Table 1.** Well Location Coordinates and Screen Depths Below the Water Table[a]

| Well | $x$ (m) | $y$ (m) | $z_0$ (m) | $z_1$ (m) |
|------|---------|---------|-----------|-----------|
| R-11 | 499882.61 | 539296.05 | 5.57 | 12.55 |
| R-13 | 500174.36 | 538579.77 | 36.73 | 55.14 |
| R-15 | 498442.06 | 538969.46 | 0 | 15.04 |
| R-28 | 499563.69 | 538995.82 | 13.15 | 20.41 |
| PM-1 | 502229.40 | 538920.57 | 44.55 | 512.92 |
| PM-2 | 498865.40 | 536571.86 | 33.48 | 423.29 |
| PM-3 | 500661.43 | 539352.74 | 47.81 | 528.99 |
| PM-4 | 498537.89 | 537892.75 | 49.02 | 535.69 |
| PM-5 | 497467.13 | 538822.39 | 53.59 | 551.82 |
| O-4 | 499060.43 | 540408.91 | 88.06 | 540.28 |

[a]$z_0$, screen top; $z_1$, screen bottom.

observed data in Figure 6 are: $\rho_{R_{11}} = 0.991$, $\rho_{R_{13}} = 0.976$, $\rho_{R_{15}} = 0.986$, and $\rho_{R_{28}} = 0.978$. The correlation coefficients are close to 1 demonstrating the high quality of the NMF$k$ representation of the observed data. The estimated standard Signal-to-Noise-Ratio (SNR$_i$) criteria for each observation well based on equation (10) are (rounded to two significant digits): SNR$_{R_{11}} \simeq 16,000$, SNR$_{R_{13}} \simeq 10,000$, SNR$_{R_{15}} \simeq 32,000$, and SNR$_{R_{28}} \simeq 5000$. The SNR criteria are substantially greater than 1 demonstrating the NMF$k$ estimated signals are substantially larger (orders of magnitude) than the estimated noise in the data.

The mixing matrix **A** representing contribution of the two source signals in the observed pressures at the monitoring wells is listed in Table 2 (note that mixing components for each well add up to 1). The estimated mixing components are not explicitly related to physical properties representing the propagation of the source signals through the subsurface flow medium. The estimated mixing components only characterize how the unknown source (forcing) signals are observed at the monitoring wells. It is important to emphasize that in the presented analysis, we do not estimate the source locations and the magnitude of the transients at the source locations; we also do not estimate how the unknown source transients (forcings) changed during their propagation form their initiation points (the water-supply wells for pumping effect, and the ground surface for the barometric effects) through the subsurface flow medium (the aquifer for pumping effects, and the vadose zone and the aquifer for the barometric effects) before reaching the observation points (the monitoring wells). The NMF$k$ algorithm only characterizes the manifestation of the source transients in the observed data. It also estimates the relative magnitudes (proportions) of the source transients at the observation points.

We expected the two source signals to be related to barometric pressure effects and the water-supply pumping, respectively. Figure 7 shows a comparison between the first reconstructed source signal and measured barometric pressure fluctuations at the LANL site (barometric pressure data are collected at the "TA-54" barometric station; data were downloaded from http://weather.lanl.gov). The figure demonstrates a strong correlation between the first source signal and barometric pressure fluctuations. The high correlation
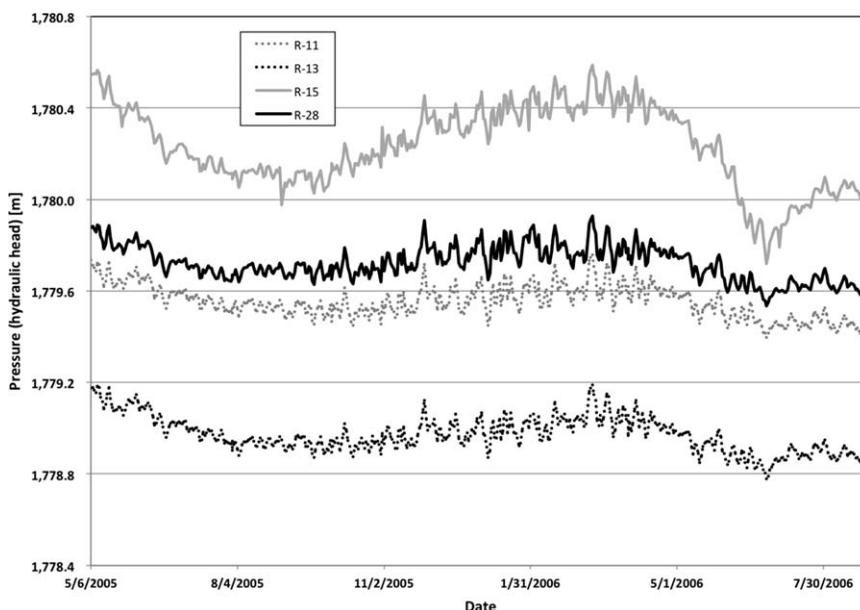


