7-19-2011

# Assessing Uncertainty Associated with Groundwater and Watershed Problems Using Fuzzy Mathematics and Generalized Regression Neural Networks

Bree R. Mathon
*University of Vermont*

Follow this and additional works at: http://scholarworks.uvm.edu/graddis

# ASSESSING UNCERTAINTY ASSOCIATED WITH GROUNDWATER AND WATERSHED PROBLEMS USING FUZZY MATHEMATICS AND GENERALIZED REGRESSION NEURAL NETWORKS

A Dissertation Presented

by

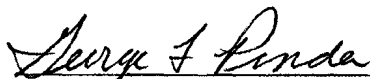Bree R. Mathon

to

The Faculty of the Graduate College

of

The University of Vermont

In Partial Fulfillment of the Requirements
for the Degree of Doctor of Philosophy
Specializing in Civil and Environmental Engineering

January, 2011

Accepted by the Faculty of the Graduate College, The University of Vermont, in partial fulfillment of the requirements for the degree of Doctor of Philosophy, specializing in Civil and Environmental Engineering.

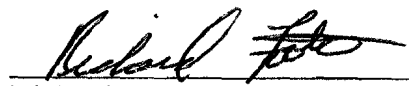Dissertation Examination Committee:

George F. Pinder, Ph.D.      Co-Advisor

Donna M. Rizzo, Ph.D.      Co-Advisor

Richard Foote, Ph.D.

Lori Stevens, Ph.D.      Chairperson

Domenico Grasso, Ph.D.      Dean, Graduate College

Date: September 23, 2010

# ABSTRACT

When trying to represent an environmental process using mathematical models, uncertainty is an integral part of numerical representation. Physically-based parameters are required by such models in order to forecast or make predictions. Typically, when the uncertainty inherent in models is addressed, only aleatory uncertainty (irreducible uncertainty) is considered. This type of uncertainty is amenable to analysis using probability theory. However, uncertainty due to lack of knowledge about the system, or epistemic uncertainty, should also be considered. Fuzzy set theory and fuzzy measure theory are tools that can be used to better assess epistemic, as well as aleatory, uncertainty in the mathematical representation of the environment.

In this work, four applications of fuzzy mathematics and generalized regression neural networks (GRNN) are presented. In the first, Dempster-Shafer theory (DST) is used to account for uncertainty that surrounds permeability measurements and is typically lost in data analysis. The theory is used to combine multiple sources of subjective information from two expert hydrologists and is applied to three different data collection techniques: drill-stem, core, and pump-test analysis. In the second, a modification is made to the fuzzy least-squares regression model and is used to account for uncertainty involved in using the Cooper-Jacob method to determine transmissivity and the storage coefficient. A third application, involves the development of a GRNN to allow for the use of fuzzy numbers. A small example using stream geomorphic condition assessments conducted in the state of Vermont is provided. Ultimately, this fuzzy GRNN will be used to better understand the relationship between the geomorphic and habitat conditions of stream reaches and their corresponding biological health. Finally, an application of the GRNN algorithm to explore links between physical stream geomorphic and habitat conditions and biological health of stream reaches is provided. The GRNN proves useful; however, physical and biological data collected concurrently is needed to enhance accuracy.

# CITATIONS

Material from this dissertation has been published in the following form:

Material from this dissertation has been published in the following form:

# ACKNOWLEDGEMENTS

So many people have helped me along this path, to begin I would like to thank the incredible faculty of the Civil and Environmental Engineering program for welcoming me (even when I was a graduate student in the math department) and teaching me how to think more like an engineer. I would also like to thank the faculty of the Mathematics and Statistics Department for the wonderful education they provided me. A special thanks to the faculty who served on my comprehensive exam committee: Arne Bomblies, Britt Holmén, George Pinder, Donna Rizzo, and Mun Son.

Of course this dissertation would not have occurred without the constant support and enthusiasm from my co-advisors George Pinder and Donna Rizzo. George has been a continuous source of support, guidance, and advice throughout my years at UVM, especially when I was trying to figure out exactly what research path was right for me. He is an unending source of knowledge and experience; I am appreciative for the time I had to work with him. Donna is an amazing human being, advisor, and teacher. She always made the time to listen and offer advice even before she became my co-advisor. I would like to thank Donna for the opportunity to work with her and all that she taught me about artificial neural networks, coding, writing and geomorphology. She has helped me more than she will ever know.

I met Lori Stevens when I started working on the watershed projects. During the time we worked together, I developed an appreciation for oligochaetes and macroinvertebrates. I thank Lori for sharing her time and knowledge, and for her

I thank all my friends that have continued to believe in me and support me over the years. A special thanks to some of my dearest friends Kristin, Zoi, Helen, and Kirsten who were always there to listen and offer advice through the good and the occasional bad times.

My family deserves a huge thank you and hug for all the support and encouragement they have provided. My parents' love and belief in me means so much and I am forever grateful for that.

And finally, I am incredibly grateful to my husband Greg. He has been so understanding and patient during the longest "three more years" of his life. I want to thank him for always knowing what I needed (at times even when I didn't) and making sure I was fed during the crazy times. As this one journey ends, I look forward to our future journeys together…

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# CHAPTER 1
# INTRODUCTION AND OBJECTIVES

## 1.1 Introduction

When trying to represent an environmental process using mathematical models, uncertainty is an integral part of numerical representation. Physically-based parameters are required by such models in order to forecast or make predictions. For example, in subsurface hydrology, soil permeability must be specified in equations descriptive of groundwater flow. Most permeability measurements are assumed to represent the area immediately surrounding the measurement. However, due to the subsurface heterogeneity, these measurements say very little about porous medium as you move further away from the measurement location. Many deterministic models use the observations at hand and ignore the matter of the uncertainty.

Typically, when the uncertainty inherent in models is addressed, only aleatory uncertainty (irreducible uncertainty) is considered. This type of uncertainty is amenable to analysis using probability theory. However, uncertainty due to lack of knowledge about the system, or epistemic uncertainty, should also be considered. The theory of fuzzy mathematics is a tool that allows incorporation of epistemic uncertainty into the mathematical representation.

### 1.1.1 Overall Goal and Specific Objectives

The overall goal of this dissertation is to apply nontraditional mathematical tools, i.e., fuzzy set theory, fuzzy measure theory and a generalized regression neural network, to

better assess epistemic (as well as aleatory) uncertainty. These methods are used in two environmental application areas.

1. Groundwater applications:

    a. Dempster-Shafer theory (DST, Dempster, 1967 and Shafer, 1976) is used to account for uncertainty associated with soil permeability measurements.

    b. A modified fuzzy least-squares regression (in place of linear regression) and the Cooper-Jacob method (Cooper and Jacob, 1946) are used to determine subsurface transmissivity and storage coefficient membership functions.

2. Watershed applications:

    a. A generalized regression neural network (GRNN) is used to explore linkages amongst physical geomorphic and habitat condition, and biological health.

    b. From the above GRNN work, the algorithm is modified to accommodate the use of fuzzy numbers as inputs and outputs since the assessments conducted on geomorphic and habitat condition contain subjective information.

## 1.2 Dissertation Overview

Chapter 1 continues with a literature review on the use of Dempster-Shafer theory (DST), fuzzy least-squares regression (FLSR), and generalized regression neural networks (GRNN) in environmental applications. Chapter 2 applies DST combination rules to join field-measured permeability data (quantitative data) with hydrogeologists'

expert opinions (subjective information) to examine uncertainty. Three data sets consisting of permeability (k) values measured in the Dakota Sandstone within the Denver Basin (Belitz and Bredehoeft, 1988) were analyzed. Each data set has a different collection method: well-water pump-test, core analysis, and drill-stem analysis. Dempster's rule of combination (chosen to combine the two forms of information), which has received criticism in the literature (Zadeh, 1986; Yager, 1987), was compared to two alternative combination methods.

Chapter 3 discusses the development of a modified fuzzy least-squares regression (MFLSR) method that allows the use of imprecise pumping-test data to obtain fuzzy intercept and slope values that are then used in the Cooper-Jacob method. Fuzzy membership functions for the soil transmissivity and storage coefficient are then calculated using the extension principle. The supports of the fuzzy membership functions incorporate the transmissivity and storage coefficient values that would be obtained using ordinary least-squares regression and the Cooper-Jacob method. The MFLSR coupled with the Cooper-Jacob method allows for the inclusion of inherent uncertainty due to a lack of knowledge regarding the heterogeneity of the subsurface. The methodology is tested on a pumping-test data set collected in an intermediate scale groundwater facility.

In Chapter 4 the focus of the application is on the ability to identify streams with high environmental risk, which is essential for a proactive adaptive watershed management approach. In efforts to describe the health and geomorphic condition of streams, environmental managers must gather and assess various forms of information - quantitative, qualitative and subjective. These geomorphic and habitat assessments used to characterize streams include some uncertainty, and fuzzy numbers can be used to

capture this. In this work, a new ANN is created by embedding the ability to calculate fuzzy numbers into the hidden (called pattern) nodes of the GRNN to leverage the uncertainty associated with field data collected by experts. The Vertex Method (Dong and Shah, 1987) is used to calculate the crisp functions in the algorithm using fuzzy numbers. This allows uncertainty from expert field assessments to be accounted for in the data analysis, typically this is not quantified. The algorithm is tested and validated on habitat and geomorphic assessment data collected by the Vermont Agency of Natural Resources (VTANR) River Management Program throughout the state.

Chapter 5 is an application of the GRNN, developed by Donald Specht (1991), to explore linkages between the geomorphic, physical habitat and biological health of stream reaches in Vermont. Since physical processes occurring in a stream form the habitat, habitat assessments look at physical ecological parameters that might help understand the relationship between fluvial processes and aquatic communities (VTANR, 2008). The GRNN is first used to predict habitat conditions for stream reaches throughout the state of Vermont using only geomorphic data. Further analysis added biological health data (fish and macroinvertebrate) into the algorithm, first as an input, then as the output.

Chapter 6 presents a summary and discussion of the projects in this dissertation. Appendices A and B contain the MATLAB (The MathWorks, 2010) codes for the GRNN algorithms used in this work.

## 1.3 Literature Review

A review of applications of fuzzy mathematics, specifically Dempster-Shafer theory (DST) and fuzzy least-squares regression (FLSR), is followed by an introduction to artificial neural networks and review of the applications of GRNNs.

### 1.3.1 Quantifying uncertainty using fuzzy mathematics

#### 1.3.1.1 Dempster-Shafer Theory

Dempster-Shafer theory (DST), also known as evidence theory (Shafer, 1976), is a branch of the theory of monotone measures, a generalization of classical measure theory (Klir, 2003), and is one of the few areas of mathematics developed to explore uncertainty due to a lack of knowledge about the system. Strengths of the DST framework include its well developed theory, ability to handle various types of evidence (consonant, consistent, arbitrary, or disjoint), ability to combine evidence from different sources, lack of any assumptions about the distribution of the data, and ability to use all available data (outliers are not discarded from the analysis).

Fields of study where DST has been applied include a study that combines fuzzy logic with DST to evaluate slope instability (Binaghi *et al.*, 1998). Agarwal *et al.* (2004) quantify uncertainty of design tools in multidisciplinary systems analysis. Cayuela *et al.* (2006) apply DST to remote sensing information along with expert opinion to more accurately classify land cover. Kriegler and Held (2005) use DST to model future climate change, which they then project to make an estimate of global mean warming. Carranza and Hale (2003) use DST to produce data-driven (instead of knowledge-driven)

maps of gold potential in the Baguio district of the Philippines. Luo and Caselton (1997) explore the use of DST to address uncertainty associated with climate change models.

The characteristic of DST of most interest in this dissertation is its ability to combine multiple sources of evidence. The original method derived to combine data, Dempster's rule of combination, has been criticized (Zadeh, 1986; Yager, 1987) for how conflicting evidence is handled because it provides counterintuitive results when the level of conflict among the evidence is high (Zadeh, 1984). Several papers discuss the different combination rules developed to overcome this difficulty (Sentz and Ferson, 2002; Smarandache, 2004; and Smets, 2005). Two other rules are compared to Dempster's rule in this work: Yager's rule and the Hau-Kashyap method, which differ in how they handle conflicting evidence. In the case where there is no or little conflict, Yager's rule and the Hau-Kashyap method produce very similar, if not identical, results to Dempster's rule, however the methods are superior when conflict is greater.

Dempster's rule of combination has found numerous applications where conflict is low. It is used to combine evidence from remote sensing information to assist in the production of accurate plant functional type maps (Sun *et al.,* 2008). In Bi *et al.* (2007), Dempster's rule of combination is used to explore the impact of combining four machine-learning methods for text categorization. In this dissertation, various combination rules are used to combine evidence on permeability data from two independent sources (Mathon *et al.*, 2010).

**1.3.1.2 Fuzzy least-squares regression**

Regression is a statistical tool that is widely used to examine the relationship between dependent and independent variables. Ordinary regression is capable of analyzing and producing models only for crisp data. In reality, data can have a fuzziness to it that when ignored, weakens the model used for prediction. Fuzzy set theory has been used to develop fuzzy regression, which can better address the uncertainties associated with the regression model of fuzzy data. Fuzzy regression was originally introduced by Tanaka *et al*. (1982) and since then several other fuzzy regression methods have been developed (Chang and Ayyub, 2001). The regression method that will be discussed in this dissertation is fuzzy least-squares regression (FLSR) developed by Savic and Pedrycz (1991). Fuzzy least-squares regression has proven to be effective when few data points are available, as traditional statistics need a large number of data points to be valid (Bardossy et al., 1990; Ozelkan and Duckstein, 2000). This lends itself nicely to hydrologic applications where data collection is often expensive, and thus sparse.

Applications of FLSR in hydrology include Bardossy *et al.* (1990) where FLSR was used in a case study that looked at the imprecise relationship between electrical resistivity and hydraulic permeability of soil. Groundwater availability was assessed by Uddameri and Honnunger (2007). Uddameri (2004) used FLSR to explore the relationship between scale and longitudinal dispersivity. Si and Bodhinayake (2005) used FLSR to determine soil hydraulic properties (e.g., hydraulic conductivity) using tension infiltrometer measurements. Ozelkan and Duckstein (2001) estimated parameters for a rainfall-runoff model using FLSR.

