# ResearchOnline@JCU

JAMES COOK
UNIVERSITY
AUSTRALIA

# Optimal Monitoring Network Design and Identification of Unknown Pollutant Sources in Polluted Aquifers

Thesis submitted by

**Om Prakash, B.Sc. Engg, M.Tech**



for the degree of Doctor of Philosophy (PhD)

in the School of Engineering and Physical Sciences

James Cook University

April 2014

# STATEMENT OF ACCESS

I, the undersigned, author of this work, understand that James Cook University will make this thesis available for use within the University Library and, via the Australian Digital Theses network, for use elsewhere.

I understand that, as an unpublished work, a thesis has significant protection under the Copyright Act and I wish this work to be embargoed until April 2015, after which date I do not wish to place any further restriction on access to this work.

Signature                                                                                          Date

# STATEMENT OF SOURCES DECLARATION

I herewith declare that I have produced this thesis without the prohibited assistance of third parties and without making use of aids other than those specified; notions taken over directly or indirectly from other sources have been identified as such. This thesis has not previously been presented in identical or similar form to any other Australian or foreign examination board.

Signature                                                                                    Date

# <u>ELECTRONIC COPY</u>

I, the undersigned, the author of this work, declare that the electronic copy of this thesis provided to James Cook University Library is an accurate copy of the print thesis submitted.

Signature                                                                 Date

# STATEMENT OF CONTRIBUTION OF OTHERS

Financial contribution towards this PhD project was received from:

Apart from the financial assistance, the following have contributed to this PhD project as specified hereunder:

Dr. Bithin Datta: Dr. Datta supervised the entire PhD project and helped conceptualise the problem, suggested several ways of solving it and provided insights on the tools and technologies used in this thesis.

# **Acknowledgements**

# Abstract

Increasing stress from various anthropogenic activities has resulted in widespread pollution of groundwater resources. Often, when the pollutant is first detected in groundwater, little is known about the pollutant sources. Identification of source characteristics in terms of locations, activity initiation times, and source flux release histories and activity durations are vital in planning effective remediation measures and determining the liability of the polluter. Groundwater pollution source characterization is an inverse and ill-posed problem. Finding a solution to this inverse problem remains a challenging task due to uncertainties in accurately predicting the aquifer response to source flux injection, generally encountered sparsity of concentration measurements in the field, and the non-uniqueness in the aquifer response to the subjected hydraulic and chemical stresses. This study presents linked simulation-optimization, and sequential monitoring network design based methodologies for identification of unknown groundwater pollution source characteristics.

Pollution in groundwater aquifers is generally first detected in an arbitrarily located water supply well or a group of wells. Often pollutants are detected much after activity at the sources may have initiated, or even after it has ceased to exist. There may be a gap of years, or even decades, between the start of source activity and detection of pollutants in an aquifer. Other important issues in accurately identifying unknown groundwater pollution source characteristics are the quality, usability and extent of pollution measurement data from the study area. Existing methodologies for unknown groundwater pollution source characterization have several limitations. Methodologies developed in

this study aim to address some of these limitations. The major limitations addressed in this study include:

i.    sparsity of pollutant concentration measurement data,

ii.   inefficient monitoring network for concentration measurements,

iii.  difficulty in identifying the source locations,

iv.   difficulty in establishing the pollutant source activity initiation time,

v.    applicability of optimal source characterization with missing observation data.

In many cases of aquifer pollution, especially in clandestine underground disposal of toxic wastes, no information is available about the number and location of such sources. Moreover, monitoring wells where pollution is first detected may not be optimally located for accurately identifying the release history of unknown pollution sources. A large number of pollutant concentration measurements spread over time and space is necessary for accurate source identification. However, long term monitoring over a large number of monitoring locations has budgetary constraints. This study presents a sequential optimal monitoring network design methodology based on geostatistical kriging, a pollutant concentration gradient based search for identification of source locations, and a Genetic Programming (GP) based optimal monitoring network design model for collecting concentration measurements for efficient source characterization.

To address the issue of unknown starting times of activity of the sources, a new methodology is developed for simultaneously identifying the starting times of the activity of the sources and their flux release history. A new optimum decision model is formulated and solved such that the starting times of the activity of the sources are directly obtained as solution. Simulated Annealing (SA) is used for solving the

optimization problem with the starting time of pollutant source activity incorporated as explicit decision variable.

Subsequent to the detection of pollution in an aquifer, a more formal methodology for source characterization is generally initiated only after large numbers of spatiotemporal concentration measurements, spaced over a sufficiently long period of time, are obtained. During this time, the spread of the pollutant continues while temporal measurements are being obtained at monitoring locations. A feedback-based sequential methodology for efficient identification of unknown pollutant source characteristics, integrating optimal monitoring network design and an optimization based source identification model, is developed. The main advantage of this methodology is that source characterization can start at the same time as when pollutant is first detected in the aquifer. In every sequence, feedback from the source identification model improves the optimal monitoring network design and vice-versa. This results in efficient and accurate source characterization, within a few sequences of source identification and monitoring network design.

The performances of the developed methodologies are evaluated for different scenarios of groundwater pollution incorporating transient flow and advective-dispersive transport in heterogeneous anisotropic conditions. The applicability of the developed methodologies is tested for a real aquifer site polluted with petrochemical waste (BTEX). These evaluation results demonstrate the potential applicability of the developed methodologies to correctly estimate the unknown source flux's magnitude, and location and source activity initiation times, while improving the accuracy of source flux identification. Results of performance evaluation of each of these methodologies indicate their potential for field application.

# Table of Contents

# List of Figures

# List of Tables

# 1. Introduction

Increasing stress from various anthropogenic activities has resulted in widespread pollution of groundwater resources. Often, when a pollutant is first detected in groundwater, little is known about the pollutant source, or the extent of pollution in the aquifer. This study presents linked simulation-optimization, and sequential monitoring network design based methodologies for identification of unknown groundwater pollution source characteristics. The optimization models are solved using a Simulated Annealing (SA) optimization algorithm. Genetic Programming is also utilized to model the relevance of potential pollutant sources to the pollutant concentration measurement at potential monitoring locations. These impact factors form the basis of optimal monitoring network design for unknown groundwater pollution source characterization. The



**Figure 1.1 Total Water and Fresh Water Distribution**

performances of the developed methodologies are evaluated for different scenarios of groundwater pollution.

Sustenance of all forms of life depends on water. Of all the water available on earth, only 2.5 percent is deemed as fresh water. Figure 1.1 shows the percentage distribution of total water and freshwater on the planet (Shiklomanov, 1993). As the world population increases, so does the demand for fresh water. As per the recent estimates published in United Nations World Water Development Report 3 (WWDR-3, 2009), global water consumption has tripled over the past 50 years to meet the demands of a growing population. Groundwater constitutes a significant part of the water resources required to meet the global water demand. As per the recent estimates for the year 2010, published in the United Nations World Water Development Report 4 (WWDR-4, 2012), 26 percent of the total water demand was met by groundwater resources. The global groundwater abstraction rate has at least tripled over the past 50 years and continues to increase at an annual rate of 1 to 2 percent.

Increasing stress from various anthropogenic factors has led to depletion and degradation of groundwater resources (WWDR-4, 2012). The biggest challenge to sustainability of groundwater resources is deterioration in groundwater quality. Groundwater resources are polluted mainly because of pollutants generated by increasing industrial activities, and due to use of chemicals in agriculture. Leaking underground pipelines or tanks, waste water lagoons, accidental spills, landfill leachate and improper disposal of chemical waste are among the common causes of groundwater pollution.

Due to the vulnerability of these groundwater resources to pollution, it is essential to develop efficient techniques for prevention, detection and remediation of groundwater

pollution. Pollutants in groundwater can stay undetected for significant periods of time. In order to develop effective and economical strategies for control and remediation of polluted groundwater aquifers, it is necessary to identify sources of pollution, determine the extent of pollution, and predict future pollution scenarios. The first step in groundwater aquifer remediation should be the detection and characterization of unknown pollutant sources. The effectiveness of remediation strategies will depend on the accuracy and reliability of information regarding the sources of pollution in terms of their location, release history and time of initial activity.

## 1.1. Unknown Groundwater Pollutant Source Characterisation

Sources of groundwater pollution can be characterized by:

i.     number of pollutant sources present,

ii.    type of sources (point or distributed),

iii.   spatial location and extent of the sources,

iv.    activity initiation time of the sources, and

v.     source flux release history as a function of time since the initiation of source activities

The process of flow and transport of a conservative pollutant in a three dimensional aquifer system is modelled using mathematical equations, referred to as the Advection Dispersion Equation (ADE). This equation can be solved using numerical techniques, especially when two-dimensional or three-dimensional processes are considered. When the above-mentioned groundwater source characteristics are precisely known and the aquifer parameters governing flow and transport, such as hydraulic conductivity, or porosity, for example are measured accurately and adequately, numerical simulation

models using ADE can be solved to predict the fate and transport of the pollutants in groundwater aquifers with respect to space and time. Hence, identification of source characteristics is helpful in understanding the fate and transport of the pollutants in groundwater aquifers, and in developing a remediation strategy accordingly.

In real world scenarios, there is seldom any information about the source characteristics. In such cases, the numerical simulation model can be solved backwards in time and space to find these unknown source characteristics, using pollutant concentration measurements from various locations in the polluted aquifer area. Theoretically, this is possible but solving inverse problems has its own limitations. Unlike forward modelling, inverse modelling does not have a unique solution, as different combinations of source characteristics may result in similar pollutant plume. Groundwater systems are generally non-linear, and inverse method of solution may lack in stability.

Characterization of unknown groundwater pollution sources is a challenging task. A large number of pollutant concentration measurements spread over space and time is necessary for reliable characterization of unknown groundwater pollution sources. However, in reality, most of the time only limited amounts of relevant data are initially available to understand the ongoing pollution in the aquifer. Even this pollutant concentration data has inherent measurement errors. Aquifer parameters governing flow and transport, (such as hydraulic conductivity, or porosity, for example) and boundary conditions are only average estimates and contain varying degrees of uncertainty.

## 1.2.    Linked Simulation-Optimization Approach

One way to solve the unknown pollutant source characterization problem is to simulate the physical process of flow and transport in the groundwater system for different combinations of source characteristics and try to match the results with the observed data. Since the number of possibilities with different combinations of source characteristics is almost infinite, such a comparison without the use of an optimum decision model is exhaustive and computationally infeasible.

An alternative is the use of numerical simulation of flow and transport processes in conjunction with an optimum decision model to solve this problem. The solution of these numerical simulation models represents the effect of pollutant sources on groundwater quality in space and time. This method is known as the "linked simulation-optimization approach". Earlier implementations of this approach used linear programming and response matrix along with forward simulations. However, with the advent of evolutionary algorithms such as Genetic Algorithm (GA) and Simulated Annealing (SA), source characterization using linked simulation-optimization has become computationally efficient and feasible.

## 1.3.    Sequential Monitoring Network Dedicated to Source Characterization

Monitoring wells are pivotal in understanding groundwater hydraulics and aquifer pollution. Pollution in groundwater aquifers is monitored by installing sets of monitoring wells at different locations in a polluted aquifer area often termed as a "monitoring network". Monitoring networks are installed with different underlying objectives. However, long term monitoring over multiple monitoring locations is expensive and often

governed by budgetary constraints. Also, the initially observed pollutant concentration data is often sparse in space and in time. Earlier attempts at the design of optimal monitoring networks were aimed at early detection of pollutants, redundancy reduction, compliance monitoring and so on, using different optimization based decision models. An overview of these optimal monitoring network design models can be found in Amirabdollian and Datta (2013).

Observed pollutant concentration measurements from monitoring networks are used for characterization of unknown groundwater pollution sources using the linked simulation-optimization approach. The source characterization problem is solved by minimizing the difference between the observed and simulated concentration measurements at the monitoring locations. However, not all monitoring locations are ideally suited for source characterization using the linked simulation-optimization approach. Hence, an optimally designed monitoring network is necessary to improve the accuracy of source characterization, using concentration measurements from such a network.

The effect of pollution sources varies in time and in space in a polluted aquifer. Any monitoring network, with an underlying objective, may seem to be suitable at one time but prove to be redundant at other times. Hence, it is important that the monitoring networks are designed sequentially in time to adapt to these changes.

Given the sparsity of pollutant concentration measurement data, it is imperative that the monitoring network be designed such that it improves the accuracy of source identification. Sequential optimal monitoring networks can provide vital feedback information in identifying unknown source locations, and in efficient characterization of unknown groundwater pollution sources.

## 1.4. Research Objectives

Existing methodologies for unknown groundwater pollution source characterization have several limitations. Some of the main limitations are:

   i.    sparsity of pollutant concentration measurement data,

   ii.   inefficient monitoring networks for concentration measurements,

   iii.  difficulty in identifying source locations,

   iv.   difficulty in establishing pollutant source activity initiation time,

   v.    applicability of optimal source characterization with missing observation data.

This study is aimed at developing optimization based methodologies for characterization of unknown groundwater sources. In the absence of any information about source locations and source activity initiation times, source characterization is even more difficult. Using limited amounts of observed pollutant concentration data from randomly located monitoring wells further increases the difficulty. Real life scenarios may have multiple sources, each having different activity initiation times, and missing measurements from the observation data.

This study presents a methodology for optimal sequential monitoring network design for identification of unknown pollution source locations. Also, a methodology for dedicated monitoring network design to improve the accuracy of source characterization is presented in this study. A new linked simulation-optimization method for source characterization is developed to overcome the limitations posed due to difficulty in estimating the initial activity initiation times. Finally, a feedback based optimal sequential monitoring network design model is combined with a linked simulation-optimization model for efficient source characterization.

The specific objectives of this study are to:

i.   develop a sequential monitoring network design methodology for identification of unknown pollution source locations;

ii.  develop a dedicated monitoring network design methodology for improving the accuracy of source characterization in a linked simulation optimization approach;

iii. develop a linked simulation-optimization methodology for simultaneous optimal identification of unknown groundwater pollution source fluxes and source activity initiation times, considering source activity initiation times as explicit decision variables;

iv.  develop feedback based methodology for efficient identification of unknown pollutant source characteristics linking sequential monitoring networks with a linked simulation-optimization based source characterization model;

v.   extend the developed methodologies to incorporate measurement errors;

vi.  evaluate the performances of the developed methodologies using illustrative example sites; and

vii. evaluate the applicability of developed methodologies for a real contaminated aquifer site in New South Wales polluted with BTEX.

## 1.5.  Organization of the Thesis

This thesis comprises of seven chapters including the introduction. Chapter 2 of the thesis discusses the state-of-art of various techniques used in this study.

Chapter 3 presents a linked simulation-optimization methodology for simultaneous estimation of unknown pollutant source characteristics in terms of their location, magnitude, duration and unknown activity initiation time. Simulated annealing is used as

the optimization algorithm. Performance of the methodology is tested for different complex scenarios, consisting of multiple active sources and no information about the source activity initiation times. The proposed methodology is extended to include complex scenarios with multiple sources, each having different activity initiation times, missing observation data and unsteady state flow conditions.

Chapter 4 presents a sequential optimal monitoring network design model for identification of source locations using concentration gradient information, and data interpolation using kriging. A separate dedicated monitoring network design methodology is presented for efficient characterization of unknown pollutant sources using monitoring impact factors. Genetic programming based impact factors and frequency factors are used for assessing the impact of a potential source on a potential monitoring location. Efficiency of the developed methodology is demonstrated by comparing the performance of the source identification model, when utilizing concentration measurements obtained from the designed network, with those obtained using measurements from arbitrary monitoring networks.

Chapter 5 presents a feedback based methodology for efficient identification of unknown pollutant source characteristics, integrating sequential monitoring network design with a source identification model. The linked simulation-optimization model is solved using observation data from a sequentially designed optimal monitoring network, such that the feedback information in the nature of measured concentration data from the designed network improves the source characterization results. A gradient based search technique is used for the optimal monitoring network design. Performance of this method is evaluated using synthetic aquifer data.

Chapter 6 presents an application of the methodologies developed in Chapter 3, Chapter 4 and Chapter 5 to a real life managed polluted groundwater aquifer site in New South Wales, Australia.

Chapter 7 presents a brief discussion of salient features of this study and conclusions. Some of the limitations of the methodologies developed are also mentioned.

# 2. Review of Literature

This chapter provides a brief overview of the literature relevant to solving the unknown groundwater pollutant source characterisation problem. Variations of the source characterization problem, in terms of the different unknown source characteristics that were addressed and the solution methodologies that were adopted, are discussed. Various approaches for monitoring network design for solving the unknown groundwater pollutant source characterisation problem are also presented. An overview of the tools and techniques utilized for simulation, optimization and data interpolation, are discussed in this chapter.

## 2.1. Unknown Groundwater Pollution Source Characterisation

Pollution in groundwater can remain undetected for long periods of time. Pollution in groundwater aquifers is generally first detected by an arbitrarily located water supply well or a group of wells. Many times pollutants are detected much after the activity at the sources may have started, or may even have ceased to exist. Identification of pollutant source characteristics in terms of pollutant flux magnitude, location, and activity duration from sparse pollutant concentration measurements, belongs to a category of inverse problems which are generally ill-posed (Yeh, 1986). It is like finding the cause from limited knowledge of the effects. The ill-posedness is characterized by non-uniqueness and instability in the solution, as different combinations of source characteristics can result in similar pollutant plume, and small errors in input can result in large estimation errors.

Various methodologies for solving this ill-posed inverse problem have been suggested. These methods can be broadly classified as follows:

i.   heat transport inversion based approach

ii.  direct approach

iii. analytical and regression based approach

iv.  probabilistic and geostatistical simulation approach

v.   optimization based approach

A detailed review of these methodologies can be found in Atmadja and Bagtzoglou (2001b); Bagtzoglou and Atmadja (2005); Michalak and Kitanidis (2004a, b); Neupauer et al., (2000); and Sun et al. (2006a, b).

The governing equation for heat conduction is similar to the transport equation, discounting the advection phenomenon. Owing to the similarity between heat transfer and groundwater flow and transport, a solution approach based on inversion of heat transfer (Cannon, 1966), is often applied to solve the source characterization problem. However, such approaches require that the parameters used in mathematical models are precisely known. This may be possible in the case of heat transfer as the medium is mostly homogeneous and the parameter values can be accurately measured. The same does not hold true in groundwater systems as they are highly heterogeneous and flow parameters such as hydraulic conductivity, and porosity, are average approximations for the entire study area. Hence, the use of heat transport inversion methods has been limited.

A more direct approach (Skaggs & Kabala, 1994, 1995; Liu & Ball, 1999) used a technique called Tikhonov regularization to transform the ill-posed algebraically indeterminate inverse problem into a minimization problem to get a unique solution.

However, even small errors in the input data significantly reduced the accuracy of the solution. Other direct approaches, namely Marching-Jury Backward Beam Equation (MJBBE) (Atmadja & Bagtzoglou 2000; 2001a; 2003), Minimum Relative Entropy (MRE) inversion (Woodbury & Ulrych, 1996; Woodbury et al. 1998), particle-based censored random walk (Ababou 2010), optimization based inverse method (Bagtzoglou & Baun 2005), and reversed time particle tracking (Bagtzoglou 2003), have been applied for solving the source characterization problem.

An analytical and regression based approach has been applied in characterization of unknown groundwater pollution sources. Some of the significant studies using these methods are Sidauruk et al. (1998); Alapati and Kabala (2000); and Ala and Domenico (1992). However, these analytical solutions work only for a limited number of cases where the aquifer is homogeneous with simple geometry and flow conditions.

Bagtzoglou et al. (1991, 1992) presented a probabilistic frame work for solving the source characterization problem. They used random walk method to solve the transport ADE equation backwards in time. Neupauer and Wilson (1999, 2005) used an adjoint method to find backward-in-time location and travel time probabilities. Snodgrass and Kitanidis (1997) also used a probabilistic approach combining Bayesian theory and geostatistical techniques. Though these methods could handle heterogeneity in the aquifer system, they required extensive computation.

One of the early methods of source characterization was to run forward simulations for different combinations of source characteristics and try to match the results with the observed effect. Due to the non-uniqueness of the solution and the infinite number of plausible combinations, optimization based methods were evolved for finding the best

possible combination of source characteristics. Various techniques involving different optimization algorithms have been evolved to solve the source characterization problem.

First attempts in this regard were made by Gorelick et al. (1983). They formulated the source characterization problem as forward-time simulations coupled with an optimization model using linear programming and multiple regressions. Solute transport model was treated as a constraint and presented as a concentration response matrix in this formulation.

Datta et al. (1989) developed an expert-system embedding pattern-recognition technique for pollution-source characterization. The pattern recognition technique was based on stochastic dynamic programming, which minimized the loss due to recognition error. This pattern recognition model was utilized as a screening model for optimization based source characterization model.

Wagner (1992) developed an optimization based methodology for simultaneous estimation of model parameters along with source characterization. In his attempt, Wagner used an inverse model as a non-linear maximum likelihood estimation problem. Estimates of hydrogeological and source parameters were based on measurements of hydraulic head and pollutant concentration.

Mahar and Datta (1997; 2001) were the first to combine optimal source characterization with the design of a groundwater quality monitoring network to improve the efficiency of the source characterization process. An embedding technique was used in which the simulation model is embedded as binding constraints in the optimization model. They applied their method to a hypothetical 2-D homogeneous, isotropic and saturated aquifer

with a conservative pollutant plume in a two-step process. In the first step, optimization model was used to identify an unknown pollution source based on observation data. In the second step, different realizations of pollutant plumes were simulated using perturbed sources. These realizations of the pollutant plumes were used in an integer programming to determine the optimal locations of the monitoring wells. Pollutant concentration measurements from these monitoring well locations are used in the non-linear optimization model to obtain a more accurate estimation of sources. Mahar and Datta (2000) were also able to estimate the magnitude, location and duration of pollutant sources using a non-linear optimization technique. Datta et al. (2009a) used an optimal dynamic monitoring network design for identification of unknown groundwater pollution sources.

Aral et al. (2001) formulated a pollutant source characterization problem as a nonlinear optimization model, in which pollutant source locations and release histories are defined as explicit unknown variables. The optimization model minimized the residual error between observed and simulated pollutant concentrations at the observed locations. Simulated concentrations were implicitly embedded in the formulation through the simulation models. As repeated solution of these models is a computationally intensive but a necessary feature of the optimization process, a progressive genetic algorithm (PGA) was applied to solve the optimization problem.

Singh et al. (2004) and Singh and Datta (2007) used a trained multilayer, feed-forward Artificial Neural Network (ANN) for simultaneous estimation of unknown pollution source characteristics and hydrogeological parameters. Universal function approximation capability of the ANN was utilized to estimate the source characteristics and flow and transport parameters. ANN was trained on data patterns which consisted of a set of source

fluxes and corresponding temporally varying simulated concentration measurements. The methodology was evaluated for varying degrees of concentration measurement error.

Mahinthakumar and Sayeed (2005) investigated and compared several hybrid optimization approaches that combine genetic algorithms with a number of local search approaches for reconstructing the release histories of the pollutant sources. The methodology was evaluated for a three dimensional, heterogeneous flow field considering multiple sources. The results indicate that hybrid optimization methods, which combine an initial global heuristic approach (for example, genetic algorithms) with a gradient-based local search approach (for example, conjugate gradients), are very effective in estimating the flux release history.

Singh and Datta (2006) used a simulation optimization approach for characterization of unknown groundwater pollution sources. The performance of the developed methodology was evaluated for combinations of source characteristics, data availability conditions, and concentration measurement error levels. The main advantage of this method was that the numerical simulation model could be externally linked to the optimization model. This solution approach enables to solve source characterization problems for complex aquifer study areas with multiple pollution sources.

Yeh et al. (2006) proposed an approach, SATS-GWT, which combines Simulated Annealing (SA), Tabu Search (TS) and the three-dimensional groundwater flow and solute transport model (MODFLOW-GWT), to estimate the source information: source location, release concentration and release period. The source location is selected by TS within the suspected source area. SA is used to optimally estimate the release concentration, and release period. Search for the optimal estimate of these unknown

source characteristics is terminated based on the best objective function value. The performance of this method was evaluated for homogeneous and heterogeneous aquifer study areas with transient flow conditions.

He et al. (2009) presented a coupled simulation–optimization approach for optimal design of petroleum contaminated groundwater remediation under uncertainty. It had the following advantages: (1) it addressed the stochasticity of the modelling parameters in simulating the flow and transport of NAPLs in groundwater, (2) it provided a direct and response-rapid bridge between remediation strategies (pumping rates) and remediation performance (pollutant concentrations) through the created proxy models, (3) alleviated the computational cost in searching for optimal solutions, and (4) it gave confidence levels for the obtained optimal remediation strategies.

Datta et al. (2009b) developed a methodology for simultaneous source identification and parameter estimation in groundwater systems in which the simulation model is externally linked to a nonlinear optimization model. The simulator defines the flow and transport processes, and serves as a binding equality constraint. The search direction is determined by the Jacobian matrix in the nonlinear optimization model, linking the groundwater flow-transport simulator and the optimization method. This addresses the limitation of embedding the discretised flow and transport equations as equality constraints in the optimization.

Ayvaz (2010) developed a linked simulation–optimization model in which the locations and release histories of the pollution sources are treated as the explicit decision variables. MODFLOW and MT3DMS packages are used to simulate the flow and transport

processes in the groundwater system. These models are then integrated with an optimization model which is based on the heuristic Harmony Search (HS) algorithm.

Datta et al. (2011) developed a methodology linking a classical nonlinear optimization model to flow and transport simulation model. The essential link between the simulator and the optimization method are the derivatives or gradient information required for the optimization algorithm. This methodology does not possess some of the computational limitations of some earlier developed methodologies, using nonlinear programming with the flow and transport process governing equations embedded as equality constraints within the optimization model.

Jha and Datta (2011) presented a linked simulation optimization method for solving the source characterization problem. Use of Adaptive Simulated Annealing (ASA) as optimization algorithm showed superior performance as compared to similar methods using GA. Jha and Datta (2012b) showed superior performance of ASA over GA in solving groundwater source characterization problems. An overview of different optimization based methodologies for solving characterization of unknown groundwater pollution sources is given in Chadalavada et al. (2011b).

## 2.2.   Dedicated Monitoring Network design for Unknown Groundwater Pollution Source Characterisation

Monitoring networks are integral to groundwater management. Design of monitoring network may have different underlying objectives, and vary as per site specific conditions and budgetary constraints. Monitoring networks are essentially installed for extracting information which would assist in achieving the underlying objectives for which the monitoring network was installed in the first place. A large body of literature exists,

dealing with the design of monitoring networks for different groundwater quality management objectives.

Massmann and Freeze (1987) discussed a compliance monitoring network design as the risk cost benefit for a landfill site. Risk is defined as the cost associated with probability of failure. Cost is that of construction and operation of the facility and benefit is the revenue generated from the operation of the facility. The probability of failure was estimated using reliability theory, where failure was designated as the breach of the containment structure resulting in transport of the pollutants to the compliance surface through the hydrogeological surface. This was relatively simple and included only advection transport.

Meyer and Brill (1988) developed a method for locating wells in a monitoring network under conditions of uncertainty. Uncertainty in estimating the simulation parameters is translated to uncertainty in pollutant concentration distribution by generating multiple realizations of the pollutant plume. The monitoring well locations are determined using a facility location model such that the probability of detection is maximised.

Loaiciga (1989) proposed an optimization based methodology for groundwater quality monitoring network design. The optimization problem was solved using mixed-integer programming by minimizing the variance of estimation error subject to resource and unbiasedness constraints. The main objective was to design an optimal sampling plan defining the number, location and sampling frequency of the sampling sites.

Loaiciga and Hudak (1992) developed an optimal monitoring network design for early detection of migrating pollutant plumes. New monitoring well locations are augmented to the existing wells to optimize the probability of detection of pollutants from a waste

impoundment. Loaiciga and Hudak (1993) extended this further, using an analytical approach for the design of optimal monitoring wells in a multilayered groundwater flow system. Monitoring wells are located among sets of candidate monitoring locations to gain more information on maximum pollutant concentration and spatial extent of pollution.

Meyer et al. (1994) designed a multi objective non inferior monitoring network design under conditions of uncertainty. The main objectives considered in this design were to: (1) minimize the number of monitoring wells, (2) maximize the probability of detecting a contaminant leak, and (3) minimize the expected area of contamination at the time of detection. The multi-objective network design problem was formulated as an integer programming problem and solved using SA.

Fethi et al. (1994) designed an optimal monitoring network for jointly monitoring several variables. The method was formulated as an optimization problem in which the variance of estimation was minimized. The proposed technique is based on the geostatistical method of cokriging and the optimization model was solved using branch and bound algorithm. The method gave better estimates of the monitored parameters.

Hudak et al. (1995) designed a monitoring network for a multilayered, regional ground water flow system at risk of contamination from waste storage facilities. Monitoring weights are assigned to candidate locations in terms of the prospect of plume detection and exposure hazard criteria. The weights are used in a binary integer programming problem to select the monitoring locations.

Cieniawski et al. (1995) extended the work of Meyer and Brill (1988) on the optimal location of a network of groundwater monitoring wells under conditions of uncertainty using GAs.

Datta and Dhiman (1996) designed a groundwater quality monitoring network polluted with radioactive pollutants. A simulation model was used for prediction of radioactive pollutant transport. Chance constraint was used to formulate the problem and solved using a mixed-integer programming algorithm. The simulation model was used for prediction of radioactive pollutant transport. Nonlinearities due to the inclusion of cumulative distribution functions (CDFs) of actual spatial concentrations are accommodated in the optimization model through a piecewise linearization scheme.

Groundwater flow and pollution transport is dynamic in nature. At any given time, an optimal monitoring network designed with an underlying objective may not be adequately suitable or may even be redundant in achieving the desired objective at another time. Hence, the monitoring network should be designed in a sequential fashion at different time steps to counter the effect of changes due to the dynamic nature of flow and transport in the groundwater system. Grabow et al. (2000) designed a sequential monitoring network to gain plume information. In their attempt they used an empirical model to predict the future location of the monitoring well using two dimensional concentration data from the existing wells. The method showed a significant reduction of 50 percent in the number of wells installed to get the plume information without loss of plume information.

Montas et al. (2000) developed an optimization based sequential monitoring network design for a stochastic flow field. Optimization problem was solved by directly

21

incorporating the time dimension in the objective function. As a result they were able to find a set of monitoring well locations and a sampling schedule that minimizes plume characterization error while satisfying constraints on the maximum number of wells and allowable number of active wells.

Reed and Minsker (2004) designed a multi-objective formulation for monitoring network design considering Long Term Monitoring (LTM) applications. The specific objectives considered in this design were: (1) minimizing sampling costs, (2) maximizing the accuracy of interpolated plume maps, (3) maximizing the relative accuracy of contaminant mass estimates, and (4) minimizing estimation uncertainty. A combination of quantile kriging and Nondominated Sorted Genetic Algorithm-II (NSGA-II) was used to balance these conflicting objectives in this high order Pareto optimization problem.

Mugunthan and Shoemaker (2004) developed a sequential monitoring network design for LTM considering multiple monitoring periods and uncertain flow conditions. The main objective was to minimize the monitoring cost under the constraint to meet an acceptable level of error in the estimation of total mass for multiple contaminants simultaneously. A new Myopic Heuristic Algorithm (MS-ER), Error-Reducing Neighbourhood (SA-ER) and Genetic Algorithm (GA) were used for solving the optimization problem.

Nunes et al. (2004a) proposed three optimization models for monitoring network design: (1) one that maximizes spatial accuracy; (2) one that minimizes temporal redundancy; and (3) a model that maximizes spatial accuracy and minimizes temporal redundancy. The optimization problem was solved using SA. The third model resulted in selection of the most relevant monitoring locations. This was extended by Nunes et al. (2004b) as a redundancy reduction problem in which cost-benefit analysis was performed to determine

the number of monitoring locations to include in the new design versus loss of information. This resulted in a relative reduction in exploration costs.