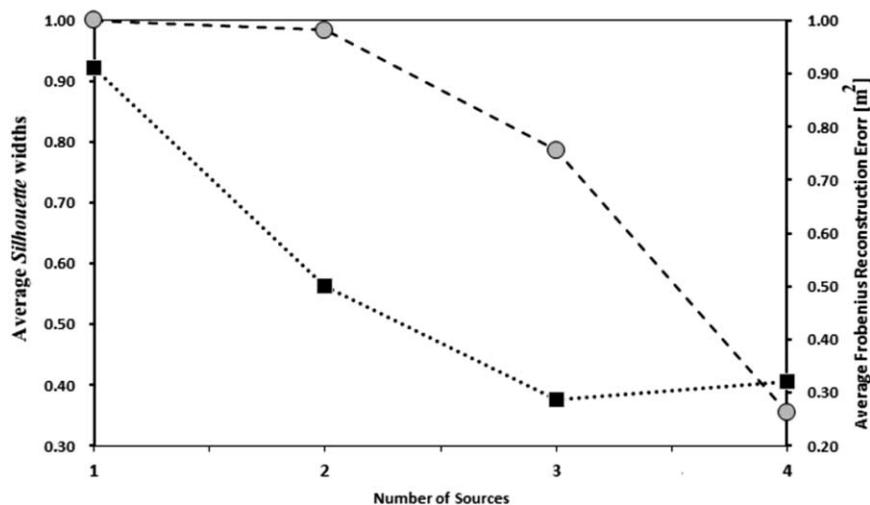**Figure 2.** Observed pressure transients at the monitoring wells.

**Figure 3.** The average Silhouette widths (gray circles; left vertical axis) and the average Frobenius reconstruction norm (black squares; right vertical axis; (m$^2$)) of the solutions.

is also demonstrated in a scatter plot between the two transients in Figure 8; the correlation coefficient is $\sim-0.85$. Figure 9 presents a comparison between the total groundwater volume pumped by the nearby six water-supply wells and the second reconstructed source signal. There are well-defined water-level minima between 80–150 days and 380–430 days that correspond to increased water-supply pumping in the summer months of 2005 and 2006 (Figure 6). Figure 9 suggests that the observed aquifer pressures are responding with a general lag time of about 20 days which is caused by the aquifer hydrodynamics. This lag is well defined for the period of high drawdowns between 380 and 430 days (Figure 9b). The correlation between the pumped volumes and pressure declines is not expected to be perfect; important is the general consistency of the observed temporal trends (with some delay). The transients in the second recon-structed source signal are also consistent with the water-supply pumping effects estimated in previous anal-yses [cf. *Harp and Vesselinov*, 2011]. Therefore, the second signal is associated with the water-supply pumping.

The observed seasonal similarities between the two signals are expected but somewhat coincidental. In general, the barometric pressures are lower during the summer months and higher in the winter. Similarly, the water consumption is higher in the summer and lower in the winter.

From a hydrogeologic perspective, it is interesting that based on our analyses, we conclude that four moni-toring wells spread over relatively large distances, and located at varying distances from multiple pumping wells (Figure 1) are responding similarly to the water-supply pumping (Figure 6). All the monitoring wells are responding to the same signal (Source Signal 2 in Figure 5) only with different scaling coefficients as listed in Table 2. The lack of spatial dependence (due to differences in the relative distance between individ-ual monitoring and pumping wells) in the observed pumping signal is potentially caused by the aquifer conditions. The monitoring wells are screened in the top portion of the regional aquifer while the municipal wells are pumping at depth. Based on the available hydrogeological information, it is expected that the aquifer will behave as a leaky, confined aquifer where the vertical propagation of the pumping drawdowns

is subdued by hydrostrati-graphic layering causing verti-cal anisotropy. As a result, the observed pumping drawdowns in the monitor-ing wells are relatively simi-lar and do not exhibit strong spatial dependence. This is an important finding that has implications for
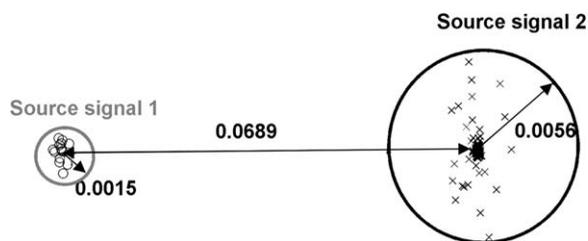