The mathematics of the FLSR method is explained in detail in Chapter 3 of this dissertation, however the general methodology is outlined here. FLSR can be used with crisp dependent, X, and independent, Y, data or with crisp X and fuzzy triangular Y data. The method proposed by Savic and Pedryz (1991) combines two steps. First, least-squares regression is carried out using the crisp X and crisp or center value of the fuzzy Y data. This provides the regression model with the center values for the regression coefficients (i.e., slope and intercept in the bivariate linear case). The second step uses optimization to obtain the halfwidth (or fuzziness) of the coefficients. The resulting coefficients take the form of symmetrical triangular functions.

One key problem for FLSR (as well as other fuzzy regression models) has been the inability to approach the ordinary regression model when the data are crisp and there is no fuzziness associated with the system (Chang and Ayyub, 2001). A hybrid FLSR was introduced by Chang (2001), however that methodology produced negative halfwidths, which was confusing, as halfwidths are defined by positive values. One part of this dissertation developed an alternative modified FLSR that reduced to the ordinary regression model in the crisp case (Mathon *et al*., 2008).

## 1.3.2 Artificial neural networks

Artificial neural networks (ANNs) are parallel, nonparametric statistical methods that can be used in pattern classification, pattern completion, function approximation, prediction, optimization, and system control applications among others (Wasserman, 1993). In general, a supervised ANN takes an input vector and maps it to either a vector

or a scalar output. The mapping relationship is defined by a set of weights that are determined during a training phase.

ANNs are trained via one of two methods: supervised or unsupervised. In supervised learning, the network is presented with a training data set that consists of input values and their corresponding output or target value(s). During training, the algorithm produces an output from the input vector, which is then compared to the target output and the weights of the algorithm are adjusted to minimize the distance between the ANN output and the target values. This process is done iteratively until a user-defined amount of error is achieved. Once the weights produce a satisfactory mapping between training inputs and outputs, they are fixed, and used to predict output from additional input vectors. In unsupervised learning algorithms, the network does not have a target output vector to compare predictions values. Instead, training is accomplished using only input vectors and the weights are adjusted to achieve similarities in the data (e.g. classifying like things together; clustering).

The most commonly used supervised algorithm is the feed-forward back-propagation network (FFBP). In fact, more than 95% of ANNs currently reported in environmental engineering literature have used either a FFBP or a radial basis function (RBF) neural network (Govindaraju and Ramachandra, 2000). In this work, however, an alternative ANN, a GRNN (Specht, 1991), is used to explore watershed management issues. In addition, a new GRNN algorithm is created to allow for the use of fuzzy numbers as inputs and outputs in order to capture expert opinion typically not captured in field geomorphic and habitat assessments.

GRNNs have been used in various modeling applications, and fuzzy mathematics has been used to preprocess data or as a comparison to GRNN results. However, the project presented in Chapter 4 appears to be the first use of fuzzy numbers in a modified GRNN algorithm.

**1.3.2.1 Generalized regression neural network**

The GRNN is a one-pass learning algorithm with a parallel structure capable of estimating continuous variables (Specht, 1991) and designed to be used on data where the functional form is unknown (i.e., a linear assumption cannot be validated and the order of the "optimal" polynomial is unknown). Due to its one-pass design, it does not require iterative training like the more widely used FFBP. The advantages of the GRNN are 1) the computational speed; 2) the ability to update easily as new information becomes available; and 3) the accuracy of prediction from sparse data sets.

The GRNN has extensive applications in the water resources and hydrological fields. Aksoy and Dahamsheh (2009) explore using a GRNN for forecasting monthly precipitation. Several studies have had success predicting leaf wetness (Chtioui *et al*., 1999a; Chtioui *et al*., 1999b) and evapotranspiration (Kim and Kim, 2008; Kisi, 2008a). Cigizoglu and Alp (2004) found the GRNN to be successful in predicting rainfall runoff and, unlike the radial basis function and multiple linear regression, did not produce negative flow estimations. There have been numerous applications of GRNNs in forecasting stream flow. Firat (2008) explored its use in daily stream flow forecasting, while Ng *et al.* (2009) use a GRNN to estimate missing observations in extreme daily streamflow records. Besaw *et al.* (2009a) use a recurrent feed-back loop on a GRNN to

predict flow in ungauged streams. Several studies found the GRNN to outperform the FFBP when forecasting intermittent stream flow (Cigizoglu, 2005a) and monthly stream flow (Cigizoglu, 2005b; Kisi, 2008b). Turan and Yurdusev (2009) tested the GRNN on the prediction of stream flow from measured upstream flow records. The GRNN was also used to predict water quality based on rainfall, surface discharge and nutrient concentration (Kim and Kim, 2007) and to estimate daily mean sea level heights (Sertel *et al.*, 2008).

The GRNN has also been used to help manage water supply. Asefa *et al.* (2007) predict groundwater levels from one to four weeks into the future for the purposes of water demand planning. Chlorine residuals in a water distribution system were predicted to ensure that the water is safe for human consumption (Bowden *et al.,* 2006), and monthly water consumption was forecasted based on several socio-economic and climatic factors (Firat *et al.*, 2009).

River sediment transport has also been modeled with GRNNs (Cigizoglu and Alp, 2006; Cobaner *et al*., 2009; Kisi *et al.*, 2008). Wang *et al.* (2009) used data collected in a weir during storm events for one year (the variables considered were turbidity, water discharge, and suspended sediment concentrations) in a GRNN to model event-based suspended sediment concentration following storm events.

Predicting environmental contamination is another area where GRNNs have been applied. Kanevski *et al.* (1999) use a GRNN for spatial prediction of surface soil contamination by radionuclides released during the Chernobyl accident in 1986. In a "what if" scenario of biological contamination of water systems, Kim *et al.* (2008) use *E. coli* transport patterns and a GRNN to locate the pathogenic release location. Ligang *et*

11

*al.* (2008) modeled the relationship between coal-fired boilers and NOx emissions into the environment. Li *et al.* (2008) predicted nitrogen concentrations in disturbed and undisturbed streams and Durdu (2009) used a GRNN to spatially predict polluted surface water.

Agricultural applications of GRNNs include predicting nitrate release from a controlled release fertilizer (Du *et al.,* 2008), which found the thickness of the polymer coating on the fertilizer was the most important factor controlling nitrate release. Sun *et al.* (2008) found a GRNN preferable over a FFBP to model air-quality near livestock production facilities.

Finally, Ustaoglu *et al.* (2008) found GRNNs compared quite well to the conventional method of multiple linear regression when forecasting daily mean, maximum, and minimum temperature time series as related to agriculture, water resources and tourism.

**1.3.2.2 Fuzzy generalized regression neural network**

The fuzzy GRNN, developed in Chapter 4, is a modified GRNN that can accommodate fuzzy numbers. A formal definition of a fuzzy number will be given in Chapter 4, but a conceptual definition follows: Given a real number $x$, a fuzzy number consists of the real numbers close to or around $x$. From fuzzy set theory, the extension principle (Zadeh, 1975) is used to fuzzify the crisp mathematical functions that are used in the GRNN. Here, the Vertex Method (Dong and Shah, 1987) is used to calculate the functions as an approximation to the extension principle. This function is embedded into the GRNN algorithm to carry out the appropriate calculations using fuzzy numbers.

Following a literature search of a GRNN and fuzzy numbers, several studies have used fuzzy clustering as a method to preprocess a large amount of training data so as to simplify the GRNN (Lee *et al*., 2004, Lee *et al*., 2006, Husain *et al*., 2004, Goulermas, *et al*., 2007, Zhao *et al*., 2007).

In the field of image processing, identifying head pose is helpful in applications such as face recognition.  Bailly and Milgram (2009) use fuzzy functional criterion as a filter to select relevant features from images and couple it with a GRNN to assist in mapping between features and corresponding head pose.  Li and Fenli (2008) propose a digital image watermarking method based on fuzzy c-mean clustering and a GRNN.

Traffic models have been constructed that use the nonlinear mapping capabilities of fuzzy systems and then pass the data to a GRNN (Gharavol *et al*., 2007).  Kumara *et al*. (2003) used fuzzy logic to cluster noisy traffic data and then a GRNN was used to predict the hazardousness of a traffic intersection.

The combination of fuzzy mathematics and GRNNs has also been found in other engineering fields as well as control systems studies.  Ravanbod (2005) uses a GRNN to predict the dimensions of pipeline corrosions.  Fuzzy decision-based neural networks are then used for the detection and classification of the corrosions.  Singh *et al*. (2007) compare an adaptive neuro fuzzy inference system (ANFIS) to several ANN algorithms, including a GRNN, on the ability to predict thermal conductivity using various physico-mechanical properties such as porosity and density.  Seng *et al*. (1998) propose an adaptive neuro-fuzzy control system where they use a radial basis function neural network as a neuro-fuzzy controller and a GRNN as a predictor.

In stream flow prediction, Firat (2008) compared an ANFIS to several ANN methods (including a GRNN).  Turan and Yurdusev (2009) compare GRNN, FFBP, and fuzzy logic methodologies independently on the ability to predict streamflow from measured upstream flow records.  Various ANN models (including a GRNN) were compared to a neuro-fuzzy model on estimation abilities of suspended sediment in rivers (Kisi *et al*., 2008).  In weather forecasting, Tham *et al*. (2002) uses fuzzy c-means clustering on satellite images to help deduce cloud cluster velocities and then use a GRNN to predict cloud velocities over the area of interest.

The literature search, however, was unable to find a GRNN algorithm that incorporated fuzzy numbers.  To our knowledge, this is the first development and application of such an algorithm.

# CHAPTER 2
# DEMPSTER-SHAFER THEORY APPLIED TO
# UNCERTAINTY SURROUNDING PERMEABILITY

## 2.1 Abstract

Typically, if uncertainty in subsurface parameters is addressed, it is done so using probability theory. Probability theory is capable of only handling one of the two types of uncertainty (aleatory), hence epistemic uncertainty is neglected. Dempster-Shafer evidence theory (DST) is an approach that allows analysis of both epistemic and aleatory uncertainty. In this paper DST combination rules are used to combine measured field data on permeability, along with the expert opinions of hydrogeologists (subjective information) to examine uncertainty. Dempster's rule of combination is chosen as the combination rule of choice primarily due to the theoretical development that exists and the simplicity of the data. Since Dempster's rule does have some criticisms, two other combination rules (Yager's rule and the Hau-Kashyap method) were examined which attempt to correct the problems that can be encountered using Dempster's rule. With the particular data sets used here, there was not a clear superior combination rule. Dempster's rule appears to suffice when the conflict amongst the evidence is low.

## 2.2 Introduction

While uncertainty is an integral part of the mathematical representation of the environment, behavior forecasting requires the use of mathematical models that depend upon the specification of physically based parameters descriptive of the environment. In

subsurface hydrology, for example, the permeability must be specified in equations descriptive of groundwater flow. Typically, when the uncertainty surrounding such parameters is addressed, only aleatory uncertainty (irreducible uncertainty) is considered. However, there is another type of uncertainty, epistemic (lack of knowledge about the system), which should also be considered when using mathematical models to represent the environment. Currently, probability theory, usually within the framework of spatial interpolation (kriging), is used in an effort to generate a random field representation of a parameter (e.g. permeability). An effort to accommodate subjective information (e.g. expert opinion) into these analyses has been limited. For example, Ross et al. (2008) has developed a fuzzy Kalman filtering approach to incorporate expert knowledge into hydraulic conductivity field approximations.

Analysis of both aleatory and epistemic uncertainty surrounding permeability measurements (since classical probability is not sufficient to handle epistemic uncertainty, Sentz and Ferson, 2002) requires other avenues for assessing the uncertainty to be considered. Before making the decision on how to combine the evidence at hand, one must assess the evidence to determine what type it is: *consonant evidence, consistent evidence, arbitrary evidence*, or *disjoint evidence* (Sentz and Ferson, 2002). Consonant evidence can be described as a nested structure of subsets, so the smallest set is included in the next larger set, which is included in the next larger set, continuing until the largest set is reached. As an example, the following permeability (md) intervals from different sources A= [0.6, 0.8], B= [0.5, 0.9], C= [0.2, 1.2] form consonant evidence. With consistent evidence there is at least one element that is shared by all subsets as is the case in the following example: A= [0.2, 1.2], B= [0.1, 0.8], C= [0.6, 1.0]. Here the interval

[0.6, 0.8] is common to all sources. For arbitrary evidence there is no one element common to all subsets, however some subsets may share elements. As an example, in A= [0.2, 0.7], B= [0.5, 0.9], C= [1.1, 1.4], while the permeability interval [0.5, 0.7] is shared by sources A and B, source C has no permeability value common to sources A or B. This is the type of evidence encountered in this paper. Finally, disjoint evidence describes the situation where any two distinct subsets in the collection of sets have no element in common.

Once the evidence was defined, the use of evidence theory (Shafer, 1976) or Dempster-Shafer theory (DST), a branch of the theory of monotone measures (a generalization of classical measure theory) (Klir, 2003), was chosen to explore the uncertainty surrounding permeability. The main focus of this paper will be on the application of DST to combine subjective information (expert defined uncertainty bounds) with objective permeability data sets measured in the Dakota Sandstone. The Dempster-Shafer theory framework was selected due to its well developed theory, ability to combine evidence from many different sources, ability to handle the type of evidence in this study (arbitrary), numerous applications in the sciences (Agarwal *et al*., 2004; Binaghi *et al*., 1998; Cayuela *et al*., 2006; Kriegler and Held, 2005), lack of any assumptions about the distribution of the data, and ability to use all available data to analyze permeability uncertainty (outliers are kept in the analysis).

The following sections will provide a review of Dempster-Shafer theory and introduce the three combination methods used in this paper, describe the data sets acquired from three different methods to measure permeability, discuss the results obtained from combining expert opinions on the uncertainty surrounding each

17

measurement technique to obtain a more comprehensive representation of the uncertainty surrounding the measured data, and finally, compare the results of using modified versions of the Dempster's rule of combination.