Dhar and Datta (2007) developed an optimization based sequential model for groundwater quality monitoring networks. The optimization model incorporated uncertainties in prediction of aquifer parameters like hydraulic conductivity and dispersivity. Randomly generated aquifer parameter values are used to simulate different realizations of the pollutant plume. Cumulative Distribution Functions (CDFs) of actual concentrations at different spatiotemporal locations are incorporated in the optimization problem. These CDFs are used to define chance constraints with associated reliabilities.

Chadalavada and Datta (2008) developed an optimal groundwater monitoring network design for detecting pollution in groundwater aquifers. The method not only considered uncertainty in estimation of aquifer parameters, but extended to include source uncertainty as well. Two separate objectives were considered: (1) minimize the summation of unmonitored concentrations, and (2) minimizes estimation variances of pollutant concentrations at unmonitored locations. The developed optimization models were solved using Genetic Algorithm and the variance was calculated using kriging.

Kollat et al. (2008) developed a new multi-objective evolutionary algorithm (MOEA) to solve large, long-term groundwater monitoring (LTM) design problems. A new class of probabilistic model building evolutionary algorithms called the epsilon-dominance hierarchical Bayesian Optimization Algorithm (ε-hBOA) was used. Kollat et al. (2011) used bias-aware Ensemble Kalman Filtering (EnKF) for Adaptive Strategies for Sampling in Space and Time (ASSIST) framework for improving long-term groundwater monitoring. Multi-objectives were considered in the formulation. In a laboratory-based

physical aquifer tracer experiment, the position and frequency of tracer sampling was optimized to: (1) minimize monitoring costs, (2) maximize the information provided to the EnKF, (3) minimize failures to detect the tracer, (4) maximize the detection of tracer fluxes, (5) minimize error in quantifying tracer mass, and (6) minimize error in quantifying the centroid of the tracer plume. Reed and Kollat (2012) further extended it to include groundwater flow-and-transport forecasting uncertainties and contaminant observation uncertainties.

Azghadi-Bashi and Kerachian (2010) developed a new methodology for optimally locating monitoring wells in order to identify unknown pollution sources. The method combines the capability of Monte Carlo analysis, groundwater flow and transport simulation models and Probabilistic Support Vector Machine (PSVM). The optimization model maximizes both the reliability of contamination detection and the probability of detecting an unknown pollution source. This was further modified by Azghadi-Bashi et al. (2010).

Dhar and Datta (2010) developed an optimization model for redundancy reduction in groundwater quality monitoring networks. The optimization problem was formulated as a logic-based mixed-integer linear optimization model and solved using branch and bound algorithm. Reduction in the redundancy resulted in prevention of loss of economy and overall inefficiency of the network.

Chadalavada et al. (2011a) designed an optimal monitoring network to delineate the pollutant plume using minimum number of monitoring wells. Monitoring wells were installed at locations having minimum measurement uncertainty. The uncertainty in the

study area was quantified by using concentration estimation variances at all the potential monitoring locations.

Jha and Datta (2012a) used a Dynamic Time Warping system to estimate the starting time of the source activity. The estimated starting time of the source activity was further utilized in comparing the observed pollutant concentration measurements with the simulated pollutant concentration measurements correctly in time in a linked simulation optimization model. ASA was used to solve the optimization problem.

## 2.3.    Relevant Tools and Techniques

Relevant literature on tools for modelling groundwater flow and transport processes, data interpolation technique, optimization algorithm, and regression modelling used at different stages throughout this study are discussed in this section.

## 2.4.    Flow and Transport Modelling

Mathematical models are abstractions that represent the physical processes describing the cause and effect relationship.  In groundwater systems, these models are used to simulate the process of flow and solute transport to compute the concentration of a dissolved chemical species in an aquifer at any specified time and place.

Groundwater flow is generally governed by Darcy's law and conservation of mass. The theoretical basis for the equation describing solute transport has been well documented in Bear (1979). Analytical solutions are also available (Javandel et al., 1984), but are limited to unrealistically idealised scenarios. Numerical methods, namely finite-difference and finite-element methods are commonly used for the solution of mathematical equations

used for flow and transport simulation. A comprehensive discussion on the application of these numerical methods to groundwater problems is presented by Wang and Anderson (1982); Anderson and Woessner (1992); Zheng and Bennett (1995); and Domenico and Schwartz (1998).

McDonald and Harbaugh (1988) developed a finite-difference based modular three dimensional groundwater flow model called MODFLOW. MODFLOW has been widely used in groundwater flow simulations. MODFLOW has continuously been updated and the most recent version was released in 2005.

Zheng and Wang (1999) developed a modular three dimensional transport model MT3DMS for simulation of solute transport process in groundwater system. This was an extension of model MT3D developed by Zheng (1990) to simulate the various transport processes such as advection, dispersion and chemical reactions of pollutants in groundwater systems. The MT3DMS includes a multi-component program structure which can accommodate add-on reaction packages for modelling various biological and geochemical reactions. MODFLOW and MT3DMS have been consistently used in this study.

### 2.4.1. Techniques for Geostatistical Data Interpolation

Data interpolation requires estimating the value of a variable at an unmeasured location from observed values at surrounding locations. Matheron (1963) developed a method of geostatistical data interpolation called "kriging". Kriging is a collection of generalised linear regression techniques for minimizing an estimation variance defined from a prior model for covariance (Deutsch & Journel, 1998). An overview of kriging based

geostatistical data interpolation techniques can be found in Journel and Huijbregts (1978); and Cressie (1990).

Geostatistical data interpolation was initially developed with an emphasis on solving mining related problems but has found wide application in all major engineering fields. Geostatistical kriging has been extensively applied to solve various groundwater management problems.

Reed et al. (2000) developed a long term cost-effective monitoring network design for polluted aquifer sites. The method combined a transport simulation model, plume interpolation, and a genetic algorithm to identify cost-effective sampling plans. Inverse Distance Weighting (IDW), Ordinary Kriging (OK) and a hybrid method that combines the two approaches were used for plume interpolation.

Wu et al. (2005) extended the work by Reed et al. (2000) by introducing the first and second moments of a three-dimensional pollutant plume as new constraints in the optimization formulation. Application of geostatistical kriging in groundwater monitoring network design can also be found in Yeh et al. (2006); and Feng-guang et al. (2008).

A MATLAB open source code, mGstat version 0.99 (Hansen, 2004) and kriging packages from the Geostatistical Software Library (GSLIB) (Deutsch & Journel, 1998) are used for data interpolation in this study.

### 2.4.2. Optimization Algorithm: Simulated Annealing

Characterization of unknown groundwater pollution sources is often formulated as an optimization problem and can be solved using different optimization algorithms. Choice

of optimization algorithm largely depends on the type of problem to be solved. In this study, SA is used as the solution algorithm to solve the optimization problem.

Objective functions for solving unknown groundwater pollution source characterization are complex multi-variate optimization problems. Such formulations are highly non-linear, containing several local and global optima. Simulated annealing is a meta-heuristic search algorithm capable of escaping from local optima. Its use of hill-climbing moves to escape local optima makes SA efficient in solving non-convex optimization problems. Its ease of implementation of complex objective functions, and convergence to a global optimal solution, enhances its suitability for solving ill-posed inverse problems, as is the case with unknown groundwater pollution source characterization.

SA was first introduced by Kirkpatrick et al (1983), as an extension of the Metropolis Algorithm (Metropolis et al., 1953). The basic concept of SA is derived from thermodynamics. Cerny (1985) used a thermodynamical approach to solve the travelling salesman problem. Each step of SA algorithm replaces the current solution by a random nearby solution, chosen with a probability that depends on the difference between the corresponding function values and algorithm control parameters, (such as initial temperature, or temperature reduction factor, for example). In this study, SIMANN, a FORTRAN public domain code for SA developed by Goffe (1996) is utilized for the solution algorithm.

### 2.4.3. Genetic Programming

Genetic programming is an evolutionary optimization algorithm based on the concepts of genetics and natural selection and bears strong resemblance to GA. A GP model is

essentially a highly fit computer programs describing the relationship between output values and inputs, evolved using genetic programming (Koza, 1994).

GP is often used to perform symbolic regression. Most conventional regression algorithms optimize the parameters for a pre-specified model structure. However, with GP, the model structure and parameters are determined simultaneously. GP optimises the parameter values of a given model structure within predefined parameter space to find a highly fit computer program that produces desired output for a particular set of inputs.

GP typically codes solutions as tree structured variable length chromosomes. The first step towards development of GP was performed by Cramer (1985), in which he developed the first tree structured GAs for basic symbolic regression. Classification rules using structured GA were developed by Forsyth and Rada (1986). However, it was Koza (1994) who coded the GP algorithm in LISP. This was applied to a wide range of problems, including symbolic regression and classification.

GP is domain independent, and this flexibility renders it the capability to be used for structural optimisation of various engineering problems. However, GP has not been widely applied in groundwater resource management problems. The potential applicability of GP in groundwater problems has been advocated by Sreekanth and Datta (2012) due to the following reasons: (1) GP's ability to develop simple models with interpretability to overcome the curse of the "black box" nature of data intensive models, (2) the lesser numbers of parameters used in GP models as compared to parallel neural network architectures, and (3) GP's ability to parsimoniously identify the significance of the modelling inputs. In this study, this feature has been exploited for the design of an optimal monitoring network and is explained in Chapter 4. A professional software

package, Discipulus$^{TM}$ 5.1 (RML Technologies, Inc.) is used in this study for GP modelling.

## 2.5.    Motivation for this Study

Unknown groundwater pollution source characterization methods are designed to find the answers to the three most important questions about the pollutant sources (Pinder, 2009): (1) When was the pollutant released from the source (release history)? (2) Where is the contamination source (source location)? and (3) At what concentration was the pollutant flux coming from the source (source magnitude)? Based on the unknown source characteristics that a method tries to find, he classified the pollutant source characterization problem into the following categories:

  i.    reconstruction of source release history problems,

 ii.    identification of source location or release time of contaminant,

iii.    identification of source location and magnitude,

 iv.    identification of source location and release time of contaminant, and

  v.    identification of location, magnitude of source and release time of contaminant

In some real world scenarios of groundwater pollution there may be some information about the actual/potential source locations but the release history and the release time is often unknown. In case of clandestine disposal, the number of sources and their locations are also unknown. This makes source characterization a challenging task.

Methodologies developed so far rely on a large number of spatiotemporal concentration measurements for reliable estimation of unknown groundwater pollution source characteristics. All the existing linked simulation-optimization based approaches

implicitly assume that the starting time of the activity of the sources is precisely known, or the time span for the possible start of source activity is known with a fair degree of certainty. These methods fail to give any meaningful result for most of the real world scenarios as such information is seldom available and spatiotemporal concentration measurements are sparse and erroneous. Often, in real scenarios, some measurements are missing or the time interval between measurements is not uniform. Moreover, there may be multiple sources, each starting at different times.

In this study, a linked simulation-optimization methodology is developed for simultaneous identification of unknown groundwater pollution source fluxes and source activity initiation time. The methodology uses source activity starting time as an explicit decision variable. This method attempts to eliminate a major deficiency in the existing methods where the source activity starting time, or the time span for the possible start of source activity, is assumed to be precisely known.

One of the difficulties in accurate characterization of unknown groundwater pollution sources is the uncertainty regarding the number and the location of such sources, especially in a clandestine disposal scenario. Only when the number of source locations is estimated with some degree of certainty, can the characterization of the sources in terms of location, magnitude and activity duration be meaningful. In this study, a sequential monitoring network design that addresses the issue of identification of pollutant source locations plausible in clandestine disposal scenarios is presented.

Monitoring networks are vital for any unknown pollution source characterization problem. However, not all monitoring locations are ideally suited for source identification using linked-simulation optimization technique. The importance of a scientifically

designed monitoring network for efficient identification of unknown source characteristics has not been adequately addressed in the past. To address this issue, a dedicated monitoring network design for efficient identification of unknown source characteristics has been developed in this study. The dedicated monitoring network uses GP based monitoring impact factors and frequency factors for the optimal design.

Systematic and planned monitoring at optimal monitoring locations can provide vital feedback information for efficient source characterization. This aspect has received only limited attention in the past. This study, therefore, incorporates a sequential monitoring network design and gathered information feedback based methodology for efficient identification of unknown pollutant source characteristics. The linked simulation-optimization model is solved using observation data from sequentially designed monitoring networks, such that the feedback information from these networks improves the source characterization results.

# 3. Linked Simulation-Optimization Methodologies for Simultaneous Estimation of Unknown Pollutant Source Characteristics

Similar versions of this chapter have been published and copyrighted to appear in the following journals:

- Prakash, O., & Datta, B., (2014). "Characterization of Groundwater Pollution Sources with Unknown Release Time History". *Journal of Water Resource and Protection (JWARP)*. To be published.

- Prakash, O., & Datta, B., (2014). "Simultaneous Optimal Identification of Unknown Groundwater Pollution Source Fluxes and Source Activity Initiation Time". *Journal of Hydrologic Engineering (ASCE)*. Under review.

In this chapter, a linked simulation-optimization based methodology for simultaneous optimal identification of unknown groundwater pollution source fluxes, and source activity initiation time is discussed. Performance of the developed methodology evaluated for different real life like scenarios using a synthetic study area is also discussed.

## 3.1. Background of the Problem

In the event of detection of any pollutant in the groundwater system, the next step is to design an effective remediation strategy for reclaiming a polluted aquifer. This requires precise information of the source characteristics in terms of magnitude, location, activity

initiation time and activity duration. To achieve this, a preliminary assessment of such a polluted aquifer site along with historical information on the land use pattern is often conducted. In some cases, such assessments may provide reliable estimates on potential source locations but activity initiation times are mostly unknown.

Owing to the slow nature of flow and transport in the groundwater system, pollution in an aquifer is often detected long after the pollution sources have become active or may even have ceased to exist. There may be a gap of years, or even decades, between the start of source activity and detection of pollutants in the aquifer. Therefore, even a rough guess of the time span within which the source activity has started may be difficult.

Linked simulation-optimization based approaches are often used for solving groundwater pollution source identification problems. These existing approaches are efficient when the starting times of the activity of the sources are precisely known, or the possible time window within which sources activity actually starts is not too large and can be specified. However, in real life scenarios, the starting times of the activity of the sources is either unknown or can lie anywhere within a time window of years or decades. Absence of any prior information about the span of time window, within which the sources become active, makes existing source identification methodologies inefficient. To address this major deficiency, a methodology is developed, based on a new optimal decision model that can efficiently solve the source identification model, when it is impossible to specify a small time span within which the source activity may have started.

All existing optimization based source identification techniques implicitly assume that the starting time of source activity is precisely known, or the time span for the possible start of source activity is known with a fair degree of certainty. This is because, in the existing

34

methodologies, the starting time is not treated as an explicit decision variable, but generally estimated indirectly by solving for source flux magnitudes at defined discretised time intervals. The following scenarios represent the existing methods and their limitations.

i.  The actual starting time of the source activity is assumed to be precisely known. This is an impractical assumption and in all real world scenarios, there is seldom any clue about the actual starting time of the sources. As a result, most of these developed methods cannot be applied to such real world scenarios.

ii. In scenarios where the time span for the possible start of sources activity is known with a fair degree of certainty, a typical range being of 1 to 5 years, the precise actual starting time of the sources is estimated indirectly by estimating the source flux magnitude for the entire time span, discretised into smaller stress periods. As a result, a large number of decision variables are added to the optimization problem. Each of these decision variables represents the source flux magnitude for a potential/actual source for a given stress period. Hence, the optimization problem needs to be solved for a large number of source flux magnitudes at different time intervals. Such a solution may theoretically seem possible but with every added decision variable, the dimension of the search space increases, leading to an exponential rise in computation time. This indirect approach does not represent an efficient formulation of the optimum decision model.

iii. In the third scenario, where there is no information of time span for the possible start of source activity, solution of source identification becomes highly challenging. If the solution approach as discussed for the previous scenario is

applied to solve for the actual starting time, this would mean solving for hundreds

of source flux magnitudes for a very large number of stress periods, over a very

large time span. This would not only increase the size of the optimization

problem, but may render it computationally infeasible.

iv. The solution methodology for all the above mentioned scenarios implicitly

assumes that the actual starting time of the source will always be within a

specified time span considered in the study. If this is untrue, the existing methods

fail to give any meaningful result.

To address these limitations, a new methodology is proposed for simultaneously

identifying the starting time of the activity of the sources and their flux release history. A

new optimum decision model is formulated and solved for this purpose. SA is used for

solving the optimization problem with starting time of the sources as one of the explicit

decision variables. Performance of the proposed methodology is evaluated by solving an

illustrative example problem to demonstrate its potential applicability.

## 3.2. Methodology

The proposed linked simulation-optimization methodology for reconstructing the release

history of an unknown pollution source and simultaneously estimating the starting time of

the sources is essentially a two-step process. The first step involves the simulation of

physical processes of flow and solute transport in groundwater system using candidate

solutions. The second step involves using an optimization algorithm for finding the

optimal candidate solution.

The linked simulation optimization model simulates the physical processes of flow and

solute transport within the optimization model. The flow and solute transport simulation

models are treated as important binding constraints for the optimization model. Therefore, any feasible solution of the optimization model needs to satisfy the flow and transport simulation models. The advantage of this approach is that it is possible to link any complex numerical model to the optimization model. In this simultaneous optimal source flux and activity starting time identification model, the flow and transport simulation models are linked to the optimization model using the SA algorithm for solution. The flowchart showing the linking of the simulation model to the optimization model is shown in figure 3.1.

The proposed methodology incorporates the starting times of source activity and source fluxes as explicit unknown decision variables in the optimization model. Candidate values of these unknown decision variables are generated in the optimization algorithm. These candidate values of the unknown decision variables are used as input in the numerical groundwater transport simulation model to generate spatio-temporal pollutant concentration measurements for the entire study area. The generated pollutant concentrations are matched in space and time to the observed pollutant concentration measurements at designated monitoring locations. The difference between simulated and observed concentration is used to calculate the objective function value, which is utilized by the optimization algorithm to improve the candidate solution. The process continues until an optimal solution for the unknown decision variables is obtained. SA is used for solving the optimization problem such that the unknown source fluxes and the starting times of source activity are obtained as direct solutions of the source identification model.

**Figure 3.1 Flowchart Showing the SA based Linked Simulation-Optimization Model**

38

## 3.3. Simulation of Groundwater Flow and Solute Transport Processes

Mathematical models are sets of partial differential equations that represent the physical process of flow and solute transport in the groundwater system. These models are solved to compute the concentration of a dissolved chemical species in an aquifer at any specified time and place. A discretised numerical method is used to solve the governing equations of groundwater flow and transport in this study.

### 3.3.1. Mathematical Representation of Groundwater Flow and Solute Transport Processes Model

The fundamental principle of conservation of mass of fluid or of solute is applied for deriving mathematical equations that describe groundwater flow and transport processes. The governing principle of conservation of mass, often known as the "continuity equation", is used in conjunction with mathematical equations for the relevant process to obtain a differential equation describing flow or transport.

The process of groundwater flow is governed by Darcy's law and the conservation of mass. Darcy's law states that the flow of liquid through a porous media is related to the properties of the liquid, the properties of the porous media, and the hydraulic gradient. A general form of the equation describing the transient flow of a compressible fluid in a non-homogeneous anisotropic aquifer is derived by combining Darcy's law with the continuity equation. The partial differential equation for three dimensional groundwater flow through a porous medium (Rushton & Redshaw, 1979) is given by equation (3.1).

$$\frac{\partial}{\partial x}\left(K_{xx}\frac{\partial h}{\partial x}\right) + \frac{\partial}{\partial y}\left(K_{yy}\frac{\partial h}{\partial y}\right) + \frac{\partial}{\partial z}\left(K_{zz}\frac{\partial h}{\partial z}\right) \pm W = S_s\frac{\partial h}{\partial t} \qquad (3.1)$$

where:

$K_{xx}$, $K_{yy}$ and $K_{zz}$ represent the values of hydraulic conductivity along the x, y and z axes, respectively (LT$^{-1}$),

h is the potentiometric head (L),

W is the volumetric flux per unit volume where positive sign (+) means sources and negative sign (-) means sinks (T$^{-1}$),

$S_s$ is the specific storage of the porous material (L$^{-1}$),

t is time (T),

x, y and z are the Cartesian co-ordinates (L),

The process of solute transport in groundwater systems is a combination of different phenomena acting together:

  i.   advective transport, in which dissolved chemicals are moving with the flowing groundwater;

 ii.   hydrodynamic dispersion, in which molecular and ionic diffusion and small-scale variations in the flow velocity through the porous media cause the paths of dissolved molecules and ions to diverge or spread from the average direction of groundwater flow;

iii.   fluid sources, where water of one composition is introduced into and mixed with water of a different composition;

iv.   reactions, in which some amount of a particular dissolved chemical species may be added to or removed from the groundwater as a result of chemical, biological, and physical reactions in the water or between the water and the solid aquifer materials or other separate liquid phases.

An equation describing the transport and dispersion of a dissolved chemical in flowing groundwater is derived using the principle of conservation of mass through a controlled volume considering all the above mentioned processes. However, transport due to advection and hydrodynamic dispersion is more predominant. Hence, the transport equation is often referred to as the ADE (advection dispersion equation). The partial differential equation describing three-dimensional transport of pollutants in groundwater (Domenico & Schwartz, 1998) is given by equation (3.2).

$$\frac{\partial C}{\partial t} = \frac{\partial}{\partial x_i}\left(D_{ij}\frac{\partial C}{\partial x_j}\right) - \frac{\partial}{\partial x_i}(v_i C) + \frac{q_s}{\theta}C_s + \sum_{k=1}^{N}R_k \qquad (3.2)$$

where:

$C$ is the concentration of pollutants dissolved in groundwater ($ML^{-3}$),

$t$ is time (T),

$x_i$ is the distance along the respective Cartesian coordinate axis (L),

$D_{ij}$ is the hydrodynamic dispersion coefficient tensor ($L^2T^{-1}$),

$v_i$ is the seepage or linear pore water velocity ($LT^{-1}$). It is related to the specific discharge or Darcy flux through the relationship, $v_i = q_i/\theta$,

$q_s$ is volumetric flux of water per unit volume of aquifer representing fluid sources (positive) and sinks (negative) ($T^{-1}$),

$C_s$ is the concentration of the sources or sinks ($ML^{-3}$),

$\theta$ is the effective porosity of the porous medium (dimension less),

$\sum_{k=1}^{N} R_k$ is chemical reaction term for each of the $N$ species considered ($ML^{-3}T^{-1}$).

The transport equation is solved to compute the concentration of a dissolved chemical species in an aquifer at any specified time and place in the study area. In order to solve the transport equation, linear pore water velocity needs to be known for the study area. Hence, it becomes necessary to first calculate the hydraulic head distribution using a groundwater flow simulation model.

### 3.3.2. Numerical Models for Simulation of Groundwater Flow and Transport Processes

Analytical methods for solution of the partial differential equations describing groundwater flow and solute transport require that the properties and boundaries of the flow system be highly idealized. Analytical methods provide exact solutions, but come at the cost of over simplifying assumptions of the complex field environment. Alternatively, for problems where analytical methods are inadequate, partial differential equations are approximated numerically. Mathematical models of groundwater flow or transport assume the variables to be continuous. In order to apply numerical methods of solution, the study area is discretised into grids. Continuous variables are replaced with discrete variables defined at these grid blocks. Thus, the continuous differential equation is

replaced by a finite number of algebraic equations. In order to solve these governing equations of groundwater flow and solute transport, it is necessary to specify the boundary condition. Boundary condition is generally specified as values of head or solute concentration around a boundary (Dirichlet condition), or the flux or concentration gradient around a boundary (Neumann condition). In some cases it is also possible to specify a combination of these two boundary conditions. In order to solve for transient conditions, it is also essential to know the initial conditions. Since groundwater flow and solute transport are inherently transient, it is necessary to know the initial condition in order to solve the governing equations. The initial condition is essentially the starting head or pollutant concentration in groundwater aquifer.

Computer programs that use block-centred finite difference spatial discretization are used for solving the governing equations of groundwater flow and transport. The entire study area is discretised into smaller cuboidal finite difference cells. Continuously varying aquifer parameters, such as hydraulic conductivity, porosity, or dispersivity, are discretised by associating the aquifer parameter value to the centre of each cell of the finite difference grid. The governing equations are solved using iterative methods.

### 3.3.2.1. *MODFLOW*

MODFLOW is a computer program that numerically solves the three-dimensional groundwater flow equation through a porous medium by using a finite-difference method. It was first developed by United States Geological Survey in 1984 and is coded in FORTRAN 77. There have been four major releases of MODFLOW since its initial release in 1984: MODFLOW-88, MODFLOW-96, MODFLOW-2000, and MODFLOW-2005. MODFLOW consists of different independent modules for simulating flow due to

various hydrogeological stresses, such as flow to wells, areal recharge, evapotranspiration, flow to drains, flow through river beds etc. In this study, MODFLOW-2000 (Harbaugh et al., 2000) is used to simulate the groundwater flow. MODFLOW-2000 can simulate steady and transient flow in an irregularly shaped flow system in which aquifer layers can be confined, unconfined, or a combination of confined and unconfined. In case of transient flow conditions, the time domain is discretised into a number of stress periods of finite time length.

### 3.3.2.2.   *MT3DMS*

MT3DMS is a modular three-dimensional multispecies transport model for simulation of advection, dispersion and chemical reactions of pollutants in groundwater systems. MT3DMS was developed by Zheng and Wang (1999) at the University of Alabama. MT3DMS includes three major classes of transport solution techniques in a single code: the standard finite difference method; the particle-tracking based Eulerian-Lagrangian methods; and the higher-order finite-volume TVD method. MT3DMS can accommodate very general spatial discretization schemes and transport boundary conditions, including: (1) confined, unconfined or variably confined/unconfined aquifer layers; (2) inclined model layers and variable cell thickness within the same layer; (3) specified concentration or mass flux boundaries; and (4) the solute transport effects of external hydraulic sources and sinks such as wells, drains, rivers, areal recharge and evapotranspiration. MT3DMS is designed for use with any block-centred finite-difference flow model. In this study MT3DMS is used to compute the pollutant plume utilizing the flow field generated by MODFLOW.

### 3.3.3. Optimization Model

The basic concept of SA is derived from thermodynamics, where molten metals are slowly cooled (annealing) to achieve a low energy state. Unlike classical optimization, which assumes that a function is invariantly quadratic and has one optimum, SA can handle multiple optima. SA derives its strength from its ability to move uphill, and thus escape from local optima to find global optima. SA capability for convergence to a global optimal solution in complex, multivariate problems involving higher degree non-linear functions makes it an ideal choice for solving unknown groundwater pollution source identification problems.

SA, first introduced by Kirkpatrick et al. (1983), is an extension of the Metropolis Algorithm (Metropolis et al., 1953). Each step of the SA algorithm replaces the current solution by a random nearby solution, chosen with a probability that depends on the difference between the corresponding function values and algorithm control parameters (initial temperature, temperature reduction factor, for example). In this study, SA is used as a solution algorithm to solve the optimization problem. SIMANN a FORTRAN public domain code for SA developed by Goffe (1996) is utilized for the solution algorithm.

### 3.3.4. Formulation of Optimization Model for Simultaneous Identification of Unknown Source Flux and Source Activity Initiation time

The proposed methodology incorporates the starting time of the activity of the sources as the explicit unknown decision variables in the optimization model. SA is used for solving the optimization problem with an objective of minimizing the absolute difference between the simulated and measured pollutant concentrations at the observed locations.

The unknown source fluxes and the starting times of source activity are obtained as direct solutions of the source identification model.

In source identification problems where the starting time of the activity of the sources is known, temporal pollutant source fluxes from all the potential sources, represented by the term $q_s C_s$ in the transport equation (3.2) are the only explicit decision variables. Source flux identification using linked simulation-optimization is solved by minimizing the difference between the simulated concentration measurements and the observed concentration measurements in space and time. Hence, source identification problems need observed concentration measurement data at different locations and times. Actually, typical concentration measurements at a given observation location represents only a



**Figure 3.2 Details of Variables with respect to an Observed Breakthrough Curve at a Monitoring Location**

small portion of the entire concentration versus time plot (breakthrough curve) as shown in figure 3.2. If multiple sources are active at different locations and times, the concentration measurements at an observation location over a period of time represent a portion of the combined breakthrough curve at that location.

If the starting time of the activity of the sources is assumed to be known, the simulated concentration measurements can be correctly mapped in time with respect to observed concentration measurements to calculate the residual error at the observed locations. However, in a real life scenario, the starting times of source activity is mostly unknown. Hence, it becomes unclear which temporal concentration measurement values on the simulated breakthrough curve correspond in time to the observed concentration measurements. So that the time matched concentration values can be used to calculate the residual error at the observed locations. This procedure can be understood from the illustrative example shown below.



**Figure 3.3 Snapshot of the Pollutant Plume at 1200 days after the start of Source Activity**

**Figure 3.4 Snapshot of the Pollutant Plume at 1400 days after the start of Source Activity**



**Figure 3.5  Snapshot of the Pollutant Plume at 1600 days after the start of Source Activity**

**Figure 3.6 Snapshot of the Pollutant Plume at 1800 days after the start of
Source Activity**

Figures 3.3 to 3.6 show instantaneous snapshots of pollutant plumes in the polluted aquifer at 1200, 1400, 1600 and 1800 days since the start of source activity. It is evident from these figures, that the pollutant plume is dynamic in nature. This implies that the observed concentration measurement at an observation location will vary in magnitude when taken at different times. However, if the starting time of the sources is unknown, it cannot be determined if the resulting plumes and the corresponding concentration measurements are for 1200, 1400, 1600 and 1800 days since the start of source activity, or for a different set of days. A meaningful comparison between the observed pollutant plumes and the simulated pollutant plumes is only possible when they correspond to the same time. Therefore, a correct estimate of the source fluxes is only possible when the actual and estimated plumes are compared correctly in time and space. Furthermore, the time lag between the start of source activity and the observed concentration measurement needs to be ascertained in order to estimate the source activity starting time.

To overcome the limitations of the methodologies proposed earlier, an explicit decision variable, time lag $\Delta T$ as explained below, is defined in the source identification model. $\Delta T$ is the time between the first concentration measurement and the actual source activity initiation time. Figure 3.2 shows a typical breakthrough curve at any observation location *iob* with the source activity plotted on the x axis. The x axis represents the time axis, the primary y axis represents the source flux (g/s) and the secondary y axis represents the concentration of the dissolved pollutant in groundwater at observation location *iob* (g/lit). Only a small portion of the breakthrough curve, shown by a solid line in figure 3.2, represents concentration measurements of the pollutant at observation location *iob*. The following variables are used in the formulation of the source identification model:

$T^{act}$ = a calendar date representing the actual starting time of the activity of the sources;

$T^{sim}$ = a calendar date representing the starting time of the simulation for any particular candidate solution in the optimization model, which also represents a candidate starting time of the source activity;

$t^m$ = a calendar date representing the first concentration measurement obtained from the site;

$\Delta t$ = time interval between concentration measurements (T),

N = number of concentration measurements available (assumed same for all measurement locations) such that n = {0, 1, 2, 3......... N-1},

$cobs_{iob}^{t^a}$ = observed concentration measurement at observation location *iob* at time $t^a$ (ML$^{-3}$),

$cest_{iob}^{t^s}$ = corresponding simulated concentration at observation location *iob* at time $t^s$

$(ML^{-3})$,

The objective function is defined as:

$$MinimizeF = \sum_{n=0}^{N-1} \sum_{iob=1}^{nob} Abs(cest_{iob}^{t^s} - cobs_{iob}^{t^a}) \qquad (3.3)$$

$$t^s = \Delta T + n\Delta t \qquad (3.4)$$

$$t^a = t^m + n\Delta t \qquad (3.5)$$

subject to

$$cest_{iob}^t = f(q_s, C_s, t) \qquad (3.6)$$

where

$f(q_s, C_s, t)$ represents the simulated concentration obtained from the transport simulation model at an observation location at time *t* and source flux $q_s C_s$,

*nob* is the total number of concentration measurement wells,

*Abs* (…….) represents the absolute value.