**Figure 4.** Schematic representation of the relative distances between the two clusters.
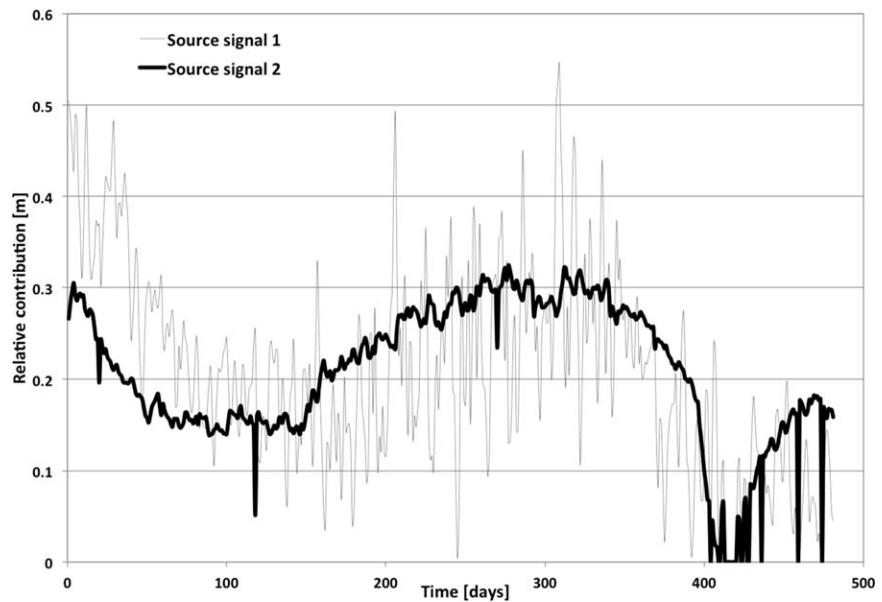
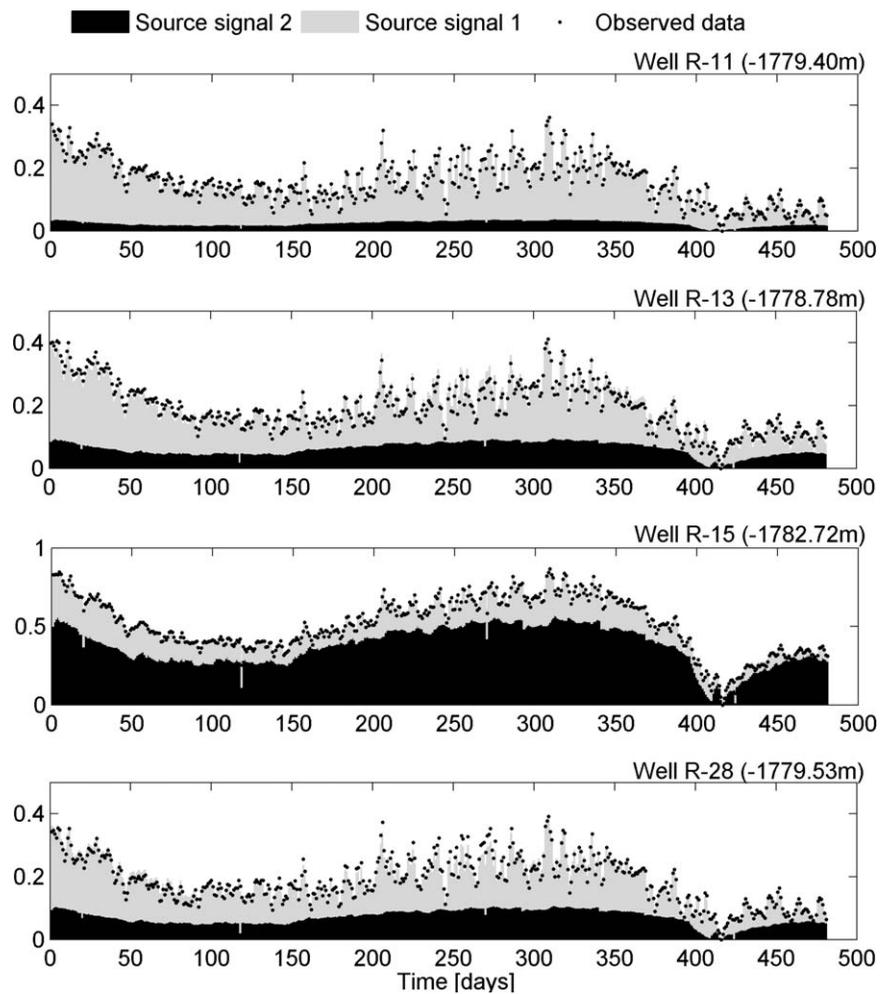**Figure 5.** The NMF*k* identified two source signals.