## 2.3 Theory

### 2.3.1 Dempster-Shafer Theory (or Evidence Theory)

The Dempster-Shafer theory (DST) used today was originally introduced by Arthur Dempster (1967) then later expanded upon by Shafer (1976). The theory is based on belief measures, (*Bel*) and plausibility measures, (*Pl*). A common interpretation of belief and plausibility measures is as bounds of the unknown probability of the (permeability) set of interest. These two measures are equal in the case of pure probabilistic information (Klir, 2003). To further explain these measures, let $X$ be a universal set (frame of discernment) e.g., all the possible permeability values in the data set, and *P(X)* denote the set of all subsets of $X$; or all possible intervals of permeability. The degree of belief, *Bel*($A$), is defined for all $A$ in $P(X)$ and it quantifies the total amount of 'justified specific' support given to the claim that the unknown permeability value is in $A$. The term 'justified' means that $B$ supports $A$, thus $B$ is contained in $A$, and the term 'specific' means that $B$ does not support any permeability outside of $A$. Similarly, the degree of plausibility, *Pl*($A$), is defined for all $A$ in $P(X)$ and it quantifies the maximum amount of 'potential specific' support that could be given to the claim that the unknown permeability value is in $A$. The term 'potential' means that $B$ might come to support $A$ without supporting any permeability values outside of $A$ if a further piece of evidence is taken into consideration, thus the intersection of $A$ and $B$ is nonempty.

18

Belief and plausibility measures can be characterized by the basic mass (probability) assignment function:

$$m : P(X) \rightarrow [0,1] \qquad \text{where } m(\varnothing) = 0 \text{ and } \sum_{A \in P(X)} m(A) = 1, \tag{2.1}$$

using the following relations:

$$Bel(A) = \sum_{B|B \subseteq A} m(B)$$

$$Pl(A) = \sum_{B|A|\ B \neq \varnothing} m(B). \tag{2.2}$$

Mass assignments, $m(A)$, characterize the degree of evidence that the unknown permeability value of interest belongs exactly to the set $A$ and not to any of its subsets. For example, suppose there is evidence that permeability k belongs to the set containing values between 20 and 50 md. Say that the degree of membership of k to this set is 0.8 ($m([20,50]) = 0.80$). The evidence that is associated with this set says nothing about k belonging to a smaller subset of the interval [20, 50], i.e., the degree of membership of k to the set [30, 40] is not known. For every set $A$ contained in $P(X)$, such that $m(A)$ is greater than zero, is a focal element.

The original method derived to combine multiple sources of evidence, Dempster's rule of combination, has been criticized (Zadeh, 1986; Yager, 1987) for how it handles conflict among the evidence and, therefore, provides counterintuitive results when the level of conflict among the evidence is high (Zadeh, 1984). Several papers discuss the different combination rules that have developed over the years in response to this criticism (Sentz and Ferson, 2002; Smarandache, 2004; and Smets, 2005). Table 2.1 provides a summary of some of the major points brought forward in these papers. It should be noted that in the case where there is no or little conflict, Yager's rule, Inagaki's

rule, Zhang's rule and the Hau-Kashyap method produce very similar, if not identical, results as Dempster's rule. Since the evidence in this research is independent, not highly conflicting, and the sources of information are assumed to be very reliable, Dempster's rule of combination, along with two other similar combination rules, Yager's rule and the Hau-Kashyap method, have been chosen to analyze the data. The theory of these combination rules is discussed in the following subsections. Sentz and Ferson (2002) provide nice examples on the use of Dempster's rule of combination and Yager's rule of combination. For an example on the use of the Hau-Kashyap method, the reader is referred to the original paper (Hau and Kashyap, 1990).

### 2.3.3.1 Dempster's Rule of Combination

There exist numerous ways to combine evidence under Dempster-Shafer theory. The first technique derived, and the most widely used, is Dempster's rule of combination, which is used to combine evidence obtained from two or more independent sources.

$$m_{1,2}(J) = \frac{\sum\limits_{B \cap C = J} m_1(B) m_2(C)}{1 - T} \qquad \text{for all J} \neq \varnothing \text{, where}$$
$$T = \sum\limits_{B \cap C = \varnothing} m_1(B) m_2(C)$$

(2.3)

Here, $J$ is simply the resulting joint focal element formed from the nonempty intersections of the expert focal elements. The symbol $m_{1,2}(J)$ is referred to as a joint basic mass assignment and represents the degree to which the combined evidence supports the premise that the unknown permeability value belongs exactly to the set $J$. The variable $T$ represents the mass associated with conflict in the combined evidence. In other words, the denominator acts as a normalization factor since the mass assignments of

**Table 2.1: Comparison of various well-known combination rules.**

| Combination rule | Highlights | Weaknesses |
|---|---|---|
| **Dempster's rule of combination** (Dempster, 1967) | Most widely used rule; Easy to implement | Counter-intuitive results can occur when the conflict is high |
| **Yager's rule** (Yager, 1987) | Based on Dempster's rule; Removes normalization term; Assigns conflict to mass of universe in order to get more intuitive results when conflict is high | Total ignorance can grow rapidly implying a lack of knowledge even when there is knowledge about the case at hand |
| **Hau and Kashyap method** (Hau and Kashyap, 1990) | Based on Dempster's rule; Conflict assigned to union of conflicting sets | Creates "new" focal elements for each set of conflicting evidence, this can become a computational burden |
| **Inagaki's unified combination rule** (Inagaki, 1991) | Encompasses both Dempster's rule and Yager's rule; Incorporates a restriction that makes the rule only applicable to situations where nothing is known about the reliability of the sources | Normalization factor must be determined by the user, no well justified procedure has been developed to determine this value |
| **Zhang's center combination rule** (Zhang, 1994) | Allows for two frames of discernment; Considers the degree of intersection of sets | Degree of intersection can be defined in many ways, hence, so can the combination rule |
| **Dubois & Prade's disjunctive consensus rule** (Dubois and Prade, 1986, 1992) | Calculates mass assignments by taking the union of sets; No conflict is encountered and no information from the sources is rejected | Results can be very imprecise |
| **Discount & Combine method** (Shafer, 1976) | For use when evidence is highly conflicting; Can apply a discounting rate to belief functions; Uses an averaging function to combine information | Analyst would need to be qualified to determine how reliable the sources of information are |
| **Mixing or averaging** (Ferson and Kreinovich, 2002) | Uses an averaging function to combine information; Mass assignments are weighted | In cases of extreme conflict, analyst must consider appropriateness of an averaged result that was not originally suggested as a viable outcome by the sources |
| **Smets' TBM rule** (Smets and Kennes, 1994) | Unknown quantity is not restricted to be in the frame of discernment. | When high conflict exists, the mass of the empty set is large, loss of information |
| **Dezert-Smarandache classic rule** (Smarandache and Dezert, 2004) | Does not consider conflict (defined on free Dedekind's lattice) | If there exists a *Bel*=0, the result of the combination rule is automatically 0; Newer theory, has not been widely used |
| **Dezert-Smarandache hybrid rule** (Smarandache and Dezert, 2004) | Extension of Dubois & Prade's rule; Considers conflict in that the user forces elements to be empty based on model constraints | Difficult to compute; Newer theory has not been widely used |

the focal elements must sum to one. In this approach it is assumed that the unknown value is within the universal set. This is different from the approach used by Smets and Kennes (1994) (i.e., Transferable Belief Model-TBM) where one considers the possibility that the unknown value is not in the universal set. Typically, admitting a nonzero basic mass assignment for the empty set does this. In the case of combination, it is reflected by the lack of a normalization factor, whereas the normalization factor in Equation (4) ensures that the total mass is unity and $m_{1,2}(\varnothing) = 0$.

### 2.3.3.2 Yager's Rule

Yager (1987) proposes an alternative combination rule that has become known in the literature as Yager's rule:

$$
\begin{aligned}
m_{1,2}(J) &= \sum_{B \cap C = J} m_1(B) m_2(C) \qquad \text{for all } J \neq \varnothing \\
m_{1,2}(X) &= m_1(X) m_2(X) + \sum_{B \cap C = \phi} m_1(B) m_2(C).
\end{aligned}
\tag{2.4}
$$

The main differences between Yager's rule and Dempster's rule are the removal of the normalization term from the definition of the joint mass assignment for *J* and the assignment of the conflict to the mass of the universe *X*. Yager's thought is that since conflict represents the portion of the universe about which nothing is known, it makes more sense to distribute it among all the elements instead of only those focal elements about which there is information (Yager, 1987).

### 2.3.3.3 Hau-Kashyap (H-K) Method

Yager's rule does provide more intuitive joint mass assignments and belief values than Dempster's rule when applied to conflicting evidence (Yager, 1987). As the conflict

increases, however, the plausibility value of each focal element increases. This in turn yields a large Belief-Plausibility range, which can artificially imply a lack of knowledge in focal elements where, in fact, something is known about them (Hau and Kashyap, 1990). An alternative approach to Yager's rule is proposed by Hau and Kashyap (1990) where the mass associated with conflict is assigned to the union of the sets whose intersection is empty, instead of to the entire set of the universe,

$$
\begin{aligned}
m_{1,2}(J) &= \sum_{B \cap C = J} m_1(B) m_2(C) \\
m_{1,2}(B \cup C) &= m_1(B) m_2(C) \quad \text{if } B \cap C = \varnothing.
\end{aligned}
\tag{2.5}
$$

Here, the term $m_{1,2}(B \cup C)$ represents the conflict associated with the particular sets $B$ and $C$. Hau and Kashyap (1990) argue that instead of "eliminating" or "erasing" the conflict as is done using Dempster's or Yager's rule, they seek compromise among the conflicted and choose to resolve the conflicts until after more information becomes available.

## 2.4 Data Sets

The data sets that are analyzed in this paper are permeability (k) values measured in the Dakota Sandstone within the Denver Basin (Belitz and Bredehoeft, 1988). There are three data sets that are considered independent of each other and each set was determined via a different technique; water-well pumping test, core analysis, and drill-stem analysis. Though a previous statistical study of this data (Ricciardi, 2002) found cause to remove several outliers in each data set, this analysis included all data points; no outliers were removed from any of the sets. All values have units of millidarcies (md).

Water-well pump-test data were compiled from state water reports in the regions of South Dakota, southwestern Kansas, and southeastern Colorado.  The sandstone here is a source of water and the measurements are taken at relatively shallow depths, less than 3000 feet.  There are 74 points in this set.  In the second set, there are 161 core permeability data values that were compiled from state petroleum reports and other literature pertaining to regions of northeastern Colorado, southeastern Wyoming, and the Nebraska panhandle.  Here the sandstone is primarily used as a source of oil, so the measurements are taken at depths from approximately 3,200 feet to 8,400 feet.  The final data set consists of drill-stem data that were interpreted by Belitz and Bredehoeft (1988) using data from the USGS Petroleum Library in Denver.  The data were obtained from the regions of northeastern Colorado, southeastern Wyoming, and the Nebraska panhandle.  This was the largest data set at 453 data points.

The methodology described in this paper is applicable for vertically averaged sections.  Each of the three data sets analyzed here provided only the depth measurements along with corresponding permeability values, therefore, spatial attributes could not be considered in these data sets.  In order to determine whether there is a depth dependency for the permeability values, plots of depth versus permeability on a log scale were created (Figure 2.1 (A)-(C)) and correlation coefficients were calculated.  The correlation coefficients for the water-well pump-test, core, and drill-stem data are 0.004, -0.42, and -0.15, respectively.  The small values for the water-well pump-test data and the drill-stem data suggest no linear relationship between depth and permeability in these data sets.  The core data, however, does exhibit a negative linear trend.  A trend such as this could increase the range between belief and plausibility, inflating the representation of

24

**Figure 2.1: Depth versus permeability plots for (A), pump-test, (B), core, and (C), drill-stem data.**

uncertainty found in the measurements. One possible approach to correcting for this would be to detrend the data by subtracting the least-squares fit. Investigation of the most appropriate way to handle data trends using the methodology presented here is a topic for further exploration and is not considered in this paper.

## 2.5 Results

### 2.5.1 Random Intervals to Probability Boxes

The field-measured permeabilities needed to be converted into structures that could be used in the Dempster-Shafer theory framework. Random sets are noted as being mathematically isomorphic to Dempster-Shafer bodies of evidence (Joslyn and Booker, 2004). A random set can be thought of as a random variable that has sets as its values rather than points. A finite random set, $P$, can be defined as (Joslyn and Ferson, 2004)

$$P \equiv \{(A_j, m(A_j)) : m(A_j) > 0\} \tag{2.6}$$

where $A_j \subseteq X$ and $1 \leq j \leq N$. A finite random interval, denoted $Q$, follows as a finite random set on $X = \Re$ for which the focal elements can be denoted as intervals $I_j$ such that $F(Q) = \{I_j\}$, $1 \leq j \leq N$. The finite random interval is a random left-closed interval of the reals, [a,b). In Joslyn and Booker (2004) it is noted that random intervals are important to engineering reliability studies due to their ability to incorporate randomness and imprecision or nonspecificity in one mathematical structure.

Though the domain that is considered here is the entire real line, the data can be represented as finite random intervals. These are in turn examples of Dempster-Shafer structures (Joslyn and Booker, 2004), which can prove to be difficult to represent, manipulate, and interpret. Typically, therefore, these structures are approximated by

26

simpler mathematical structures; one example is probability boxes (p-box) from which one can obtain an equivalence class of random intervals that are consistent with the p-box (Joslyn and Ferson, 2004). In the case of a piecewise constant p-box, one can construct a random interval in a canonical way as it is done in this paper. From the p-box, one can discretize it into rectangles, and then the width of a specific rectangle defines a focal element. Their corresponding basic mass assignments are the step sizes on the ordinate or the height of a rectangle (Ferson *et al*., 2002).

Since the data are in the form of finite random intervals that have a finite number of focals, their representation is not computationally restrictive; hence all the focal elements can be used. Therefore, the p-box is less of an approximation but more of an exact representation of a Dempster-Shafer structure. In fact, the p-boxes here are equivalent to the cumulative belief and plausibility distributions created using the intervals in each data set as the focal elements.