The constraint set in equation (3.6) essentially represents the linking of the optimization algorithm with the numerical groundwater flow and transport simulation model through the decision variables. The optimization algorithm searches for optimal values of unknown pollutant source fluxes $q_s C_s$ and the lag time $\Delta T$ by generating candidate

solutions of these decision variables in the optimization algorithm. Candidate values of fluxes $q_sC_s$ are used as input for simulations of flow and transport models.

This optimal search process with different candidate values of unknown variables results in an optimal solution that minimizes residual error between the simulated and observed pollutant concentrations. It is to be noted that $\Delta T$ is a decision variable whose value is to be optimally determined to ascertain the time of initiation of the pollutant source fluxes. At a given search iteration of the optimization algorithm, starting time of the source fluxes also represents the starting time of the transport simulation. The optimal value of lag time $\Delta T$ is the best estimate of the actual source flux starting time $T^{act}$ obtained as solution.

$\Delta T$ obtained as the direct solution of the optimization problem invariably represents optimal time lag between the first concentration measurement and starting time of the simulation $T^{sim}$ for optimal candidate solution in the optimization model. In the objective function formulation in equation (3.3) to equation (3.6) all of the sources are assumed to start at the same time. It is also assumed that at any given search iteration of the optimization algorithm, starting time of the source fluxes also represents the starting time of the transport simulation $T^{sim}$.

**Figure 3.7 Details of Variables with respect to an Observed Breakthrough Curve at a Monitoring Location with Source Starting at Different Times**

However, in scenarios having multiple sources and starting at different times, the objective function formulation in equation (3.4) is modified to solve for such scenarios. An unknown decision variable $\delta T(i)$ for each of the potential sources S($i$) is introduced in the formulation to find the actual source activity initiation time as shown in figure 3.7. If $\Delta T$ be the lag time between the first concentration measurement and a candidate transport simulation start time $T^{sim}$, then $\delta T(i)$ represents the time delay between the candidate transport simulation start time $T^{sim}$ and the start of the source activity of a potential source S($i$). In this case the objective function formulation in equation (3.3) to equation (3.6) remains the same, except equation (3.4) is modified to accommodate the new variable $\delta T$ and is rewritten as shown in equation (3.7).

$$t^{s} = \Delta T - \delta T(i) + n\Delta t \tag{3.7}$$

where:

$\delta T(i)$ is the lag time between the start of the simulation $T^{sim}$ and the start of activity for source $i$.

## 3.4.    Performance Evaluation of Developed Methodology

The performance of the developed methodology is evaluated for an illustrative polluted aquifer study area as shown in figure 3.7. The study area is specified as comprising of heterogeneous, anisotropic and confined aquifer. This study area has a total dimension of 2100 metres in the $x$ direction and 1950 metres in the $y$ direction. The entire study area is



**Figure 3.8 Plan view of the illustrative study area**

discretised into smaller grids of size $\Delta x$, $\Delta y$ and $\Delta z$ in $x$, $y$ and $z$ direction respectively. The study area contains three different hydrogeological zones with different values of hydraulic conductivity $K_{xx}$ and effective porosity $\theta$.   Groundwater flow and solute transport processes are simulated using the value for saturated thickness of the aquifer $b$, longitudinal dispersivity $\alpha L$, transverse dispersivity $\alpha T$ and horizontal anisotropy as given in table 3.1. In the discretised study area, cells marked with a red star represent the grid

locations containing a pollutant source S($i$) where $i$ represents the source number. Cells marked with green and purple circles are the grid locations containing an observation well.

| Parameter | Unit | Value |
|---|---|---|
| Maximum length of study area | m | 2100 |
| Maximum width of study area | m | 1950 |
| Saturated thickness, b | m | 30 |
| Grid spacing in x-direction, $\Delta x$ | m | 50 |
| Grid spacing in y-direction, $\Delta y$ | m | 50 |
| Grid spacing in z-direction, $\Delta z$ | m | 30 |
| Set 1: Kxx (Zone 1, Zone 2, Zone 3) | m/d | 20, 18, 17 |
| Set 2: Kxx (Zone 1, Zone 2, Zone 3) | m/d | 20, 30, 15 |
| $\theta$ (Zone 1, Zone 2, Zone 3) | dimensionless | 0.3, 0.28, 0.25 |
| Longitudinal Dispersivity, $\alpha L$ | m/d | 20 |
| Transverse Dispersivity, $\alpha T$ | m/d | 10 |
| Horizontal Anisotropy | dimensionless | 1.5 |
| Initial contaminant concentration | g/lit | 0 |
| Source Grid Location S1, S2 and S3 | | S1(11,28), S2(15,22) S3(10,16) |

**Table 3.1 Hydrogeological Parameters for Study Area**

### 3.4.1. Performance Evaluation Scenarios

The performance of the developed methodology is evaluated for different real life like scenarios. These evaluations are carried out by varying source activity initiation times, varying the degree of spatial heterogeneity in the study area, considering transient flow conditions, missing observation data and different combinations of observation locations.

#### 3.4.1.1. Case 1: Activity Initiation Times

In the first case, the performance is evaluated by varying the source activity starting time. Three different sources are considered. However, it is assumed that all the three sources

start activity at the same time. The activity duration of the sources is divided into three equal stress periods of 500 days. The pollutant flux from the sources is assumed to be constant over a stress period. The pollutant flux from each of the sources is represented as S($i$)($j$), where $i$ represents the source number and $j$ represents the stress period number. A total of nine source fluxes (S11, S12, S13, S21, S22, S23, S31, S32, S33) and $\Delta T$ are considered as explicit unknown variables in the optimization problem. Concentration measurements *Cobs* from monitoring well locations marked by purple circles (figure 3.7) are utilized in calculating the objective function (equation 3.3). Four different scenarios for having different starting times $T^{act}$ are evaluated (table 3.2). The hydraulic conductivity $K_{xx}$ for all four scenarios in zone 1, zone 2 and zone 3 in the aquifer system are kept as 20 m/d, 18 m/d and 17 m/d respectively.

| Scenario | Actual Time Lag between the start of the source and first concentration measurement $t^m - T^{act}$ (d) | Hydraulic Conductivity (Zone1, Zone2, Zone 3) $K_{xx}$ (m/d) | Observation Wells Grid Locations for Scenarios 1 to 4 |
|---|---|---|---|
| Scenario 1 | 800 | 20, 18, 17 | W1(15, 16), W2(16, 28), W3(18, 17), W4(19, 21), W5(20, 27), W6(22, 24) |
| Scenario 2 | 1000 | 20, 18, 17 | |
| Scenario 3 | 1200 | 20, 18, 17 | |
| Scenario 4 | 1400 | 20, 18, 17 | |

**Table 3.2 Test Scenarios for Case 1**

### 3.4.1.2.    *Case 2: Different Activity Initiation Times for Different Sources and Missing Observation Measurement Data*

In scenario 5, the performance of the developed methodology is evaluated for a more realistic case. Unlike the previous case, all three sources are assumed to have different activity initiation times (table 3.3) and an activity duration of 900 days. The entire

activity duration of the sources is divided into three equal stress periods of 300 days. All the other aquifer parameters and the monitoring well locations for pollutant concentration measurement are kept the same as in the previous case. A total of nine source fluxes (S11, S12, S13, S21, S22, S23, S31, S32, S33) and $\Delta T$ are considered as explicit unknown variables in the optimization problem. Since each of the sources starts at a different time, three additional variables $\delta T$ (*1*), $\delta T$ (*2*) and $\delta T$ (*3*) are introduced into the optimization problem for each of the potential sources S1, S2 and S3 respectively. Equation (3.7) is used to find the actual starting time of each of these sources. The scenario is further complicated by considering irregular monitoring for pollutant concentration measurements and missing observation measurement data *Cobs*.

| Scenario 5 | Actual Time Lag between the start of the source and first concentration measurement $t^m - T^{act}$ (d) | Activity Duration (d) |
|---|---|---|
| Source 1 | 500 | 900 |
| Source 2 | 200 | 900 |
| Source 3 | 800 | 900 |

**Table 3.3 Test Scenario for Case 2**

### 3.4.1.3.    *Case 3: Effect of Higher Degree of Heterogeneity in Hydraulic Conductivity and Monitoring Well Locations*

In scenarios 6 and 7, performance evaluation is carried out to see the effect of a high degree of heterogeneity in aquifer parameters. All the source characteristics and aquifer parameters, except hydraulic conductivity, are kept the same as in case 1. Specified larger values of hydraulic conductivity in this case (compared to case 1 and 2) ensure that the pollutants travel faster, resulting in a greater degree of overlapping/mixing of the pollutant plumes from individual sources. Thus the observed concentration measurement

would be a combined effect of all of the polluting sources, such that the effect due to a single source cannot be separately quantified. This would increase the degree of non-uniqueness in the solution, making source identification more challenging. The value of hydraulic conductivity in one of the aquifer zones is comparatively larger (table 3.4), as the idea was to test the performance for a highly heterogeneous scenario. This case also evaluates the effect of a different monitoring network for concentration measurement.

| Scenario | Actual Time Lag between the start of the source and first concentration measurement $t^m - T^{act}$ (d) | Hydraulic Conductivity (Zone1, Zone2, Zone 3) $K_{xx}$ (m/day) | Number of Observation Locations and their Respective Grid Locations |
|---|---|---|---|
| Scenario 6 | 1600 | 20, 30, 15 | For scenario 6 W1(15, 16), W2(16, 28), W3(18,17), W4(19, 21), W5(20, 27),W6(22, 24) |
| Scenario 7 | 1600 | 20, 30, 15 | For scenario 7 W1(21, 17), W2(21, 27), W3(24,19), W4(24, 25), W5(26, 22),W6(22, 22) |

**Table 3.4 Test Scenarios for Case 3**

### 3.4.1.4.    *Case 4: Transient Flow Condition*

Since all flow and pollutant transport in groundwater system is inherently transient, the performance of the developed methodology was evaluated for transient flow and solute transport conditions. Three different sources are considered and it is assumed that all the three sources start activity at the same time. The activity duration of the sources is divided into three equal stress periods of 364 days. The pollutant flux from the sources is assumed to be constant over a stress period. A total of nine source fluxes (S11, S12, S13, S21, S22, S23, S31, S32, S33) and $\Delta T$ are considered as explicit unknown variables in the

optimization problem. Table 3.5 gives the aquifer parameters, source locations and monitoring well locations used in the evaluation of scenario 8.

| Scenario 8 Source Locations | Actual Time Lag between the start of the source and first concentration measurement $t^m - T^{act}$ (d) | Hydraulic Conductivity (Zone1, Zone2, Zone 3) $K_{xx}$ (m/day) | Number of Observation Locations and their Respective Grid Locations |
|---|---|---|---|
| S1 (11,28) S2 (15,22) S3 (10,16) | 1638 | 20, 30, 15 | W1(21, 17), W2(21, 27), W3(24,19), W4(24, 25), W5(26, 22),W6(22, 22) |

**Table 3.5 Test Scenario for Case 4**

### 3.4.2. Simulation of Observed Concentration

In actual field application, the concentration measurements are to be obtained from field data. However, for performance evaluation purposes, these measurements are synthetically generated by simulation for assumed actual pollutant sources. The observed aquifer responses are simulated by numerical simulation models, MODFLOW and MT3DMS in GMS7.0. Initial and boundary conditions (initial heads and boundary heads) are specified in the numerical simulation models. Actual source fluxes are utilized to simulate these observed concentration measurement data at specified measurement locations.

The start time of the numerical simulation model for synthetically generating the observed concentration measurements is designated as the initial time $T^{act}$. This initial time can be defined with respect to a calendar date, which in practice would correspond to the actual starting time of the sources. The time lag between the start of the numerical

simulation model for synthetically generating the observed concentration measurements and the first concentration measurement used for source identification is also specified.

In the different evaluation scenarios considered, the values of time lag between the time of first pollutant concentration measurement $t^m$ used in the source identification model and the actual starting time of the source activity $T^{act}$ at the observation locations are given in table 3.6. Concentration measurements are taken every 200 days for case 1, case 2 and case 3, and every 182 days for case 4, starting from the first pollutant concentration measurement time $t^m$. A total of four temporal pollutant concentration measurements from each of the six observation wells are utilized.

| Case | Scenario | Actual Time Lag between the start of the source and first concentration measurement $t^m - T^{act}$ (d) | Concentration Measurement Intervals (d) |
|---|---|---|---|
| Case 1 | Scenario 1 | 800 | 200 |
| | Scenario 2 | 1000 | 200 |
| | Scenario 3 | 1200 | 200 |
| | Scenario 4 | 1400 | 200 |
| Case 2 | Scenario 5 | S1=500, S2=200, S3=800 | 200 (missing data) |
| Case 3 | Scenario 6 | 1600 | 200 |
| | Scenario 7 | 1600 | 200 |
| Case 4 | Scenario 8 | 1638 | 182 |

**Table 3.6 Actual Time Lag between the start of the source and first concentration measurement and Concentration Measurement Intervals**

### 3.4.3. Evaluation of Methodology using Erroneous Concentration Measurement Data

In order to reflect real life conditions, where the contamination measurements are erroneous, the numerically simulated concentrations were perturbed to incorporate measurement errors. The observed concentrations generated using MODFLOW and MT3DMS were perturbed by adding error terms to the simulated measurement data to

represent the effect of random measurement errors. These errors were added in order to incorporate realistic measurement errors. The observed pollutant concentration data is perturbed with a random measurement error with maximum deviation of 10 percent of the measured concentration value as shown in equation 3.8.

$$^{Pert}cobs_{iob}^{t^a} = cobs_{iob}^{t^a}(1+err)$$
(3.8)

$$err = \mu per \times rand$$
(3.9)

where:

$^{Pert}cobs_{iob}^{t^a}$ is the perturbed simulated erroneous concentration measurement at location *iob* at time $t^a$,

*err* is error term,

*μper* is maximum deviation expressed as percentage,

*rand* is a random fraction between -1.0 and +1.0 generated using a latin hypercube distribution.

Latin hypercube distribution is chosen for generating random error data evenly distributed across all class intervals, thus eliminating any clustering of sample data in few of the class intervals.

### 3.4.4. Solution Procedure

The proposed methodology for source flux identification and estimation of source activity initiation time is evaluated using synthetically generated observed concentration measurements as explained in the previous section. The source identification model simulates the aquifer response for a period of 4000 days, and concentration measurements are recorded for different cases as specified in table 3.6 at specified locations. Since the actual starting time $T^{act}$ of the activity of the source is not known to the source identification model, it assumes the starting time $T^{sim}$ in the source identification model. Time $T^{sim}$ can be any calendar date before or after the actual starting time $T^{act}$ of the sources.

Candidate values of unknown source fluxes and time lag are generated within the optimization model. The generated flux values are used for simulation of flow and transport processes as a part of the linked simulation-optimization model. Value of time lag $\Delta T$ obtained as a solution of the source identification model determines the temporal spacing between the assumed starting time $T^{sim}$ of the sources, and first concentration measurement $t^s$. This concentration measurement will lie on the breakthrough curve, a portion of which is used for calculating the residual error. Optimal source identification is evaluated for all the cases, first using error-free measurement data and then using measurement data perturbed with random error. SA parameters used to optimally estimate the source fluxes and the starting time are kept same in all the scenarios and is as described in table 3.7. These performance evaluation results are discussed in the following section.

| SA Parameter | Parameter Value |
| --- | --- |
| Temperature reduction factor | 0.85 |
| Initial Temperature | 10E8 |
| Error tolerance for termination | 1.0E-8 |
| Number of cycles | 20 |
| Number of final function values used to decide upon termination | 4000 |
| Number of iterations before temperature reduction | 10000 |
| Maximum number of iteration | 10E8 |

**Table 3.7 SA Parameters used in Source Identification Model**

## 3.5. Discussion of Solution Results

The evaluation results of all the scenarios are presented in the following sections. The solution results of identification using error-free data and perturbed-error data are compared to the actual source flux and lag times for all scenarios as shown in figure 3.9 to figure 3.17. Each of the unknown source flux variables $S(i)(j)$ and lag time $\Delta T$ is marked on the x axis, having three corresponding bars. The first bar is the actual value. The second bar represents the estimated values using error-free concentration measurements. The third bar represents the estimated values using concentration measurements with perturbed error.

Results of source flux identification using error-free data closely match with the actual source flux values for all source fluxes in all scenarios. The time lag $\Delta T$ between the start of source activity and the first pollutant concentration measurement is precisely identified by the model in all eight scenarios. This estimated time lag can be used to estimate the starting times of the activity of the sources using equation (3.4) and (3.5). Even while using perturbed concentration measurement data in the identification model, the time lag between the start of source activity and the first pollutant concentration measurement is

accurately estimated, leading to accurate identification of the starting time of the activity of the sources.

The time lag $\Delta T$ was estimated correctly irrespective of the assumed starting time in the source identification model. This was one of the drawbacks in the earlier models where the starting time of the sources in the simulation model $T^{sim}$ was implicitly assumed to be earlier than the actual starting time of the sources $T^{act}$. The earlier models failed to give any meaningful result if this condition was not met. These evaluation results for scenario 1, scenario 2, scenario 3 and scenario 4 can be seen in figure 3.9, figure 3.10, figure 3.11 and figure 3.12 respectively.

Simultaneous Source flux and Starting Time Identification: Scenario 1



**Figure 3.9 Identification Results for Scenario 1**

**Figure 3.10 Identification Results for Scenario 2**



**Figure 3.11 Identification Results for Scenario 3**

Simultaneous Source flux and Starting Time Identification: Scenario 4

**Figure 3.12 Identification Results for Scenario 4**

Evaluation results for scenario 5 show that the source activity initiation times for all three sources are estimated correctly. $\delta T$ (1), $\delta T$ (2) and $\delta T$ (3) representing the time delay between the candidate transport simulation start time $T^{sim}$ and the start of the source activity of potential sources S1, S2 and S3 respectively, are the direct outputs. Activity initiation time of source S1, S2 and S3 are estimated by using equation (3.7). $\delta T(1)$, $\delta T(2)$ and $\delta T(3)$ are subtracted from $\Delta T$ to find lag time between the start of the source activity and the first concentration measurement respectively. From figure 3.13 it is evident that source S1 starts at the same time as the start of the simulation ($\delta T(1) = 0$). Source S2 starts activity 300 days after the simulation start time $T^{sim}$ ($\delta T(2) = 300$) and source S3 starts activity 600 days after the simulation start time $T^{sim}$ ($\delta T(3) = 600$). Estimates of starting times for all sources present in the scenario are estimated accurately, using error-free data and perturbed-error data and missing concentration measurement data.

Simultaneous Source flux and Starting Time Identification: Scenario 5



**Figure 3.13 Identification Results for Scenario 5**

Results from scenario 6 and 7 show that the actual starting time is estimated correctly in spite of the high degree of heterogeneity shown in figures 3.14 and 3.15 respectively. To understand the effect of change of observation well locations, results of source flux identification from scenario 6 and scenario 7 are compared in figure 3.16. These two scenarios are identical except that they use concentration measurements from different observation locations for solving the identification problem. The source fluxes estimated in scenario 6 show large deviations in estimating S11, S21 and S31using perturbed-error data as shown in figure 3.14. This shows that the location of observation wells strongly impacts the identification results.

Simultaneous Source flux and Starting Time Identification: Scenario 6



**Figure 3.14 Identification Results for Scenario 6**

Simultaneous Source flux and Starting Time Identification: Scenario 7



**Figure 3.15 Identification Results for Scenario 7**

Comparison of Observation location change on Simultaneous Source
flux and Starting Time Identification: Scenario 6 vs Scenario 7



**Figure 3.16 Effect of Observation Location on Source Identification Result**

It was also interesting to note that result of source flux identification in scenario 7 performed better than scenario 6, despite the fact that the observation well locations were relatively further away from the sources in scenario 7 compared to scenario 6. This may be contrary to intuition. Therefore, concentration measurements from the observation locations used in scenario 6 were analysed. It was found that, after a lag time of 1600 days, as was the case in scenario 6 and 7, most of the pollutant plume had already passed over the observation well locations used in scenario 6. As a result, only the tail end of the breakthrough curve was captured in the observed pollutant concentrations. It is also highly probable that the observed pollutant concentration from the tail end of the breakthrough curve was due to the later part of source activity, and the effect of the initial source activity was not captured by the tail end of the breakthrough curve. This could

explain the reason for large deviations in estimating the flux values from all three sources during stress period one using erroneous concentration measurements.

It is seen that observation well locations can significantly affect the results of source fluxes estimation in a linked simulation-optimization problem. Spatial locations of these observation wells determine what part of the breakthrough curve will be captured as observed pollutant concentration. Observed pollutant concentration will represent different parts of the same breakthrough curve for different lag times. Hence, it is important to choose the observation well locations such that concentration measurements from these locations improve the accuracy of source identification results. This issue has been explained in detail in the next chapter of this study. The chapter specifically deals with optimal design of dedicated monitoring networks for efficient identification of unknown groundwater pollution sources incorporating genetic programming based monitoring.

Evaluation results for scenario 8, having transient condition, is shown in figure 3.17. It can be seen that the source activity starting time is estimated accurately using error-free data and perturbed-error data even in transient condition. The source fluxes estimated in scenario 8 shows large deviation in estimating S11, S21 and S31 using perturbed error data. The same explanation as in case of scenario 6 can be used, since measurement data from the same monitoring wells is used for identification of source flux and source activity starting time in scenario 6 and scenario 8. Also, the lag time in scenario 6 and scenario 8 is nearly same.

Simultaneous Source flux and Starting Time Identification: Scenario 8



**Figure 3.17 Identification Results for Scenario 8**

These results show that this proposed methodology is capable of overcoming major limitations in the earlier methods. The main limitations of these earlier proposed source identification methods were: (1) it assumed that the starting time of the activity of the sources is precisely known; (2) in scenarios where the time span for the possible start of the sources activity is known with a fair degree of certainty, actual starting time was estimated indirectly by solving for source flux magnitudes; and (3) in scenarios where there was no information of the time span for the possible start of the sources activity, solving indirectly for a large number of source flux magnitudes over a very large time span rendered the earlier methods computationally infeasible.

The solution results for source flux estimates using erroneous concentration measurements data show large errors of estimation in comparison to error-free measurement data. These deviations between the actual and the estimated value of the

source fluxes show the effect of errors in concentration measurement data, which accounts for random measurement errors in a real world scenario.

## 3.6. Summary and Conclusions

A methodology is developed for simultaneous identification of the source fluxes and the starting time of the source activity. In this linked simulation-optimization based methodology the starting time of the activity of the source is an explicit unknown decision variable which is estimated as a part of the optimal solution result.

This developed methodology for simultaneous identification of source fluxes and their starting times appears to perform satisfactorily in estimating the unknown groundwater pollution source fluxes and their starting times for the illustrative example. The tested scenario assumed there is very little information initially available regarding the time span within which the pollutant sources became active. The performance evaluation results show that the developed methodology is successful in estimating the source flux values and the starting time of the sources correctly, even when a high degree of heterogeneity is introduced in the aquifer system. However, the result of source identification is highly affected by the spatial location of the observation wells used in the source identification. The spatial locations of the observation wells determine what part of the breakthrough curve will be captured in the form of concentration measurements for a given lag time. In other words, some observation locations may be efficient in accurate source flux estimation for a given lag time, but may not be efficient for a different lag time.

This method overcomes one of the critical limitations in the methods proposed earlier using linked simulation-optimization. In these earlier proposed methods the actual starting times of the sources were assumed to be precisely known (impractical assumption). In scenarios where the starting time was totally unknown, the simulation component of the optimal search procedure needed to start from a much earlier point in time. This was necessary to cover the actual starting times of the sources within the range of discretised time intervals considered. These existing approaches resulted in increasing the number of decision variables substantially, making the optimal search algorithm inefficient, if not infeasible, in some cases. Also, the convergence to a correct optimal solution may be hindered by inclusion of such a large number of decision variables, each representing an unknown source flux magnitude for each of the discretised time intervals for possible source activity.

The proposed methodology is applicable to real world problems of source identification in polluted aquifer sites, where no information is available about the starting time of the source activity and the uncertainty spans over a large time period. This methodology in its current form can estimate the activity starting time for all the potential sources, both in steady state and transient conditions. It can also handle multiple sources having different source activity initiation times and missing observation data.

The main advantage of the proposed methodology is its capability to treat the starting time of source activity as an explicit decision variable. This capability has the potential to render many real life unknown groundwater pollution source identification problems computationally feasible. This will ultimately enhance the capability of addressing the groundwater pollution remediation issue in complex and large scale polluted aquifer

study areas, where there is very little prior knowledge of source location, magnitude, and source activity initiation time.

# 4. Methodologies for Monitoring Network Design for Efficient Source Characterization

Similar versions of this chapter have been published and copyrighted in the following journals:

- Prakash, O., & Datta, B., (2012). "Sequential optimal monitoring network design and iterative spatial estimation of pollutant concentration for identification of unknown groundwater pollution source locations". *Environment Monitoring Assess. (EMAS)*. DOI: 10.1007/s10661-012-2971-8

- Datta, B., Prakash, O., Campbell, S., & Escalada, G., (2013). "Efficient Identification of Unknown Groundwater Pollution Sources using Linked Simulation-optimization incorporating Monitoring Location Impact Factor and Frequency Factor". *Water Resources Management*. DOI: 10.1007/s11269-013-0451-8

- Prakash, O., & Datta, B., (2013). "A Multi-Objective Monitoring Network Design for Efficient Identification of Unknown Groundwater Pollution Sources Incorporating Genetic Programming Based Monitoring". *Journal of Hydrologic Engineering, (ASCE)*. To be published.

- Prakash, O., & Datta, B., (2014). "Optimal monitoring network design for efficient identification of unknown groundwater pollution sources". *Int. J. of GEOMATE*, March, 2014, Vol. 6, No. 1 (Sl. No. 11), pp. 785-790.

- Datta, B., Prakash, O., & Sreekanth, J., (2014). "Application of Genetic Programming Models Incorporated in Optimization Models for Contaminated Groundwater Systems Management". *A bridge between Probability, Set Oriented Numerics and Evolutionary Computations VI,* (Advances in Intelligent and Soft Computing, Springer Series). Book chapter to be published.

In this chapter, two separate monitoring network design models are discussed. The first monitoring network design model deals with identification of unknown groundwater pollution source locations. The second is a dedicated monitoring network design model for efficient identification of unknown groundwater pollution sources. Both monitoring network design models use SA as the solution algorithm to solve the optimization problem.

## 4.1.    Sequential Optimal Monitoring Network Design for Identification of Unknown Groundwater Pollution Source Locations

One of the difficulties in accurate characterization of unknown groundwater pollution sources is the uncertainty regarding the number and location of such sources.  Only when the number of source locations is estimated with some degree of certainty can the characterization of the sources in terms of location, magnitude and activity duration be meaningful. A fairly good knowledge of potential source locations can substantially decrease the degree of non-uniqueness in the set of possible aquifer responses to subjected geochemical stresses.

In a source characterization problem, often the number of potential sources, and their respective locations, are assumed to be known. In some cases of aquifer pollution, the number and locations of polluting sources may be evident or can be guessed using the

available information. In such scenarios, this assumption holds true. However, in most cases, especially in clandestine underground disposal of toxic wastes, no information is available about the number and locations of such sources. In the absence of preliminary information, such assumptions may often lead to inaccurate source characterization results. Hence, the first step in the process of source characterization should be identification of number and locations of such sources.

A large amount of observed concentration measurement data spread over time and space is necessary for detection of unknown groundwater pollution source locations. However, long term monitoring over a large number of sampling locations has budgetary constraints. Contrary to the above mentioned-requirement, pollution in groundwater aquifers is generally first detected by an arbitrarily located water supply well or a group of wells. Pollution concentration measurement data from these wells is often sparse. Moreover, these wells may not be optimally located to identify the number of potential source locations.

To address this limitation, a new monitoring network design methodology is developed for the estimation of unknown groundwater pollution source locations starting from very sparse and arbitrary pollutant concentration observation data. The methodology incorporates concentration gradient information and a sequence of monitoring network implementation, for concentration feedback information. The methodology combines the technique of optimization with the geostatistical data interpolation technique for design of a groundwater monitoring network for more efficient identification of unknown source locations. Once the potential locations of the unknown sources are identified, this information can be used as an important input for optimal pollutant source characterization in terms of location, magnitude and activity duration.

### 4.1.1. Methodology for Sequential Optimal Monitoring Network Design for Source Locations Identification

The proposed methodology consists of two steps for the sequential design of the monitoring network. In the first step, a geostatistical data interpolation technique called kriging is used to interpolate pollutant concentration data for the entire aquifer study area. Initially, this data consists of sparsely available observed concentration from randomly located existing wells. In the second step, an SA-based optimization model is solved to obtain the locations of the next sequence of monitoring wells to be implemented. Once a sequence of monitoring wells is implemented, data subsequently collected from these wells and pre-existing wells is together used in the kriging model to interpolate the concentration data for the study area in the next iteration. This forms the input for the optimization model for finding the optimum locations of the next sequence of monitoring wells to be implemented. Thus, the observed concentration data from the designed and implemented monitoring network is used iteratively as feedback information for identification of potential groundwater pollution source locations, and to estimate the number of sources present in the aquifer study area.

#### 4.1.1.1. Data Interpolation Model: Geostatistical Kriging

Almost all practical scenarios of groundwater pollution lack adequate information about the sources or spatiotemporal pollutant concentration data throughout the affected aquifer region. The observed concentration data is only available at a few sparsely and arbitrarily located wells. Data interpolation techniques are often employed to overcome this problem. Data interpolation requires estimating the value of a variable at an unmeasured location from observed values at surrounding locations.

Geostatistics offers various deterministic and statistical tools for modelling spatial variability. Kriging is a geostatistical interpolation technique which provides optimal unbiased estimates of unknown data points taking into account the distance and the degree of variations between known data points. Deutsch and Journel (1998) defined kriging as a collection of generalised linear regression techniques for minimizing an estimation variance defined from a prior model for covariance.

The first step involved in kriging is constructing an experimental semivariogram using equation (4.1), and fitting it to a standard model. The semivariogram describes the relationship between the variance in data values and the distance between data points by plotting the variance against distance. The basic technique in kriging is to use a weighted sum of neighbouring sample values to estimate the unknown value at a given location. These weights are optimized using the semivariogram model, the location of the samples and all of the relevant inter-relationships between known and unknown values.

$$\gamma(h) = \frac{1}{2N(h)} \sum_{i=1}^{N(h)} [z(x_i) - z(x_i + h)]^2 \qquad (4.1)$$

where:

$\gamma(h)$ is the estimated value of the semivariance for lag $h$,

$N(h)$ is the number of experimental pairs separated by vector $h$,

$z(x_i)$ and $z(x_i + h)$ are the values of variable $z$ at $x_i$ and $x_i + h$, respectively,

$x_i$ and $x_i + h$, are the positions in two dimensions separated by a lag distance of $h$.

The experimental semivariogram estimated using equation (4.1) is fitted to various standard semivariogram models such as the spherical model (equation 4.5 and 4.6), exponential model (equation 4.7) and Gaussian model (equation 4.8), for example, using the weighted least square method. The standard semivariogram model with the lowest error is then chosen for further analysis. A typical experimental semivariogram is shown in figure 4.1.



**Figure 4.1 Typical Semivariogram Structure**

In figure 4.1, sill is the amount of semivariance achieved at the plateau of the curve often represented by variable $c$. It is equivalent to the variance of the data. Range is defined as the lag distance beyond which the data is no longer correlated. It is represented by variable $a$. The nugget is the semivariance at $h = 0$. This basically represents the noise in the data. Range marks the lag distances such that data within the range are correlated. Only data within the range is used for making predictions. The values of sill and range are calculated by fitting a model to the experimental semivariogram. However, different theoretical semivariogram models yield different values for the sill and range.

In kriging based interpolation, the interpolated value is expressed as the weighted sum of the measured values (equation 4.2)

$$z(x_0) = \sum_{i=1}^{N} \lambda_i z(x_i) \qquad i = 1, 2, 3 \text{ .....................} N \tag{4.2}$$

where:

$\lambda_i$ = weight for the observation $z$ at location $x_i$.