**Figure 6.** Reconstruction of the observed pressure transients (m) by the two identified source signals (note the difference in the vertical axis in the third subplot for R-15).

**Table 2.** Mixing Matrix **A** Representing Contribution of the Two Source Signals in the Observed Pressures at the Monitoring Wells

| Source Signal | R-11 | R-13 | R-15 | R-28 |
|---|---|---|---|---|
| #1 (barometric effects) | 0.84 | 0.67 | 0.27 | 0.62 |
| #2 (pumping effects) | 0.16 | 0.33 | 0.73 | 0.38 |

model development and conceptualization of the site conditions. The similarity in the water-supply response may be also caused because the nearby pumping wells had a relatively similar pumping regime within the analyzed time period [cf. *Harp and Vesselinov*, 2011].

Additional analyses that include data from more monitoring wells and over longer time periods may be capable to discriminate between the pumping effects of the individual water-supply wells.

Another important conclusion from the analyses is that there are apparently no strong signal transients in the observed pressure data caused by other factors such as variability in infiltration recharge or surface/subsurface water flow. The lack of strong signal caused by infiltration recharge is caused by the properties of the vadose zone. The relatively large vadose zone thickness (~300m) and the existing perched saturated-groundwater horizons within the vadose zone delay, subdue, and diffuse the impact of infiltration events on the water levels in the regional aquifer. This impacts the conceptualization of the hydrogeological conditions at the site.

The estimated mixing matrix $\tilde{A}$ defines the relative contribution of the two components (Table 2). The mixing matrix components associated with the barometric pressure also allow us to compute the barometric efficiencies of the monitoring wells. The barometric efficiencies are computed based on the estimated contribution of the barometric pressure changes on the observed water levels. Here, the water-level fluctuations observed at R-11, R-13, R-15, and R-28 range 0.33, 0.42, 0.9, and 0.38 m, respectively (this is the range between the highest and lowest water levels within observation period for each well; Figure 2). Taking into account the estimated contribution of the barometric pressure changes (Table 2), the barometric contributions on the water-level changes at R-11, R-13, R-15, and R-28 are 0.27, 0.28, 0.24, and 0.24 m, respectively. Assuming that (1) each well is affected by the same barometric pressure fluctuations (i.e., ignoring small differences in the ground-surface elevations at the well head and local impacts of the topography on the barometric pressure distribution on the ground surface), and (2) the highest barometric impact (0.28 at R-13) is representative for 100% efficiency, the barometric efficiencies of R-11, R-15, and R-28 are 99%, 86%, and 84%, respectively. Previous analyses at the site estimated that all four wells have barometric efficiency of ~100% [*Koch and Schmeer*, 2008]. The analyses presented here suggest that only two of the wells have efficiency close to 100% (R-11 and R-13). The other two wells (R-15 and R-28) have much lower barometric
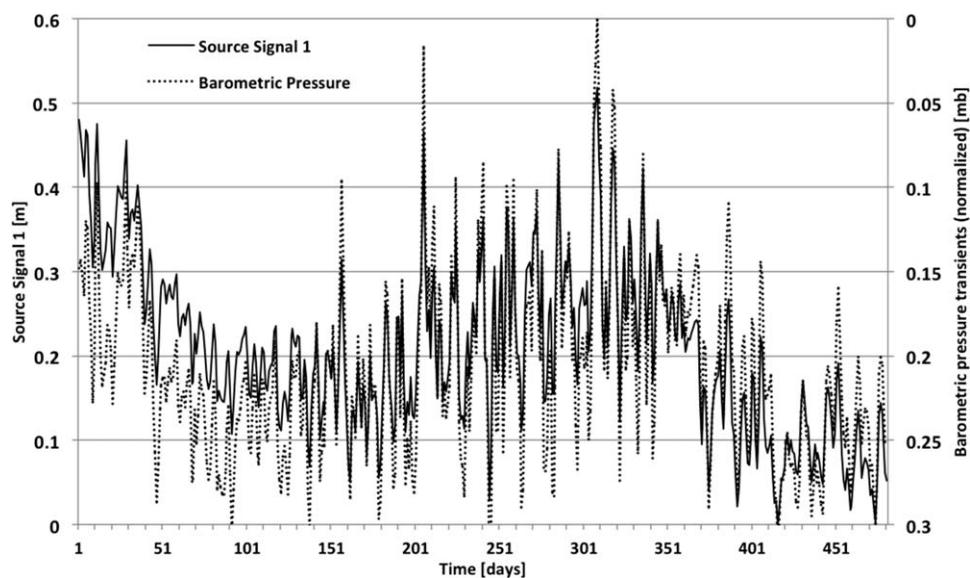