In order to construct the p-boxes used in this paper, two experts in the field of hydrogeology and familiar with the Denver Basin were asked to provide a range of uncertainty for each of the three methods. Neither expert had knowledge of the others responses. The values are given in Table 2.2. The uncertainty values were then used to create two p-boxes for each data set, one for each expert (Figs. 2.2-2.4). The resulting

**Table 2.2: Expert assigned uncertainty to the three different methods for measuring hydraulic conductivity.**

|  | Water-well Pump-test | Core Analysis | Drill Stem Analysis |
|---|---|---|---|
| **Expert 1** | +/- 1 order of magnitude | +/- 2 orders of magnitude | +/- 0.75 orders of magnitude |
| **Expert 2** | +/- 0.5 orders of magnitude | +/- 1 order of magnitude | +/- 0.5 orders of magnitude |

**Figure 2.2: Probability box constructed from water-well pump-test data and a measurement uncertainty of +/-1 order of magnitude assigned by Expert 1 and +/- 0.5 orders of magnitude assigned by Expert 2.**



**Figure 2.3: Probability box constructed from core data with a measurement uncertainty of +/-2 orders of magnitude assigned by Expert 1 and +/-1 order of magnitude assigned by Expert 2.**

focal elements are consistent with the definition of arbitrary evidence. The lognormal

cumulative distribution, the distribution typically used to analyze permeability values of

the respective data set is also plotted to observe how well it is contained in the p-box. It

should be noted here that the resulting p-boxes in Figures 2.2, 2.3, and 2.4 only display

part of the plot in order to show detail.

## 2.5.2 Combination Rules

Once all the focal elements and corresponding mass assignments were determined,

the calculations necessary to combine the information were performed. Analysis using

Dempster's rule yielded conflict values for the pump-test, core, and drill-stem data of

$T=6.57 \times 10^{-2}$, $3.09 \times 10^{-4}$, and $4.00 \times 10^{-1}$, respectively. Even though the conflict values



**Figure 2.4: Probability box constructed from drill-stem data with a measurement uncertainty of +/-0.75 orders of magnitude assigned by Expert 1 and +/- 0.5 orders of magnitude assigned by Expert 2.**

for the pump-test and core data are low, the other combination methods were explored to see if there was a marked difference in this type of application. In order to compare the results of the three combination rules within the different measurement techniques, plots of cumulative belief and plausibility were examined. Ultimately, a decrease in space between the lower (belief) and upper (plausibility) bounds upon combination of the information would suggest a decrease in the uncertainty for permeability.

### 2.5.2.1 Pump-test Data

The results of the application of Dempster's rule of combination to the pump-test data can be seen in Figure 2.5. This combination yielded 662 joint focal elements. Upon comparison to Figure 2.2, the distance between the bounds (or uncertainty) for the permeability is obviously decreased. This is due to the overall decrease in the size of the permeability intervals that form the focal elements of the random interval representing evidence on permeability upon combination of information from the two experts.

Next, the results of combination via Yager's rule, yields 662 joint focal elements plus the set of the universe, $X$, to which is assigned an additional mass equal to the conflict between the experts, $T = 6.57 \times 10^{-2}$. The uncertainty range is similar to that obtained using Dempster's rule (Fig. 2.5), however, due to the addition of the mass assignmentassociated with the universal set, plausibility values are inflated resulting in a wider gap between the cumulative belief and plausibility plots, i.e., greater uncertainty at higher permeability values.

Finally, the Hau-Kashyap (H-K) method produces a total of 750 joint focal elements (this includes the joint focals that are created by taking the union of the sets that conflict,

**Figure 2.5: Plot of cumulative belief and plausibility for water-well pump-test data. These values are calculated using the joint focal elements obtained through Dempster's rule of combination, Yager's rule of combination, and the Hau-Kashyap method.**

i.e., the intersection is empty). The H-K method, like Dempster's and Yager's rule, appears to reduce the uncertainty upon combination (Fig. 2.5). However, unlike the combination from Dempster's or Yager's rule, it considers permeability values greater than 46,802 md. Also, unlike Yager's rule, the H-K method does not appear to inflate the plausibility values, as convergence to one is achieved by both the cumulative belief and plausibility. In all cumulative belief and plausibility plots for water-well pump-test data, the lognormal curve fits within the bounds of the "box" created by the uncertainty.

## 2.5.2.2 Core Data

The second data set to be analyzed, core data, produced 2,115 joint focal elements when combined using Dempster's rule. Yager's rule produces 2,115 joint focal elements

plus the universal set, which is assigned a mass of 3.09 x $10^{-4}$. The H-K method produces 2,123 joint focal elements. The resulting cumulative belief and plausibility plots for these combination methods appear identical to each other except for the inclusion of the larger permeability values (up to 12,990 md ) when the H-K method is used (Fig. 2.6). In this case, due to the extremely low conflict among the evidence, there does not appear to be any significant differences between the combination methods. Compared to the experts' probability boxes (Fig. 2.3), the bounds decrease upon combination (more closely resembling the belief of Expert 2). This is a result of how the joint focals are created and analyzed. Again the lognormal curve fit's within the box for all cases for the core data.



**Figure 2.6: Plot of cumulative belief and plausibility for core data. These values are calculated using the joint focal elements obtained through Dempster's rule of combination, Yager's rule of combination, and the Hau-Kashyap method.**

32

## 2.5.2.3 Drill-stem Data

The final data set to be analyzed was the drill-stem data. Recall that this data set had the most conflict of the three data sets, $T=4.00 \times 10^{-1}$. Dempster's rule applied to the evidence provided by the experts resulted in 55,473 joint focal elements. Note that the lognormal curve clearly fits into the experts' uncertainty opinion (Fig. 2.4). However, in looking at the results of Dempster's rule of combination (Fig. 2.7), the lognormal curve violates the bounds that are established with this method, suggesting either that the lognormal distribution may not be the best distribution in this case or that Dempster's rule may not be the best combination rule to choose with this level of conflict.



**Figure 2.7: Plot of cumulative belief and plausibility for drill-stem data. These values are calculated using the joint focal elements obtained through Dempster's rule of combination, Yager's rule of combination and the Hau-Kashyap method.**

33

Yager's rule also yields 55,473 joint focal elements plus the universal set, which as before, is assigned a mass equal to the conflict present amongst the data, 4.00 x $10^{-1}$. Since this data set has a relatively high conflict, it becomes more apparent how Yager's rule can inflate the plausibility (Fig. 2.7) when compared to the other two data sets (Figs. 2.5 and 2.6).

The H-K method produces 102,877 joint focal elements. Examining the results of the H-K method (Fig. 2.7) the cumulative belief and plausibility plots provide more uncertainty than Dempster's rule, yet less than Yager's rule. In neither Yager's rule nor the H-K method does the lognormal curve violate the bounds.

## 2.6 Conclusions

In this paper the use of Dempster-Shafer theory is examined as an alternative way to assess the uncertainty surrounding permeability measurements. The benefits of DST include not having to choose a distribution that may or may not be a best fit for the data and all available data can be used. Here it is shown that field measured permeability data can be joined with expert subjective data and then the different sources of evidence can be combined. Being able to incorporate multiple sources of evidence would, theoretically, provide a better representation of the uncertainty surrounding permeability.

The second matter considered here is the comparison of combination processes, i.e., Dempster's rule of combination, and its two modified versions, Yager's rule and the Hau-Kashyap method. Yager's rule appears to err on the side of caution by applying the conflict to the mass of the universe. This results in inflated plausibility values which, in particular for cases of higher conflict, results in wide uncertainty ranges (Figs. 2.5 and

2.7).  Proceeding with too much caution can actually lead to a lack of knowledge across the entire universe.  It can overshadow the areas where much is known, resulting in the loss of important information.  Based on the study here, it appears that if there is little conflict amongst the data (as in the pump-test and core data) and the data sources are reliable, Dempster's rule is sufficient.  If the level of conflict is questionably high (the drill-stem data case), then it may be safer to choose an alternative combination method such as that proposed by Hau-Kashyap.

## 2.7 Acknowledgements

## 2.8 References

Agarwal, H., Renaud, J.E., Preston, E.L., and Padmanabhan, D., 2004, Uncertainty quantification using evidence theory in multidisciplinary design optimization: Reliability Engineering and System Safety, v. 85, p. 281-294.

Belitz, K. and Bredehoeft, J. D., 1988, Hydrodynamics of Denver Basin - Explanation of subnormal fluid pressures:  American Association of Petroleum Geologists Bulletin, v. 72, no. 11, p. 1334-1359.

Binaghi, E., Luzi, L., Madella, P., Pergalani, F., and Rampini, A., 1998, Slope instability Zonation: a comparison between certainty factor and fuzzy Dempster-Shafer approaches: Natural Hazards, v. 17, p. 77-97.

Cayuela, L., Golicher, J.D., Salas Rey, J., and Rey Benayas, J.M., 2006, Classification of a complex landscape using Dempster-Shafer theory of evidence: International Journal of Remote Sensing, v. 27, no. 10, p. 1951-1971.

Dempster, A. P., 1967, Upper and lower probabilities induced by a multivalued mapping: Annals of Mathematical Statistics, v. 38, p. 325-339.

Dubois, D. and Prade, H., 1986, A set-theoretic view on belief functions: logicaloperations and approximations by fuzzy sets: International Journal of General Systems, v. 12, p. 193-226.

Dubois, D. and Prade, H., 1992, On the combination of evidence in various mathematical frameworks, *in* Flamm , J. and Luisi, T., eds., Reliability Data Collection and Analysis, Kluwer Academic Publishers, Brussels, p. 213-241.

Ferson, S. and Kreinovich, V., 2002, Representation, Propagation, and Aggregation of Uncertainty: Sandia National Laboratories, Technical Report, Albuquerque, New Mexico.

Ferson, S., Kreinovich, V., Ginzburg, L., Myers, D.S., and Sentz, K., 2002, Constructing Probability Boxes and Dempster-Shafer Structures: Sandia National Laboratories, Technical Report SAND2002-4015, Albuquerque, New Mexico, Available at: http://www.sandia.gov/epistemic/Reports/SAND2002-4015.pdf.

Hau, H.Y. and Kashyap, R.L., 1990, Belief combination and propagation in a lattice-structured inference network: IEEE Trans. on Systems, Man, and Cybernetics, v. 20, no. 1, p. 45-57.

Inagaki, T., 1991, Interdependence between Safety-Control Policy and Multiple-Sensor Schemes Via Dempster-Shafer Theory: IEEE Transactions on Reliability, v. 40, no. 2, p. 182-188.

Joslyn, C. and Booker, J.M., 2004, Generalized Information Theory for Engineering Modeling and simulation, *in* Nikolaidid, E., Ghiocel, D., and Singhal, S., eds., Engineering Design Reliability Handbook, CRC Press, p. **9**-1 - **9**-40.

Joslyn, C. and Ferson, S., 2004, Approximate representations of random intervals for hybrid uncertainty quantification in engineering modeling, *in* Hanson, K.M., and Hemez, F.M., eds., Sensitivity Analysis of Model Output (SAMO04), LANL, Los Alamos, p. 453-469 http://library.lanl.gov/cgi-bin/getdoc?event=SAMO2004&document=samo04-83.pdf.

Klir, G.J., 2003, Uncertainty: Encyclopedia of Information Systems, v. 4, p. 511-521.

Kriegler. E., and Held, H., 2005, Utilizing belief functions for the estimation of future climate change: International Journal of Approximate Reasoning, v. 39, p. 185-209.

Ricciardi, K.L., 2002, Optimal groundwater remediation design subject to uncertainty: Ph.D dissertation, University of Vermont, USA, p. 50-66.

Ross, J., Ozbek, M., and Pinder, G.F., 2008, Kalman filter updating of possibilistic hydraulic conductivity:  Journal of Hydrology, v. 354, p. 149-159.

Sentz, K. and Ferson, S., 2002, Combination of Evidence in Dempster-Shafer Theory: Sandia National Laboratories, Technical Report SAND2002-0835, Albuquerque, New Mexico, Available at: http://www.sandia.gov/epistemic/Reports/SAND2002-0835.pdf.

Shafer, G., 1976, A mathematical theory of evidence: Princeton University Press, Princeton, New Jersey, 312 p.

Smarandache, F., 2004, An in-depth look at information fusion rules and the unification of fusion  theories: arXiv electronic archives, available at: http://xxx.lanl.gov/ftp/cs/papers/0410/0410033.pdf.

Smarandache, F. and Dezert, J., eds., 2004, Applications and Advances of DSmT for Information Fusion: American Research Press, Rehoboth, NM: http://www.gallup.unm.edu/~smarandache/DSmT-book1.pdf.

Smets, P., 2005, Analyzing the combination of conflicting belief functions: available at: http://iridia.ulb.ac.be/%7epsmets/Combi_Confl.pdf.

Smets, P. and Kennes, R., 1994, The transferable belief model: Artificial Intelligence, v. 66, p. 191-234.

Yager, R.R., 1987, On the Dempster-Shafer framework and new combination rules: Information Sciences, v. 41, p. 93-138.

Zadeh, L. A., 1984, Review of Books: A Mathematical Theory of Evidence: The AI Magazine, v. 5, no. 3, p. 81-83.

Zadeh, L. A., 1986, A simple view of the Dempster-Shafer theory of evidence and its implication for the rule of combination: The AI Magazine. v. 7, p. 85-90.

Zhang, L., 1994, Representation, independence, and combination of evidence in the Dempster-Shafer theory, *in* Yager, R.R., Kacprzyk, J., and Fedrizzi, M., eds., Advances  in the Dempster-Shafer Theory of Evidence: John Wiley & Sons, Inc., New York, p. 51-69.

# CHAPTER 3
# TRANSMISSIVITY AND STORAGE COEFFICIENT ESTIMATION BY COUPLING THE COOPER-JACOB METHOD AND MODIFIED FUZZY LEAST-SQUARES REGRESSION

## 3.1 Abstract

Traditionally the Cooper-Jacob equation is used to determine the transmissivity and the storage coefficient for an aquifer using pump test results. This model, however, is a simplified version of the actual subsurface and does not allow for analysis of the uncertainty that comes from a lack of knowledge about the heterogeneity of the environment under investigation. In this paper, a modified fuzzy least-squares regression (MFLSR) method is developed that uses imprecise pump test data to obtain fuzzy intercept and slope values, which are then used in the Cooper-Jacob method. Fuzzy membership functions for the transmissivity and the storage coefficient are then calculated using the extension principle. The supports of the fuzzy membership functions incorporate the transmissivity and storage coefficient values that would be obtained using ordinary least-squares regression and the Cooper-Jacob method. The MFLSR coupled with the Cooper-Jacob method allows the analyst to ascertain the uncertainty that is inherent in the estimated parameters obtained using the simplified Cooper-Jacob method and data that are uncertain due to lack of knowledge regarding the heterogeneity of the aquifer.