Weights $\lambda_i$ are estimated such that the variance of estimation is minimized. Two different forms of kriging, Simple Kriging (SK) and Ordinary Kriging (OK) are often applied. In SK the weights are calculated using equation (4.3).

$$\gamma(x_i, x_0) = \sum_{j=1}^{N} \lambda_j \gamma(x_i, x_j) + \mu \qquad i = 1, 2, 3 \text{ .....................} N \tag{4.3}$$

where:

$\mu$ is a Lagrange multiplier,

$\gamma(x_i, x_j)$ is semivariogram between two points $x_i$ and $x_j$.

SK requires a prior knowledge of the mean and the covariance matrix. This is often unrealistic. However, in the case of OK, the kriging weights are summed to unity shown in equation (4.4). Simple kriging uses the average of the entire data set while ordinary kriging uses a local average. As a result, simple kriging can be less accurate than ordinary kriging.

$$\sum_{j=1}^{N} \lambda_j = 1 \tag{4.4}$$

Kriging is a geostatistical estimation method that has wide application in all major engineering fields. In this study, Geostatistical kriging uses pollutant concentration data

from sparsely and arbitrarily located wells to interpolate the pollutant concentration for the entire study area. A MATLAB open source code, mGstat version 0.99 (Hansen, 2004) is used for kriging. Some of the commonly used variogram structures are given in equation (4.5) to equation (4.8).

$$\gamma(h) = cSph\left(\frac{h}{a}\right) = c\left[1.5\frac{h}{a} - 0.5\left(\frac{h}{a}\right)^3\right] \qquad \text{for } h \leq a \qquad (4.5)$$

$$\gamma(h) = cSph\left(\frac{h}{a}\right) = c \qquad \text{for } h \geq a \qquad (4.6)$$

$$\gamma(h) = cExp\left(\frac{h}{a}\right) = c\left[1 - \exp\left(-\frac{3h}{a}\right)\right] \qquad (4.7)$$

$$\gamma(h) = c\left[1 - \exp\left(-\frac{(3h)^2}{a^2}\right)\right] \qquad (4.8)$$

where $a$ represents the value of the range and $c$ represents the sill value (figure 4.1) .

### 4.1.1.2.   Optimization Model for Monitoring Network Design

The sequential monitoring network design model is developed to find optimal monitoring well locations with the aim of finding well locations with high concentrations of pollutant in every sequence, eventually isolating the source locations. The model for optimal monitoring network design for efficient identification of unknown pollution source locations is defined by equation (4.5). This optimal design model is solved at every iteration in the sequence of network design. The optimization model maximizes the objective function value, subject to the constraint that the maximum number of

monitoring wells that can be selected in any design iteration is limited. The objective function defined by equation (4.9) maximizes the summation of the product of estimated concentration gradients, and the kriged concentration at all the monitoring locations chosen as optimal in particular design iterations.

$$F_{obj} = Max \sum f_{i.j} C_{i,j}^* \left\{ \frac{\Delta C_{i,j}^*}{\Delta X} + \frac{\Delta C_{i,j}^*}{\Delta Y} \right\} \quad \forall \quad i,j \qquad (4.9)$$

These concentration gradients are computed based on the interpolated spatial concentrations as obtained from the implemented monitoring network at the beginning of the design iteration. The concentration gradient along $x$ axis and $y$ axis at any grid location $i, j$ is given by equation (4.10) and (4.11) respectively.

$$\frac{\Delta C_{i,j}^*}{\Delta X} = \frac{\left| C_{i-1,j}^* - C_{i,j}^* \right| + \left| C_{i+1,j}^* - C_{i,j}^* \right|}{\Delta X} \qquad (4.10)$$

$$\frac{\Delta C_{i,j}^*}{\Delta Y} = \frac{\left| C_{i,j-1}^* - C_{i,j}^* \right| + \left| C_{i,j+1}^* - C_{i,j}^* \right|}{\Delta Y} \qquad (4.11)$$

where:

$C_{i,j}^*$ is the concentration value obtained from kriging at the grid $i, j$,

$\Delta X$ and $\Delta Y$ is size of the grid in the $i, j$ direction respectively.

The objective function is maximized subject to the following constraints:

$$\varepsilon_{max} \geq C_{i,j}^* \geq \varepsilon_{min} \quad \forall \quad i,j \qquad (4.12)$$

$\varepsilon_{\min}$ is the average of the measured concentration from the initial and implemented well locations given by equation (4.13).

$$\varepsilon_{\min} = (\sum_{n=1}^{m} (C_{i,j})_n) / m \qquad (4.13)$$

The value of $\varepsilon_{\min}$ changes with every iteration as new monitoring wells are implemented in the study area with every sequence of the design.

$\varepsilon_{\max}$ is the upper bound pollutant concentration value in the aquifer region.

The following constraint in equation (4.14) essentially represents the imposed limit on the total number of permissible monitoring wells at the current design iteration.

$$\sum f_{i,j} \leq k \quad \forall \; i,j \qquad (4.14)$$

where:

$f_{i,j}$ represents the binary decision variable to place or not to place a monitoring well at grid location $i, j$. $f_{i,j} \equiv \{0, 1\}$ such that when $f_{i,j}$ value equal to 1 representing monitoring well to be placed at grid $i, j$, and zero otherwise,

$m$ is the total number of monitoring wells already existing before each design iteration,

$k$ is the maximum permissible number of monitoring wells that can be placed in the study area including the existing ones in the current design stage,

$C_{i,j}$ is the concentration measurement value at the initial wells and the implemented wells in the study area,

The solution of the model provides locations for the placement of new monitoring wells, having larger concentrations of pollutant compared to previous iteration, in every sequence. The observed concentration values from these new monitoring wells and from the pre-existing monitoring wells is used to krig the values of concentration for the next iteration, thereby every time increasing the value of $\varepsilon_{\min}$ . The higher $\varepsilon_{\min}$ value confines the search to locations closer to the actual source locations, thereby finding the potential candidates to be tested for actual sources.

### *4.1.1.3.    Pollution Source Location Identification*

Pollution source locations are identified by implementing monitoring well locations sequentially. Every sequence of the pollution source locations identification model comprises various steps as explained below. The entire methodology for source location identification is schematically described in figure 4.2:

1. The method starts with observed pollutant concentration data, $Cobs$, from sparse, randomly located, existing wells showing some pollutant concentration. Total number of initially available wells is $M^{int}$.

2. The total number of monitoring wells available in the field at the beginning of a current design iteration $ITR$ is denoted by $M_{ITR}$ which is the addition of the initially available wells $M^{int}$ and the total number of monitoring wells implemented from iteration one through iteration $ITR-1$.

$$M_{ITR} = M^{int} + \sum_{itr=1}^{ITR-1} M_{itr}^{imp} \qquad\qquad (4.15)$$

where:

$ITR$ is the current design iteration number and $itr$ is the iteration variable; the

number of implemented monitoring wells in every iteration is denoted by $M^{imp}$;

In the first iteration: $ITR = 1$, $M_{ITR} = M^{int}$ as there are no monitoring wells

implemented in the field at the beginning of the first iteration.

3. The average concentration computation at the beginning of the current iteration

$ITR$ is $Cavg_{M_{ITR}}^{ITR}$ and is given by equation (4.16).

$$Cavg_{M_{ITR}}^{ITR} = \frac{1}{M_{ITR}} \sum_{k=1}^{M_{ITR}} Cobs_k^{ITR} \qquad\qquad (4.16)$$

where:

$Cobs^{ITR}$ is the observed pollutant concentration value at all available monitoring

wells in the study area measured at the beginning of the current iteration $ITR$.

4. The current observed pollutant concentration data, $Cobs^{ITR}$, from the total

number of wells available in the field $M_{ITR}$ is used to interpolate by kriging the

concentration values, $Ckrig_{i,j}^{ITR}$, at all the grid locations $i, j$ for current iteration

$ITR$. $Ckrig_{i,j}^{ITR}$ value corresponds to the $C_{i,j}^*$ variable value described in the

optimization model for monitoring network design in equation (4.11).

5. The variance of the distribution of the local probability density function for the

current iteration of interpolated concentration, $Ckrig_{i,j}^{ITR}$, is calculated for all the

nodes $i, j$ in the aquifer study area and is denoted by $Var_{i,j}^{ITR}$. Data uncertainty is

defined as the variance of Gaussian noise distribution associated with each data measurement.

6. The interpolated concentration is set to zero for those grid locations $i, j$ where the variance of the local probability density function is greater than or equal to one.

   Set $Ckrig_{i,j}^{ITR} = 0$ if $Var_{i,j}^{ITR} \geq 1$.

7. At the beginning of any current iteration $ITR$, constraint $\varepsilon_{\min}$ in equation (2) is updated and is equal to the computed value of $Cavg_{M_{ITR}}^{ITR}$. The optimization model described by equation (4.11) is solved to find new monitoring well locations $i, j$.

8. The new monitoring wells are implemented at grid locations given by the optimization model. Concentration measurement data $Cobs^{ITR}$ from these newly implemented wells and already existing wells in the study is collected and used as input at the beginning of the next iteration.

9. The average concentration value $Cavg_{M_{ITR}}^{ITR}$ at the beginning of current iteration $ITR$ is calculated as in step 3 and compared to the average concentration value from the previous iteration: $ITR - 1$ given by $Cavg_{M_{ITR-1}}^{ITR-1}$.

   If $Cavg_{M_{ITR}}^{ITR} > Cavg_{M_{ITR-1}}^{ITR-1}$ proceed to step 4.

10. If $Cavg_{M_{ITR}}^{ITR} \leq Cavg_{M_{ITR-1}}^{ITR-1}$ then the interpolated concentration $Ckrig_{i,j}^{ITR}$ is plotted over the entire grid locations $i, j$ for the iteration, $ITR$ to obtain a concentration contour profile.

11. The concentration contour profile is visually analysed to determine the plausibility of another unidentified source location in the polluted aquifer region.

| Step 1 | Collect pollutant concentration data from sparsely available observed concentration from randomly located, existing wells. |
|---|---|

| Step 2 | Calculate the average of observed concentrations from these existing wells. |
|---|---|

| Step 3 | Interpolate the pollutant concentration data for the entire aquifer region by kriging the observed pollutant concentration data obtained from existing wells and implemented wells. |
|---|---|

| Step 4 | Calculate the error variance for the interpolated pollutant concentration. |
|---|---|

| Step 5 | Set the interpolated concentration values to zero where the variance of the distribution is greater than or equal to one. |
|---|---|

| Step 6 | Run the optimal monitoring design model to choose locations of the monitoring wells for next sequence, setting the constraint to choose well locations having concentration value greater than the average concentration. |
|---|---|

| Step 7 | Implement the monitoring wells in the field and collect pollutant concentration data from the implemented monitoring wells and randomly located, existing wells. |
|---|---|

| Step 8 | Calculate the new average value of observed concentrations from randomly located, existing wells and the implemented monitoring wells. |
|---|---|

| Step 9 | Compare the new average concentration values with the average value of the previous iteration. If the new average is greater than the old average go to step 3 or else go to step 10. |
|---|---|

| Step 10 | Plot the kriged pollutant concentration profile for the entire aquifer region and look for traces of another possible source. |
|---|---|

| Step 11 | If no trace of another source is evident, terminate the process, or else eliminate all wells that have no impact on the traces of new source by looking at the concentration contour plot. |
|---|---|

| Step 12 | Recalculate the new average value of observed concentrations for the remaining wells after elimination process. |
|---|---|

| Step 13 | Interpolate the pollutant concentration data for the entire aquifer region by kriging the observed pollutant concentration data obtained from these remaining wells and go to step 4. |
|---|---|

**Figure 4.2 Flowchart of the Methodology for Pollution Source Locations Identification**

12. If no additional source location is apparent, the method is terminated. Otherwise, the average concentration value, $Cavg_{M_{ITR}}^{ITR}$, is recalculated after eliminating wells from the observation locations which do not have any bearing on the traces of another possible source. The number of wells eliminated in this process is denoted by $M^{eli}$, and is based on the visual interpretation of the plume contours.

13. The total number of wells to be considered for interpolating the concentration, $Ckrig_{i,j}^{ITR}$, at all the grid locations $i$, $j$, and the average concentration value, $Cavg_{M_{ITR}}^{ITR}$, in the beginning of current iteration is given by

$$M_{ITR} = M^{int} + \sum_{itr=1}^{ITR-1} M_{itr}^{imp} - \sum M^{eli} \qquad (4.17)$$

14. Steps 3 to 13 are repeated until the method terminates as no other source location appears possible.

### 4.1.2. Performance Evaluation of Developed Methodology

The performance of the developed methodology is evaluated for an illustrative study area 2100 metres wide and 1600 metres long. Five different scenarios are evaluated as shown in table 4.1. The scenarios differ in terms of number of sources and the initial number of pollutant observation wells. These scenarios represent various degrees of complexity in locating the pollution sources. The source location problem in general increases in complexity as the number and proximity of sources increases. Figures 4.3 to 4.7 show the plan view of the illustrative study area used in scenarios 1 to 5 respectively. Grid cells marked with a red star represent the grid locations containing the actual pollutant source. Grid cells marked with a blue star are the grid locations of existing monitoring wells with

non-zero pollutant concentration measurement. Pollution in the aquifer is first detected in these arbitrarily located monitoring wells.



**Figure 4.3 Plan View of Study Area for Scenario 1**

| Scenario | Number of Sources of pollution | Plume Overlapping | Initial Number of Wells |
|---|---|---|---|
| Scenario 1 | 1 | NA | 3 |
| Scenario 2 | 2 | Low | 6 |
| Scenario 3 | 2 | Low | 6 |
| Scenario 4 | 2 | High | 6 |
| Scenario 5 | 4 | Very High | 11 |

**Table 4.1 Test Scenarios for Source Location Identification**

**Figure 4.4 Plan View of Study Area for Scenario 2**



**Figure 4.5 Plan View of Study Area for Scenario 3**

91

**Figure 4.6 Plan View of Study Area for Scenario 4**



**Figure 4.7 Plan View of Study Area for Scenario 5**

In scenario 1, only one actual source is present with three initial arbitrary pollutant concentration measurement locations. Figure 4.3 shows the plan view of the study area for scenario 1. In scenarios 2, 3 and 4 two sources of pollutants and six initial arbitrary pollutant concentration measurement locations are present in each of the scenarios. In scenario 2, shown in figure 4.4, and scenario 3, shown in figure 4.5, the pollution source locations are the same but the initial arbitrary pollutant concentration measurement locations differ. In scenario 4 the sources of pollution are much closer to each other, resulting in a higher degree of overlapping of the individual pollutant plumes originating from the respective sources (figure 4.6). In scenario 5, four actual sources of pollutant are present with a very high degree of overlapping of the pollutant plumes resulting from the individual sources considered. Eleven arbitrary pollutant concentration measurement locations are initially present in this scenario. Figure 4.7 shows the plan view of the study area for scenario 5.

### *4.1.2.1.    Solution Procedure for Source Location Identification*

Concentration measurement data from the initial arbitrary locations is first interpolated by Geostatistical kriging to obtain the pollutant concentration data at all the other grid locations in the study area. An open source MATLAB code mGstat, version 0.99 (Hansen, 2004) is used in this study for kriging computations. The kriged concentration data is set to zero for those grid locations where the variance of the distribution of the local probability density function is greater than or equal to one. A lower variance value for the local probability density function represents a higher reliability of the interpolated concentration data, and vice-versa. Setting the values of interpolated concentrations to zero for grid locations, where the variance of the local probability density function is

greater than or equal to one, ensures that grid locations with higher reliability of estimation are incorporated in the optimization model. SA is used as an optimization algorithm to solve the optimization problem. The optimization algorithm is solved using these interpolated concentration values to find locations for placing the new monitoring wells.

In the first iteration, the optimization problem is solved with the constraint $\varepsilon_{min}$, as specified in equation (4.13), as the average concentration value of the initial arbitrary pollutant concentration measurements. The new sequence of monitoring wells is implemented in the field and pollutant concentration is measured. The observed concentration measurement values from these new monitoring wells and from the pre-existing monitoring wells is used to spatially krig the values of concentrations for the next iteration. As every iteration of the optimization algorithm searches for locations with a higher concentration of pollution, the value of $\varepsilon_{min}$ increases with every iteration. This process is iteratively repeated, comparing the value of $\varepsilon_{min}$ for the current iteration with that of the previous iteration.

Whenever the value of $\varepsilon_{min}$ for the current iteration decreases as compared to the previous iteration, the kriged concentration contour profile is plotted and analysed to verify the possibility of another source location. If another source location is apparent, the wells having no impact on the other possible source are eliminated from the next iteration. $\varepsilon_{min}$ value and the interpolated concentration values are again calculated, using concentration measurement data from the remaining wells. This process is repeated until no other source location appears plausible.

### 4.1.2.2.    *Summary of Results of Source Location Identification*

In scenario 1, one actual source and three initial arbitrary observation wells are present. The grid location coordinates and the observed concentration ($i, j, Cobs$) are (15, 23, 589 mg/l), (14, 21, 8.759 mg/l) and (13, 20, 17.9 mg/l).  Figure 4.8 shows the kriged concentration contour profile plotted for the entire study area once the source location steps are completed and no further source locations appear evident. The model is able to identify the source location correctly, although only one of the initial wells shows some significant pollutant concentration. The concentration from the other two wells is negligible. The solution results correctly show that there is only one source and no possibility exists of individual plume due to separate sources.



**Figure 4.8 Kriged Concentration Contour Profile for Scenario 1**

In scenario 2, two actual sources are considered. The overlapping zone of the plumes from the individual sources is in the concentration range of approximately 100mg/l to 200mg/l.  The location coordinates and the observed concentration of the initial arbitrary

observation wells are (15, 24, 253 mg/l), (16, 15, 249 mg/l), (16, 26, 359 mg/l), (16, 28, 214 mg/l), (17, 17, 368 mg/l) and (17, 19, 332 mg/l). Although the distances of the initial wells from the actual sources are in the range of 350m to 600m, the proposed method is still able to identify the two sources accurately. Figure 4.9 shows the kriged concentration contour profile plotted for the entire study area once the source location steps are completed and no further source locations appear evident.



**Figure 4.9 Kriged Concentration Contour Profile for Scenario 2**

Scenario 3 is same as scenario 2, except that the initial observation wells are different. The location coordinates and the observed concentration of the initial arbitrary observation wells are (14, 25, 412 mg/l), (17, 14, 109 mg/l), (15, 27, 398 mg/l), (17, 27, 213 mg/l), (16, 18, 402 mg/l) and (18, 18, 321 mg/l). Figure 4.10 shows the kriged concentration contour profile plotted for the entire study area once the source location steps are completed and no further source locations appear evident. The results show that both of the source locations are identified, however one of the identified source locations is offset by one grid location in the *i* direction from the actual source location. It is found

that the number of monitoring wells required to identify the same sources as in scenario 2, starting from different initial arbitrary observation wells, differ.



**Figure 4.10 Kriged Concentration Contour Profile for Scenario 3**

Scenario 4 is designed to test a high degree of overlapping. The overlapping zone of the plumes from the individual sources is in the concentration range of approximately 500mg/l to 650mg/l. The location coordinates and the observed concentration of the initial arbitrary observation wells are (7, 26, 78.8 mg/l), (11, 23, 49.63 mg/l), (10, 27, 1520 mg/l), (12, 21, 899 mg/l), (13, 18, 267 mg/l) and (13, 26, 460 mg/l). The developed methodology is able to identify the two sources correctly. Figure 4.11 shows the kriged concentration contour profile plotted for the entire study area once the source location steps are completed and no further source locations appear evident.

**Figure 4.11 Kriged Concentration Contour Profile for Scenario 4**

In Scenario 5, four actual sources and eleven initial arbitrary observation wells are present. The location coordinates and the observed concentration of the initial arbitrary observation wells are (7, 27, 115 mg/l), (8, 26, 269 mg/l), (10, 24, 175 mg/l), (10, 25, 340 mg/l), (10, 30, 89.71 mg/l), (11, 29, 260 mg/l), (12, 26, 554 mg/l), (13, 20, 604 mg/l), (13, 24, 546 mg/l), (14, 15, 221 mg/l)  and (17, 20, 232 mg/l). Figure 4.12 shows the kriged concentration contour profile plotted for the entire study area once the source location steps are completed and no further source locations appear evident. The developed methodology is able to identify three of the four actual source locations present. The solution results show that there are only three source locations present and the possibility of a fourth source is not indicated. The actual source not identified by using the solution results is the source with the lowest magnitude of strength and is located close to the plumes from other sources having a higher magnitude of strength. The reason for missing the low strength source location may be that the concentrations resulting from this source are dominated by concentrations of larger magnitude resulting from the other sources.

**Figure 4.12 Kriged Concentration Contour Profile for Scenario 5**

| | Actual sources present | | Source locations identified by the model | | Total Number of monitoring wells implemented |
|---|---|---|---|---|---|
| | Source No. | Co-ordinate of the source $(i, j)$ | Source No. | Co-ordinate of the source $(i, j)$ | |
| Scenario 1 | 1 | 10, 19 | 1 | 10, 19 | 19 |
| Scenario 2 | 1 | 7, 16 | 1 | 7, 16 | 50 |
| | 2 | 8, 28 | 2 | 8, 28 | |
| Scenario 3 | 1 | 7, 16 | 1 | 8, 16 | 40 |
| | 2 | 8, 28 | 2 | 8, 28 | |
| Scenario 4 | 1 | 10, 19 | 1 | 10, 19 | 30 |
| | 2 | 8, 28 | 2 | 8, 28 | |
| Scenario 5 | 1 | 10, 19 | 1 | ? | 64 |
| | 2 | 8, 28 | 2 | 8, 28 | |
| | 3 | 12, 22 | 3 | 12, 22 | |
| | 4 | 13, 16 | 4 | 13, 16 | |

**Table 4.2 Summary of Source Location Identification Results**

The summary of the results in table 4.2 shows that in scenario 1, scenario 2, scenario 3 and scenario 4 all of the source locations are identified. In scenario 5, only 3 source

locations out of the actual 4 source locations are identified by the developed methodology. The accuracy of the source locations identified in scenario 1, scenario 2, scenario 4 and scenario 5 are 100 percent. In scenario 3 one of the identified source locations is offset by one grid in the *i* direction from the actual source location.

## 4.2. Optimal Monitoring Network Design Models for Efficient Identification of Unknown Groundwater Pollution Sources Incorporating Genetic Programming Based Monitoring

In this section, three separate methodologies for optimal monitoring network design are presented. The aim is to improve the accuracy of groundwater pollution source identification using concentration measurements from a designed optimal monitoring network. The proposed methodology combines the capability of GP, and linked simulation-optimization for recreating the flux history of the unknown conservative pollutant sources with a limited number of spatiotemporal pollution concentration measurements.

The most common problem encountered in remediation of a polluted aquifer is the accurate characterization of pollution sources in terms of location, magnitude and activity duration, utilizing a limited set of spatiotemporal pollutant concentration measurements. In scenarios where potential source locations and activity duration are known with a fair degree of certainty, a linked simulation-optimization based approach is often applied for recreating the flux release history of the sources. A large amount of observed pollutant concentration data spread over time and space is necessary for accurate source identification. However, long term monitoring over a large number of monitoring locations has budgetary constraints. Often, monitoring networks consist of arbitrarily

located single water supply wells, or a group of arbitrarily placed wells where pollution in the aquifer is first detected. Moreover, these monitoring locations may not be optimally located for accurately identifying the release history of unknown pollution sources. This study aims to address this aspect of efficient source identification, using designed monitoring networks.

In a real world problem the number of monitoring wells to be implemented is governed by budgetary constraints. Therefore, it is important that the monitoring locations are chosen such that the concentration measurements from these locations, when utilized in a source identification model, improve the accuracy of source identification results. To achieve this, three separate methodologies are proposed for the design of an optimal monitoring network aimed at improving the accuracy of source identification. The proposed methodologies use trained GP models to calculate the impact factors and the frequency factors of the sources on the candidate monitoring locations. These impact factors and frequency factors are utilized as design criteria to formulate three optimal monitoring network design models: (1) a heuristic design model using GP based impact factors, (2) a heuristic design model using GP based frequency factors, and (3) an impact factor based multi-objective optimal monitoring network design model. Concentration measurements from the designed monitoring networks can reduce the possibility of missing an actual source and decrease the degree of non-uniqueness in possible aquifer responses by utilizing the monitored data.

### 4.2.1. Genetic Programming Models for Impact Factor Assessment and Frequency Factor Assessment

Genetic programming is an evolutionary optimization algorithm based on the concepts of genetics and natural selection. A GP model is essentially a computer program that represents the mathematical relationship between dependent variables (output) and independent variables (input). GP optimises the parameter values of a given model structure within predefined parameter space to find a highly fit computer program that produces desired output for a particular set of inputs. In this study, highly fit computer programs describing the relationship between output values (pollutant concentration at candidate monitoring locations at any monitoring time step) and input (flux values of pollutant at potential pollutant source locations) are evolved using genetic programming.

Genetic programming has not been widely applied in groundwater resource management problems. However, potential applicability of GP in groundwater problems has been proposed due to the following reasons: (1) GP's ability to develop simple models with interpretability to overcome the curse of "black box" nature of data intensive models, (2) lesser number of parameters are used in GP models as compared to parallel neural network architectures, and (3) GP's ability to parsimoniously identify the significance of the modelling inputs. In this study, this feature has been exploited for the design of an optimal monitoring network. The impact (significance) of potential source fluxes on the resulting concentration measurement at a potential monitoring location is obtained using GP.

GP modelling starts with an initial population of randomly generated computer programs composed of functions and terminals consisting of arithmetic operations, programming operations, logical functions or domain specific functions. Principles of reproduction,

crossover and mutation are imitated to create offspring computer programs from the initial generation of programs. These randomly generated individual programs are then tested for a fitness measure in terms of how well they perform in the problem environment. The computer programs with the better fitness measure values are allowed to survive by passing over to the new generation. After performing reproduction, crossover and mutation on the population, the parent population is replaced by the offspring population. Each program in the new population is evaluated against the fitness measure and the process is repeated over many generations to obtain the best individual program.

Each GP model is ranked based on the $R^2$ fitness value. A chosen subset of best fitting GP model ($\mu$) is used to compute the impact factor of each input variable (for example, flux values of pollutant at potential pollutant source locations). The impact factor is described as a measure of how much an input variable accounts for the output result; a factor by which the result would differ if the variable was removed. This essentially implies that, if by removing a variable from the mathematical function (GP model) the output differs highly, then the removed variable has a high impact on the output and hence the impact factor of that variable will be high.

The impact factor of a potential source at any given monitoring location at any sampling time step is given by equation (4.18).

$$IF_{iob}^{S} = \sum_{t=1}^{nt} (I_{iob}^{St}) \qquad (4.18)$$

where:

$IF_{iob}^{S}$ is the impact factor of source $S$ on monitoring well location $iob$,

$I_{iob}^{St}$ is the impact factor of source $S$ on monitoring well location $iob$ at stress period $t$,

$nt$ is the total number of stress periods.

In order to compute the impact for an entire monitoring time horizon, the impact factor at a given monitoring location obtained from the GP model (equation 4.18) for each time step is summed over all the monitoring time steps. The total impact factor of a potential source at any given monitoring location for all sampling time steps is given by equation (4.19).

$$SumIF_{iob}^{S} = \sum_{k=1}^{nk}(IF_{iob}^{S})^{k} \tag{4.19}$$

where:

$SumIF_{iob}^{S}$ is the sum of the impact factors of a potential source $S$ at any given monitoring location $iob$ for $nk$ sampling steps,

$IF_{iob}^{S}$ is the impact factor of source $S$ on monitoring well location $iob$,

$nk$ is the total number of monitoring time steps.

The normalised sum of impact factor due to all the potential sources at any monitoring location for all sampling time steps is given by Equation (4.20).

$$SumIF_{iob}^{norm} = \sum_{S=1}^{nS}\frac{SumIF_{iob}^{S}}{\dfrac{1}{nob}\displaystyle\sum_{iob=1}^{nob}SumIF_{iob}^{S}} \tag{4.20}$$

where:

$SumIF_{iob}^{norm}$ is the normalised sum of impact factor at any monitoring location *iob* due to all the potential sources *nS* for all *nk* monitoring time steps,

$\dfrac{1}{nob}\displaystyle\sum_{iob=1}^{nob}SumIF_{iob}^{S}$ is the average impact factor due to a source *S* at all monitoring well locations *nob*,

*nob* is the total number of monitoring well locations.

The relative impact factor is defined as the difference between the impact factor of the source having the maximum impact factor and the sum of residual contributions from the remaining sources. If *nS* is the total number of sources then the relative impact factor is given by equation (4.21). Monitoring well locations having higher values of relative impact factor signifies that the monitoring well location is predominantly influenced by one potential source and the influence of the other potential sources is significantly low. This reduces the non-uniqueness in the solution that may arise due to the presence of multiple sources of pollution.

$$^{\mathrm{Re}l}SumIF_{iob} = Max\{SumIF_{iob}^{S}\} - ((\sum_{S=1}^{nS}(SumIF_{iob}^{S})) - Max\{SumIF_{iob}^{S}\}) \qquad (4.21)$$

where:

$^{\mathrm{Re}l}SumIF_{iob}$ is the relative impact factor as a measure of the impact factor of the maximum contributing potential source relative to the combined impact factor of the rest of the potential sources at a given monitoring well location *iob*,

*nS* is the total number of potential sources.

The normalised relative impact factor for a monitoring well location *iob* is given by equation (4.22).

$$^{\mathrm{Re}l}SumIF_{iob}^{norm} = \frac{^{\mathrm{Re}l}SumIF_{iob}}{\dfrac{1}{nS}\displaystyle\sum_{S=1}^{nS}SumIF_{iob}^{S}} \tag{4.22}$$

where:

$^{\mathrm{Re}l}SumIF_{iob}^{norm}$ is the normalized relative impact factor at monitoring well location *iob* for all potential sources,

$\dfrac{1}{nS}\displaystyle\sum_{S=1}^{nS}SumIF_{iob}^{S}$ is the average impact factor at monitoring well location *iob* for all potential sources.

The frequency factor describes the percentage of subset of GP models ($\mu$) that incorporate the input variable. This simply means that, in a subset of best GP models (chosen as per the $R^2$ fitness value), if an input variable is incorporated in most of the GP models, then that variable has more influence on output than the other input variables.

The frequency factor of a potential source at any given monitoring location at any sampling time step is given by equation (4.23).

$$FF_{iob}^{S} = \sum_{t=1}^{nt}(F_{iob}^{St}) \tag{4.23}$$

where:

$FF_{iob}^{S}$ is the frequency factor of source *S* on monitoring well location *iob*,

$F_{iob}^{St}$ is the frequency factor of source $S$ on monitoring well location $iob$ at stress period $t$,

$nt$ is the total number of stress periods.

Total frequency factor of a potential source at any given monitoring location for all sampling time steps is given by equation (4.24).

$$SumFF_{iob}^{S} = \sum_{k=1}^{nk} (FF_{iob}^{S})^{k} \qquad (4.24)$$

where:

$SumFF_{iob}^{S}$ is the sum of the frequency factors of a potential source $S$ at any given monitoring location $iob$ for $nk$ sampling steps,

$FF_{iob}^{S}$ is the frequency factor of source $S$ on monitoring well location $iob$,

$nk$ is the total number of monitoring time steps.

The normalised sum of frequency factors due to all of the potential sources at any monitoring location for all sampling time steps is given by equation (4.25).

$$SumFF_{iob}^{norm} = \sum_{S=1}^{nS} \frac{SumFF_{iob}^{S}}{\dfrac{1}{nob} \sum_{iob=1}^{nob} SumFF_{iob}^{S}} \qquad (4.25)$$

where:

$SumFF_{iob}^{norm}$ is the normalised sum of frequency factor at any monitoring location $iob$ due to all the potential sources $nS$ for all $nk$ monitoring time steps,

$\dfrac{1}{nob} \displaystyle\sum_{iob=1}^{nob} SumFF_{iob}^{S}$ is the average frequency factor due to a source $S$ at all

monitoring well locations $nob$,

$nob$ is the total number of monitoring well locations.