**Figure 7.** Comparison between the first reconstructed source signal and measured barometric pressure fluctuations (mb) at the LANL site (note that the right axis is reversed).
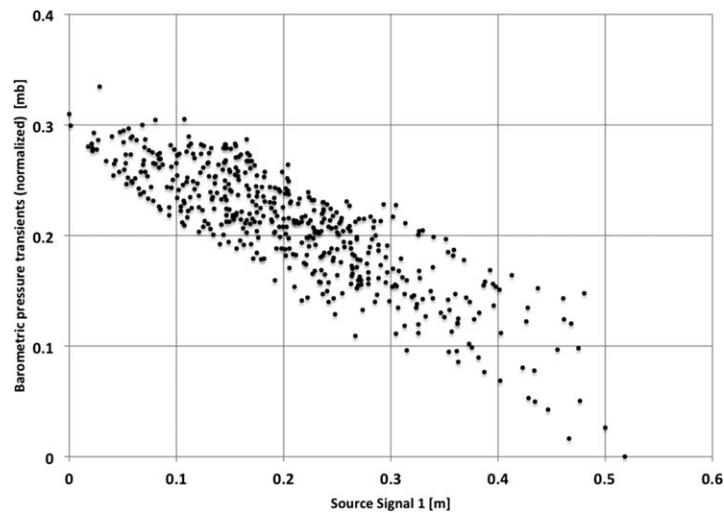
**Figure 8.** Scatter plot between the first reconstructed source signal and measured barometric pressure fluctuations (*mb*) at the LANL site; the correlation coefficient is ∼−0.85.

efficiency (∼85%) The deviations from 100% barometric efficiency suggest a more complicated regime of propagation of the barometric pressures through the vadose zone than previously expected.

All the analyses were performed using MatLab code implementing the NMF*k* algorithm; the simulations were executed in parallel using eight 3.0GHz Intel-Xeon processors and the final simulations were completed after about 12 h. In comparison, the analyses using physics-based models [*Harp*
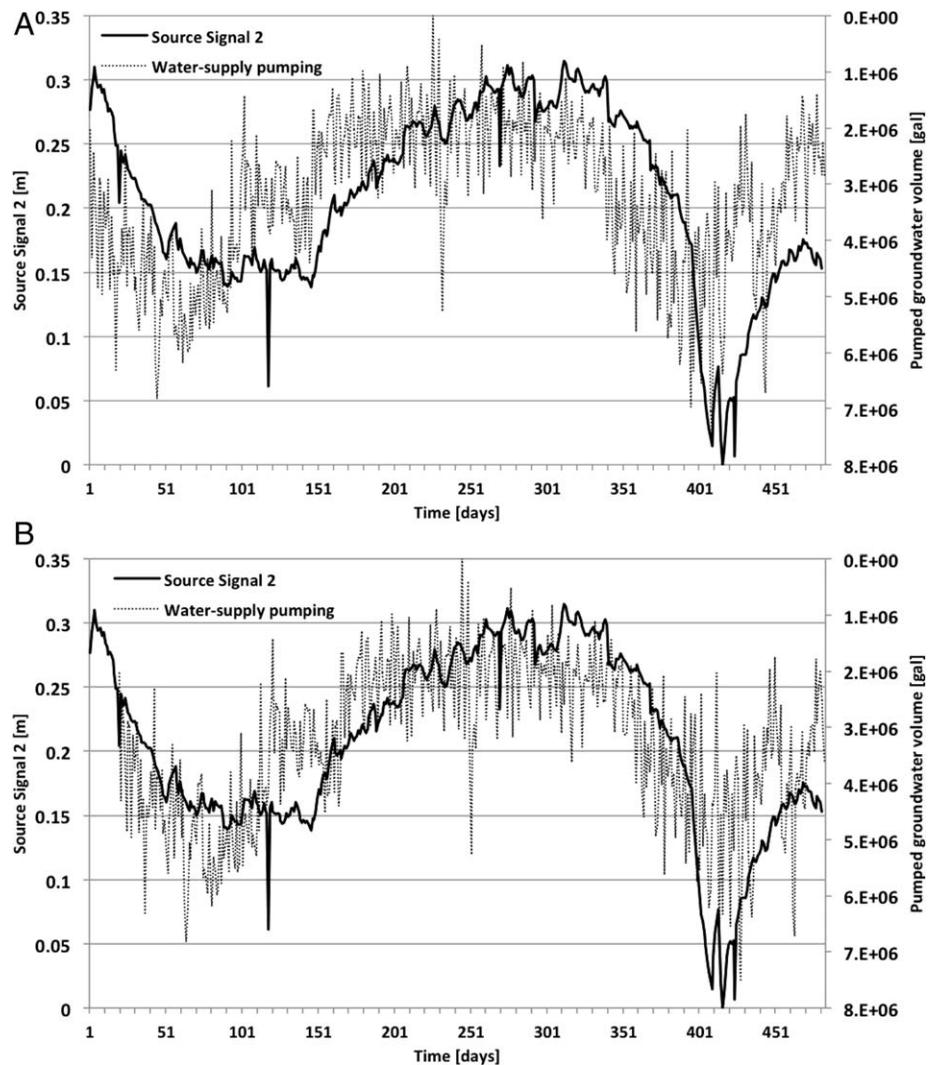