## 3.2 Introduction

For decades, water well pump tests have been used to predict the characteristics of the subsurface. While, in the conduct of a pumping test, the water-level measurements and their location are relatively crisp with small measurement error, the nature of the porous medium with which one identifies these measurements is uncertain. Due to heterogeneity there will be variability in the material properties in the neighborhood of the observation well. Adding more observation wells would provide a more detailed picture of the subsurface, however this can be costly and impractical. Hence, in the absence of additional observation wells, the question is `to what degree do the changes in the water levels in the observation wells measured during a pumping test reflect the heterogeneous nature of the properties in the neighborhood of the well?'

Cooper and Jacob (1946) proposed the 'straight-line' method, built on the theory introduced by Theis (1940), for obtaining the transmissivity and the storage coefficient through a simplified analysis of pump test results. The measured water level values represent the solution to an equation that includes unknown parameters that reflect the heterogeneities in the volume of the geologic formation that is identified with the pumping test; that is, the region that is impacted by the pumping test in the specified test period. Denote this solution, or observation, as $h_{obs}$. The Theis (1940) solution, $h_{Theis}$, on the other hand generates a water level time profile at a specific point in response to a specified pumping rate that assumes a homogeneous aquifer. If the 'straight-line' method is used to determine transmissivity and storage coefficient values, the parameters extracted are not those of the heterogeneous aquifer, but a surrogate homogeneous formation. Using crisp water-level measurements, the 'straight-line' method will provide

crisp values of the transmissivity and the storage coefficient.  The parameter identification process is silent on the matter of the uncertainty with which these parameters represent the heterogeneity in the media in the neighborhood of the observation wells.

If the values of transmissivity and the storage coefficient identified via the 'straight-line' method are substituted into the physically correct, but unknown, equation for the aquifer (the equation that generated $h_{obs}$) a new water level value, would be generated. The difference between the values $h_{obs}$ and $h_{Theis}$ could be considered as the model error; that is the error committed when using the Cooper-Jacob equation rather than the physically correct equation to represent the actual heterogeneous aquifer.  This error will be denoted by $\varepsilon$.

Uncertainty due to a lack of knowledge, such as $\varepsilon$ noted above (rather than randomness), is called epistemic uncertainty.  Unlike aleatory uncertainty that is associated with irreducible uncertainty and amenable to analysis using probability theory, epistemic uncertainty is not easily analyzed using probability theory and is more appropriately analyzed using other mathematical tools.  Fuzzy sets constitute such a tool. In this paper it is shown how to incorporate epistemic uncertainty in the 'straight-line' method of pump test analysis to examine the impact of model uncertainty on transmissivity and the storage coefficient.

Traditionally, the 'straight-line' method employs ordinary linear regression in an attempt to fit a slope and intercept to water levels measured over time at specified well locations or over a series of wells at a specified time.  Ordinary linear regression can be used to analyze aleatory uncertainty due to observation errors.  However, such errors are

generally small relative to model errors and do not reflect the inherent uncertainty in the estimated coefficients attributable to heterogeneity.

Fuzzy linear regression, introduced by Tanaka *et al*. (1982), is an approach that will allow for the accommodation of epistemic uncertainty attributable to lack of knowledge. Recent application of fuzzy regression in hydrology can be seen in the work of Bardossy *et al*. (1990), Ozelkan and Duckstein (2001), Uddameri (2004), Si and Bodhinayake (2005), and Uddameri and Honnungar (2007). Many fuzzy linear regression methods exist and Chang and Ayyub (2001) provide a nice review of some of these.

Due to the limitations, as summarized by Ozelkan and Duckstein (2001), of the fuzzy regression (FR) method originally proposed by Tanaka *et al*. (1982), and due to the ease of implementation of the fuzzy least-squares regression (FLSR) method as proposed by Savic and Pedrycz (1991) the latter was chosen for use in this analysis. With a slight modification to the technique, the proposed modified fuzzy least-squares regression (MFLSR) method improved the results, which were found to be similar to those obtained using the hybrid fuzzy-least squares regression outlined by Chang (2001).

In the following sections, the theory of FLSR will be introduced and the reasons behind the modification will be discussed. The hybrid method will also be discussed briefly. The paper will conclude with the results and a discussion of the analysis that was conducted by using the MFLSR method in conjunction with the Cooper-Jacob method to determine transmissivity and the storage coefficient of an aquifer.

## 3.3 Fuzzy least-squares regression

A typical bivariate regression model could be represented by:

$$\tilde{Y} = \tilde{A}_0 + \tilde{A}_1 X \qquad (3.1)$$

where $\tilde{A}_0$ and $\tilde{A}_1$ are the fuzzy intercept and fuzzy slope coefficients, respectively, and are assumed to have symmetrical triangular membership functions (Figs. 3.1 and 3.2).

Data identified with $X$ (the independent variable) is crisp and the output $\tilde{Y}$ (or dependent variable) is either crisp or a fuzzy number. The fuzzy coefficients can be represented for the case of a symmetrical triangular basis function using a center point $m_j$ and a spread (or halfwidth) $c_j$, i.e. $\tilde{A}_j = (m_j, c_j)$. The fuzzy coefficients can be determined by



**Figure 3.1: Fuzzy intercept coefficient ($A_0$ term) membership function for 10% and 50% (epistemic) uncertainty cases.**

42

**Figure 3.2: Fuzzy slope coefficient ($A_1$ term) membership function for 10% and 50% (epistemic) uncertainty cases.**

solution of the optimization problem defined by the following objective function and constraints derived by Tanaka et al. (1982):

$$\text{Minimize} \quad \sum_{j=0}^{m} \sum_{i=1}^{n} c_j \, |x_{ij}| \tag{3.2}$$

which, for the bivariate case, simplifies to

$$\text{Minimize} \quad nc_0 + c_1 \sum_{i=1}^{n} |x_i| \tag{3.3}$$

subject to the following constraints:

$$\sum_{j=0}^{1} m_j x_{ij} + (1-b) \sum_{j=0}^{1} c_j \, |x_{ij}| \; \ge y_i + (1-b)e_i \quad \text{for i} = 1 \text{ to n,} \tag{3.4}$$

43

$$\sum_{j=0}^{1} m_j x_{ij} - (1-b)\sum_{j=0}^{1} c_j \,|\, x_{ij} \,| \;\; \le y_i - (1-b)e_i \quad \text{for } i = 1 \text{ to } n,$$

$$c_0 > 0; \; c_1 \ge 0$$

$(3.5)$

where n is the number of data points, $x_{ij}$ is the independent variable (in this case $x_{i0} = 1$ and $x_{i1}$ is the input variable from the given data set), $y_i$ is the center of the fuzzy dependent (output) variable, $e_i$ is the spread of the fuzzy dependent variable, and $b \in [0,1]$ is a degree of compatibility which can be viewed as a measure of fit between the regression model and the actual data. This measure, $b$, imposes a threshold on the model to express the fact that the fuzzy model result should contain all the (crisp) observed data $y_i$ to a certain degree, and it is of the form (Savic and Pedrycz, 1991):

$$\mu_Y(y_i) \ge b \quad \text{for } i = 1,2,...,n$$

$(3.6)$

where $\mu_Y$ is the membership function for Y. The choice of $b$ influences the widths $c_j$. In particular, Chang and Ayyub (2001) have shown that as $b$ approaches 1 the fuzziness of the model increases. In several cases (Tanaka *et al*., 1982; Bardossy *et al*., 1990; Savic and Pedrycz, 1991; Uddameri, 2004; Si and Bodhinayake, 2005; Uddameri and Honnungar, 2007) $b$ values of 0.5 to 0.75 have been used.

The FLSR model is a two-step process (Savic and Pedrycz, 1991). First the fuzzy coefficient centers $m_j$ are determined from ordinary least-squares regression, i.e., the $e_i$ are considered to equal zero. Once these center values are obtained the values are substituted in Eqs. (3.4) and (3.5) above and the optimization problem (Eqs. (3.2) − (3.5)) is solved in order to obtain the halfwidth values, $c_j$.

One of the limitations of FR and FLSR is that as $b$ and the $e_i$ tend to zero the results of the regression do not converge to those of ordinary regression as would be expected.

In reviewing the above method, upon examination of the case where $b$ and $e_i$ are set equal to zero, the constraints in Eqs. (3.4) and (3.5) reduce to:

$$\sum_{j=0}^{1} m_j x_{ij} + \sum_{j=0}^{1} c_j \, | x_{ij} | \geq y_i \quad \text{for } i=1 \text{ to n,}$$

(3.7)

$$\sum_{j=0}^{1} m_j x_{ij} - \sum_{j=0}^{1} c_j \, | x_{ij} | \leq y_i \quad \text{for } i=1 \text{ to n.}$$

(3.8)

Now by letting $y_i = y_{i,observed}$, recalling that $x_{i0} = 1$, and expanding Eqs. (3.7) and (3.8) yields

$$m_0 + m_1 x_{i1} + c_0 + c_1 \, | x_{i1} | \geq y_{i,observed} \quad \text{for } i=1 \text{ to n,}$$

(3.9)

$$m_0 + m_1 x_{i1} - c_0 - c_1 \, | x_{i1} | \leq y_{i,observed} \quad \text{for } i=1 \text{ to n.}$$

(3.10)

Recalling that the $m_j$ are obtained using least-squares regression, it can be written that

$$m_0 + m_1 x_{i1} = y_{i,calculated}.$$

(3.11)

Substituting Eq. (3.11) into Eqs. (3.9) and (3.10) yields

$$c_0 + c_1 \, | x_{i1} | \geq y_{i,observed} - y_{i,calculated} \quad \text{for } i=1 \text{ to n,}$$

(3.12)

$$-c_0 - c_1 \, | x_{i1} | \leq y_{i,observed} - y_{i,calculated} \quad \text{for } i=1 \text{ to n.}$$

(3.13)

The right hand side of the constraints in Eqs. (3.12) and (3.13) introduces a difference between the observed and calculated $y_i$ based on the least-squares regression to determine the $m_j$. Because of this difference, the use of $y_{i,observed}$ in Eqs. (3.4) and (3.5) introduces an artificial fuzziness into the model. This fuzziness is manifested in $c_j$ that can take non-zero values even for crisp observed data, which is not desirable when assessing the effect of non-crisp observed data on calculated model results. Therefore, in the MFLSR method, $y_{i,calculated}$ is used, and in doing so the model also converges to crisp results as desired when observed data are crisp.

45

In the cases where there exist non-crisp observed data (i.e., $e_i > 0$), it can be easily shown that the use of $y_{i,calculated}$ cancels the effect of the measure of compatibility $b$ in Eq. (3.6) by setting its effective value to 0:

$$\mu_Y(y_i) \geq 0 \quad \text{for } i = 1, 2, ..., n.$$

(3.14)

From a fuzzy set theoretic point of view, the comparison of Eqs. (3.6) and (3.14) requires a distinction between uncertain and imprecise model results. The membership function in Eq. (3.6) refers to an imprecise (or fuzzy) value that is certain to a degree of $1 - b$ (i.e., imprecise and uncertain), whereas the membership function in Eq. (3.14) refers to an imprecise value that is certain to a degree of 1 (i.e., imprecise and certain) (Dubois *et al.*, 1988). It is an objective in this paper to quantify and propagate the imprecision in the observed data, therefore the resulting effective value of $b = 0$ due to the use of $y_{i,calculated}$ is consistent with the application. The above considerations modify the FLSR method to the MFLSR approach in the following form:

$$\text{Minimize} \quad nc_0 + c_1 \sum_{i=1}^{n} |x_i|$$

(3.15)

subject to the following constraints:

$$-\sum_{j=0}^{1} c_j |x_{ij}| \leq -e_i \quad \text{for } i = 1 \text{ to } n,$$

$$c_0 \geq 0; c_1 \geq 0.$$

(3.16)

The hybrid method (Chang, 2001), which uses weighted fuzzy arithmetic and the least-squares fitting-criterion, was an alternative approach that was considered as it addressed the issue of convergence upon crisp results given crisp data. The method was used as a comparison for the results obtained during the analysis with MFLSR. For details on the method the reader is referred to the original paper, Chang (2001).

46

## 3.4 Cooper-Jacob Equation

A typical representation of the Cooper-Jacob method (Cooper and Jacob, 1946) is given as follows:

$$s \approx \frac{Q}{4\pi T}(-0.5772 - \ln(u))$$  (3.17)

for sufficiently small $u$, where $u = \frac{r^2 S}{4Tt}$, $r$ is the distance from the pumping well to the observation well (L), $S$ is the storage coefficient (dimensionless), $T$ is the transmissivity (L$^2$/time), $t$ is time, $Q$ is the pumping rate (L$^3$/time), and $s$ is the drawdown (L). Expanding Eq. (3.17) and substituting for $u$ provides the following equation

$$s \approx -\frac{Q}{4\pi T}\left[0.5775 + \ln\frac{r^2 S}{4T}\right] - \frac{Q}{4\pi T}\ln(1/t)$$  (3.18)

which is the equation of a straight line where $s$ can be viewed as a function of $1/t$.

Therefore, the slope is $-\frac{Q}{4\pi T}$ and $T$ can be solved for directly. Typically, this method solves for storativity, $S$, by extrapolating the line to where it intercepts the time axis (where $s = 0$). This is denoted as $t_0$. Through some manipulation, $S$ can be obtained from the following relation

$$S = \frac{2.25Tt_0}{r^2}.$$  (3.19)

This method was not used to solve for the storage coefficient in this paper since it is not clear how to define a fuzzy zero drawdown value in order to extrapolate the regression to determine $t_0$. A more direct approach was chosen instead and the details are explained in section 3.5.3.