The relative frequency factor is defined as the difference between the frequency factor of the source having the maximum frequency factor and the sum of residual contributions from the remaining sources. Here, the term relative is used to show how one input variable performs with respect to the other input variables, in terms of influencing the outcome of the result. The same analogy has been applied to show the influence of one source with respect to the other sources on the pollutant concentration at a monitoring location. If $nS$ is the total number of sources then the relative frequency factor is given by equation (4.26). A monitoring well location having a higher value of relative frequency factor signifies that the monitoring well location is predominantly influenced by one potential source and the influence of the other potential sources is significantly low. This reduces the non-uniqueness in the solution that may arise due to the presence of multiple sources of pollution.

$$^{\mathrm{Re}l} SumFF_{iob} = Max\{SumFF_{iob}^{S}\} - ((\sum_{S=1}^{nS}(SumFF_{iob}^{S})) - Max\{SumFF_{iob}^{S}\}) \quad (4.26)$$

where:

$^{\mathrm{Re}l} SumFF_{iob}$ is the relative frequency factor as a measure of the frequency factor

of the maximum contributing potential source relative to the combined frequency

factor of the rest of the potential sources at a given monitoring well location $iob$,

$nS$ is the total number of potential sources.

The normalised relative frequency factor for a monitoring well location *iob* is given by equation (4.27).

$$^{\mathrm{Re}l}SumFF_{iob}^{norm} = \frac{^{\mathrm{Re}l}SumFF_{iob}}{\dfrac{1}{nS}\displaystyle\sum_{S=1}^{nS}SumFF_{iob}^{S}} \tag{4.27}$$

where:

$^{\mathrm{Re}l}SumFF_{iob}^{norm}$ is the normalized relative frequency factor at monitoring well location *iob* for all potential sources,

$\dfrac{1}{nS}\displaystyle\sum_{S=1}^{nS}SumFF_{iob}^{S}$ is the average frequency factor at monitoring well location *iob* for all potential sources.

### 4.2.2. Methodology for Optimal Monitoring Network Design Models for Source Identification

The impact factor and frequency factor are numerical representations of the influence of a potential pollutant source on a candidate monitoring location. Impact factor and frequency factor are used as separate design criteria for design of optimal monitoring networks for source identification. The aim of the monitoring network design model is to reduce the possibility of missing a pollution source. At the same time the designed monitoring network decreases the degree of non-uniqueness in the set of possible aquifer responses to subjected geochemical stresses.

The optimal monitoring network design model chooses: (1) monitoring well locations where combined influence of all the potential sources is high, and (2) monitoring well locations where the relative influence (with respect to other potential sources) of a

potential source is high over a chosen observation period. This is achieved by maintaining the right balance between the monitoring well locations having maximum normalized total impact/frequency factor from all the potential sources, and monitoring well locations having maximum normalized relative impact/frequency factor from an individual potential source. Choosing monitoring well locations with maximum normalized total impact/frequency factors ensures choosing well locations where the expected overlapping of plumes due to all potential sources is maximal. This reduces the possibility of missing any actual source. This also reduces the likelihood of choosing monitoring well locations where the influence of potential sources is significantly low. Choosing well locations with maximum normalized relative impact/frequency factor ensures that the influence of one of the potential sources is predominantly higher than the rest of the potential sources, and therefore reduces the degree of non-uniqueness in the solution. Hence, uncertainty of source location is addressed by objective (1) above, and non-uniqueness is addressed by objective (2).

Three separate design models are considered in this study. The first design model uses impact factor as a design criterion for selecting optimal monitoring well locations. In the second design model, frequency factor is used as a design criterion for selecting optimal monitoring well locations. Both design models select monitoring well locations where the combined influence of all the potential sources is high, and monitoring well locations where the relative influence (with respect to other potential sources) of a potential source is high. A heuristic approach is used in the two design models to choose the number of wells of each type. In the third design model, a multi-objective formulation for optimal monitoring network design is presented. The design model uses impact factor as the

design criterion. The Pareto-optimal solutions obtained from the two objective model is utilized to design a set of Pareto-optimal monitoring networks.

The proposed methodologies consist of two steps. In the first step, GP models are trained against a large set of data patterns comprising possible source flux history for all the potential sources as input, and corresponding aquifer responses at all potential monitoring locations and different monitoring time steps as the output. Based on the $R^2$ fitness value of the GP models, a subset of best GP models ($\mu$) is selected for computing the impact factor of a potential source, on a potential monitoring location at any monitoring time step. The impact factor and frequency is calculated for all candidate monitoring locations at each monitoring time step. These impact factors and frequency factors directly obtained from the GP models are further utilized to calculate the relative impact factor, relative frequency factor, total impact factor and total frequency factor.

In the second step, a linked simulation-optimization model for source identification is solved. SA is used as a solution algorithm for solving the optimization problem, which minimizes the difference between the simulated and measured pollutant concentrations at optimally chosen monitoring locations. The source identification model is solved using concentration measurements from different optimal monitoring networks. The schematic diagram illustrating the steps involved in the three methodologies is shown in figure 4.13.

**Figure 4.13 Schematic Diagram Illustrating Salient Steps in the Proposed Methodologies**

112

### 4.2.2.1. Heuristic Design Models using GP based Impact Factor and Frequency Factor

The heuristic model for designing a monitoring network for accurate identification of unknown pollution sources using the impact factor and frequency factor are given by equations (4.28) and (4.29), respectively.

$$Maximize \sum_{iob=1}^{n_1} SumIF_{iob}^{norm} f_{iob} + Maximize \sum_{iob=1}^{n_2} {}^{\mathrm{Re}l} SumIF_{iob}^{norm} f_{iob} \qquad (4.28)$$

$$Maximize \sum_{iob=1}^{n_1} SumFF_{iob}^{norm} f_{iob} + Maximize \sum_{iob=1}^{n_2} {}^{\mathrm{Re}l} SumFF_{iob}^{norm} f_{iob} \qquad (4.29)$$

where:

$$nob = n_1 + n_2 \qquad (4.30)$$

$nob$ is the total number of monitoring well locations.

$n_1$ is the number of monitoring wells where the total impact/frequency from all the sources are found to be maximum, and $n_2$ is the number of wells where the relative impact/frequency due to a source is found to be maximum. $n_1$, $n_2$ are heuristically determined integer values. These values are pre-determined and subject to user judgement and may vary as per site specific conditions.

$f_{iob}$ is the binary decision variable to select a monitoring well location. $f_{iob} \equiv \{1, 0\}$ such that when $f_{iob}$ value equal to 1 representing monitoring well to be selected, and zero otherwise.

### *4.2.2.2.    Multi-Objective Design Model using GP based Impact Factor*

The optimal monitoring network design model based on impact factor, as determined by a chosen subset of best fitting GP model ($\mu$), finds monitoring well locations with the following objectives: (1) finding well locations with maximum normalized total impact from all the potential sources, and (2) finding well locations with maximum normalized relative impact from an individual potential source over a chosen observation period. Finding well locations with maximum normalized total impact factor, is conflicting with the other objective of finding well locations with maximum normalized relative impact factor from an individual potential source.

A multi-objective optimization model is formulated for the design of an optimal monitoring network with above stated conflicting objectives. One of the objectives is traded off to improve the other objective and vice-versa. The two-objective optimization model is solved by maximizing one of the objectives subject to the other objective being defined as an implicit constraint. The number of monitoring wells to be selected is essentially governed by budgetary constraints. The two objectives of the multi-objective optimization model for optimal monitoring network design for accurate identification of unknown pollution sources are defined by equations (4.31) and equation (4.32) respectively.

$$F1 = Maximize \sum_{iob=1}^{nob} SumIF_{iob}^{norm} f_{iob} \qquad (4.31)$$

Objective function *F1* maximizes the normalised sum of impact factors due to all the potential sources at any monitoring location.

$$F2 = Maximize \sum_{iob=1}^{nob} {}^{\text{Re}l}SumIF_{iob}^{norm} f_{iob} \qquad (4.32)$$

Objective function *F2* maximizes the normalised relative impact factor due to a source at any monitoring location subject to constraint defined in equation (4.33).

$$\sum_{iob=1}^{nob} f_{iob} \leq \alpha \qquad (4.33)$$

where:

$\alpha$ is integer constant representing the maximum number of wells that can be chosen,

$f_{iob}$ is the binary decision variable to select a monitoring well location. $f_{iob} \equiv \{1, 0\}$, such that when $f_{iob}$ value is equal to 1 representing the monitoring well to be selected, and zero otherwise.

The two objective optimization model is solved using the constrained method. In the constrained method, one of the objective functions (*F1*) is maximized, constraining the minimum level of satisfaction of the second objective function (*F2*) as shown in equation (4.34).

$$\sum_{iob=1}^{nob} {}^{\mathrm{Re}l}SumIF_{iob}^{norm} f_{iob} - \lambda \geq 0 \qquad (4.34)$$

where:

$\lambda$ is the minimum level of satisfaction of the second objective function *F2* also termed as the trade-off constant.

Therefore, the resulting model can be solved iteratively as a single objective optimization model for different satisfaction levels of $\lambda$; thus, a Pareto-optimal solution set is generated. The second objective function can be specified as a new implicit constraint. The upper limit of $\lambda$ is defined by the new constraint of the maximum value of the second

objective function *F2* when solved as a single objective optimization (equation 4.35). The lower limit of $\lambda$ is the value of the second objective function *F2* corresponding to the maximum value of the first objective function *F1*, when the optimization model is solved as a single objective model with *F1* as the only objective (equation 4.36).

$$MaxF2 \geq \lambda \qquad (4.35)$$

$$F2_{MaxF1} \leq \lambda \qquad (4.36)$$

where:

$F2_{MaxF1}$ is the value of the objective function *F2* corresponding to the maximum value of the first objective function *F1* when solved as a single objective model.

All solutions obtained on a Pareto-optimal front correspond to a different monitoring network.

### 4.2.3. Source Identification Model

Source identification, in terms of reconstructing the release history of an unknown pollution source, is solved using a linked simulation-optimization approach. The linked simulation optimization model simulates the physical process of flow and solute transport within the optimization model. The flow and solute transport simulation models are treated as important binding constraints for the optimization model. Therefore, any feasible solution of the optimization model needs to satisfy the flow and transport simulation model. A three-dimensional numerical model MODFLOW and a three-dimensional modular pollutant transport model MT3DMS, given by equations (4.37) and equation (4.38), respectively are used for simulation. The flow simulation model

MODFLOW and solute transport simulation model MT3DMS are explained in detail in Chapter 3, sections 3.3.2.1 and 3.3.2.2 respectively. The advantage of this approach is that it is possible to link any complex numerical model to the optimization model. In this study, flow and transport simulation models are linked to the optimization model using the SA algorithm for solution.

$$\frac{\partial}{\partial x}\left(K_{xx}\frac{\partial h}{\partial x}\right)+\frac{\partial}{\partial y}\left(K_{yy}\frac{\partial h}{\partial y}\right)+\frac{\partial}{\partial z}\left(K_{zz}\frac{\partial h}{\partial z}\right)\pm W = S_s\frac{\partial h}{\partial t} \qquad (4.37)$$

where:

$K_{xx}$, $K_{yy}$ and $K_{zz}$ represent the values of hydraulic conductivity along the $x$, $y$ and $z$ axes, respectively ($LT^{-1}$),

$h$ is the potentiometric head (L),

$W$ is the volumetric flux per unit volume where positive sign (+) means sources and negative sign (-) means sinks ($T^{-1}$),

$S_s$ is the specific storage of the porous material ($L^{-1}$),

$t$ is time (T),

$x$, $y$ and $z$ are the Cartesian co-ordinates (L).

$$\frac{\partial C}{\partial t} = \frac{\partial}{\partial x_i}\left(D_{ij}\frac{\partial C}{\partial x_j}\right)-\frac{\partial}{\partial x_i}(v_i C)+\frac{q_s}{\theta}C_s+\sum_{k=1}^{N}R_k \qquad (4.38)$$

where:

$C$ is the concentration of pollutants dissolved in groundwater ($ML^{-3}$),

$t$ is time (T),

$x_i$ is the distance along the respective Cartesian coordinate axis (L),

$D_{ij}$ is the hydrodynamic dispersion coefficient tensor ($L^2T^{-1}$),

$v_i$ is the seepage or linear pore water velocity ($LT^{-1}$); it is related to the specific discharge or Darcy flux through the relationship, $v_i = q_i/\theta$,

$q_s$ is volumetric flux of water per unit volume of aquifer representing fluid sources (positive) and sinks (negative) ($T^{-1}$),

$C_s$ is the concentration of the sources or sinks ($ML^{-3}$),

$\theta$ is the effective porosity of the porous medium (dimension less),

$\sum_{k=1}^{N} R_k$ is chemical reaction term for each of the $N$ species considered ($ML^{-3}T^{-1}$).

In the source identification problem, temporal pollutant source fluxes from all potential sources, represented by the term $q_sC_s$ in the transport equation (4.38), are the explicit unknowns. The strategy for estimating these unknown pollutant source fluxes is to generate candidates of these unknown variables in an optimization algorithm, use these values for simulations of flow and transport models, compute the difference between simulated and observed pollutant concentrations, and finally obtain an optimal solution that minimizes the difference between observed and simulated values. The model for optimal identification of the unknown polluted sources is represented by the following objective function and constraint:

$$MinimizeF = \sum_{k=1}^{nk} \sum_{iob=1}^{nob} Abs(cest_{iob}^{k} - cobs_{iob}^{k}) \tag{4.39}$$

subject to

$$cest_{iob}^{k} = f(q_s, C_s), \text{ where } f(q_s, C_s, cest_{iob}^{k}) \quad \forall \, iob, k \tag{4.40}$$

The set of constraints essentially represents the linking of the optimization algorithm with the numerical groundwater flow and transport simulation model through the decision variable,

where:

$q_s C_s$ is the pollutant source fluxes ($ML^{-3}T^{-1}$),

$q_s$ is the volumetric flux of water per unit volume of aquifer ($T^{-1}$),

$C_s$ is the concentration of the sources or sinks ($ML^{-3}$),

*Abs* (.......) represents the absolute value,

$cest_{iob}^{k}$ is the concentration estimated by the transport simulation model at observation monitoring location *iob* and at the end of time period $k$ ($ML^{-3}$),

*nk* is the total number of monitoring time steps,

*nob* is the total number of observation wells,

$cobs_{iob}^{k}$ is the observed concentration at monitoring location *iob* and at the end of time period $k$ ($ML^{-3}$).

### 4.2.4. Performance Evaluation of Developed Methodologies for Efficient Source Flux Identification

To evaluate the performance of the proposed monitoring network design methodology for accurate identification of the pollution sources using linked simulation-optimization, a hypothetical homogeneous, isotropic and saturated aquifer is utilized as an illustrative example as shown in figure 4.14. Cells marked with a red star represent the grid locations containing a potential pollutant source S($i$), where $i$ represents the source number. Cells marked with a green circle are the grid locations containing a potential monitoring well. A groundwater flow and solute transport model is simulated with hydrogeological parameters as given in table 4.3. The synthetic concentration measurement data used for the specified polluted aquifer facilitates evaluation of the



**Figure 4.14 Plan view of the Illustrative Study Area**

methodology without having to account for the unknown reliability of any field data. The performance evaluation procedure for the developed methodology is schematically shown in figure 4.15.



**Figure 4.15 Schematic Illustration of the of Performance Evaluation Process**

| Parameter | Unit | Value |
|---|---|---|
| Maximum length of study area | m | 2100 |
| Maximum width of study area | m | 1950 |
| Saturated thickness, $b$ | m | 30 |
| Grid spacing in $x$-direction, $\Delta x$ | m | 50 |
| Grid spacing in $y$-direction, $\Delta y$ | m | 50 |
| Grid spacing in $z$-direction, $\Delta z$ | m | 30 |
| Hydraulic conductivity, $K$ | m/d | 20 |
| Effective porosity, $\theta$ | | 0.3 |
| Longitudinal Dispersivity, $\alpha L$ | m/d | 20 |
| Transverse Dispersivity, $\alpha T$ | m/d | 10 |
| Horizontal Anisotropy | | 1 |
| Initial contaminant concentration | g/lit | 0 |
| Diffusion Coefficient | | 0 |
| Contaminant Flux | g/s | 0-100 |
| Source Grid Location S1, S2 and S3 | | S1(8,28), S2(12,22) S3(7,16) |

**Table 4.3 Hydrogeological Parameter Values Specified for Study Area**

The activity duration of the sources is specified for three equal duration stress periods of 500 days each. The pollutant flux from each of the sources is assumed to be constant over a stress period. The pollutant flux from each of the sources is represented as S($i$)($j$), where $i$ represents the source number and $j$ represents the stress period number. A total of nine source fluxes (S11, S12, S13, S21, S22, S23, S31, S32 and S33) are considered as explicit variables in the optimization problem, representing three potential source locations and three active stress periods. The concentration measurements are simulated for 4000 days after the start of the simulation. Pollutant concentration measurements at the potential location start at time $t$ = 1600 days and are taken after every 200 days at all the potential monitoring locations until $t$ = 2200 days. Figure 4.16 and figure 4.17 show the piezometric head profile and the pollutant concentration profile in the study area respectively.



**Figure 4.16 Piezometric Head Profile for the Study Area**

**Figure 4.17 Pollutant Concentration Profile for the Study Area**

### 4.2.4.1. *Genetic Programming Impact Factor and Frequency Factor Evaluation*

The impact factor and the frequency factor are calculated for three sources at 25 potential monitoring well locations (W1 to W25 as shown in figure 4.14). The input data consists of sets of flux values for each of the nine source fluxes (S11, S12, S13, S21, S22, S23, S31, S32 and S33) representing three sources and three active stress periods. The corresponding output data consists of the resulting pollutant concentration measurement due to these source fluxes at all 25 potential monitoring wells at times $t = 1600$, $t = 1800$, $t = 2000$ and $t = 2200$ days. 3000 data patterns comprising inputs and the corresponding outputs are used in the GP models. Out of total data patterns 50 percent are used for training, 40 percent for validation, and the remaining 10 percent for testing. A Latin Hypercube distribution (MATLAB R2010b) was used for generating the random flux

values ranging between 0g/s and 100g/s, as the input. The corresponding output data was simulated using MT3DMS code.

Discipulus$^{TM}$ 5.1 (RML Technologies, Inc.) is used for training, validation and testing for GP models. In this performance evaluation, based on $R^2$ fitness value, top 30, ($\mu = 30$) GP models are used for computing the impact factor and frequency factor. The impact factor and frequency factor for all potential monitoring locations at every sampling time step is calculated likewise. These impact and frequency factors directly obtained from the GP models is used to calculate the normalized relative impact/frequency factor $^{\mathrm{Re}l} SumIF_{iob}^{norm} / ^{\mathrm{Re}l} SumFF_{iob}^{norm}$ and normalised sum of impact/frequency factor due to all potential sources $SumIF_{iob}^{norm} / SumFF_{iob}^{norm}$ at any monitoring location (equations 4.18 to 4.27). The normalised sum of impact/frequency factor and the relative impact/frequency factor due to all the sources is calculated for all the potential monitoring locations. The normalised sum of impact factor and frequency factor and the relative impact factor and frequency factor due to all sources, calculated for all potential monitoring locations is shown in table 4.4.

| Well ID | Grid Location Y | Grid Location X | Impact Factor | | Frequency Factor | |
|---|---|---|---|---|---|---|
| | | | Normalised Sum of Impact Factor | Normalized Relative Impact Factor | Normalised Sum of Frequency Factor | Normalized Relative Frequency Factor |
| W1 | 10 | 21 | 3.55 | 3.26 | 9.88 | 1.72 |
| W2 | 12 | 23 | 4.62 | 0.24 | 11.9 | 0.86 |
| W3 | 13 | 19 | 4.24 | 3.94 | 12.19 | 4.51 |
| W4 | 13 | 13 | 4.99 | 2.85 | 9.7 | 3.28 |
| W5 | 13 | 32 | 3.96 | 3.33 | 11.28 | 2.84 |
| W6 | 14 | 24 | 4.25 | 3.52 | 14.28 | -1.74 |
| W7 | 15 | 27 | 4.13 | 4.13 | 11.26 | 6.48 |

| | | | | | |
|---|---|---|---|---|---|
| W8 | 15 | 16 | 3.92 | 3.92 | 11.11 | 4.41 |
| W9 | 17 | 19 | 4.90 | 4.82 | 17.29 | 3.25 |
| W10 | 17 | 24 | 5.13 | 1.81 | 21.53 | -4.53 |
| W11 | 17 | 29 | 4.29 | 4.18 | 14.39 | 4.39 |
| W12 | 18 | 22 | 5.31 | 0.07 | 22.51 | -2.51 |
| W13 | 19 | 26 | 4.70 | 4.55 | 13.77 | 6.29 |
| W14 | 19 | 15 | 5.30 | 4.90 | 20.87 | 0.67 |
| W15 | 19 | 30 | 5.59 | 5.59 | 14.69 | 7.03 |
| W16 | 20 | 20 | 6.80 | 3.52 | 20.1 | 2.84 |
| W17 | 21 | 23 | 5.04 | 3.68 | 31.63 | -10.37 |
| W18 | 22 | 16 | 6.19 | 6.07 | 16.1 | 6.5 |
| W19 | 22 | 19 | 6.19 | 6.07 | 20.26 | 1.46 |
| W20 | 22 | 26 | 6.24 | 5.92 | 20.51 | 2.15 |
| W21 | 24 | 23 | 6.29 | 5.21 | 22.49 | -1.57 |
| W22 | 24 | 30 | 5.61 | 5.56 | 14.71 | 5.55 |
| W23 | 25 | 27 | 5.01 | 4.94 | 13.02 | 3.14 |
| W24 | 26 | 18 | 4.08 | 3.74 | 14.54 | 0.52 |
| W25 | 26 | 22 | 5.79 | 3.25 | 17.18 | 0.94 |

**Table 4.4 Impact and Frequency Factor Values for Potential Monitoring Locations**

### 4.2.4.2. *Designed Monitoring Network and Source Identification Evaluation*

The heuristic monitoring design model is solved using normalised impact factor and frequency factor values as shown in table 4.4. In the heuristic design approach, a total number of six monitoring wells (*nob*) is implemented with values of $n_1$ and $n_2$ as 3 for both. This is achieved by maximizing the sum of impact/frequency factors due to all potential sources $SumIF_{iob}^{norm}$ / $SumFF_{iob}^{norm}$ for $n_1$ monitoring wells and maximizing the normalized relative impact/frequency factors $^{Rel}SumIF_{iob}^{norm}$ / $^{Rel}SumFF_{iob}^{norm}$ for $n_2$ monitoring wells, respectively. The optimization model defined by equation (4.28) and equation (4.29) is solved and two designed monitoring networks, DMNIF and DMNFF, are obtained using impact factor and frequency factor respectively, as the design criteria.

The optimal monitoring network design model is solved using normalised relative impact factor values and a normalised sum of impact factor values as inputs (table 4.4). Twelve different Pareto-optimal solutions are obtained as different pairs of *F1* and *F2* values. The value of the minimum satisfaction level of the second objective function *F2* varies from a minimum -1.7 to a maximum of 8.06 (equation 4.31 and equation 4.32). Each of the 12 solutions on the Pareto-optimal front represents different Pareto-optimal monitoring networks represented as MN1 to MN12 respectively, for the corresponding values of objective functions *F1* and *F2* (table 4.5). A total of 6 monitoring wells are selected for each Pareto-optimal solution.

| Monitoring Network | Objective Function Value F1 | Objective Function Value F2 |
|---|---|---|
| MN1 | 8.57 | -1.75 |
| MN2 | 8.49 | -1.0 |
| MN3 | 8.23 | 0.0 |
| MN4 | 8.04 | 1.0 |
| MN5 | 7.65 | 2.0 |
| MN5 | 7.37 | 3.0 |
| MN7 | 7.01 | 4.0 |
| MN8 | 6.63 | 5.0 |
| MN9 | 6.02 | 6.0 |
| MN10 | 5.59 | 7.0 |
| MN11 | 5.08 | 8.0 |
| MN12 | 5.01 | 8.07 |

**Table 4.5 Pareto-optimal Monitoring Networks**

A linked simulation-optimization model as represented by objective function (equation 4.39) and constraint set (4.40) is solved to identify the pollution sources and to evaluate the performance of the proposed methodologies. The two heuristically designed monitoring networks DMNIF and DMNFF, based on impact factor and frequency factor respectively, and 12 Pareto-optimal monitoring networks (MN1to MN12) are used as sampling locations. Concentration measurements from each of the monitoring networks

are used to estimate the pollution sources' flux release history. These evaluation results using concentration observations from the 12 Pareto-optimal monitoring networks are compared to find the most efficient monitoring network design based on the trade-off between the two objectives.

The observed aquifer responses are simulated by solving MODFLOW (equation 4.37) and MT3DMS (equation 4.38) in GMS7.0, along with appropriate initial and boundary conditions. The resulting concentrations are then perturbed to represent the effects of sampling measurement errors. The observed pollutant concentration data is perturbed with random measurement error with a maximum deviation of 10 percent of the actual observed concentration *cobs* as shown in equation (4.41).

$$^{Pert}cobs_{iob}^{k} = cobs_{iob}^{k}(1+err) \tag{4.41}$$

$$err = \mu per \times rand \tag{4.42}$$

where:

$^{Pert}cobs_{iob}^{k}$ is the perturbed numerically simulated concentration value,

$cobs_{iob}^{k}$ is the numerically simulated concentration value,

*err* is the error term,

*μper* is the maximum deviation expressed as a percentage,

*rand* is a random fraction between -1 and +1generated using Latin Hypercube distribution.

To show the improved efficiency of the source identification model, when using concentration measurements from designed monitoring networks over arbitrary networks, the linked simulation-optimization model is solved using concentration measurement data from designed monitoring networks and ten arbitrary monitoring networks. Ten arbitrary monitoring networks represented as ARMN1 to ARMN10, are chosen from all of the 33 (W1 to W33 as shown in figure 4.14) potential monitoring locations. A total of six monitoring wells are selected in each of the arbitrary monitoring networks. The wells are numbered from 1 to 33 and a random number generator is used for selecting the wells in these arbitrary monitoring networks. In all scenarios, source fluxes are first estimated using error-free concentration measurement data, and then with concentration measurement perturbed with random errors.

### 4.2.5. Discussion of Evaluation Results for the Developed Methodologies for Efficient Source Flux Identification

Source flux identification, using a linked simulation-optimization model for heuristically designed monitoring networks DMNIF and DMNFF and twelve Pareto-optimal monitoring networks (MN1 to MN12) obtained as a solution of the multi-objective monitoring network design model and ten arbitrary monitoring networks ARMN1 to ARMN10, is evaluated. The evaluation results of souse flux identification obtained using concentration measurements from designed monitoring networks and arbitrary monitoring networks is compared. These evaluation results using error-free and erroneous concentration measurement data from designed and arbitrary monitoring networks are discussed in this section.

### 4.2.5.1. *Source Flux Identification Results: Heuristically Designed Optimal Monitoring Networks*

Source flux identification, using a linked simulation-optimization model for two designed monitoring networks DMNIF and DMNFF, and ten arbitrary monitoring networks ARMN1 to ARMN10, is evaluated. These evaluation results using error-free and erroneous concentration measurement data from designed and arbitrary monitoring networks are shown in figures 4.18 and 4.19 respectively. Each of the unknown fluxes (S11, S12, S13, S21, S22, S23, S31, S32 and S33) is marked on the *x* axis. Each of the bars corresponding to an unknown source flux shows the actual flux value, estimated flux values using arbitrary monitoring networks, and estimated flux values using designed monitoring networks.

Source Flux Identification using Designed and Arbitrary Monitoring Networks



Source Fluxes for all the sources for different stress Periods

**Figure 4.18 Source Flux Identification Results using Error-free Concentration Measurements**

Figure 4.18 and figure 4.19 show that the performance of the source identification methodology is highly impacted by the monitoring locations utilized for concentration measurements. For example, data obtained from the arbitrary monitoring network ARMN2 appears to result in more accurate source identification when compared to those obtained using arbitrary monitoring network ARMN10. However, as can be seen in figure 4.18, the designed monitoring networks DMNIF and DMNFF result in more accurate source identifications using error-free concentration measurement data. As expected, where the concentration measurements are erroneous, the source identification becomes relatively less accurate, as shown in figure 4.19. This result is consistent for both designed monitoring networks, as well as arbitrary monitoring networks.

Source Flux Identification using Designed and Arbitrary Monitoring Networks with Erroneous Data



**Figure 4.19 Source Flux Identification Results using Erroneous Concentration Measurements**

The average of the absolute differences between the actual source fluxes and the estimated source fluxes for arbitrary networks ARMN1 to ARMN10 is compared with the

absolute difference between the actual source fluxes and the estimated source fluxes for designed monitoring networks DMNIF and DMNFF, for both error-free concentration data and erroneous data as shown in figures 4.20 and 4.21 respectively.



Figure 4.20 Comparison of Average Absolute Error for Arbitrary Networks vs. Absolute Error for Designed Monitoring Networks in case of Error-free Data

Figure 4.21 Comparison of Average Absolute Error for Arbitrary Networks vs. Absolute Error for Designed Monitoring Networks in case of Erroneous Data

It is evident from figure 4.20, that source flux identification errors, when using observed concentration data from the designed monitoring network DMNIF and DMNFF are smaller than those obtained using arbitrary monitoring networks ARMN1 to ARMN10. It can be seen that the results of identification using monitoring network DMNIF and DMNFF are consistently better for all of the source flux estimations in the case of error-free concentration measurement data. However, in the case of erroneous concentration measurement data, estimated source flux shows a large deviation from the actual source flux values in the case of designed and arbitrary monitoring networks (figure 4.21). It is also seen that for some of the arbitrary networks, few of the source flux estimates are better than those obtained using designed networks.

While using erroneous data, the average absolute difference between all actual source fluxes and estimated fluxes obtained using the arbitrary networks is generally larger

compared to that of designed monitoring networks, except for source one. These comparisons using erroneous data show that, in the case of flux estimates for source one (S11, S12 and S13), designed network DMNIF performs worse than the arbitrary networks. This is because the relative impact of source one is comparatively lower than that of the other two sources. As a result of this, no monitoring well is chosen in the vicinity of source one that can nullify the effect of non-uniqueness arising due to error in concentration measurement data. In all other cases designed monitoring networks show less error in estimation than the arbitrary networks.

### 4.2.5.2.    *Source Flux Identification Results: Multi-Objective Optimal Monitoring Networks*

The two-objective optimal monitoring network design model is solved and the first objective function values *F1* (equation 4.31) are plotted against the minimum satisfaction level of the second objective function value *F2* (equation 4.32), as shown in figure 4.22. Twelve Pareto-optimal monitoring networks (MN1to MN12) are chosen for different values of objective function *F2* (table 4.5).

## Pareto-Optimal Solutions



**Figure 4.22 Pareto-Optimal Solution Front**

The non-inferior solutions show the conflicting nature of the two objective functions. The Y axis represents the value of the objective function *F1* and the X axis represents value of objective function *F2*. It is seen that, as the value of *F2* (objective function two) decreases the value of *F1* (objective function one) increases and vice-versa. This essentially shows the conflicting nature of the two objective functions and their trade-off. Larger objective function values for the first objective function *F1* (equation 4.31) increase the likelihood of choosing monitoring locations where the combined impact of potential sources is high. It also reduces the possibility of missing an actual source as it chooses those locations where summation of impact of potential sources is large. Higher values of the second objective function *F2* (equation 4.32) increase the likelihood of choosing monitoring locations where relative impact from an individual potential source with respect to other sources is higher. This essentially results in reducing the non-

134

uniqueness due to overlapping of different pollutant plumes resulting from individual sources.

The results of source flux identification using a linked simulation-optimization model is compared for all the 12 Pareto-optimal monitoring networks (MN1 to MN12) obtained as solutions, using error-free and perturbed concentration measurement data. This comparison is shown in figures 4.23 and 4.24 respectively. Each of the unknown fluxes (S11, S12, S13, S21, S22, S23, S31, S32 and S33) is marked on the X axis. Each of the bars corresponding to an unknown source flux shows the actual flux value and estimated flux values using Pareto-optimal monitoring networks. The results of source flux identification for all 12 Pareto-optimal monitoring network (MN1 to MN12) designs are very close to the actual flux when solved using error-free measurement data. However, when the concentration measurement data is perturbed with random errors, the results of source flux identification show a greater amount of deviation from the actual flux values in all of the 12 Pareto-optimal monitoring networks (MN1 to MN12), when compared to the results with error-free concentration measurement data.