**Figure 9.** Comparison between the total groundwater volume (gal) pumped by the nearby six water-supply wells (note that the right axis is reversed) and the second reconstructed source signal: (a) shows the original data; and (b) shows the pumping data shifted by 20 days forward in time.

*and Vesselinov*, 2011] were faster requiring a computational time on the order of 1 h. However, *Harp and Vesselinov* [2011] did not estimate the barometric pressure effects. As already discussed, the barometric efficiencies were assumed to be 100%.

The analyses performed using NMF*k* and physics-based inversion have their pros and cons. The physics-based inversion can be faster if substantial simplifications were made about the governing processes and their properties. For example, *Harp and Vesselinov* [2011] applied the Theis equation which can be considered a substantial simplification of the aquifer conditions. They also did not account for complicated propagation of the barometric pressures through the vadose zone (as suggested by the analyses presented above). The NMF*k* analyses do not make any assumptions about the physical process. In our case, this allowed us to make conclusions about the complexity in the propagation of (1) the barometric pressures through the vadose zone (deviations from 100% barometric efficiency) and (2) pumping drawdowns through the regional aquifer (relatively uniform pumping drawdowns in all the analyzed monitoring wells). These conclusions can be utilized in the future to develop of physics-based site models.

## 5. Conclusions

Our analyses demonstrate the applicability of a new inverse method for analysis of pressure transients based on a Nonnegative Matrix Factorization (NMF) algorithm applied to Blind Source Separation (BSS). The algorithm is called NMF*k* and is based on previous work by [*Alexandrov et al.*, 2013b]. The unknown sources are identified from a set of mixed signals observed at monitoring wells without any information about (1) the sources, (2) their location, (3) their transients, (4) the physical processes impacting the signal propagation through the subsurface, and (5) the properties associated with these process. In presented analyses, we extract (deconstruct) the source signals as observed at the monitoring wells. The only assumption that we make is that the source signals are proportionally manifested at each observation well. The proportionality factor implicitly depends on the processes that occur in the flow medium and their properties. The NMF*k* algorithm allows for decoupling (deconstruction) of the source (forcing) signals as recorded at the observation points (wells) without models and model assumptions. The solved inverse problem is underdetermined (ill-posed). To address this, the NMF*k* algorithm thoroughly explores the plausible inverse solutions, and seeks to narrow the set of possible solutions by estimating the optimal number of source signals needed to robustly and accurately characterize the observed data. The NMF*k* results identify two unique sources causing the observed transients: these are barometric pressure and water-supply pumping effects. Other potential sources such as infiltration recharge and surface-water stages were not identified. The detected sources appear to be proportionally manifested at the observation wells. This allows us to estimate the barometric pressure efficiencies. Future work will include (1) probabilistic analyses of the uncertainty associated with multiple acceptable solutions identified by NMF*k*, (2) estimation of the source locations and spatial properties of the flow medium, and (3) coupling of the NMF*k* algorithm in physics-based inverse analyses. The possible applications of the developed NMF*k* algorithm are not limited to hydrogeology problems; NMF*k* can be applied to any real problem where temporal system behavior is observed at multiple locations and mixing of an unknown number of physical sources are causing these fluctuations.