## 3.5 Calculations

To test the impact of imprecision in the dependent variable (water levels), attributable to model error, on the uncertainty of the computed transmissivity and storage coefficient, a data set created in an intermediate scale groundwater facility was used in conjunction with fuzzy least-squares regression. The data set consisted of change in pressure values in observation wells in response to a pumping test. The pressure changes were measured using a pressure transducer connected to a continuous recording device. The values so measured are assumed to have negligible measurement error. The data set contains 220 points, where the pumping rate was 3.75 cm$^3$/s with a distance of 133.78 cm between the observation and pumping well, Table 3.1. While in this particular case there is a large number of data available for analysis, the same methodology can be applied to data sets with fewer points. To explore the sensitivity of this, a second analysis was conducted where the number of data points was reduced by a factor of two. The results of the second analysis proved to be very similar to the results presented in this section for the full data set.

As noted earlier, while the values observed are crisp their interpretation in terms of the mathematical model underlying the 'straight-line' method of analysis contains the model error $\varepsilon$. In other words, it should be expected that the water level values when

**Table 3.1: Pumping well data**

| Data points in the set | Discharge rate, Q (cm$^3$/s) | Distance between observation and pumping well, r (cm) |
|:---:|:---:|:---:|
| 220 | 3.75 | 133.78 |

used in the 'straight-line' method are imprecise and that imprecision is reflected in imprecision in the resulting parameters. The degree of uncertainty is a function of the inconsistency between the simplified mathematical model used and physical model producing the measurement values. The higher the degree of heterogeneity and flow complexity, the less confident the analyst is that the head values observed represent the set of values that, via the 'straight-line' method, would produce an accurate volume averaged heterogeneous hydraulic conductivity. The analyst is faced with determining through observation of the properties that constitute the reservoir and his/her professional experience, the level of confidence they have that the head values are consistent with the set of values which, if placed in the 'straight-line' model would provide the best volume averaged heterogeneous hydraulic conductivity.

The data set considered is, by design of the intermediate scale facility, relatively homogeneous. Thus, the imprecision in the head values, in the context of the above, is relatively small. This small imprecision is reflected in the form of the membership functions that exhibit 10% (epistemic) uncertainty. In a field situation, it would be anticipated that a greater degree of imprecision would be assigned to the data and in an effort to illustrate this situation an extreme value of 50% is also assumed. The following results show how the imprecision in the water-level values impacts the uncertainty in the estimated parameters.

### 3.5.1 Optimization Results

The optimization problems were solved using the 'linprog' function in MATLAB (The MathWorks). Solving the inverse problem using the MFLSR method, the following fuzzy coefficients were obtained:

$$Y_{10} = (0.067207, 0.006721) + (-0.046013, 0.004601) \ X \tag{3.20}$$

for the 10% (epistemic) uncertainty case and

$$Y_{50} = (0.067207, 0.033603) + (-0.046013, 0.023007) X \tag{3.21}$$

for the 50% (epistemic) uncertainty case. Since the fuzzy centers were determined via ordinary least-squares regression they are identical to the intercept and slope of a least-squares regression. In this study since the spread of the data is a percentage applied to the entire data set, the halfwidths are this same percentage of the center. Based on the linear equation from the Cooper-Jacob method, since the intercepts and the slopes have nonzero halfwidths, Figs. 3.1 and 3.2, the storage coefficient and transmissivity must be fuzzy numbers, the calculation of these values follows in the next two sections. Figs. 3.3 and 3.4 show how these regression results relate to the observed data.

For comparison, the results of the FLSR method as proposed by Savic and Pedrycz (1991) and the hybrid method (Chang, 2001) are listed in Table 3.2. For these methods the independent data were the same, but the dependent data were the observed drawdown values, not the calculated values as is used in the MFLSR. All three methods have the same fuzzy centers and they are equivalent to the coefficients of the ordinary least-squares regression. The MFLSR method produces halfwidths smaller than the FLSR except for the slope halfwidth in the 10% case that is essentially zero. This is not an uncommon occurrence with the FLSR method (Savic and Pedrycz, 1991; Uddameri,

50

**Figure 3.3**: **This plot shows how the observed data is "bounded" by the results of the modified fuzzy least-squares regression results for the 10% (epistemic) uncertainty case.**



**Figure 3.4: This plot shows how the observed data is "bounded" by the results of the modified fuzzy least-squares regression results for the 50% (epistemic) uncertainty case.**

51

**Table 3.2: Results of various fuzzy regression methods**

| | 10% Uncertainty | | | | 50% Uncertainty | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | $\tilde{A}_0$ | | $\tilde{A}_1$ | | $\tilde{A}_0$ | | $\tilde{A}_1$ | |
| | $m_0$ | $c_0$ | $m_1$ | $c_1$ | $m_0$ | $c_0$ | $m_1$ | $c_1$ |
| Savic and Pedrycz FLSR | 0.06721 | 0.06450 | -0.04601 | 1.75E-14 | 0.06721 | 0.04500 | -0.04601 | 0.02430 |
| Hybrid FLSR | 0.06721 | 0.00672 | -0.04601 | -0.00460 | 0.06721 | 0.03360 | -0.04601 | -0.02301 |
| MFLSR | 0.06721 | 0.00672 | -0.04601 | 0.00460 | 0.06721 | 0.03360 | -0.04601 | 0.02301 |

2004; Si and Bodhinayake, 2005; and Uddameri and Honnungar, 2007). In this case, the halfwidth is determined to be zero by FLSR due to the minimization of the objective function in Eq. (3.3). The slope halfwidth will have a small (near zero) value when the sum of the $|x_i|$ is greater than *n*. The technique of MFLSR addresses this issue by allowing the fuzziness to be distributed over both the intercept and the slope halfwidths. The halfwidths of the MFLSR are equivalent to the absolute value of the halfwidths of the hybrid method. The hybrid method, however, has the tendency to produce counterintuitive negative halfwidths. Spreads or halfwidths are commonly defined as positive values.

## 3.5.2 Transmissivity

Recall from the Cooper-Jacob method that the slope of the linear equation is used to solve for transmissivity. Since the numbers being dealt with are no longer crisp, in order to solve for transmissivity, the following equation would have to be solved:

$$T_{10} = -\frac{Q}{4\pi(\text{-}0.046013,\ 0.004601)} \tag{3.22}$$

52

for the 10% case and for the 50% case the equation would be:

$$T_{50} = -\frac{Q}{4\pi(\text{-}0.046013, 0.023007)} \qquad (3.23)$$

where the values contained in the parentheses are the fuzzy slopes obtained from the $\tilde{A}_1$ term in Eqs. (3.20) and (3.21). In order to solve these equations, a Fortran 90 based program, ExtFUZZ, written by Ozbek and Pinder (2005), was used to perform the calculations. The program ExtFUZZ implements the n-dimensional form of the extension principle (Zadeh, 1975). From Dubois and Prade (1991) the extension principle can be written as:

$$\mu_{f(F_1,...,F_n)}(y) = \sup\{\min(\mu_{F_1}(x_1), \mu_{F_2}(x_2),..., \mu_{F_n}(x_n)) \mid y = f(x_1, x_2,..., x_n)$$

$$(3.24)$$

where $\mu_{f(F_1,...,F_n)}(y)$ represents the membership function of the fuzzy result $f(F_1,..., F_n)$ of the state variable $y$ and $\mu_{F_1}$ denotes the membership function of the fuzzy set associated with input parameter $i$.

Here, the implementation of the extension principle applied to fuzzy sets is based on a linear approximation of $f$. Given the inquiry on $\mu_{f(F_1,...,F_n)}(y)$ which represents the degree of membership of $y$ as the value of the model state variable, it proceeds in three steps:

<u>Step 1:</u>

A Delaunay tessellation of the *n*-dimensional parameter space is constructed. This results in a number of simplices. In an *n*-dimensional problem each simplex will have *n*+1 vertices. For example, in the two-dimensional case, if *f* is evaluated at only four vertices, this will result in two triangular simplices. The tessellation of the parameter is

then followed by the identification of simplices $X^j$, $j = 1,...,K$ that contain $y$. A simplex contains $y$ if

$$\min_i f_j^i \leq y \leq \max_j f_j^i \tag{3.25}$$

holds where $f_j^i$ denotes the function value at vertex $i$ of simplex $X^j$.

Finally, a trial function $\hat{f}^j$ is constructed within each of the $K$ simplices. The trial function gives the exact function value at the vertices and uses a linear approximation of $f$ within the simplex:

$$\hat{f}^j(x_1,...,x_n) = \sum_{i=1}^n a_j^i x_i + a_{n+1}^j. \tag{3.26}$$

Step 2:

An optimization (linear programming) is performed within each simplex $X^j$ of Step 1:

$$\max \quad a^j \tag{3.27}$$

subject to:

$$\mu_{F_i}(x_i) \geq a^j \quad i = 1,2,...,n$$
$$\hat{f}^j(x_1,x_2,...,x_n) = y \tag{3.28}$$

for $j = 1,2,...,K$.

Step 3:

Using $a^j$ of Step 2, $\mu_{f(F_1,...,F_n)}(y)$ is determined as:

$$\mu_{f(F_1,...,F_n)}(y) = \max_j a^j. \tag{3.29}$$

The code is written such that the results are sent to a MATLAB m-file from which the membership functions of the $n$ inputs, and the newly calculated transmissivity membership function can be plotted. The resulting membership functions for

54

transmissivity for the MFLSR method are shown in Fig. 3.5.  The supports of the fuzzy

membership functions for the 10% and 50% uncertainty cases are 5.62 cm$^2$/s to 7.59

cm$^2$/s and 4.12 cm$^2$/s to 13.65 cm$^2$/s, respectively.  The transmissivity value with a

membership degree of one in both cases coincides with the transmissivity value that

would be obtained via the standard Cooper-Jacob method, 6.49 cm$^2$/s.

### 3.5.3  Storage Coefficient

Since both the intercept and slope are fuzzy numbers, instead of using the typical

Cooper-Jacob method explained earlier in this paper, a more direct calculation for the

storage coefficient uses Eq. (3.18) and the fact that the value of the first term is known, it

is the value of the intercept from the regression optimizations, namely $\tilde{A}_0$ . Therefore,



**Figure 3.5: Transmissivity membership functions for the 10% and 50% (epistemic) uncertainty cases using modified fuzzy least-squares regression.**

$$\tilde{A}_0 = -\frac{Q}{4\pi\tilde{T}}\left[0.5772 + \ln\frac{r^2\tilde{S}}{4\tilde{T}}\right] \tag{3.30}$$

where $\tilde{T}$ and $\tilde{S}$ are the fuzzy transmissivity and fuzzy storage coefficient values, respectively. Upon rearranging,

$$\tilde{S} = \frac{4\tilde{T}}{r^2}\exp\left[\frac{-4\pi\tilde{T}}{Q}[\tilde{A}_0] - 0.5772\right]. \tag{3.31}$$

In general, to solve the above equation, two known fuzzy numbers, $\tilde{T}$ and $\tilde{A}_0$, must be used to calculate the fuzzy storativity membership function, $\tilde{S}$. ExtFUZZ is called upon again in order to solve Eq. (3.31) for the storativity membership function.

Using the $\tilde{A}_0$ terms from Eqs. (3.20) and (3.21) for the 10% and 50% (epistemic) uncertainty cases, respectively, and the corresponding transmissivity membership functions obtained in 3.5.2, results in the storativity membership functions as seen in Fig. 3.6. The supports of the storativity membership functions for the 10% and 50% uncertainty cases are 1.26 x 10$^{-4}$ to 2.40 x 10$^{-4}$ and 1.11 x 10$^{-5}$ to 4.32 x 10$^{-4}$, respectively. Similar to the transmissivity case, the storativity value with a membership degree of one in both cases coincides with the storativity of this data set calculated using the standard Cooper-Jacob method, 1.89 x 10$^{-4}$.

**Figure 3.6: Storativity membership functions for the 10% and 50% (epistemic) uncertainty cases using modified fuzzy least-squares regression.**

## 3.6 Disscussion and Conclusions

This paper looks at how to incorporate epistemic uncertainty (different from the well-studied aleatory uncertainty) into the 'straight-line' method, developed by Cooper and Jacob (1946), for pump test analysis. More specifically, the impact of model uncertainty on transmissivity and the storage coefficient is examined.

Since, traditionally, ordinary linear regression is used to solve for the transmissivity and storage coefficient via the Cooper-Jacob method, fuzzy least-squares regression seems an appropriate method to examine the epistemic uncertainty associated with these resulting parameters. In this paper a modified fuzzy least-squares regression was used instead of one of the pre-existing fuzzy least-squares regression methods because it had the following benefits, beyond those of FLSR and the hybrid approach: 1) the limitation

on fuzzy regression and fuzzy least-squares regression, that as the data approach crisp values the regression solution does not converge to the ordinary least-squares regression, is removed, 2) the optimization problem is simple to solve, 3) the distribution of the model fuzziness is more evenly distributed amongst the regression coefficients (the halfwidths on the coefficients are nonzero), and 4) compared to the hybrid method, the absolute values of the halfwidths are nearly identical, however the negative halfwidths that can be encountered using the hybrid method are avoided. A negative halfwidth traditionally has no meaning since typically a halfwidth is defined as a positive value.

The technique of using the MFLSR combined with the Cooper-Jacob method as described in this paper allowed for the incorporation of an uncertainty that has previously been neglected. By assigning an imprecision around the measured data, traditionally treated as crisp values, the MFLSR method produces a fuzzy linear regression relationship that, when used in place of ordinary linear regression results in the Cooper-Jacob equation, transmissivity and storage coefficient ranges, or membership functions, can be determined which better describe the uncertainty around those numbers. For example, in the 50% case a membership function is obtained that has a transmissivity value with a membership degree of one at the 6.49 $cm^2$/s, which is exactly what would be obtained using the standard Cooper-Jacob approach. With the approach presented here, transmissivity is allowed to have varying values of degree of membership that increase with transmissivity values from 4.12 $cm^2$/s to 6.49 $cm^2$/s then decrease with transmissivity values from 6.49 $cm^2$/s to 13.65 $cm^2$/s. Representing this epistemic uncertainty in transmissivity and storage coefficient values will allow for a better understanding of the heterogeneous subsurface.