Source Flux values as Identified by Pareto-Optimal Monitoring
Networks with Non-Erroneous Concentration Measurement Data

**Figure 4.23 Source Identification Results using Error-free Measurement data from Pareto-Optimal Monitoring Networks**



Source Flux values as Identified by Pareto-Optimal Monitoring
Networks with Erroneous Concentration Measurement Data

**Figure 4.24 Source Identification Results using Erroneous Measurement data from Pareto-Optimal Monitoring Networks**

To choose the most efficient monitoring network out of the 12 Pareto-optimal monitoring networks (MN1 to MN12), the absolute difference between actual source fluxes and the estimated source fluxes for all 12 Pareto-optimal monitoring networks is calculated. Figures 4.25 and 4.26 show the absolute difference between actual source fluxes and the estimated source fluxes. These figures also show the average of the absolute differences for all unknown source fluxes plotted for each of the Pareto-optimal monitoring networks (MN1 to MN12) using error-free and erroneous measurement data. The absolute difference of actual source fluxes and estimated source fluxes for all 12 Pareto-optimal monitoring networks (MN1 to MN12) using error-free and erroneous measurement data show similar trends. The average of this difference shows a decreasing trend as the value of the second objective function *F2* is first increased and is minimum for monitoring network 5 (MN5). A further increase in the value of *F2* increases the average absolute difference between actual source fluxes and estimated source fluxes.



Figure 4.25 Absolute Difference between Actual Fluxes and Estimated Fluxes using Error-free Measurement Data from Pareto-Optimal Monitoring Networks

Absolute difference of Actual and Estimated fluxes using Pareto-Optimal Monitoring Locations for Erroneous Concetration Measurement Data

**Figure 4.26 Absolute Difference between Actual Fluxes and Estimated Fluxes using Erroneous Measurement Data from Pareto-Optimal Monitoring Networks**

These results show that, for efficient identification of unknown source fluxes, it is important to have the right balance between monitoring well locations where combined impact of individual potential sources are high and at the same time have monitoring well locations that reduce the non-uniqueness in the response of the aquifer system. These evaluations show that Pareto-optimal monitoring network 5 (MN5) results in a minimum average absolute difference between the actual source flux and the estimated source flux, both in case of error-free data, and erroneous measurement data. In any polluted aquifer with polluted sources, where individual polluted plumes from each of the sources do not overlap, absence of any monitoring well within any of these pollutant plumes will result in failure to identify the relevant sources. Hence, choosing monitoring locations with high potential source impact reduces the chance of missing a possible source.

The results of source flux identification using concentration measurements from ten arbitrary monitoring networks (ARMN1 to ARMN10) and monitoring network 5 (MN5) with error-free data and erroneous data are shown in figures 4.27 and 4.28 respectively. The estimated source flux values using the arbitrary networks are averaged and compared with the actual flux values and estimated flux values obtained using monitoring network 5 (MN5). Both error-free and perturbed measurement data are used, as shown in figures 4.27 and 4.28 respectively. It is seen that the results of source flux estimates when using monitoring network 5 (MN5) are better for all source fluxes except S22, when using erroneous measurement data.

Source Flux Identification using Arbitrary Monitoring Networks using Error Free Concentration Measurement Data



**Figure 4.27 Source Identification Results using Error-free Measurement Data from Arbitrary Monitoring Networks**

**Figure 4.28 Source Identification Results using Erroneous Measurement Data from Arbitrary Monitoring Networks**



**Figure 4.29 Comparison of Source Identification Results using Error-free Measurement Data**

Comparison of Source Flux Identification: Actual, Average of Arbitrary Monitoring Networks and MN5 using Erroneous Concentration Measurement Data

**Figure 4.30 Comparison of Source Identification Results using Erroneous Measurement Data**

The absolute difference between the actual flux and the estimated flux using the arbitrary networks is averaged and compared with the absolute difference between actual flux and the estimated flux from monitoring network 5 (MN5). Theses comparisons, using error-free and erroneous measurement data, are shown in figures 4.29 and 4.30 respectively. It is seen that the averaged absolute difference between actual source fluxes and estimated fluxes, obtained using arbitrary networks, is larger compared to that for monitoring network 5 (MN5). The estimated source flux for S22, obtained using erroneous measurement data has a larger averaged absolute error when using data from arbitrary monitoring networks. However, the average source flux estimates appear better with the arbitrary networks. This is because the negative and positive differences between the actual source flux and estimated source flux for S22, obtained using concentration measurements from the 10 arbitrary networks, cancel out each other, showing smaller apparent average estimation errors as seen in figure 4.30. The absolute differences between the actual source flux and estimated source flux for S22 are shown in figure

4.31. It is clear from figure 4.31 that arbitrary monitoring networks have a larger error of source flux estimation when compared with those obtained using the formally designed monitoring network, MN5.



Comparison of Absolute Difference between Actual and Estimated Source Fluxes: MN5 vs Averaged Absolute Difference for Arbitrary Networks with Error-Free Concentration Measurement Data

**Figure 4.31 Comparison of Averaged Absolute Error for Arbitrary Networks vs. Absolute Error for Pareto-Optimal Monitoring Network 5 (MN5) using Error-free Measurement Data**

Comparison of Absolute Difference between Actual and Estimated Source
Fluxes: MN5 vs Averaged  Absolute Difference for Arbitrary Networks
with Erroneous Concentration Measurement Data

**Figure 4.32 Comparison of Averaged Absolute Error for Arbitrary Networks vs. Absolute Error for Pareto-Optimal Monitoring Network 5 (MN5) using Erroneous Measurement Data**

Comparing the results of source flux estimation using concentration measurements from an optimally designed monitoring network and an arbitrary network may not justify the merits of such a designed monitoring network. However, when the same identification results are averaged for different arbitrary monitoring networks and compared, then the utility of such optimally designed monitoring networks becomes clear. Such designed monitoring networks are even more important in the case of aquifers polluted by multiple sources. Placing monitoring wells arbitrarily in such situations may yield good results in the case of some source fluxes, but the estimates may be very poor in the case of other sources. This is because the pollutant concentration measurements at these arbitrary observation locations are not representative of all of the sources present in the aquifer system, and the observed measurements may represent the effect of a subset of the sources only.

## 4.3. Conclusions

In this study, two methodologies aimed at optimal source identification are presented; (1) optimal monitoring network design models for source location identification, and (2) optimal monitoring network design models for efficient identification in terms of location, magnitude and time of activity of unknown groundwater pollution sources.

### 4.3.1. Optimal monitoring network design model for source location identification

In all real life groundwater pollution scenarios, the source characterization process generally starts with data collected from a few randomly located wells. These limited amounts of data are not sufficient for solving the source characterization problem. An accurate source characterization requires reliable knowledge of pollution source locations. Only then can a more formal source identification method be applied to characterise the source magnitude and activity duration. To address this problem a methodology has been developed for preliminary estimation of possible pollution source locations. This methodology can be utilized to improve the source characterization results.

The developed methodology, based on the sequential design of an optimal monitoring network and collecting concentration measurement data from such an implemented network, appears to perform satisfactorily for source location identification. In scenarios where overlapping of the pollutant plume from the individual sources is low, the developed methodology is able to identify all pollution source locations successfully. It is also found that the number of monitoring wells required increases with an increase in the distance of the source from the initial set of wells where the pollution is first detected.

The number of monitoring wells implemented also varies depending on the location of the initial set of wells where the pollution is first detected. This method also performs satisfactorily when the number of actual sources of pollution increases. However, larger numbers of monitoring wells need to be implemented for such scenarios.

The specific advantage of this method is that the flow of groundwater and transport of the pollutant in the groundwater system need not be modelled. This method can be applied to various groundwater pollution scenarios where some pollutant concentration is detected. This method appears to work well even with very small amount of initial pollution concentration measurement data. The solution of the methodology can provide information regarding plausible groundwater pollution source locations, which can be utilized by an optimal source characterization model to accurately estimate locations, magnitude and duration of activity of the unknown sources.

The proposed methodology in its current form is limited to pollutant sources that are continuous in time, although it may vary in magnitude with respect to time. Further work is necessary to incorporate scenarios with sources that are not continuous in time. The developed methodology may be further refined to reduce the number of monitoring wells required, and to integrate it with an optimal source characterization methodology to improve source characterization accuracy.

### 4.3.2. Optimal monitoring network design models for efficient identification of unknown groundwater pollution sources

Unknown groundwater pollution source identification models utilize spatiotemporal pollutant concentration measurement data for source flux identification. Accuracy of source flux identification results depends on the number of pollutant concentrations and

the spatiotemporal locations at which they are measured. However, in all real life scenarios, the number of such spatiotemporal locations at which pollutant concentrations are measured is limited due to budgetary constraints, and often have measurement errors.

In groundwater pollution source identification, the location and timing of pollutant concentration measurement data have a direct bearing on the accuracy of source identification results. Not all monitoring wells are ideally located for accurate identification of source fluxes. An optimal monitoring network design for source flux identification is a complex multi-objective problem. It requires the right balance between well locations where the impact of all potential sources is significantly high, reducing the possibility of missing an actual source, and well locations where non-uniqueness due to overlapping of pollutant plumes from the individual sources is less. These two conflicting goals are combined to form a two-objective optimal monitoring network design model solved using: (1) a heuristic design approach and (2) a multi-objective design approach. Pollutant concentration measurements from these monitoring locations, when used in source flux identification, can improve the accuracy of source identification results.

The developed heuristic design methodology uses a GP based monitoring location impact factor and frequency factor as separate design criteria for design of an optimal dedicated monitoring network. The developed methodology, based on GP impact factor and frequency factor for the design of an optimal dedicated monitoring network for improving source identification efficiency, appears to perform satisfactorily for efficient identification of unknown groundwater pollution source fluxes. However, from the limited solution results in the illustrative example problem, it cannot be concluded if impact factor or the frequency factor is a better design criterion for dedicated monitoring network design. A more exhaustive study is required to conclude which of the two design

criteria, impact factor or frequency factor, is  better suited for a monitoring network design for given site specific conditions. The solution results also show variation in the accuracy of source flux identification results with varying monitoring well locations.

The developed multi-objective design methodology based on GP impact factor for design of an optimal monitoring network appears to perform satisfactorily for efficient identification of unknown groundwater pollution source fluxes. The solution results in the illustrative example problem show that the accuracy of source flux identification varies when using pollutant concentration measurement data from different monitoring locations. The designed monitoring network results in better source identification compared to other arbitrary networks, both with error-free and erroneous measurement data.

The proposed methodology is shown to improve source identification efficiency as compared to arbitrary measurements. It shows that there is a trade-off in selecting optimal monitoring locations in terms of isolating the impact of individual potential sources and choosing locations which reduce the possibility of missing the impact of any of the potential sources. However, the ideal levels of trade-off needs to be studied and may depend on site specific conditions. The proposed methodology need to be expanded to incorporate parameter uncertainty. Also, the computational effort increases with the number of potential sources as well as with an increase in the number of candidate monitoring locations.

In all real world problems of source identification, the degree of uncertainty in terms of source locations and aquifer response to subjected geochemical stress is high. Increased monitoring data can reduce some of these uncertainties. Moreover, the number of

monitoring wells to be implemented for concentration measurement data is governed by budgetary constraints. This method can be applied in such polluted aquifer sites. This method can decrease such uncertainties by using a limited number of monitoring wells which otherwise would have to be reduced by implementing a large number of monitoring wells, resulting in capital loss. This method can increase the accuracy of source identification with concentration measurement data from a limited number of monitoring wells, as it reduces the non-uniqueness in the aquifer response to subjected geochemical stress by incorporating the impact factor. This would increase the efficiency of implementing a monitoring network by reducing the costs associated with arbitrary and partially redundant monitoring networks.

# 5. Feedback based Methodology for Efficient Identification of Unknown Pollutant Source Characteristics Integrating Sequential Monitoring Network and Source Identification Model

In this chapter, a sequential approach integrating a monitoring network design model with a source identification model is presented. This sequential methodology uses feedback information from a designed monitoring network to improve the accuracy of source identification results in every sequence of implementation. Performance of the developed methodology is evaluated for a real life like scenario using a synthetic study area. The efficiency of this method is demonstrated by comparing the solution results obtained using concentration measurements from arbitrarily chosen monitoring networks and a sequentially designed optimal monitoring network.

## 5.1. Integration of Source Identification Model and Sequential Monitoring Network Design Model

When pollution in a groundwater aquifer system is first detected, very little is known about the pollution source characteristics. Often, the pollutant is first detected in a random water supply well or a group of wells. The first concentration measurements showing detection of pollutants are insufficient for estimating the source characteristics. Moreover, these water supply wells may not be ideally located for efficient source identification. Unknown groundwater pollution source characterization requires large numbers of

spatiotemporal concentration measurements from ideally located monitoring wells for efficient source characterization.

The need is twofold: a large number of temporal concentration measurements spaced over a sufficiently long time period is required from ideally located monitoring wells. Taking several temporal concentration measurements at regular time intervals, spaced over a sufficiently long period, ensures that at least a part of the breakthrough curve (concentration versus time at monitoring location) is captured so that the source identification process can be initiated. Since the groundwater flow and pollutant transport is dynamic in nature, the spread of the pollutant continues while temporal measurements are being obtained at monitoring locations.

However, if the monitoring locations are not ideally suited for source characterization, the concentration measurements from such monitoring wells may have limited utility in accurate source characterization. Hence, it might not be worth waiting to get the temporal concentration measurements from a set of monitoring wells which may later prove to be less than optimal in identifying the source characteristics accurately. This process of late realization would lead to a loss of time, during which a much larger part of the aquifer will get polluted. This methodology aims to address this deficiency by integrating the source identification model and optimal monitoring network design model to optimize the efficiency of source identification.

This sequential methodology uses the source identification model to estimate source flux characteristics. These estimated source characteristics are used in the forward simulation of the pollutant transport model to predict the future pollutant concentration distribution. Figure 5.1 shows a schematic chart of this sequential methodology. This

**Figure 5.1 Flow Chart of Methodology Linking of Source Identification with Monitoring Network Design**

pollutant concentration distribution is utilized to obtain concentration gradients at different locations. The concentration gradient information is utilized to find optimal locations for implementing monitoring wells for concentration measurements at the next sampling time step. In the next sequence all available pollutant concentration measurements from the current sampling time step and previous sampling time steps are used again in the source identification model. With every sequence, more concentration measurements are obtained based on the prediction of the polluted plume as per previously estimated source characteristics. This sequential process adds additional targeted monitoring data resulting in sequential improvements in source flux estimates. Improved estimates of source flux values result in improved prediction of future pollutant concentration distribution, thus enhancing the utility of the designed monitoring network. Feedback information, in terms of observed concentration measurements from optimal monitoring wells, improves the accuracy of the source flux estimation obtained as a solution of the source identification model. Vice versa, a better source flux estimate improves the subsequent monitoring network design. This sequential methodology of source identification and optimal monitoring network design provides the feedback information to improve source characteristics estimates, at the same time optimizing the monitoring well network, ensuring that the observed concentration measurement data used in the source identification is more efficient in source flux identification.

### 5.1.1. Source Identification Model

In this methodology the source identification model is immediately applied upon detection of any pollution in the aquifer. A linked simulation-optimization approach, as discussed in Chapter 4 section 4.2.3., is utilized for reconstructing the source flux release

history of an unknown pollution source. The process of flow and solute transport is simulated within the optimization model such that the flow and solute transport simulation models are treated as important binding constraints for the optimization model. Therefore, any feasible solution of the optimization model needs to satisfy the flow and transport simulation models. A three-dimensional numerical model MODFLOW, and a three-dimensional modular pollutant transport model MT3DMS, incorporating the governing equations (5.1) and (5.2) respectively, are used for simulation. The flow simulation model MODFLOW and solute transport simulation model MT3DMS is explained in detail in Chapter 3, sections 3.3.2.1 and 3.3.2.2 respectively. SA is used as the solution algorithm to solve the optimization problem.

$$\frac{\partial}{\partial x}\left(K_{xx}\frac{\partial h}{\partial x}\right) + \frac{\partial}{\partial y}\left(K_{yy}\frac{\partial h}{\partial y}\right) + \frac{\partial}{\partial z}\left(K_{zz}\frac{\partial h}{\partial z}\right) \pm W = S_s\frac{\partial h}{\partial t} \qquad (5.1)$$

where

$K_{xx}$, $K_{yy}$ and $K_{zz}$ represent the values of hydraulic conductivity along the $x$, $y$ and $z$ coordinate axes ($LT^{-1}$),

$h$ is the potentiometric head (L),

$W$ is the volumetric flux per unit volume where positive sign (+) means sources and negative sign (-) means sinks ($T^{-1}$),

$S_s$ is the specific storage of the porous material ($L^{-1}$),

$t$ is time (T),

$x$, $y$ and $z$ are the Cartesian co-ordinates (L).

$$\frac{\partial C}{\partial t} = \frac{\partial}{\partial x_i}\left(D_{ij}\frac{\partial C}{\partial x_j}\right) - \frac{\partial}{\partial x_i}(v_i C) + \frac{q_s}{\theta}C_s + \sum_{k=1}^{N} R_k \qquad (5.2)$$

where:

$C$ is the concentration of pollutants dissolved in groundwater (ML$^{-3}$),

$t$ is time (T),

$x_i$ is the distance along the respective Cartesian coordinate axis (L),

$D_{ij}$ is the hydrodynamic dispersion coefficient tensor (L$^2$T$^{-1}$),

$v_i$ is the seepage or linear pore water velocity (LT$^{-1}$); it it is related to the specific discharge or Darcy flux through the relationship, $v_i = q_i / \theta$,

$q_s$ is volumetric flux of water per unit volume of aquifer representing fluid sources (positive) and sinks (negative) (T$^{-1}$),

$C_s$ is the concentration of the sources or sinks (ML$^{-3}$),

$\theta$ is the effective porosity of the porous medium (dimension less),

$\sum_{k=1}^{N} R_k$ is chemical reaction term for each of the $N$ species considered (ML$^{-3}$T$^{-1}$).

In the source identification problem, temporal pollutant source fluxes from all potential sources, represented by the term $q_s C_s$ in the transport equation (5.2), are the explicit unknowns. The strategy for estimating these unknown pollutant source fluxes is to generate candidates of these unknown variables in an optimization algorithm, use these values for simulations of flow and transport models, compute the difference between

simulated and observed pollutant concentrations, and finally obtain an optimal solution that minimizes the difference between observed and simulated values. The optimization task is performed by the SA optimization algorithm (Chapter 2, section 2.4.2). The model for optimal identification of the unknown polluted sources is represented by the following objective function and constraint.

$$Minimize F = \sum_{k=1}^{nk} \sum_{iob=1}^{nob} Abs(cest_{iob}^{k} - cobs_{iob}^{k}) \qquad (5.3)$$

subject to:

$$cest_{iob}^{k} = f(q_s, C_s), \text{ where } f(q_s, C_s, cest_{iob}^{k}) \quad \forall \; iob, k \qquad (5.4)$$

The set of constraints essentially represents the linking of the optimization algorithm with the numerical groundwater flow and transport simulation model through the decision variable (source fluxes),

where:

$q_s C_s$ is the pollutant source fluxes ($ML^{-3}T^{-1}$),

$q_s$ is the volumetric flux of water per unit volume of aquifer ($T^{-1}$),

$C_s$ is the concentration of the sources or sinks ($ML^{-3}$),

*Abs* (.......) represents the absolute value,

$cest_{iob}^{k}$ is the concentration estimated by the transport simulation model at observation monitoring location *iob* and at the end of time period $k$ ($ML^{-3}$),

*nk* is the total number of monitoring time steps,

*nob* is the total number of observation wells,

$cobs_{iob}^{k}$ is the observed concentration at monitoring location *iob* and at the end of time period $k$ (ML$^{-3}$).

## 5.1.2. Design of Gradient based Optimal Monitoring Network

The monitoring network design model is developed to find optimal monitoring well locations with the aim of improving the accuracy of source identification results. Pollutant concentration measurements from such optimally designed monitoring networks are used in the source identification model to estimate the source flux release history. The optimal monitoring network design model for identification of the unknown pollution source fluxes is defined by equation (5.5). The optimization model maximizes the objective function value, subject to the constraint that the maximum number of monitoring wells that can be selected in any design sequence is limited.

The objective function (equation 5.5) maximizes the summation of the product of estimated concentration gradients and the estimated concentration at all of the monitoring locations for each sequence of monitoring network design.

$$F_{obj} = Max \sum f_{i.j} C_{i,j}^{*} \left\{ \frac{\Delta C_{i,j}^{*}}{\Delta X} + \frac{\Delta C_{i,j}^{*}}{\Delta Y} \right\} \ \forall \ i,j \qquad (5.5)$$

where

$C_{i,j}^{*}$ is concentration value obtained from forward simulation of the transport model at the grid $i,j$;

$\Delta X$ and $\Delta Y$ is size of the grid in the $i,j$ direction respectively.

These concentration gradients are computed based on the spatial concentrations predicted by the forward simulation using the transport simulation model (equation 5.2), and using the estimated source flux values in the current sequence of the source flux identification model. The concentration gradient along $x$ axis and $y$ axis at any grid location $i, j$ is given by equation (5.6) and (5.7) respectively.

$$\frac{\Delta C^*_{i,j}}{\Delta X} = \frac{\left|C^*_{i-1,j} - C^*_{i,j}\right| + \left|C^*_{i+1,j} - C^*_{i,j}\right|}{\Delta X} \tag{5.6}$$

$$\frac{\Delta C^*_{i,j}}{\Delta Y} = \frac{\left|C^*_{i,j-1} - C^*_{i,j}\right| + \left|C^*_{i,j+1} - C^*_{i,j}\right|}{\Delta Y} \tag{5.7}$$

The objective function is maximized subject to the constraint in equation (5.8), which essentially represents the imposed limit on the total number of permissible monitoring wells in the current design sequence. Monitoring well locations with high concentration gradient values has a higher chance of capturing the point of inflexion in a breakthrough curve. Temporal or spatial changes in source flux characteristics may result in multiple peaks in a breakthrough curve. These peaks can be distinguished from the others at these points of inflexion. Choosing monitoring wells with steep gradient values results in capturing these peaks, which in effect are caused by changes in source flux characteristic. Often it is these changes in the source flux characteristics that are difficult to estimate using a source identification model when using concentration measurements from random well locations. The constraint defining the limit on the number of optimal monitoring wells to be obtained is given by equation (5.8).

$$\sum f_{i,j} \leq W \quad \forall \ i,j \tag{5.8}$$

where:

$f_{i,j}$ represents the binary decision variable to place or not to place a monitoring well at grid location $i$, $j$. $f_{i,j} \equiv \{0, 1\}$ such that when $f_{i,j}$ value equal to 1 representing monitoring well to be implemented at grid $i, j$, and zero otherwise;

$W$ is maximum permissible number of monitoring wells that can be placed in the study area in the current design sequence;

The solution of the design model specifies the optimal locations of the new monitoring wells. SA is used as the solution algorithm to solve the optimization model for finding optimal monitoring well locations.

### 5.1.3. Sequential Integrated Model

The sequential integrated model combines the above discussed source identification model and the pollutant concentration gradient based optimal monitoring network design model to be implemented sequentially. The optimal solution of the source identification model is utilized for obtaining a new optimal monitoring network. The new monitoring network is implemented for collecting new concentration measurement data. In the next sampling time step, pollutant concentration measurements from these newly implemented monitoring wells and at already existing monitoring wells are obtained. Subsequently, the pollutant concentration measurements from the current sampling time step and previous sampling time steps are utilized for source identification.

At the first instance of pollutant detection, only one temporal pollutant concentration measurement $cobs_{iob}^{1}$ from a few random monitoring locations $nob^{k}$ ($nob^{1} = 3$, for example) are available for solving the source identification model. Therefore, in equation

(5.3) the total number of monitoring time steps $nk = 1$, where $k$ is the sampling time step number or sequence number. The estimated source flux value $^{k}q_sC_s$ obtained by using one set of temporal pollutant concentration measurements $cobs_{iob}^{1}$ is used in the transport simulation model (equation 5.2) to predict the spatial pollutant concentrations $^{\text{Pr}ed}cobs_{iob}^{k+1}$ at monitoring time step two ($k = 2$).

The predicted pollutant concentration $^{\text{Pr}ed}cobs_{iob}^{k+1}$ is used to estimate the pollutant gradient using equations (5.6) and (5.7). For any given sampling time step:

$$C_{i,j}^{*} = \ ^{\text{Pr}ed}cobs_{iob}^{k+1} \tag{5.9}$$

where:

$i, j$ is the grid location of a monitoring well $iob$.

The objective function (equation 5.5) is evaluated for all potential well locations along with the constraint set (equation 5.8) to choose the $W$ optimal well locations to be implemented in the given sequence. Thus, the number of monitoring wells available for sampling in the next sequence (next sampling time step) is given by equation (5.10):

$$nob^{k+1} = nob^{k} \ + W \tag{5.10}$$

In the next sequence concentration measurements $cobs_{iob}^{k+1}$ from $nob^{k+1}$ for sampling time step $k+1$ and concentration measurements $cobs_{iob}^{k}$ from $nob^{k}$ for sampling time step $k$ are used again in the source identification model to estimate the value of source flux $^{k+1}q_sC_s$. This sequential process is repeated until there is almost negligible change in the estimated source flux values in any two consecutive sequences of implementation, given by

equation (5.11), where $\xi$ is a small value. The change in estimated source fluxes is given by equation (5.11), where $\xi$ is a small value:

$$|^{k+1}q_sC_s - {}^kq_sC_s| \leq \xi \qquad (5.11)$$

## 5.2. Performance Evaluation of Developed Methodology

The performance of the developed methodology is evaluated for an illustrative polluted aquifer study area, as shown in figure 5.2, comprising heterogeneous, anisotropic and confined aquifer. This study area has a total dimension of 2100 metres in the *x* direction and 1950 metres in the *y* direction. The entire study area is discretised into smaller grids



**Figure 5.2 Plan View of the Study Area**

of size $\Delta x$, $\Delta y$ and $\Delta z$ in $x$, $y$ and $z$ directions respectively. The study area contains different hydrogeological zones with different values of hydraulic conductivity $K_{xx}$ and effective porosity $\theta$. Groundwater flow and solute transport processes are simulated using the value for saturated thickness of the aquifer $b$, longitudinal dispersivity $\alpha L$, transverse dispersivity $\alpha T$ and horizontal anisotropy as given in table 5.1. In the discretised study area, cells marked with a red star represent the grid locations containing a potential pollutant source S($i$), where $i$ represents the source number. Out of the four pollutant sources, source S3 is a dummy source. The dummy source represents zero source flux and is included to test the performance of the source identification methodology. Cells marked with green circle are the grid locations containing an observation well where the pollutant is first observed. Cells marked with a blue cross represent the grid cell with a potential monitoring well location.

| Parameter | Unit | Value |
|---|---|---|
| Maximum length of study area | m | 2100 |
| Maximum width of study area | m | 1950 |
| Saturated thickness, $b$ | m | 30 |
| Grid spacing in x-direction, $\Delta x$ | m | 50 |
| Grid spacing in y-direction, $\Delta y$ | m | 50 |
| Grid spacing in z-direction, $\Delta z$ | m | 30 |
| Hydraulic Conductivity $K_{xx}$ | m/d | Between 15 and 30 |
| Effective Porosity $\theta$ | dimensionless | Between 0.25 and 0.3 |
| Longitudinal Dispersivity, $\alpha L$ | m/d | 20 |
| Transverse Dispersivity, $\alpha T$ | m/d | 10 |
| Horizontal Anisotropy | dimensionless | 1.5 |

**Table 5.1 Hydrogeological Parameters for Study Area**

In the performance evaluation scenario the activity duration of the sources is divided into three equal stress periods of 500 days. The pollutant flux from the sources is assumed to be constant over a stress period. The pollutant flux from each of the sources is

represented as S($i$)($j$), where $i$ represents the source number and $j$ represents the stress period number. A total of twelve source fluxes (S11, S12, S13, S21, S22, S23, S31, S32, S33, S41, S42 and S43) are considered as explicit unknown variables in the source identification optimization problem.

### 5.2.1. Solution Procedure

For performance evaluation purposes, observed concentration measurements are synthetically generated for assumed actual pollutant sources. The observed aquifer responses are simulated by numerical simulation models MODFLOW and MT3DMS in GMS7.0. Initial and boundary conditions (initial heads and boundary heads) are specified in the numerical simulation models. For this performance evaluation purpose, the specified actual source fluxes are utilized to simulate concentration measurement at specified measurement locations.

It is assumed that the pollutant is first observed at three random monitoring locations (figure 5.2) after 1600 days from when the source activity started. The source identification model is initiated using observed concentration measurements from these three monitoring wells. The source fluxes (S11, S12, S13, S21, S22, S23, S31, S32, S33, S41, S42 and S43) estimated using concentration measurements from three observation wells at 1600 days are used to predict the pollutant concentration distribution after 1800 days from the start of source activity. Concentration gradient information from predicted concentration distribution is used to design the monitoring network.

A total of 57 potential monitoring wells, including the initial three well locations, are considered as candidate locations for monitoring well design. A monitoring network design model is used to choose the optimal monitoring well locations with an upper limit

of three wells ($W = 3$) to be installed in every design sequence, allowing for repetition of monitoring wells from the previous design sequence. Monitoring wells are installed and concentration measurements are taken after 1800 days from all the existing monitoring wells. Observed concentration measurements at 1600 days and 1800 days are now used in the next sequence to solve the source flux identification model. Concentration measurement data is available every 200 days starting at 1800 days from specified time, $t$ = 0.

### 5.2.2. Evaluation using Erroneous Concentration Measurement Data

In order to reflect real life conditions, where the pollution measurements are erroneous, the numerically simulated concentrations were perturbed to incorporate measurement errors. The observed concentrations generated using MODFLOW and MT3DMS were perturbed by adding error terms to the simulated measurement data to represent the effect of random measurement errors. These errors were added in order to incorporate realistic measurement errors. The observed pollutant concentration data is perturbed with random measurement error containing a maximum deviation of 10 percent of the measured concentration value as shown in equation (5.12).

$$^{Pert}cobs_{iob}^{k} = cobs_{iob}^{k}(1 + err)$$

(5.12)

$$err = \mu per \times rand$$

(5.13)

where:

$^{Pert}cobs_{iob}^{k}$ is the perturbed simulated erroneous concentration measurement at location *iob* at sampling time step *k*,

*err* is error term,

*µper* is maximum deviation expressed as percentage,

*rand* is a random fraction between -1.0 and +1.0 generated using a latin hypercube distribution.

A Latin hypercube distribution is chosen for generating random error data evenly distributed across all class intervals, thus eliminating any clustering of sample data in a few of the class intervals. The procedure as stated in section 5.2.1. is followed using erroneous concentration measurements. To show the improved efficiency of the source identification when using a sequential source identification model integrated with an optimal monitoring network design model, the same problem is solved utilizing concentration measurements from five static arbitrary monitoring networks. All scenarios are solved, first using error-free concentration measurement data, and then using concentration measurements perturbed with random errors.

## 5.3.    Discussion of Solution Results

The solution results of pollution source flux identification using an integrated source identification model linked with a sequential monitoring network design model is presented in figures 5.3 to 5.9. The source flux identification model is solved using equation (5.3). Each of the unknown source flux variables (S11, S12, S13, S21, S22, S23, S31, S32, S33, S41, S42 and S43) is marked on the *x* axis having three corresponding bars. The first bar is the actual value of the source flux. The second bar represents the estimated flux value using error-free data and the third bar represents the estimated flux while utilizing erroneous concentration measurement. Since the methodology involves

multiple sequences of implementation, the source flux identification results are presented at every sequence. Figures 5.3 to 5.5 represent the source flux identification results for sequences 1 to 3 respectively.