## References

Alexandrov, L. B., et al. (2013a), Signatures of mutational processes in human cancer, *Nature*, *500*, 415–421.
Alexandrov, L. B., S. Nik-Zainal, D. C. Wedge, P. J. Campbell, and M. R. Stratton (2013b), Deciphering signatures of mutational processes operative in human cancer, *Cell Rep.*, *3*(1), 246–259.
Amari, S.-I., A. Cichocki, H. H. Yang (1996), A new learning algorithm for blind signal separation, in *Advances in Neural Information Processing Systems*, edited by D. Touretzky, M. Mozer and M. Hasselmo, *8*, pp. 757–763, MIT Press, Cambridge, Mass.
Berg, S. J., and W. A. Illman (2011a), Three-dimensional transient hydraulic tomography in a highly heterogeneous glaciofluvial aquifer-aquitard system, *Water Resour. Res.*, *47*, W10507, doi:10.1029/2011WR010616.
Berg, S. J., and W. A. Illman (2011b), Capturing aquifer heterogeneity: Comparison of approaches through controlled sandbox experiments, *Water Resour. Res.*, *47*, W09514, doi:10.1029/2011WR010429.
Broxton, D. E., and D. T. Vaniman (2005), Geologic framework of a groundwater system on the margin of a rift basin, pajarito plateau, north-central New Mexico, *Vadose Zone J.*, *4*(3), 522–550.
Cardiff, M., W. Barrash, P. Kitanidis, B. Malama, A. Revil, S. Straface, and E. Rizzo (2009), A potential-based inversion of unconfined steady-state hydraulic tomography, *Ground Water*, *47*(2), 259–270.
Carrera, J., A. Alcolea, A. Medina, J. Hidalgo, and L. J. Slooten (2005), Inverse problem in hydrogeology, *Hydrogeol. J.*, *13*(1), 206–222.
Cichocki, A., R. Zdunek, A. H. Phan, and S.-I. Amari (2009), *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*, 500 pp., Wiley, Hoboken, N. J.
Fischler, M. A., and R. A. Elschlager (1973), The representation and matching of pictorial structures, *IEEE Trans. Comput.*, *100*(1), 67–92.

Forootan, E., and J. Kusche (2012), Separation of global time-variable gravity signals into maximally independent components, *J. Geod.*, *86*(7), 477–497.

Forootan, E., and J. Kusche (2013), Separation of deterministic signals using independent component analysis (ICA), *Stud. Geophys. Geod.*, *57*(1), 17–26.

Frappart, F., G. Ramillien, P. Maisongrande, and M.-P. Bonnet (2010), Denoising satellite gravity signals by independent component analysis, *IEEE Geosci. Remote Sens. Lett.*, *7*(3), 421–425.

Frappart, F., G. Ramillien, M. Leblanc, S. O. Tweed, M.-P. Bonnet, and P. Maisongrande (2011), An independent component analysis filtering approach for estimating continental hydrology in the grace gravity data, *Remote Sens. Environ.*, *115*(1), 187–204.

Golub, G. H., and C. F. Van Loan (1980), An analysis of the total least squares problem, *SIAM J. Numer. Anal.*, *17*(6), 883–893.

Halford, K. J., C. A. Garcia, J. Fenelon, and B. Mirus (2012), Advanced methods for modeling water-levels and estimating drawdowns with seriessee, an excel add-in, *U.S. Geol. Surv. Tech. Meth. Rep.*, *4–F4*, 29 pp.

Harp, D. R., and V. V. Vesselinov (2011), Identification of pumping influences in long-term water level fluctuations, *Ground Water*, *49*(3), 403–414.

Hartigan, J. A., and M. A. Wong (1979), Algorithm as 136: A k-means clustering algorithm, *J. R. Stat. Soc. Ser. C*, *28*(1), 100–108.

Herault, J., and C. Jutten (1986), Space or time adaptive signal processing by neural network models, Neural networks for computing, *151*, 206–211, AIP Publishing, doi:10.1063/1.36258.

Hyvärinen, A., and E. Oja (2000), Independent component analysis: algorithms and applications, *Neural Netw.*, *13*(4), 411–430.

Illman, W. A., X. Liu, S. Takeuchi, T.-C. J. Yeh, K. Ando, and H. Saegusa (2009), Hydraulic tomography in fractured granite: Mizunami underground research site, Japan, *Water Resour. Res.*, *45*, W01406, doi:10.1029/2007WR006715.

Izenman, A. J. (2008), *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning*, pp. 731, Springer.

Jutten, C., and J. Herault (1991), Blind separation of sources, Part I: An adaptive algorithm based on neuromimetic architecture, *Signal Process.*, *24*(1), 1–10.

Keylock, C. (2012), A resampling method for generating synthetic hydrological time series with preservation of cross-correlative structure and higher-order properties, *Water Resour. Res.*, *48*, W01406, doi:10.1029/2007WR006715.

Koch, R., and S. Schmeer (2008), Groundwater level status report for 2008, *Tech. Rep. LA-14397-PR*, Los Alamos National Laboratory, Los Alamos, N. M.

Lee, D. D., and H. S. Seung (1997), Unsupervised learning by convex and conic coding, in *Advances in Neural Information Processing Systems*, pp. 515–521.