58

## 3.7 Acknowledgements

## 3.8 References

Bardossy, A., Bogardi, I. and Duckstein, L., 1990. Fuzzy regression in hydrology. Water Resources Research, 26(7): 1497-1508.

Chang, Y.H.O., 2001. Hybrid fuzzy least-squares regression analysis and its reliability measures. Fuzzy Sets and Systems, 119(2): 225-246.

Chang, Y.H.O. and Ayyub, B.M., 2001. Fuzzy regression methods - a comparative assessment. Fuzzy Sets and Systems, 119(2): 187-203.

Cooper, H.H. and Jacob, C.E., 1946. A generalized graphical method for evaluating formation constants and summarizing well field history. American Geophysical Union Transactions, 27: 526-534.

Dubois, D., Martin-Clouaire, R. and Prade, H., 1988. Practical computing in fuzzy logic. In: M.M. Gupta and T. Yamakawa (Editors), Fuzzy Computing. North-Holland, Netherlands.

Dubois, D. and Prade, H., 1991. Random Sets and Fuzzy Interval-Analysis. Fuzzy Sets and Systems, 42(1): 87-101.

Ozbek, M.M. and Pinder, G.F., 2005. ExtFUZZ: A Fortran 90 program to implement the (ext)ension principle in (fuzz)y set theory. (Draft), Version 1.

Ozelkan, E.C. and Duckstein, L., 2001. Fuzzy conceptual rainfall-runoff models. Journal of Hydrology, 253(1-4): 41-68.

Savic, D.A. and Pedrycz, W., 1991. Evaluation of fuzzy linear-regression models. Fuzzy Sets and Systems, 39(1): 51-63.

Si, B.C. and Bodhinayake, W., 2005. Determining soil hydraulic properties from tension infiltrometer measurements: Fuzzy regression. Soil Science Society of America Journal, 69(6): 1922-1930.

Tanaka, H., Uejima, S. and Asai, K., 1982. Linear-regression analysis with fuzzy model. Ieee Transactions on Systems Man and Cybernetics, 12(6): 903-907.

Theis, C.V., 1940. The source of water derived from wells - essential factors controlling the response of an aquifer to development. American Society of Civil Engineers, 10: 277-280.

Uddameri, V., 2004. Relationships of longitudinal dispersivity and scale developed from fuzzy least-squares regressions. Environmental Geology, 45(8): 1172-1178.

Uddameri, V. and Honnungar, V., 2007. Interpreting sustainable yield of an aquifer using a fuzzy framework. Environmental Geology, 51(6): 911-919.

Zadeh, L.A., 1975. The concept of a linguistic variable and its applications in approximate reasoning. Information Science, 8: 199-251, 301-357; 9, 43-80.

# CHAPTER 4
# FUZZY GENERALIZED REGRESSION NEURAL NETWORK METHODOLOGY

## 4.1 Background

A new artificial neural network was developed that combines fuzzy sets with generalized regression to address the relationships between physical habitat and the geomorphic condition of Vermont streams. The focus is on using fuzzy numbers to capture expert information typically lost. The Vermont Agency of Natural Resources (VTANR) River Management Program (RMP) has developed protocols, based on well-known stream classification methods, to assess both the geomorphic condition (Rapid Geomorphic Assessment – RGA; Kline *et al*., 2007) and the physical habitat (Reach Habitat Assessment – RHA; Schiff *et al*., 2008) of a stream reach. Both of these assessments involve expert-based field observations. For example, in the RGA, experts assign a score between 0 (poor) and 20 (reference) to assess the four adjustment processes (i.e. degradation, aggradation, widening, and planform change) associated with stream geomorphic condition. The sum of these four scores provides a total RGA score between 0 and 80, which is used subsequently to classify the stream reach into one of four overall condition categories: that is *poor*, *fair*, *good*, and *reference*. Assigning individual scores to a stream adjustment process, and ultimately an entire reach, relies not only on physical measurements, but also on expert opinion. Figure 4.1 illustrates a small portion of the field assessment form for the RGA (Kline *et al*., 2007 – Appendix A). As an example, protocol requires the expert to assign (circle) an integer score to the adjustment process for channel degradation (7.1 on the form), while determining the

overall (categorical) condition for the stream reach.  In Figure 4.1, the "x's" indicate an expert's field observations for a particular reach.  Once the evidence for channel degradation has been evaluated, the expert must choose, to the best of their knowledge, which of the four categories best describes the reach being studied and assign an integer score to this process.  Protocol advises the assessor to give greater weight to the channel and floodplain geometry changes (rows 2-4 under this particular adjustment process) than the human induced changes (the lower rows in the adjustment process).  In this example, the expert assigned a score of 12 to the adjustment process degradation.  However, in the field, experts express difficulty in assigning a single score.  It is common to hear, "the score could be as high as 14 and as low as 11" (personal communication, Kristen Underwood).

Fuzzy numbers may provide a means to capture information that is lost when assigning a crisp number to a process that uses subjective information and expert opinion. A fuzzy number can capture the opinion that the process score is "around 12."  In this work, Specht's (1991) generalized regression neural network (GRNN) is modified to allow the use of fuzzy numbers to capture the imprecision of the assessor's opinion.  A new predictive fuzzy algorithm is developed.  The Vertex Method (Dong and Shah, 1987), an approximation to the Extension Principle (Zadeh, 1975), is implemented to solve the resulting fuzzy equations.  A small example shows how the new methodology is designed to capture the imprecision associated with assigning RGA scores to stream reaches.  As a result, one may account for information typically lost during expert assessments.  Knowing the imprecision associated with expert assessments may be more important than knowing a crisp number, when flagging reaches for further study.

| Adjustment Process | Condition Category | | | |
|---|---|---|---|---|
| | Reference | Good | Fair | Poor |
| **7.1 Channel Degradation**<br><br>•Exposed till or fresh substrate in the stream bed or exposed infrastructure<br><br>•New terraces or recently abandoned flood prone areas<br><br>•Headcuts, or nickpoints significantly steeper bed segment and comprised of smaller bed material than typical steps<br><br>•Freshly eroded, vertical banks<br><br>•Alluvial sediments imbricated high in bank<br><br>•Tributary rejuvenation, observed through the presence of nickpoints at or upstream of tributary mouth<br><br>•Depositional features with steep faces (usually on downstream end) | Little evidence of localized slope increase or nickpoints ✖ | Minor localized slope increase or nickpoints | Sharp change in slope, head cuts present, and/or tributaries rejuvenating | Sharp change in slope and/or multiple head cuts present, tributaries rejuvenating |
| | Incision ratio ≥ 1.0 < 1.2<br>and ✖<br>Where channel slope < 4%<br>Entrenchment ratio >1.4<br>Where channel slope ≥ 4%<br>Entrenchment ratio >1.2 | Incision ratio ≥ 1.2 < 1.4<br>and<br>Where channel slope < 4%<br>Entrenchment ratio >1.4<br>Where channel slope ≥ 4%<br>Entrenchment ratio >1.2 | Incision ratio ≥ 1.4 < 2.0<br>and<br>Where channel slope < 4%<br>Entrenchment ratio >1.4<br>Where channel slope ≥ 4%<br>Entrenchment ratio >1.2 | Incision ratio ≥ 2.0<br>and<br>Where channel slope < 4%<br>Entrenchment ratio >1.4<br>Where channel slope ≥ 4%<br>Entrenchment ratio >1.2 |
| | Step-pool systems have all expected bed features, steps complete with coarser sediment (≥ D80) | Step-pool systems have full complement of expected bed features, steps mostly complete ✖ | Step-pool systems with incomplete (eroded) steps dominated by runs | Step-pool bed features eroded and replaced by plane bed features |
| | No significant human-caused change in channel confinement | Only minor human-caused change in channel confinement ✖ | Significant human-caused change in channel confinement but no change in valley type | Human caused change in valley type |
| | No evidence of historic/present channel straightening, dredging, and/or channel avulsions | Evidence of minor historic dredging, and/or channel avulsion | Evidence of significant historic channel straightening, dredging, or gravel mining, and/or channel avulsion | Extensive historic channel straightening, commercial gravel mining, and/or recent channel avulsions ✖ |
| | No known flow alterations (i.e., increase in flow and/or decrease in sediment supply) ✖ | Some increase in flow and/or minor reduction of sediment load | Major historic flow alterations, greater flows and/or reduction of sediment load | Major existing flow alterations, greater flows and/or reduction of sediment load |
| Score: | 20  19  18  17  16 | 15  14  13  ⑫  11 | 10  9  8  7  6 | 5  4  3  2  1 |

**Figure 4.1: Channel degradation section of the VTANR Rapid Geomorphic Assessment field form found in Appendix A of Kline *et al*., 2007.**

The next section provides a brief introduction to the GRNN algorithm and the Vertex Method to facilitate the subsequent development of the fuzzy GRNN algorithm. An example calculation is presented to demonstrate the performance of the algorithm. The chapter concludes with a discussion of the strengths and challenges associated with utilization of this algorithm.

## 4.2 Methodology

### 4.2.1 Generalized Regression Neural Network

The GRNN introduced by Donald Specht (1991) is a parallel, one-pass algorithm designed to perform least-squares generalized regression. The network does not require iterative training like the more popular feed-forward backpropagation networks. The training data are used to set the network weights. What makes this algorithm unique, aside from it's parallel computational nature, is that it does not require *a priori* knowledge of the function that best fits the data. Figure 4.2 shows the structure of the GRNN algorithm as applied to the prediction of stream RHA scores using the four adjustment processes that comprise the RGA (degradation, aggradation, widening, and planform change) as inputs. These input variables are equivalent to the independent variables associated with traditional regression techniques.

To begin, the algorithm needs both training data and testing/prediction data. The training data set consists of $k$ training patterns. A single pattern is defined as one set of $n$ input variables, $X = \{x_{i=1}, x_{i=2}, ..., x_{i=n}\}$ and the corresponding output (dependent) variable, $y_j$. In this algorithm, the training input variables are also the network weights; thus $X = W_j = \{w_{1j}, w_{2j}, ...w_{nj}\}$. The prediction data set consists of additional input patterns (each

64

**Figure 4.2: GRNN structure showing the four components of the RGA as inputs used to predict the total Legacy RHA score.**

comprised of *n* input variables) for which the user would like predictions.

The GRNN network consists of four nodal layers. The first *Input Layer* simply passes the *n* input variables, $X = \{x_{i=1}, x_{i=2}, ..., x_{i=n}\}$, to the weights of the next network layer. The training weights, $w_{ij}$, connect the *Input Layer* to the *Pattern Units* layer (e.g. $w_{12}$ connects input node $x_{i=1}$ with pattern unit node $I_{j=2}$). These weights are set by the training data and do not update as in other artificial neural network (ANN) algorithms. Each $j^{th}$ training pattern weight, $w_{ij}$, contains a value (e.g., degradation, aggradation, widening, planform change) for which there is a corresponding output (RHA score). The RHA score is stored in the weights, $y_j$, associated with node A of the *Summation Units* layer (Figure 4.2). The *Pattern Units* layer has one node for each of the *j* training

65

patterns and calculates a distance metric (e.g., the Euclidean distance) between all sets of

training weights and the current input pattern for which a prediction is desired (Eqn. 4.1):

$$I_j = \sqrt{\sum_{i=1}^{n}(w_{ij} - x_i)^2} \; , \qquad\qquad (4.1)$$

where $x_i$ refers to the $i^{\text{th}}$ input parameter, $w_{ij}$ are weights associated with the $i^{\text{th}}$ input

variable and the $j^{\text{th}}$ training pattern. The resulting Euclidean distance, $I_j$, is passed

through an exponential activation function (Eqn. 4.2):

$$f(I_j) = \exp\left(\frac{-I_j}{2\sigma^2}\right), \qquad\qquad (4.2)$$

where $\sigma$ is a smoothing parameter explained in greater detail below.

The third layer, *Summation Units*, calculates the dot product of the output from the

*Pattern Units* (Eqn. 4.2) and, for node A, the corresponding output training weights, $y_j$.

The weights associated with node B are set equal to 1; node B calculates the dot product

between the output from the *Pattern Units* and the weights set equal to 1. The final

output is the result of dividing the nodes in the *Summation Units*:

$$\hat{y}(X) = \frac{\sum_{j} y_j \bullet f(I_j)}{\sum_{j} 1 \bullet f(I_j)} = \frac{A}{B}. \qquad\qquad (4.3)$$

Recall that the weights are fixed in the GRNN algorithm. Thus, $\sigma$ (Eqn. 4.2) is the only

parameter that may be adjusted by the user and is used to optimize the GRNN output. As

$\sigma$ approaches zero, the predicted network output, $\hat{y}$, tends to overfit the training data.

When $\sigma$ is large, $\hat{y}$ is smoothed and assumes the value of the sample mean. For details,

the reader is referred to Specht (1991). The GRNN algorithm described in this paper was written in MATLAB 7.10.0 (R2010a).

### 4.2.1.1 Example of GRNN Calculation

To illustrate how the GRNN works, an example is presented using one expert's assessed RGA components (degradation, aggradation, widening, and planform change) to predict a total RHA score (the response variable) for a particular stream reach since a correlation between the two parameters has been previously shown (Chapter 5, Figure 5.4). Training patterns from 20 stream reaches in Vermont along the Lewis Creek and prediction patterns from 38 Middlebury River reaches that have expert assessed RGA and RHA data are provided in Tables 4.1 and 4.2, respectively. These reaches have been selected for use in past ANN studies (Doris, 2006; Besaw *et al.*, 2009) due to the similarities in watershed size and land cover type.