Source Flux Identification Result using Initial Observed Concentration Measurents



**Figure 5.3 Source Flux Identification Results using Initial Observed Concentration Measurements**

Source Flux Identification Result after First Design Sequence



**Figure 5.4 Source Flux Identification Results after First Design Sequence**

Source Flux Identification Result after Second Design Sequence



**Figure 5.5 Source Flux Identification Results after Second Design Sequence**

Figure 5.3 shows the source flux identification results utilizing concentration measurements when the pollutants are first detected in the aquifer. Large errors can be seen in all source flux estimates, as only one temporal observed concentration measurement (at 1600 days) from three random locations is utilized for source identification. After implementing one sequence (at 1800 days) of monitoring network design, the improvement in the source flux identification results can be seen in figure 5.4. However, the dummy source S3 (not actual source) is still not identified accurately. The source flux identification result after implementing the second sequence of monitoring network design (at 2000 days) is shown in figure 5.5. It can be noted that the source flux estimates now closely match the actual source flux values. The dummy source S3 is accurately identified as source flux values corresponding to the dummy source (S31, S32 and S33) are close to zero.

In order to establish efficiency in the source flux identification process using the sequentially designed monitoring network, these source identification results are

compared with the source flux identification results obtained utilizing concentration measurement from five arbitrary monitoring networks. The average of the source flux estimation errors for the five arbitrary monitoring networks is compared to estimation errors obtained using concentration measurements from the designed sequential optimal monitoring networks.

Comparison of Source Flux Identification Result: Optimal Vs. Arbitrary
Networks Utilizing Error-Free Measurements



**Figure 5.6 Source Flux Identification Results for Arbitrary Monitoring Networks Utilizing Error-free Measurement**

Figure 5.6 shows the source flux identification results when utilizing error-free concentration measurements from arbitrary monitoring networks. The first bar represents the source flux value corresponding to the actual source flux, and the second bar until the sixth bar are the estimated source flux obtained utilizing concentration measurements from the designed optimal monitoring networks. The estimation error of source flux

identification is given by the absolute difference between the actual source flux value and the estimated source flux value. The estimation error for the five arbitrary monitoring networks is averaged and compared with the estimation error for the optimal dynamic monitoring network. This comparison, using error-free concentration measurements, is shown in figure 5.7.

Comparison of Source Flux Estimation Error Utilizing Error-Free Measurements



**Figure 5.7 Comparison of Source Flux Estimation Error Utilizing Error-free Measurements**

The source flux estimation errors obtained by using arbitrary monitoring networks and the sequential optimal monitoring network are small when utilizing error-free measurement data. However, it can be seen in figure 5.7 that the source flux estimation error is smaller when sequential optimal monitoring network based concentration measurements are utilized, compared to the average estimation error for arbitrary networks. Although these results may not show much difference in estimation error for

arbitrary networks and the sequentially designed optimal monitoring network, the difference is more pronounced when using erroneous concentration measurements (figures 5.8 and 5.9).
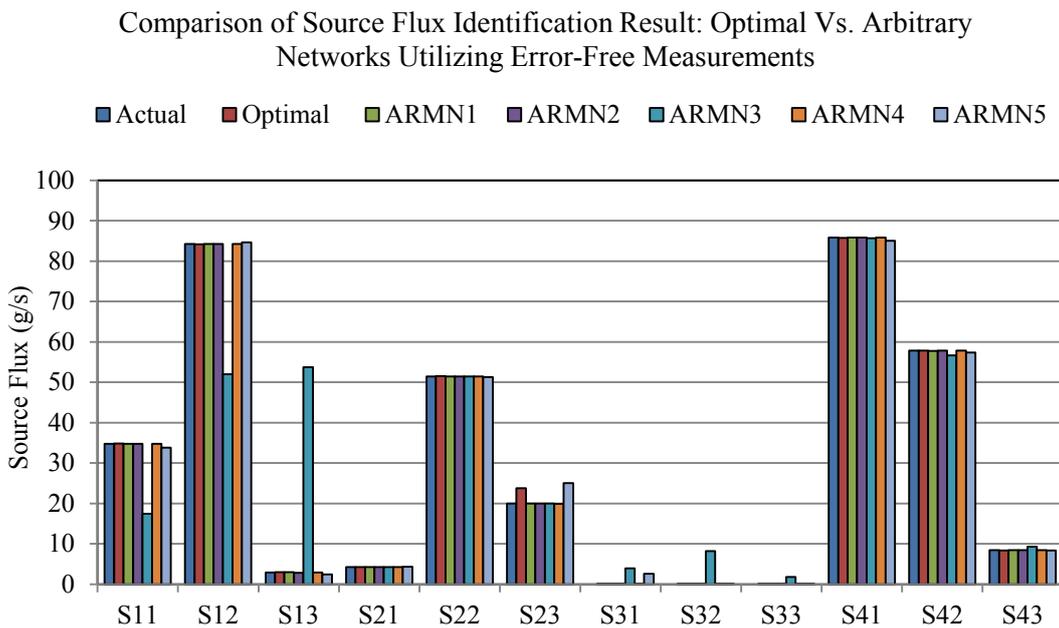


**Figure 5.8 Source Flux Identification Results for Arbitrary Monitoring Networks Utilizing Erroneous Measurement**



**Figure 5.9 Comparison of Source Flux Estimation Error Utilizing Erroneous Measurements**

It can be seen in figure 5.9 that the source flux estimation error is much larger when measurements from arbitrary monitoring networks are used, except for source fluxes S11 and S21. This discrepancy could be due to the fact that not enough monitoring data reflecting the impact of S11 and S21 are present, and only the tail end of the breakthrough curve is captured by the implemented monitoring wells, and the initial part of the source flux activity at sources S1 and S2 is not captured. This problem can be addressed by one more sequential design using the methodology. However, the cost versus benefit implications should be considered in such scenarios. Also, the source fluxes (S31, S32 and S33) corresponding to the dummy source S3 are not identified accurately while utilizing concentration measurements from arbitrary monitoring networks. Table 5.2 shows the grid locations where monitoring wells are implemented in each sequence of the sequential optimal monitoring network design.

| Design Sequence Number | Number of Wells | Grid Location of Monitoring Wells ($i, j$) |
|---|---|---|
| Initial Wells (Error Free) | 3 | (22, 24),(25, 33),(29,30) |
| Initial Wells (Erroneous) | 3 | (22, 24),(25, 33),(29,30) |
| Design Sequence 1(Error Free) | 6 | (22, 24),(25, 33),(29,30),(19,16),(21,16),(23,16) |
| Design Sequence 1(Erroneous) | 5 | |
| Design Sequence 2(Error Free) | 8 | (22, 24), (25, 33),(29,30),(21,23),(24,25) |
| Design Sequence 2(Erroneous) | 8 | (22, 24),(25, 33),(29,30),(19,16), (21,16),(23,16),(21,19),(21,23) (22, 24), (25, 33),(29,30),(21,23) (24,25),(19,16),(19,29),(23,21) |

**Table 5.2  Grid Location of Implemented Monitoring Wells**

## 5.4.  Conclusions

In all real life scenarios of groundwater pollution, the source characterization process requires a large number of spatiotemporal concentration measurements. Designed

collection of monitoring data can reduce some of these uncertainties. Significant time is generally lost between the first detection of any pollutant and collection of sufficient temporal concentration measurement data required for accurate identifications of source characteristics. Moreover, the locations from which these temporal concentration measurements are obtained may not be ideal for accurate source identification.

This sequential methodology integrating a source identification model with an optimal monitoring network design model performs satisfactorily in addressing this issue. Performance evaluation results of source flux identification using this methodology appear to result in more accurate source characterization, even when starting with just three observed concentration measurements at arbitrarily located observation wells. In most of the traditional methods, source identification starts only when sufficient temporal measurements are recorded without having any clue about the merit of the locations from which these temporal concentration measurements are recorded. The proposed methodology addresses this deficiency.

One of the biggest advantages of using this methodology is that source identification can start at the very instance when pollutants are first detected in the aquifer. The feedback based methodology helps the source identification model to rectify its estimates sequentially. Even while starting with very few concentration measurements from a random set of monitoring wells, this method is able to come up with good estimates within a couple of design sequences. Improved estimates of the source fluxes in a given sequence are passed on as important feedback information to the monitoring network design model in the form of the predicted pollutant distribution over space for the next sampling time step. This information improves the design of the new monitoring network,

such that the concentration measurement from such locations improves the accuracy of source flux estimates in the source identification model.

The proposed methodology is applicable to real world problems of source identification in polluted aquifer sites where very little or no pollutant concentration measurements are available for source identification. This method can increase the accuracy of source identification and reduce the time for implementing a source identification model after the first detection of pollutants in an aquifer.

# 6. Application of Developed Methodologies for Pollution Source Characterization in a Polluted Aquifer in NSW

In this chapter, the source characterization methodologies developed in the previous chapters are applied to a real, polluted aquifer site in the state of New South Wales, Australia. Due to confidentiality requirements, some details pertaining to the identification of this site are not included.

## 6.1. Background of the Problem

The polluted aquifer region forms a part of the Upper Macquarie Groundwater Management Area in New South Wales, Australia. Overlying the polluted aquifer is a suburban town. Several complaints of BTEX vapour emanating from building basements in the locality led to the investigation of the polluted aquifer. However, there is no record of when the event of vapour emanating from the building basement was first recorded in the area. The investigation records show BTEX concentration as high as 320mg/l in October 2009 in one of the monitoring wells. The extent of the pollution was roughly estimated to be over an area of 1km$^2$.

During the investigation spanning from October 2006 until July 2011 (concluded based on the recorded measurements), seventy-four monitoring wells were installed in the polluted aquifer region (approximately 1km$^2$) to monitor the pollutant concentration and groundwater hydraulic head. These monitoring wells were installed intuitively at different times during the investigation period to locate the pollutant plume and to understand the

plume movement in order to tackle the spread of pollutants. Nineteen of these wells were used as injection wells to inject neutralizer to contain the spread of pollutant and eventually reclaim the polluted aquifer.

Although no set pattern was followed in installing the monitoring wells, wells were more densely installed around the potential source location as compared to regions further away from the source. As per the investigation, the origin of pollutants was traced back to a leaking underground storage tank at a nearby gas station. However, the starting time of leakage and the release time history of the pollutants coming out of the tank were not ascertained in the investigation.

The main aim of this study is to identify the unknown source characteristics (source location, starting time of activity of the source, and source flux release history) in the polluted aquifer. A source identification model with unknown source starting time as an explicit decision variable (section 3.3.4.) is integrated with a gradient based dynamic optimal monitoring network design model (section 5.1.3.) for estimating the unknown source characteristics. Data collected during the five years of investigation is obtained and used to confirm the applicability of the developed methodologies for source identification and optimal monitoring network design.

## 6.2.    Polluted Aquifer Site Description

The polluted aquifer site constitutes a small part of the Upper Macquarie Groundwater Management Area. The Macquarie River forms the western boundary of the aquifer study area and runs from south to north as shown in figure 6.1. Due to confidentiality requirements the exact location of the polluted aquifer site is not disclosed.

**Figure 6.1 Plan Views of the Study Area and the Impacted Area**

### 6.2.1. Study Area Boundary

The previous investigation estimated the extent of pollutant spread roughly over an area of 1km$^2$. The Macquarie River formed a natural boundary on the western side of the study area. However, there were no distinct geological formations to be used as natural boundaries on the other three sides of the polluted aquifer area. Hence, a much larger area than the actual polluted aquifer region was considered in this study. The study area measuring 2.1871km by 2.4256km was considered in the study such that all hydrogeological conditions impacting the actual polluted aquifer region are accounted for in the model. In this study the actual polluted aquifer region is referred to as the "impacted area" and the total aquifer region considered in this study is called the "aquifer study area", as shown in figure 6.2.

**Figure 6.2 Total Area Considered in the Study**

### 6.2.2.    Topography

The ground topography generally slopes from the south east towards the river in the west. The ground elevation in the study area ranges from 292mAHD towards the river to 251mAHD on the north-eastern side. A majority of the area considered in the study area consists of roads, houses and pavements, with some agricultural land, parkland and playing grounds towards the river boundary (figure 6.1).

### 6.2.3. Stratigraphy of the Study Area

Based on the geological information and the bore-hole logs available at the site, the stratigraphy of the study area can be broadly divided into three distinct layers. The top layer is comprised of tertiary alluvium, the middle layer is comprised of quaternary alluvium and impervious bedrock forms the third layer. The thickness of these layers varies from one point to the other. However, due to sparse bore-holes, the layer thickness had to be interpolated at some points in the aquifer study area. Figure 6.3 shows the stratigraphical details of the aquifer study area.



**Figure 6.3 Cross-sectional view showing Layers in the Study Area**

### 6.2.4. Extraction and Recharge

The main sources of recharge to the aquifer are from rainfall and from the river. The region receives moderate to low rainfall with a long term average of 583mm/year during

the wet season, running from November until February. The Macquarie River is a major source of groundwater recharge in the aquifer study area.

Water extracted from the aquifer through pumping wells, is mainly used for city water supply and irrigation. Extraction of groundwater in the study area is mainly through wells for the purpose of drinking water supply and agriculture. The pumping rate has varied over the years due to changes in the volumetric town water extraction limits from the groundwater system, and due to voluntary groundwater extraction limits in the year 2010 (Marsden Jacob, 2011). Another source of loss of water from the aquifer is through evapotranspiration, which peaks to 260mm/month during the dry season (Puech, 2010).

## 6.3.    Groundwater Flow Modelling of the Aquifer Study Area

The groundwater flow in the aquifer study area is modelled as an unconfined aquifer with specified head boundaries on all sides. A Layer Property Flow (LPF) package in GMS is used for modelling the flow in the study area. The western boundary, represented by the river, is a specific head boundary where the head at the boundary is given by the average stage in the river. A groundwater flow model of the entire Upper Macquarie Groundwater Management Area was developed by Puech (2010). Based on the information available from this report hydraulic heads at the other boundaries are estimated. Groundwater flow in the study area was modelled from 1 January, 1995 until 31 December, 2012. The entire study time horizon was divided into 18 stress periods of 1 year each.

Rainfall recharge in the study area is calculated based on three rain gauge stations around the study area. The rainfall measurements from the three rain gauging stations are averaged and the average annual rainfall is calculated for the entire study area. The

rainfall is considered to be uniform throughout the study area and constant over a stress period. As the study area is largely suburban, only 10 percent of the entire rainfall is assumed to constitute the infiltration recharge. The remaining 90 percent is considered as surface runoff.

Extraction rates from the wells vary from one well to the other. The long term extraction rates are derived from Puech (2010) and the total annual extraction rate from all wells in the study area is supplied from the study area city council. The long term extraction rates are proportionally adjusted such that the sum of total extraction from all the wells matches with the total annual extraction rate. The extraction rates used in the flow model are provided in table 6.1 and are assumed to be constant for a given stress period.

| Stress Period | Well 1 $m^3$/day | Well 2 $m^3$/day | Well 3 $m^3$/day | Well 4 $m^3$/day | Well 5 $m^3$/day | Well 6 $m^3$/day | Well 7 $m^3$/day | Well 8 $m^3$/day |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.14 | 1.14 | 11.58 | 5.98 | 1283.90 | 983.80 | 1308.06 | 1454.54 |
| 2 | 0.17 | 1.34 | 13.58 | 7.02 | 1506.12 | 1154.07 | 1534.46 | 1706.29 |
| 3 | 0.22 | 1.78 | 18.03 | 9.32 | 1999.92 | 1532.45 | 2037.56 | 2265.73 |
| 4 | 0.24 | 1.94 | 19.72 | 10.20 | 2187.57 | 1676.24 | 2228.73 | 2478.32 |
| 5 | 0.23 | 1.82 | 18.52 | 9.57 | 2054.24 | 1574.08 | 2092.90 | 2327.27 |
| 6 | 0.22 | 1.80 | 18.25 | 9.44 | 2024.61 | 1551.37 | 2062.71 | 2293.70 |
| 7 | 0.24 | 1.92 | 19.54 | 10.10 | 2167.82 | 1661.10 | 2208.61 | 2455.94 |
| 8 | 0.26 | 2.11 | 21.46 | 11.09 | 2380.16 | 1823.81 | 2424.94 | 2696.50 |
| 9 | 0.25 | 2.03 | 20.66 | 10.68 | 2291.27 | 1755.70 | 2334.38 | 2595.80 |
| 10 | 0.21 | 1.65 | 16.74 | 8.65 | 1856.72 | 1422.72 | 1891.66 | 2103.49 |
| 11 | 0.18 | 1.41 | 14.29 | 7.39 | 1585.12 | 1214.61 | 1614.95 | 1795.80 |
| 12 | 0.17 | 1.36 | 13.85 | 7.16 | 1535.74 | 1176.77 | 1564.64 | 1739.86 |
| 13 | 0.18 | 1.46 | 14.87 | 7.69 | 1649.32 | 1263.80 | 1680.35 | 1868.53 |
| 14 | 0.17 | 1.36 | 13.85 | 7.16 | 1535.74 | 1176.77 | 1564.64 | 1739.86 |
| 15 | 0.15 | 1.22 | 12.42 | 6.42 | 1377.72 | 1055.69 | 1403.65 | 1560.84 |
| 16 | 0.16 | 1.32 | 13.40 | 6.93 | 1486.36 | 1138.93 | 1514.33 | 1683.91 |
| 17 | 0.16 | 1.32 | 13.36 | 6.90 | 1481.42 | 1135.15 | 1509.30 | 1678.32 |
| 18 | 0.15 | 1.17 | 11.85 | 6.13 | 1314.76 | 1007.45 | 1339.50 | 1489.51 |

**Table 6.1 Extraction rate from the Wells in the Study Area**

**Figure 6.4 Isometric View of the Discretised Study Area**

In the three dimensional simulation models, the study area is discretised into small grids of size 21.87m by 21.08m in the $x$ and $y$ directions respectively, as shown in figure 6.4. The size of the grid in the $z$ direction varies to match with the layer thickness.

### 6.3.1.  Flow Model Calibration

The flow in the groundwater aquifer was simulated using an LPF package. The hydrogeological properties, such as hydraulic conductivity, porosity, specific storage and specific yield, were obtained from previous studies conducted on this study area by Puech (2010) and Jha and Datta (2012a). These hydrogeological properties are listed in table 6.2.

| Parameter | Unit | Value |
|---|---|---|
| Maximum length of study area | m | 2187.1 |
| Maximum width of study area | m | 2425.6 |
| Saturated thickness, $b$ | m | Variable |
| Number of layers in z-direction | | 3 |
| Grid spacing in x-direction, $\Delta x$ | m | 21.87 |
| Grid spacing in y-direction, $\Delta y$ | m | 21.08 |
| Grid spacing in z-direction, $\Delta z$ | m | Variable |
| $Kxx$ (Layer 1, Layer 2, Layer 3) | m/d | 12.37, 16.24, 0.001 |
| $Kyy$ (All Layers) | m/d | 0.2 |
| $\theta$ ( All Layers ) | dimensionless | 0.27 |
| Longitudinal Dispersivity, $\alpha L$ | m/d | 12 |
| Transverse Dispersivity, $\alpha T$ | m/d | 6 |
| Horizontal Anisotropy | dimensionless | 1.5 |
| Specific Yield $Sy$ (All Layers) | dimensionless | 0.1 |
| Specific Storage $Ss$ (All Layers) | dimensionless | 0.000006 |
| Initial pollutant concentration | g/l | 0.00 |

**Table 6.2 Hydrogeological Properties used in Flow Modelling of the Study Area**

The developed groundwater flow model was calibrated using hydraulic head measurement data from 31 observation locations spread across the impacted area. The recorded hydraulic head data used for model calibration was recorded for the duration starting from October 2006 to July 2011, at discrete time intervals. Calibration targets were set to be within one metre intervals of the observed hydraulic head value with a confidence level of 90 percent. The model boundary conditions were manually adjusted to achieve the calibration targets. Figures 6.5 to 6.8 show model calibration results as box plots, at different observation time steps, in the year 2011.

**Figure 6.5 Flow Model Calibration Results 1**



**Figure 6.6 Flow Model Calibration Results 2**

**Figure 6.7 Flow Model Calibration Results 3**



**Figure 6.8 Flow Model Calibration Results 4**

183

The calibration results for the groundwater flow model are represented by box plots (Figures 6.5 to 6.8). Box plots in green indicate that the absolute difference between the estimated heads and observed heads were well within the chosen target interval. Similarly, a yellow colour on the box plot shows slight over/under estimation of the heads such that the absolute difference between the estimated heads and observed heads is slightly beyond the target interval. Red colour on the box-plot shows larger estimation errors and poor calibration at those locations.

These calibration results show that the groundwater flow model performs satisfactorily for these calibration time steps only in the year 2011. However, with the same initial and boundary conditions calibration results in the initial part of the calibration time horizon are poor. The observed hydraulic head data shows a steep rise of about 1.5 metres towards the end of year 2010 in all the observation locations. To explain this steep rise in the groundwater head value the rainfall data and well extraction data for this period were analysed (figure 6.9).



**Figure 6.9 Rainfall and Groundwater Extraction for the Calibration Time Period**

Although there was a high rainfall incidence of 1117.4mm in the year 2010, considering only 10 percent infiltration was not enough to cause a rise of 1.5m in the head value.

Similarly, there was no drastic reduction in groundwater extraction that could have led to a rise of 1.5m of hydraulic head value in the entire aquifer region. This led to the conclusion that there might be some systematic errors or inconsistencies in the observed hydraulic head measurements. Hence, during calibration more emphasis was given to observed data taken at a later period of the calibration time horizon rather than the earlier periods. Therefore, the box plots representing the calibration errors shown in figures 6.5 to 6.8 can be considered to represent satisfactory calibration.

Once the entire study area is modelled and calibrated the flow model for the actual impacted area is derived from the calibrated model. The GMS7.0 feature, Regional to Local, is used to interpolate the starting head and layer thickness values for the impacted area from the entire study area model. Figure 6.10 shows the actual impacted area where the pollutant is estimated to be present. The grid sizes are refined further in the flow



**Figure 6.10 Actual Impacted Area**

model for the impacted area. All of the boundaries are considered as time varying specified head boundary conditions. The value of the time varying specified heads at the

boundary of the impacted area are extrapolated from the calibrated model for the entire study area. All of the other hydrogeological flow parameters are kept the same as in table 6.2.

## 6.4.    Pollutant Transport Simulation in the Impacted Area

A three-dimensional transient transport simulation model was developed to study the fate and transport of the petrochemical pollutant BTEX originating from a specified point source. For the purpose of implementation, the pollutant is assumed to be conservative in nature and the pollutant plume boundary is assumed to be contained within the boundary of the impacted area. The transport model uses the flow field generated by the flow model to predict the movement of the pollutants in the impacted area of the aquifer over time. The initial concentration of BTEX in the aquifer at the start of the transport simulation is assumed to be zero. All the other relevant transport parameters used in the transport model are shown in table 6.2.

## 6.5.    Performance Evaluation of the Applied Methodologies

Once the groundwater flow model for the impacted area is developed, the next step is to identify the unknown source characteristics of the pollutant. However, in this study though the location of the source is approximately known but the starting time of the activity of the source and the pollutant source flux (BTEX) release history is unknown. Also, to evaluate the performance of the source identification methodology, the exact locations of the pollutant sources are assumed to be unknown, with two different possible locations. One of these needs to be identified as not an actual or dummy source. The two potential sources in the study area are shown in figure 6.11. Methodologies developed for

simultaneous identification of unknown source flux release history and source activity starting times (chapter 3) are utilized for recreating the source flux release history and the source activity initiation time. For efficient identification of unknown source characteristics, the methodology is integrated with a sequential optimal monitoring network design (chapter 5). The points marked in yellow circles are the grid locations containing the potential sources, and the red dots are the observation wells where the concentration of BTEX is observed. A total of 24 observation locations are present in the study area.



**Figure 6.11 Discretised Plan View of the Study Area**

### 6.5.1. Simultaneous Source Flux Release History and Source Activity Initiation Time Identification

In the simultaneous source flux release history and source activity initiation time identification, the simulation model starts from 1 January, 1995. However, the starting time of the activity of the sources is unknown and can start anywhere between 1 January, 1995 and 31 December, 2011. The activity duration of the sources is assumed to be 10 years divided into 10 equal stress periods of 1 year each. The pollutant flux from the sources is assumed to be constant over a stress period. The pollutant flux from each of the sources is represented as $S(i)(j)$, where $i$ represents the source number and $j$ represents the stress period number. In this case S1 is the actual source and S2 is the dummy source. A total of twenty source fluxes (S11, S12, S13, S14, S15, S16, S17, S18, S19, S10, S21, S22, S23, S24, S25, S26, S27, S28, S29 and S20) are considered as explicit unknown variables in the optimization problem.

To determine the starting time of the sources an additional time lag variable $\Delta T$ (chapter 3) is introduced in the optimization problem. It is assumed that both the actual source S1 and the dummy source S2 start at the same time. Pollutant concentration measurements starting 22 January, 2009 are used in the simultaneous source flux release history and source activity initiation time identification model.

### 6.5.2. Sequential Optimal Monitoring Network for Efficient Source Characterization

To demonstrate the efficiency of source characterization using concentration measurements in an optimal monitoring network, the above mentioned simultaneous source flux release history and source activity initiation time identification model is integrated with a dynamic optimal monitoring network design model (chapter 5).

The simultaneous source flux release history and source activity initiation time methodology is started using concentration measurements from three wells. The three wells are randomly chosen out of the twenty-four monitoring wells in the aquifer study area. Concentration measurements on 22 January, 2009 from these three wells are used initially to estimate the source flux release history and the source activity initiation time. This implies that pollution in the aquifer was first detected at the three wells on 22 January, 2009. These initial source flux estimates are then used to predict the pollutant (BTEX) concentration at the next sampling time step, 30 April, 2009. Concentration gradient information from the predicted pollutant (BTEX) concentration is used to find the optimal monitoring locations for the next sampling time step, 30 April, 2009.

Three optimal monitoring locations as chosen by the optimal monitoring network design model are implemented. All the twenty-four observation locations are treated as potential monitoring locations in the optimal monitoring network design model. Once a new monitoring network is designed and implemented the concentration measurements from all wells in the monitoring network obtained on 30 April, 2009, together with measurements obtained from the pre-existing three arbitrary wells taken on 22 January, 2009, are utilized in the source identification model. This approach is repeated for the subsequent sampling time steps, until changes in the source flux and starting time estimates are negligible.

### 6.5.3. Performance Evaluation Criteria

In order to evaluate the performance of the sequential approach integrating an optimal monitoring network design model with simultaneous source flux release history and the source activity initiation time model, only three initial observation locations were

assumed to exist in the study area on 22 January, 2009. Also, new monitoring wells are added in every sequential monitoring network design time step. Concentration measurements in the implemented monitoring well starts only in the next sampling time step after the implementation of the monitoring wells in the new monitoring network. All concentration measurements available are used in the source identification model.

Since the actual source flux release history or the source activity initiation time are not known, the estimated source flux magnitudes cannot be validated. The performance of the methodology can only be evaluated in terms of how well the methodology is able to identify the dummy source. The utility of the sequential optimal monitoring network design model can be demonstrated by comparing the numbers of observation wells installed in order to identify unknown source characteristics. The efficiency of the designed methodology can be demonstrated by comparing the monitoring time horizon required over which the observation wells are monitored to obtain pollutant concentration measurements. This comparison will be done in terms of the number of spatiotemporal concentration measurements required for identifying the unknown source characteristics. Long-term monitoring over a large number of monitoring wells leads to higher costs.

## 6.6.    Performance Evaluation Results

The performance evaluation results of source location identification, source flux release history and source activity initiation times, using an integrated source identification model linked with a sequential monitoring network design model, are presented in figures 6.12 and 6.13. Each of the unknown source flux variables (S11, S12, S13, S14, S15, S16, S17, S18, S19, S10, S21, S22, S23, S24, S25, S26, S27, S28, S29, S20) and time lag variable $\Delta T$ are marked on the $x$ axis. The source flux magnitude in gram per second is

shown on the primary *y* axis and the lag time in days is shown on the secondary *y* axis. Since the methodology involves multiple sequences of implementation, the source flux identification results are presented after every sequence of monitoring well implementation.

Simultaneous Source Location, Source Flux Release History
and Source Activity Starting Time Identification: Initial Estimates



**Figure 6.12 Source Identification Result using Initial Observed Concentration Measurements**

Figure 6.12 shows the source flux and starting time estimates using concentration measurements taken on 22 January, 2009 from three initial arbitrary well locations. The lag time estimated by the methodology shows that source activity started in the year 1999. From the estimated source flux values it cannot be concluded which of the two sources is the dummy source. Therefore, to find the accurate source location, the next sequence of the optimal monitoring network was implemented for the next monitoring time step, dated 30 April, 2009. In the next sequence, concentration measurements from all wells in the monitoring network obtained on 30 April, 2009, and measurements obtained from the pre-existing three arbitrary wells taken on 22 January, 2009, are

utilized in the source identification model. Figure 6.12 shows the source flux and the starting time estimates using concentration measurements from all wells in the monitoring network obtained on 30 April, 2009, and measurements obtained from pre-existing three arbitrary wells taken on 22 January, 2009.

Simultaneous Source Location, Source Flux Release History
and Source Activity Starting Time Identification: Initial Estimates



**Figure 6.13 Source Identification Results after One sequence of Implementation of Monitoring Networks**

The source flux magnitudes for source two (S2) are estimated close to be zero, thus showing no contribution of pollutants from source two. This confirms that source two is a dummy source and source one (S1) is the actual source. This is validated by the fact that the location of source one coincides with the location of the gas station, although it cannot be ascertained if the source flux estimates for source one (actual source) are correct. The lag time estimates show that the source activity started in 1999.

The estimated source flux magnitude for S10 shows a steep jump in source flux value. As the lag time $\Delta T$ estimate shows that the source activity started in the year 1999, source

**192**

flux S10 represents the source flux magnitude for 2008. All the concentration measurements used in the identification of source characteristics are from the beginning of the year 2009 (22 January, 2009 and 30 April, 2009). It seems that the source flux magnitude S10 may not have impacted the concentration measurements taken on 22 January, 2009 and 30 April, 2009, as there is always a lag time between the source flux activity and resulting changes in the concentration of the pollutants at a monitoring location. Since only a small portion of the entire breakthrough curve is utilized, it seems that the impact due to source flux magnitude S10 is not captured in the utilized concentration measurements taken on 22 January, 2009 and 30 April, 2009. Thus, the steep jump in source flux value for S10 is a result of numerical redundancy arising due to insufficient representation of the source flux in the utilized concentration measurements, and hence can be ignored.
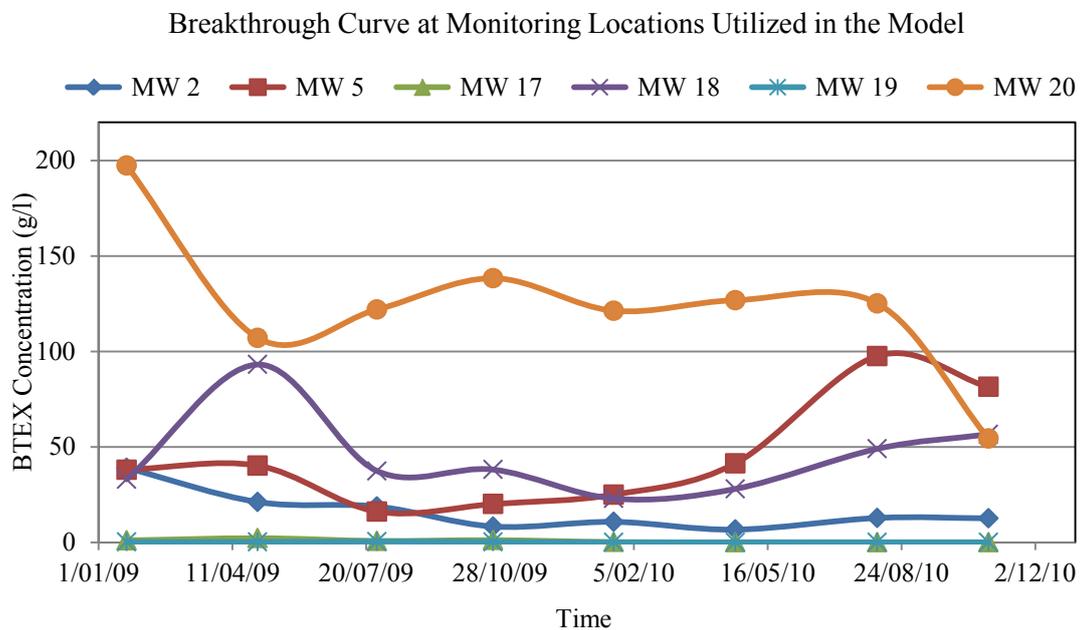


**Figure 6.14 Breakthrough Curves at Monitoring Locations Utilized in the Source Identification Model**

This also seems to be intuitively true based on the concentration breakthrough curves obtained from the monitoring wells utilized in the source identification. As can be seen in figure 6.14 that there is no steep rise in concentration measurements from any of the monitoring wells used in the source identification. This suggests that the high value for source flux S10 is a result of numerical redundancy of the optimization algorithm in estimating the value. The increase in the concentration value for MW5 is due to the late impact of source activity at this location as this location is farthest from the source as compared to the other monitoring well locations. Hence, there is a considerable time gap between source activity and the resulting concentration measurement.