Lee, D. D., and H. S. Seung (1999), Learning the parts of objects by non-negative matrix factorization, *Nature*, *401*(6755), 788–791, Nature Publishing Group.

McDermott, J. H. (2009), The cocktail party problem, *Curr. Biol.*, *19*(22), R1024–R1027.

Ni, C.-F., and T.-C. J. Yeh (2008), Stochastic inversion of pneumatic cross-hole tests and barometric pressure fluctuations in heterogeneous unsaturated formations, *Adv. Water Resour.*, *31*(12), 1708–1718.

Nuzhny, A., E. Saveleva, S. Kazakov, and S. Utkin (2008), Virtual sources for spatio-temporal monitoring in data analysis, paper presented at iEMSs Fourth Biennial Meeting: International Congress on Environmental Modeling and Software, (iEMSs 2008), edited by M. Sànchez-Marrè, J. Bèjar, J. Comas, A. E. Rizzoli, and G. Guariso, vol. 3, pp. 1734–1741, Barcelona, Catalonia.

Paatero, P., and U. Tapper (1994), Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values, *Environmetrics*, *5*(2), 111–126.

Pang-Ning, T., M. Steinbach, V. Kumar (2006), *Introduction to Data Mining*, 769 pp., Addison-Wesley, Boston, Mass.

Parra, L. C., C. Spence, P. Sajda, A. Ziehe, and K.-R. Müller (1999), Unmixing hyperspectral data, in *Advances in Neural Information Processing Systems*, *12*, pp. 942–948, MIT Press.

Rasmussen, T. C., and L. A. Crawford (1997), Identifying and removing barometric pressure effects in confined and unconfined aquifers, *Ground Water*, *35*(3), 502–511.

Rousseeuw, P. J. (1987), Silhouettes: A graphical aid to the interpretation and validation of cluster analysis, *J. Comput. Appl. Math.*, *20*, 53–65.

Seung, D., and L. Lee (2001), Algorithms for non-negative matrix factorization, in *Advances in Neural Information Processing Systems*, *13*, 556–562, MIT Press.

Srivastava, S., S. Niebling, B. Küstner, P. Wich, C. Schmuck, and S. Schlücker (2008), Characterization of guanidiniocarbonyl pyrroles in water by ph-dependent uv raman spectroscopy and component analysis, *Phys. Chem. Chem. Phys.*, *10*(45), 6770–6775.

Tapley, B. D., S. Bettadpur, J. C. Ries, P. F. Thompson, and M. M. Watkins (2004), Grace measurements of mass variability in the earth system, *Science*, *305*(5683), 503–505.

Theis, C. V. (1935), The relation between the lowering of the piezometric surface and the rate and duration of discharge of a well using groundwater storage, *Trans. AGU*, *16(2)*, 519–524, doi:10.1029/TR016i002p00519.

Toll, N. J., and T. C. Rasmussen (2007), Removal of barometric pressure effects and earth tides from observed water levels, *Ground Water*, *45*(1), 101–105.

Vasco, D., H. Keers, and K. Karasaki (2000), Estimation of reservoir properties using transient pressure data: An asymptotic approach, *Water Resour. Res.*, *36*(12), 3447–3465.

Vesselinov, V. V. (2004), An alternative conceptual model of groundwater flow and transport in saturated zone beneath the pajarito plateau, *Tech. Rep. LA-UR-05–6741*, Los Alamos National Laboratory, Los Alamos, N. M.

Vesselinov, V. V., S. P. Neuman, and W. A. Illman (2001a), Three-dimensional numerical inversion of pneumatic cross-hole tests in unsaturated fractured tuff: 1. Methodology and borehole effects, *Water Resour. Res.*, *37*(12), 3001–3017.

Vesselinov, V. V., S. P. Neuman, and W. A. Illman (2001b), Three-dimensional numerical inversion of pneumatic cross-hole tests in unsaturated fractured tuff: 2. Equivalent parameters, high-resolution stochastic imaging and scale effects, *Water Resour. Res.*, *37*(12), 3019–3041.

Westra, S., C. Brown, U. Lall, and A. Sharma (2007), Modeling multivariable hydrological series: Principal component analysis or independent component analysis?, *Water Resour. Res.*, *43*, W06429, doi:10.1029/2006WR005617.

Xu, W., X. Liu, and Y. Gong (2003), Document clustering based on non-negative matrix factorization. In Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval, pp. 267–273, ACM, Toronto, Canada, doi:10.1145/860435.860485.

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

Zarzoso, V., and A. K. Nandi (1999), Blind source separation, in *Blind Estimation Using Higher-Order Statistics*, edited by A. K. Nandi, pp. 167–252, Kluwer Academic Publishers.