The starting time estimates do not seem to change from the initial design sequence to the next design sequence. It is also evident that the dummy source is identified correctly and that there is no way to validate the estimated source flux for the actual source (S1); therefore, the methodology is terminated. These results of source flux identification and selected optimal monitoring locations are presented in table 6.3.

| Source | Grid Location ($i, j, k$) | Flux Values |
|---|---|---|
| Source 1 (Actual) | (17, 29, 1) | 0.44, 30.40, 26.94, 79.28, 23.46, 14.83, 0.31, 0.01, 0.01, 100.00 |
| Source 2 (Dummy) | (16, 24, 1) | 2.2E-04, 1.6E-04, 3.3E-04, 8.7E-05, 9.1E-05, 2.0E-05, 3.0E-05, 1.2E-06, 1.6E-05 |

| Design Sequence Number | Number of Wells | Grid Location of Monitoring Wells ($i, j, k$) |
|---|---|---|
| Initial Wells | 3 | (20, 25, 1), (22, 26, 1), (21, 28, 1) |
| Design Sequence 1 | 6 | (20, 25, 1), (22, 26, 1), (21, 28, 1), (18, 23, 1), (17, 26, 1), (20, 30, 1) |

**Table 6.3 Results of Source Identification**

## 6.7. Conclusions

The feedback based approach integrating a sequential source characterization model with a sequential optimal monitoring network performs satisfactorily in a real groundwater pollution scenario. Identification of source flux release history, source activity starting time and accurate source locations demonstrates the potential applicability of the developed methodology to all real groundwater pollution scenarios.

An earlier study initiated by the local council for site remediation installed seventy-four observation wells in the polluted aquifer site. Concentration measurements from these seventy-four observation wells were used to understand the problem, predicting the flow and transport, and for site remediation. However, no formal flow and transport modelling of the aquifer was conducted to estimate the unknown source characteristics. Due to confidentiality requirements the source of this information cannot be disclosed.

The developed methodology shows greater efficiency in identifying the unknown source characteristics as it utilizes only six monitoring wells, as compared to the seventy-four observation wells used in the previous study. Only two temporal readings (22 January, 2009 and 30 April, 2009) were utilized in estimating the source characteristics, as compared to nineteen temporal concentration measurements spanning over a period of five years. If 22 January, 2009 is considered to be the time when pollution was first observed in the aquifer study area, then the developed methodology is able to find the unknown source characteristics by 30 April, 2009. As a result, little time is lost after pollutants are first detected and the unknown source characteristics are identified. Thus, little time is lost before starting a formal remediation process to check the spread of pollutants in the aquifer.

There is also a saving of the capital cost as only six monitoring wells are installed as compared to seventy-four monitoring wells in the previous study. The methodology is able to estimate the source characteristics in a fairly short time period (4 months), thus saving on the operation costs associated with long term monitoring (five years in the previous study). The methodology estimates the source activity starting time to the year 1999, which seems to be intuitively true based on the concentration breakthrough profile of the utilized monitoring wells.

However, the feedback based monitoring is not tested to its full potential in this scenario. This is because the methodology is able to identify the source characteristics in one design sequence of the sequential optimal monitoring network design implementation without showing the need for the second design sequence. This could be due to the presence of only one actual source that makes the source identification in this scenario relatively simple. Also, the calibration results confirmed closely with the field observation measurements set as the calibration target.

# 7. Conclusions

This chapter summarises the main findings of this study. Some of the limitations of the methodologies developed in this study and scope for future work are also highlighted. The following models related to efficient pollutant source identification and the design of optimal monitoring networks were developed in this study: (i) linked simulation-optimization models for source characterization in terms of location, source activity initiation times, source flux release history and activity duration of the sources; (ii) optimal monitoring network design models for identification of source locations; and (iii) optimal monitoring network design models for concentration measurements to increase efficiency of source characterization. A sequential feedback based methodology for efficient identification of unknown pollutant source characteristics integrating optimal monitoring network design with an optimization based source identification model is also developed in this study.

Existing methodologies for unknown groundwater pollution source characterization have several limitations. Methodologies developed in this study aim to address some of these limitations. The major limitations addressed in this study include:

i.   sparsity of pollutant concentration measurement data,

ii.   monitoring network design for concentration measurements to improve the accuracy of source identification,

iii.   difficulty in identifying source locations when locations are totally unknown,

iv.   difficulty in establishing the pollutant source activity initiation times from a large potential time span, and

v.	applicability of the optimal source characterization model with missing concentration measurement data.

Existing approaches of pollution source characterization are efficient only when the starting times of the activity of the sources are precisely known, or the possible time window within which the source activity actually starts is not too large and can be specified. To address this deficiency, an optimization based methodology is developed for simultaneous identification of the source fluxes and their starting times.

This method overcomes one of the critical limitations in the earlier proposed methods in which actual starting times of the sources were estimated indirectly by estimating the source flux magnitude for the entire time span discretised into smaller stress periods. These existing approaches resulted in increasing the number of decision variables substantially, making the optimal search algorithm inefficient, if not infeasible in some cases.

The developed methodology for simultaneous identification of source fluxes and the starting times of source activity can drastically reduce the large number of discretised source flux magnitude decision variables, and instead utilize one lag time variable for each of the potential sources. This new approach decreases the number of decision variables. Thus, the optimization algorithm can quickly converge to a correct optimal solution. This methodology in its current form can estimate the activity starting times for all potential sources, both under steady-state and transient flow conditions. It can also handle multiple sources having different source activity initiation times and missing observation data.

One of the difficulties in accurate characterization of unknown groundwater pollution sources is the uncertainty regarding the number and location of such sources. A fairly good knowledge of potential source locations can substantially decrease the degree of non-uniqueness in an optimal source characterization model. To address this problem a methodology is developed based on sequential design of an optimal monitoring network and collecting concentration measurement data from such an implemented network. The sequential design is based on iterative pollutant concentration measurement information from the sequentially designed monitoring networks. The optimal monitoring network design utilizes concentration gradient information from the monitoring network at previous sequences to define the objective function. The feedback information based sequential methodology is shown to be effective in estimating the source locations when very little source location information is initially available.

In scenarios where potential source locations and activity duration are known with a fair degree of certainty, a linked simulation-optimization based approach is often applied for recreating the flux release history of the sources. A large amount of observed pollutant concentration data spread over time and space is necessary for accurate source identification. However, long term monitoring over a large number of monitoring locations has budgetary limitations. Therefore, an optimal monitoring network design model is developed to determine optimal locations for concentration measurements for accurate source identification.

The developed methodology uses trained GP models to find ideal monitoring locations. The GP based impact factor and frequency factor are used to estimate the relative importance of one monitoring well location over the other for collecting concentration

measurements. Concentration measurements from these optimally designed monitoring networks show more accurate source characterization with relatively less spatiotemporal concentration measurements.

In the event of detection of pollution in an aquifer, generally a more formal methodology for source characterization is initiated only when a large number of temporal concentration measurements spaced over a sufficiently long period of time is gathered. During this time the spread of the pollutant continues while temporal measurements are being captured at monitoring locations. However, if the monitoring locations are not optimally suited for accurate source characterization, the concentration measurements from such monitoring wells may not result in efficient source characterization. Therefore, integrated use of an optimal sequential monitoring network design model and a source identification model is presented.

Feedback information in the nature of new concentration measurements from the designed optimal monitoring network improves source characterization. It also helps in better prediction of pollutant distribution over space and time to improve the optimal monitoring network. The main advantage of this methodology is that source characterization can start at the same time as when pollutant is first detected in the aquifer, even though insufficient concentration measurement data is available.

Applicability of the developed methodologies is demonstrated by applying it to a real life polluted aquifer site in New South Wales, Australia. The feedback based methodology, integrating a sequential source characterization model with a sequential optimal monitoring network, performs satisfactorily in identifying source flux release history, source activity starting times and accurate source locations. Although the source flux

magnitudes or source activity starting time cannot be validated, the source location is identified correctly. The source magnitude, duration and activity initiation time results obtained appear to be intuitively correct, based on subjective information.

The methodology proves to be efficient as it requires only a small number of monitoring wells, monitored over a relatively small duration of time. This results in a smaller number of monitoring locations and smaller duration of monitoring. Actually, these results show the redundancy of an unplanned large monitoring network.

Performance evaluation results for each of the methodologies indicate their potential for field application. However, there are some limitations to the methodologies developed in this study. Some of the major limitations are:

1. The linked simulation-optimization methodology developed for source characterization considers the pollutant to be conservative in nature. The methodology needs to be extended to incorporate non-conservative pollutants.

2. The methodology for source location identification in its current form is limited to pollutant sources that are continuous in time. Further work is necessary to incorporate scenarios with sources that are not continuous in time.

3. The monitoring network design methodologies based on GP impact factor and frequency factor are computationally intensive.

4. The methodologies developed in this study are sensitive to uncertainties in hydrogeological parameters and need to be expanded to explicitly incorporate these uncertainties.

5. This study assumes that groundwater flow follows Darcy's law. Fractures or cracks in the subsurface have not been incorporated. In some of the mine sites, fissures and fractures may be present. This aspect needs further consideration.

6. Some of the performance evaluations are based on the assumption that the calibrated model represents actual field conditions as closely as possible. However, the performance evaluation will depend on the accuracy of the calibration.

7. The developed methodologies are computationally intensive. It is possible to further improve computational efficiency by incorporating other optimization algorithms.

# References

Ababou, R., Bagtzoglou, A.C., & Mallet, A. (2010). Anti-diffusion and source identification with the RAW scheme: a particle-based censored random walk. *J. Environ. Fluid Mech.*. doi:10.1007/s10652-009-9153-4

Ala, K. N., & Domenico, A. P. (1992). Inverse Analytical Techniques Applied to Coincident Contaminant Distributions at Otis Air Force Base, Massachusetts. *Groundwater*. DOI: 10.1111/j.1745-6584.1992.tb01793.x

Alapati, S., & Kabala, Z. J. (2000). Recovering the release history of a groundwater contaminant using a non-linear least-squares method. *Hydrological processes*, *14*(6):1003–1016, 2000. 13

Amirabdollahian, M., and Datta, B. (2013). Identification of contaminant source characteristics and monitoring network design in groundwater aquifers: an overview. *J Environ Prot.*. doi:10.4236/jep.2013.45A004

Anderson, M. P., & Woessner, W. W. (1992). *Applied Groundwater Modeling*. Academic Press.

Aral, M.M., Guan, Ji., & Maslia, M.L. (2001). Identification of contaminant source location and release history in aquifers. *J. Hydrol. Eng., 6*(3), 225-234.

Atmadja, J., & Bagtzoglou, A. C. (2000). Groundwater pollution source identification using the backward beam equation method. In: *Computational Methods for Subsurface Flow and Transport*, pp. 397–404, A. A. Balkema, Brookfield, Vt.

Atmadja, J., & Bagtzoglou, A. C. (2001a). Pollution source identification in heterogeneous porous media. *Water Resour. Res., 37*(8), 2113-2125.

Atmadja, J., & Bagtzoglou, A. C. (2001b). State of the art report on mathematical methods to reliable of groundwater pollution source identification. *Environ. Forensics, 2*(3), 205-214.

Ayvaz, T. M. (2010). A linked simulation–optimization model for solving the unknown groundwater pollution source identification problems. *Journal of Contaminant Hydrology, 117*(2010), 46–59.

Azghadi, B. N. S, & Kerachian, R. (2010). Locating monitoring wells in groundwater systems using embedded optimization and simulation models. *Science of the Total Environment, 408*(10), 2189-2198.

Azghadi, B. N. S., Kerachian, R., Lari, M. R. B., & Solouki, K. (2010). Characterizing an unknown pollution source in groundwater resources systems using PSVM and PNN. *Expert Systems with Applications, 37*(2010), 7154–7161.

Bagtzoglou, A.C. (2003). On the non-locality of reversed time particle tracking methods. *Environ. Forensics, 4*(3), 215-225.

Bagtzoglou, A.C., & Atmadja, J. (2003). Marching-jury backward beam equation and quasi-reversibility methods for hydrologic inversion: application to contaminant plume spatial distribution recovery. *Water Resour. Res., 39*(2), 10-14. SBH 10-1. doi: 10.1029/2001WR001021

Bagtzoglou, A.C., & Atmadja, J. (2005). Mathematical methods for hydrologic inversion: the case of pollution source identification. *Chapter in Environmental Impact Assessment of Recycled Wastes on Surface and Ground Waters, In: Kassim, T.A. (Ed.), Engineering Modeling and Sustainability. The Handbook of Environmental Chemistry, Water Pollution Series*, vol. 3. Springer-Verlag, Heidelberg-New York, ISBN 3-540-00268-5, pp. 65-96.Volume 5, Part F

Bagtzoglou, A.C., & Baun, S.A. (2005). Near real-time atmospheric contamination source identification by an optimization based inverse method. *Inverse Prob. Sci. Eng., 13*(3), 241-259.

Bagtzoglou, A.C., Dougherty, D.E., & Tompson, A.F.B. (1992). Application of particle methods to reliable identification of groundwater pollution sources. *Water Resour. Manage., 6*, 15-23.

Bagtzoglou, A. C., Tompson, A. F. B., & Dougherty, D. E. (1991). Probabilistic simulation for reliable solute source identification in heterogeneous porous media. In Ganoulis, J., editor, *Water Resources Engineering Risk Assessment*, pages 189–201. Springer-Verlag, Heidelberg, 1991. 38.

Bear, J. (1979). *Hydraulics of groundwater*. New York, N.Y., 241: McGraw-Hill,

Cannon, J., R. (1966). Some numerical results for the solution of the heat equation backwards in time. In: *Numerical Solutions of Non Linear Differential Equations* (*Proc. Adv. Sympos. Madison Wis.,* 1966). pp.21-54, John Wiley and Sons, Inc., New York, 1966, MR 34, 7037

Cerny, V. (1985). Thermodynamical approach to the traveling salesman problem: An efficient simulation algorithm. *Journal of Optimization Theory and Applications,* 45, 41-51. doi:10.1007/BF00940812

Chadalavada, S., & Datta, B. (2008). Dynamic optimal monitoring network design for transient transport of pollutants in groundwater aquifers. *Water Resource Management, 22*, 651-670.

Chadalavada, S., Datta, B. & Naidu, R. (2011a). Uncertainty based optimal monitoring network design for chlorinated hydrocarbon contaminated site. *Environment Monitoring Assess. 173*, 929-940.

Chadalavada, S., Datta, B. & Naidu, R. (2011 b). Optimisation approach for pollution source identification in groundwater: an overview. *International Journal of Environment and Waste Management, 8*(1-2). pp. 40-61.

Cieniawski, S. E., Eheart, W. J. & Ranjithan, S. (1995). Using genetic algorithm to solve a multiple objective groundwater monitoring problem. *Water Resource Research. 31*(2), 399–409.

Cramer, N. L. (1985). A representation for the Adaptive Generation of Simple Sequential Programs. In: *Proceedings of an International Conference on Genetic Algorithms and the Applications*, Grefenstette, John J. (ed.), Carnegie Mellon University.

Cressie, N. (1990). The origin of kriging. *Mathematical Geology*, 22p 239-252

Datta, B., Beegle, J.E., Kavvas, M.L., & Orlob, G.T. (1989). Development of an expert-system embedding pattern-recognition techniques for pollution source identification. *Technical Report: PB-90-185927/XAB, OSTI ID:6855981,* Dept. of Civil Engineering, California Univ., Davis, CA, USA.

Datta, B., & Dhiman, D. S. (1996). Chance-constrained optimal monitoring network design for pollutants in groundwater. *Journal of Water Resource Planning & Management, 122*(3), 180–188.

Datta, B., Chakrabarty, D., & Dhar, A. (2009a). Optimal dynamic monitoring network design and identification of unknown groundwater pollution sources. *Water Resour. Manage., Springer 23*(10), 2031-2049.

Datta, B., Chakrabarty, D., & Dhar, A. (2009b). Simultaneous identification of unknown groundwater pollution sources and estimation of aquifer parameters. *J. Hydrol., 376*(1, 2), 48-57. Elsevier.

Datta, B., Chakrabarty, D., & Dhar, A. (2011). Identification of unknown groundwater pollution sources using classical optimization with linked simulation. *Journal of Hydro-environment Research, 5*(2011), 25-36. Elsevier.

Deutsch, C. V., & Journel, A. G. (1998). *GSLIB: Geostatistical Software Library and user's guide*. New York: Oxford University Press.

Dhar, A., & Datta, B. (2007). Multi-objective design of dynamic monitoring networks for detection of groundwater pollution. *Journal of Water Resource Planning and Management. 133*(4), 329–338.

Dhar, A., & Datta, B. (2010). Logic-based design of groundwater monitoring network for redundancy reduction. *Journal of Water Resource Planning and Management*, 136, 88.

Domenico, P. A., & Schwartz, F. W. (1998). *Physical and Chemical Hydrogeology*. 2nd Ed., NewYork: John Wiley & Sons, Inc.

Feng-guang, Y., Shu-you, C., Xing-nian, L., & Ke-jun, Y. (2008). Design of groundwater level monitoring network with ordinary kriging. *J. Hydrodyn., 20*(2008), pp. 339–346.

Fethi, B. J., Loaiciga, A. H., & Marino, A. M. (1994). Multivariate geostatistical design of groundwater monitoring networks. *Journal of Water Resource Planning and Management, ASCE 120*(4), 505–522.

Forsyth, R., & Roy, R. (1986). Machine Learning Applications in Expert Systems and Information Retrieval, Ellis Horwood series. In: *Artificial intelligence*, Chichester, UK.

Goffe, W. L. (1996). *SIMANN: A global optimization algorithm using Simulated Annealing, Studied in Nonlinear Dynamics and Econometrics*. Berkeley Electronic Press.

Gorelick, S.M., Evans, B., & Ramson, I. (1983). Identifying sources of groundwater pollution: an optimization approach. *Water Resour. Res., 19*(3), 779-790.

Grabow, G., Mote, R.C., & Yoder, C. D. (2000). An empirically-based sequential ground water monitoring network design procedure. *Journal of American Water Resource Association. 36*(3), 549–566.

Hansen, T. M. (2004). mGstat V 0.99:MATLAB code, http://sourceforge.net/projects/mgstat/files/

Harbaugh, A.W., Banta, E.R., Hill, M.C., & McDonald, M.G. (2000). *MODFLOW-2000, the U.S. Geological Survey modular ground-water model*. U.S. Geological Survey Open-File Report 00-92, 121 p.

He, L., Huang, G. H., & Lu, H. W. (2009). A coupled simulation-optimization approach for groundwater remediation design under uncertainty: An application to a petroleum-contaminated site. *Environmental Pollution, 157*(8-9), 2485–2492.

Hudak, P. F., Loaiciga, A. H., & Marino, A. M. (1995). Regional-scale ground water quality monitoring via integer programming. *Journal of Hydrology (Amst). 164*(1–4), 153–170.

Javandel, I., Doughty, C., & Chin-Fu, T. (1984). *Groundwater Transport: Handbook of Mathematical Models. Water Resources Monograph No. 10*. Washington, D.C.: American Geophysical Union.

Jha, M. K., & Datta, B. (2011). Simulated annealing based simulation-optimization approach for identification of unknown contaminant sources in groundwater aquifers. *Desalination and Water Treatment, 32*(1-3), 79-85.

Jha, M. K., & Datta, B. (2012a). *Linked Simulation-Optimization Based Methodologies for Unknown Groundwater Pollutant Source Identification in Managed and Unmanaged Contaminated Sites*. Chapter 4, PhD Thesis, James Cook University.

Jha, M. K., & Datta, B. (2012b). Application of Simulated Annealing in Water Resources Management: Optimal Solution of Groundwater Contamination Source Characterization Problem and Monitoring Network Design Problems. *In: Simulated Annealing- Single and Multiple Objective Problems.* In:Tech, Rijeka, Croatia, pp. 157-174, DOI: 10.5772/45871.

Journel, A. G., & Huijbregts, C. J. (1978). *Mining Geostatistics.* London: Academic Press, 600p.

Kirkpatrick, S., Gelatt, D. C., & Vecchi, P. M. (1983). Optimization by simulated annealing. *Science, 220*, 671-680.

Kollat, J. B., Reed, P. M., & Kasprzyk, J. R. (2008). A New Epsilon-Dominance Hierarchical Bayesian Optimization Algorithm for Large Multi-Objective Monitoring Network Design Problems. *Advances in Water Resources, 31*(5), 828-845, 2008.

Kollat, J. B., Reed, P. M., & Maxwell, R. (2011). Many-objective Groundwater Monitoring Network Design using Bias-Aware Ensemble Kalman Filtering, Evolutionary Optimization, and Visual Analytics. *Water Resource Research, v47*, W02529.

Koza, J.R. (1994). Genetic programming as a means for programming computers by natural selection. *Statistics and Computing*, 10.1007/BF00175355.

Liu, C., & Ball, W., P. (1999). Application of inverse methods to contaminant source identification from aquitard diffusion profiles at Dover AFB, Delaware. *Water Resour. Res., 35*(7), 1975-1985.

Loaiciga, H. A. (1989). An optimization approach for groundwater quality monitoring network design. *Water Resource Research. 25*(8), 1771–1782.

Loaiciga, H. A., & Hudak, P. F. (1992). A location modelling approach for groundwater monitoring network augmentation. *Water Resource Researce. 28*(3), 643–649.

Loaiciga, H. A., & Hudak, P. F. (1993). An optimization method for network design in multilayered groundwater flow systems. *Water Resource Research, 29*, 2835.

Mahar, P. S., & Datta, B. (1997). Optimal monitoring network and ground-water-pollution source identification. *Journal of Water Resource Planning and Management, 123*(4), 199–207.

Mahar, P. S., & Datta, B. (2000). Identification of pollution sources in transient groundwater system. *Water Resour. Manage., 14*(6), 209-227.

Mahar, P. S., & Datta, B. (2001). Optimal identification of ground-water pollution sources and parameter estimation. *J. Water Resour. Plan. Manage., 127*(1), 20-29.

Mahinthakumar, G., & Sayeed, M. (2005). Hybrid genetic algorithm-local search methods for solving groundwater source identification inverse problems. *J. Water Resour. Plan. Manage., 131*(1), 45-57.

Marsden, J. (2011). Report for Strengthening Basin Communities - Planning Study Business Case - Enhancing Dubbo's Irrigation System. *Dubbo City Council*.

Massmann., J., & Freeze, R. A. (1987). Groundwater pollution from waste management sites: the interaction between risk-based engineering design and regulatory policy. I: Methodology. *Water Resource Research, 23*(2), 351–367.

Matheron, G. (1963). Principles of Geostatistics. *Economic Geology*, 58, pp 1246-1266.

McDonald, M. G., & Harbaugh, A. W. (1988). A modular three-dimensional finite-difference groundwater flow model. *Techniques of Water-Resources Investigations of the United States Geological Survey, Book 6*, Chapter A1, 586 p.

Meyer, P. D., & Brill, E. D. Jr. (1988). A method for locating wells in a groundwater pollution monitoring network under conditions of uncertainty. *Water Resource Research, 24*(8), 1277–1282.

Meyer, P. D., Valocchi, A. J., & Eheart, J. W. (1994). Monitoring network design to provide initial detection of groundwater pollution. *Water Resource Research, 30*, 2647.

Metropolis, N., Rosenbluth, A., Rosenbluth, M., Teller, A., & Teller, E. (1953). Equation of state calculations by fast computing machines. *J. Chem. Phys., 21*, 1087–1092.

Michalak, A.M., & Kitanidis, P.K. (2004a). Estimation of historical groundwater contaminant distribution using the adjoint state method applied to Geostatistical inverse modeling. *Water Resour. Res., 40*, W08302. doi:10.1029/2004WR003214

Michalak, A.M., & Kitanidis, P.K. (2004b). Application of geostatistical inverse modeling to contaminant source identification at Dover AFB, Delaware. *IAHR J. Hydraul. Res., 42 (extra issue)*, 9-18.

Montas, H. J., Mohtar, R. H., Hassan, A. E., & AlKhal, F. A. (2000). Heuristic space-time design of monitoring wells for pollutant plume characterization in stochastic flow fields. *Journal of Contaminant Hydrology. 43*(3–4), 271–301.

Mugunthan, P., & Shoemaker, C. A. (2004). Time varying optimization for monitoring multiple pollutants under uncertain hydrogeology. *Bioremediation Journal. 8*(3–4), 129–146.

Neupauer, R.M., & Wilson, J.L. (1999). Adjoint method for obtaining backward in-time location and travel probabilities of a conservative groundwater contaminant. *Water Resour. Res., 35*(11), 3389-3398.

Neupauer, R.M., Borchers, B., & Wilson, J.L. (2000). Comparison of inverse methods for reconstructing the release history of a groundwater contamination source. *Water Resour. Res., 36*(9), 2469-2475.

Neupauer, R.M., & Wilson, J.L. (2005). Backward probability model using multiple observations of contamination to identify groundwater contamination sources at the Massachusetts military reservation. *Water Resour.Res., 41*, W02015. doi:10.1029/2003WR002974

Nunes, L. M., Cunha, M. C., & Ribeiro, L. (2004a). Groundwater monitoring network optimization with redundancy reduction. *Journal of Water Resource Planning and Management. 130*(1), 33–43.

Nunes, L. M., Cunha, M. C., & Ribeiro, L. (2004b). Optimal space-time coverage and exploration costs in groundwater monitoring networks. *Environment Monitoring Assess. 93*(1–3), 103–124.

Pinder, G. (2009). Optimal search strategy for the definition of a dnapl source. *Technical report, Defense Technical Information Cente*r, United States of America.

Puech, V. (2010). Upper macquarie groundwater model. *Technical Report VW04680*, Office of Water, NSW Government and National Water Commission, Australia.

Reed, P. M., & Kollat, J. B. (2012). Save now, pay later? Multi-period many-objective groundwater monitoring design given systematic model errors and uncertainty. *Advances in Water Resources, Volume 35*, Pages 55-68.

Reed, P. M., & Minsker, B. S. (2004). Striking the balance: long-term groundwater monitoring design for conflicting objective. *Journal of Water Resource Planning and Management. 130*(2), 140–149.

Reed, P. M., Minsker, B. S., Goldberg, D. E. (2000). Designing a competent simple genetic algorithm for search and optimization. *Water Resources Research 36*(12), 3757-3761.

Rushton, K. R., & Redshaw, S. C. (1979). *Seepage and Groundwater Flow*. New York: Wiley.

Sreekanth, J., & Datta, B. (2012). Genetic programming: efficient modeling tool in hydrology and groundwater management. In: *New Genetic Programming- New Approaches and Successful Applications*. InTech, Rijeka, Croatia, pp. 227-240.

Shiklomanov, I., (1993). *World fresh water resources, Water in Crisis: A Guide to the World's Fresh Water Resources*. New York: Oxford University Press.

Sidauruk, P., Cheng, A. H. D., & Ouazar, D. (1998). Groundwater contaminant source and transport parameter identification by correlation coefficient optimization. *Groundwater, 36*, 208-214.

Singh, R.M., & Datta, B. (2004). Groundwater pollution source identification and simultaneous parameter estimation using pattern matching by artificial neural network. *Environ. Forensics, 5*(3), 143-159.

Singh, R.M., Datta, B., & Jain, A. (2004). Identification of unknown groundwater pollution sources using artificial neural networks. *J. Water Resour. Plan. Manage., 130*(6), 506-514.

Singh, R.M., & Datta, B. (2006). Identification of groundwater pollution sources using GA-based linked simulation optimization model. *J. Hydrol. Eng., 11*(2), 101-109.

Singh, R.M., & Datta, B. (2007). Artificial neural network modeling for identification of unknown pollution sources in groundwater with partially missing concentration observation data. *Water Resour. Manage., 21*(3), 557-572.

Skaggs, T. H., & Kabala, Z. J. (1994). Recovering the release history of a groundwater contaminant. *Water Resour. Res., 30*(1), 71-79.

Skaggs, T. H., & Kabala, Z. J. (1995). Recovering the release history of a groundwater contaminant plume: method of quasi-reversibility. *Water Resour. Res., 31*(11), 2669-2673.

Snodgrass, M. F., & Kitanidis, P. K. (1997). A geostatistical approach to contaminant source identification. *Water Resour. Res., 33*(4), 537-546.

Sun, A. Y., Painter, S. L., & Wittmeyer, G. W. (2006a). A constrained robust least squares approach for contaminant source release history identification. *Water Resour. Res., 42*(4), W04414. doi:10.1029/2005WR004312

Sun, A. Y., Painter, S. L., & Wittmeyer, G. W. (2006b). A robust approach for contaminant source location and release history recovery. *J. Contam. Hydrol., 88*(3, 4), 29-44.

Wagner, B. J. (1992). Simultaneous parameter estimation and contaminant source characterization for coupled groundwater flow and contaminant transport modeling. *J. Hydrol., 135*, 275-303.

Wang, H. F., & Anderson, M. P. (1982). *Introduction to Groundwater Modeling: Finite Difference and Finite Element Methods*. San Francisco: W. H. Freeman and Company, 237 pp.

Woodbury, A. D., & Ulrych, T. J. (1996). Minimum relative entropy inversion: theory and application to recovering the release history of a groundwater contaminant. *Water Resour. Res., 32*(9), 2671-2681.

Woodbury, A. D., Sudicky, E., Ulrych, T. J., & Ludwig, R. (1998). Three dimensional plume source reconstruction using minimum relative entropy inversion. *J. Contam. Hydrol., 32*, 131-158.

Wu, J., Zheng, C., & Chien, C. C. (2005). Cost-effective sampling network design for contaminant plume monitoring under general hydrogeological conditions. *Journal of Contaminant Hydrology. 77*, 41–65.

Yeh, W. W. -G. (1986). Review of parameter identification procedure in groundwater hydrology: the inverse problem. *Water Resour. Res., 22*(2), 95-108.

Yeh, M. S., Lin, Y. P., & Chang, L. C. (2006). Designing an optimal multivariate Geostatistical groundwater quality monitoring network using factorial Kriging and genetic algorithm. *Journal of Environmental Geology, 50*, 101-121.

Zheng, C., (1990). *MT3D, A modular three-dimensional transport model for simulation of advection, dispersion and chemical reactions of contaminants in groundwater systems*. Report to the U.S. Environmental Protection Agency, 170 p.

Zheng, C., & Wang, P. P. (1999). *MT3DMS, A modular three-dimensional multi-species transport model for simulation of advection, dispersion and chemical reactions of contaminants in groundwater systems*. U.S. Army Engineer Research and Development Center Contract Report SERDP-99-1, Vicksburg, MS, 202 p

Zheng C., & Bennett G. D. (1995). *Applied Contaminant Transport Modeling: Theory and Practice*. 115 Fifth Avenue, New York, NY 10003,USA Van Nostrand Reinhold, 440 + xxiii pp.