Choosing to demonstrate the computational numerics with Middlebury River reach M01 as the prediction pattern, Eqn. 4.1 is calculated for $j = 1$ and $j = 2$ as follows:

$$
\begin{aligned}
I_1 &= \sqrt{(x_1 - w_{11})^2 + (x_2 - w_{21})^2 + (x_3 - w_{31})^2 + (x_4 - w_{41})^2} \\
&= \sqrt{(18 - 18)^2 + (11 - 17)^2 + (13 - 18)^2 + (15 - 18)^2} \\
&= 8.37 \\
I_2 &= \sqrt{(x_1 - w_{12})^2 + (x_2 - w_{22})^2 + (x_3 - w_{32})^2 + (x_4 - w_{42})^2} \\
&= \sqrt{(18 - 18)^2 + (11 - 15)^2 + (13 - 13)^2 + (15 - 16)^2} \\
&= 4.12.
\end{aligned}
$$

The remainder of the $I_j$ calculations follow similarly and the results are shown in Table 4.3. Passing $I_1$ through the activation function (Eqn. 4.2) with $\sigma = 0.55$ yields:

$$
f(I_1) = \exp\left(\frac{-I_1}{2\sigma^2}\right) = \exp\left(\frac{-8.37}{2(0.55)^2}\right) = 9.87 \times 10^{-7}
$$

**Table 4.1: Lewis Creek training data for GRNN example.**

| Reach ID | | Degradation Score | Aggradation Score | Widening Score | Planform Change Score | RHA Total Score |
|---|---|---|---|---|---|---|
| j=1 | M07 | 18 | 17 | 18 | 18 | 175 |
| j=2 | M18 | 18 | 15 | 13 | 16 | 186 |
| j=3 | T2.01 | 18 | 18 | 17 | 17 | 169 |
| j=4 | M05 | 15 | 13 | 14 | 15 | 155 |
| j=5 | M14 | 18 | 15 | 15 | 18 | 152 |
| j=6 | M15A | 18 | 10 | 13 | 8 | 135 |
| j=7 | M17B | 18 | 11 | 15 | 8 | 138 |
| j=8 | M19A | 18 | 15 | 16 | 11 | 133 |
| j=9 | M20B | 5 | 10 | 10 | 13 | 143 |
| j=10 | M03 | 18 | 13 | 16 | 13 | 123 |
| j=11 | M15B | 16 | 11 | 10 | 6 | 119 |
| j=12 | M16 | 16 | 15 | 6 | 11 | 122 |
| j=13 | M17A | 15 | 13 | 15 | 11 | 127 |
| j=14 | M17C | 10 | 15 | 11 | 13 | 128 |
| j=15 | M19B | 13 | 13 | 13 | 11 | 125 |
| j=16 | M20A | 8 | 11 | 13 | 8 | 110 |
| j=17 | M21A | 10 | 13 | 10 | 13 | 125 |
| j=18 | M21B | 8 | 11 | 13 | 6 | 111 |
| j=19 | M22 | 11 | 13 | 11 | 10 | 105 |
| j=k=20 | T4.3S6.01 | 18 | 8 | 13 | 11 | 100 |

for the first *Pattern Unit* node. The remaining $f(I_j)$ results are listed in Table 4.3.

To calculate node A, the dot product of $y_j$ (the total RHA scores, Table 4.1, last column) and the $f(I_j)$ (Table 4.3) is taken:

$$A = (175)(9.87 \times 10^{-7}) + (186)(1.10 \times 10^{-3}) + (169)(1.09 \times 10^{-6}) + ... + (100)(2.58 \times 10^{-4})$$
$$= 0.7222.$$

Similarly, node B is calculated as:

$$B = (1)(9.87 \times 10^{-7}) + (1)(1.10 \times 10^{-3}) + (1)(1.09 \times 10^{-6}) + ... + (1)(2.58 \times 10^{-4})$$
$$= 0.0048.$$

The predicted GRNN output for this stream reach is a total RHA score of:

$$\hat{y} = \frac{A}{B} = 151,$$

**Table 4.2: Middlebury River prediction data for GRNN example.**

| Reach ID | Degradation Score | Aggradation Score | Widening Score | Planform Change | RHA Total Score | GRNN Prediction |
|---|---|---|---|---|---|---|
| M01 | 18 | 11 | 13 | 15 | 131 | 151 |
| M02 | 18 | 13 | 10 | 10 | 142 | 127 |
| M03 | 16 | 13 | 11 | 13 | 142 | 145 |
| M04 | 15 | 11 | 8 | 5 | 127 | 119 |
| M05 | 16 | 10 | 11 | 5 | 137 | 119 |
| M06A | 13 | 11 | 15 | 15 | 112 | 152 |
| M06B | 6 | 8 | 8 | 3 | 122 | 111 |
| M07 | 8 | 13 | 13 | 10 | 115 | 111 |
| M08A | 3 | 13 | 13 | 15 | 145 | 143 |
| M11 | 16 | 13 | 13 | 13 | 146 | 144 |
| M12A | 10 | 13 | 13 | 10 | 133 | 110 |
| M12C | 5 | 13 | 15 | 8 | 131 | 110 |
| M13A | 13 | 8 | 11 | 6 | 138 | 119 |
| M13B | 5 | 15 | 11 | 15 | 110 | 133 |
| M14 | 18 | 11 | 15 | 10 | 137 | 135 |
| M15 | 11 | 16 | 15 | 13 | 129 | 128 |
| M16 | 16 | 10 | 12 | 8 | 147 | 131 |
| M17 | 3 | 10 | 10 | 8 | 146 | 130 |
| M18 | 16 | 10 | 11 | 10 | 159 | 121 |
| M19 | 13 | 13 | 11 | 10 | 144 | 113 |
| T3.01 | 18 | 18 | 18 | 18 | 158 | 173 |
| T3.02 | 18 | 16 | 13 | 11 | 150 | 134 |
| T3.03 | 18 | 16 | 18 | 18 | 170 | 174 |
| T3.04 | 18 | 15 | 18 | 16 | 180 | 167 |
| T3.05 | 18 | 18 | 18 | 16 | 177 | 170 |
| T3.06 | 18 | 16 | 13 | 13 | 158 | 164 |
| T3.08 | 19 | 12 | 14 | 15 | 145 | 147 |
| T3.09 | 19 | 18 | 16 | 17 | 176 | 169 |
| T3.10 | 19 | 19 | 18 | 16 | 176 | 170 |
| T4.01 | 10 | 15 | 15 | 13 | 159 | 127 |
| T4.02 | 13 | 13 | 13 | 11 | 157 | 125 |
| T4.03A | 5 | 15 | 11 | 15 | 162 | 133 |
| T4.03B | 14 | 11 | 13 | 8 | 129 | 125 |
| T4.04A | 16 | 10 | 15 | 10 | 145 | 129 |
| T4.04B | 11 | 7 | 8 | 6 | 105 | 115 |
| T4.05 | 13 | 13 | 13 | 13 | 129 | 129 |
| T4.07A | 10 | 10 | 10 | 8 | 119 | 109 |
| T4.07B | 16 | 11 | 15 | 11 | 141 | 127 |

**Table 4.3: Results of Pattern Unit calculations (Figure 4.2, Eqns. 4.1 and 4.2).**

| Pattern Unit Node Number | $I_j$ | $f(I_j)$ |
|---|---|---|
| 1 | 8.37 | 9.87E-07 |
| 2 | 4.12 | 1.10E-03 |
| 3 | 8.31 | 1.09E-06 |
| 4 | 3.74 | 2.06E-03 |
| 5 | 5.39 | 1.36E-04 |
| 6 | 7.07 | 8.40E-06 |
| 7 | 7.28 | 5.94E-06 |
| 8 | 6.40 | 2.53E-05 |
| 9 | 13.53 | 1.95E-10 |
| 10 | 4.12 | 1.10E-03 |
| 11 | 9.70 | 1.10E-07 |
| 12 | 9.22 | 2.41E-07 |
| 13 | 5.74 | 7.52E-05 |
| 14 | 9.38 | 1.85E-07 |
| 15 | 6.71 | 1.53E-05 |
| 16 | 12.21 | 1.73E-09 |
| 17 | 9.00 | 3.46E-07 |
| 18 | 13.45 | 2.20E-10 |
| 19 | 9.06 | 3.16E-07 |
| 20 | 5.00 | 2.58E-04 |

(Table 4.2, last column) which classifies as a *good* habitat condition stream. The expert assigned RHA score is 131, which is also classified as a *good* habitat condition.

## 4.2.2 Vertex Method

This section will begin with terminology definitions specific to fuzzy set theory. A *fuzzy set* can be described as a set whose elements have varying degrees of membership (e.g. the sets large, medium, and small) and the elements can have membership in more than one set. This is unlike crisp sets, where an element either belongs to a set or it doesn't. A fuzzy set is *normal* if it has at least one element whose membership degree is equal to one. The *support* of a fuzzy set is defined as the crisp set of all the elements of

**Figure 4.3: Example of a fuzzy number output from the FuzzyGRNN. The dashed vertical lines show the interval cutoff values for an α-cut at membership degree 0.6 (e.g. $^{0.6}C = [148, 152]$).**

the fuzzy set with nonzero membership degrees. An *α-cut* is a crisp set, $^{\alpha}C$, on the fuzzy set $\tilde{C}$ that contains all the elements of $\tilde{C}$ whose membership degree is greater than or equal to the α value. Figure 4.3 shows an example α-cut at membership degree 0.6 on an example fuzzy number. A *fuzzy number* is a convex, normal fuzzy set on the set of real numbers with a bounded support and every α-cut must be a closed interval. To perform function evaluations using fuzzy numbers that are traditionally used on crisp numbers, the function must be fuzzified. The principle that allows for this is known as the Extension Principle (Zadeh, 1975), which generates $\mu$, the membership function of the given fuzzy set, and is defined as follows:

$$
\mu_{\tilde{D}}(t) = \begin{cases} \sup_{t=f(s_1,s_2,...s_n)} \min \left\{ \mu_{\tilde{C_1}}(s_1), \mu_{\tilde{C_2}}(s_2),...,\mu_{\tilde{C_n}}(s_n) \right\} & \\ \qquad\qquad \text{if } \exists\, t = f(s_1,s_2,...,s_n) , \\ \qquad 0 \qquad\qquad \text{otherwise} \end{cases} \tag{4.4}
$$

where $s_i$ is an independent variable and $\tilde{C}_i$ is its fuzzy set.  Then, $t = f(s_1, s_2, ...s_n)$ is the dependent variable and $\tilde{D}$ is the fuzzy set for $t$.  For details on fuzzy set theory, Klir and Yuan (1995) provide a good introduction.

Given the challenges associated with the computational coding of the expression in Eqn. 4.4, several approximations to the Extension Principle have been adopted.  One approach is to discretize the fuzzy numbers, and then apply the Extension Principle. There is also the DSW method (Dong *et al.*, 1985) and the Vertex Method (Dong and Shah, 1987), each of which begin by representing the fuzzy numbers as a series of α-cuts. For example, the fuzzy number in Figure 4.3 may be represented using $^0C$ = [145, 155]; $^{0.5}C$ = [147.5, 152.5]; $^{1.0}C$ = [150, 150].  The DSW method uses the α-cut defined intervals to carry out the mathematical function(s) using standard interval analysis rules. The Vertex Method (the method chosen for this analysis) can reduce abnormalities in the observed output when using the Extension Principle on a discretized set (Ross, 2004), resulting from the number of discretizations.  The Vertex Method is computationally easier to implement and differs from the DSW method in that it deals only with the interval endpoints defined by the α-cuts. Ross (2004) provides a nice comparison of these three methods.  To apply the Vertex Method, the function must be continuous and monotonic.  If there are extreme values in the function, the method may omit calculating

these; therefore, extreme values are treated as possible vertices along with the interval endpoints (see Eqn. 4.5).

Letting the α-cut intervals be represented by $^{\alpha}C$ = [a,b] and the extreme value(s) (if there are any) be denoted as $E_i$, then the value of the function, $f(s)$, evaluated at the endpoints of the interval of the given α value, denoted $f(^{\alpha}C)$, is defined as:

$$f(^{\alpha}C) = [\min(f(a), f(b), E_i), \max(f(a), f(b), E_i)].$$
(4.5)

### 4.2.3 Fuzzifying the GRNN

The algorithm designed here (Appendix A) allows fuzzy numbers as input that are triangular, but not necessarily symmetrical. The current algorithm is designed to work with a function that discretizes these fuzzy numbers, along with the user-supplied training weights, prediction patterns, and discretization size of the fuzzy number (i.e. values of the α-cuts at which to evaluate the function for the Vertex Method). The discretized function assumes that the edges of the triangular membership function are linear. The algorithm can accommodate triangular membership functions without linear edges by not calling the discretization function and using user-described α-cuts.

Note that the Vertex Method is *applied to the entire algorithm* and *not at each nodal layer*. More specifically, the input variables of each α-cut are carried through all steps of the GRNN and the final step, taking the minimum and maximum, provide an output interval for the given α-cut. If one makes the mistake of performing the Vertex Method for every mathematical operation at each layer, vertices that do not exist in the initial problem are introduced resulting in orders of magnitude of spread in the final output. This can produce membership functions that are meaningless in many applications.

73

**Table 4.4:  The two training input weights (a) and training pattern weights (b) for the example in Section 4.2.3.  The prediction input variables are presented in (c).**

| (a) Training α-cut | $\tilde{w}_{11}$ Left Bound | $\tilde{w}_{11}$ Right Bound | $\tilde{w}_{12}$ Left Bound | $\tilde{w}_{12}$ Right Bound | $\tilde{w}_{21}$ Left Bound | $\tilde{w}_{21}$ Right Bound | $\tilde{w}_{22}$ Left Bound | $\tilde{w}_{22}$ Right Bound |
|---|---|---|---|---|---|---|---|---|
| 0 | 17 | 19 | 15 | 17 | 5 | 7 | 12 | 14 |
| 0.25 | 17.25 | 18.75 | 15.25 | 16.75 | 5.25 | 6.75 | 12.25 | 13.75 |
| 0.50 | 17.5 | 18.5 | 15.5 | 16.5 | 5.5 | 6.5 | 12.5 | 13.5 |
| 0.75 | 17.75 | 18.25 | 15.75 | 16.25 | 5.75 | 6.25 | 12.75 | 13.25 |
| 1 | 18 | 18 | 16 | 16 | 6 | 6 | 13 | 13 |

| (b) Training α-cut | $\tilde{y}_1$ Left Bound | $\tilde{y}_1$ Right Bound | $\tilde{y}_2$ Left Bound | $\tilde{y}_2$ Right Bound |
|---|---|---|---|---|
| 0 | 22 | 26 | 27 | 31 |
| 0.25 | 22.5 | 25.5 | 27.5 | 30.5 |
| 0.50 | 23 | 25 | 28 | 30 |
| 0.75 | 23.5 | 24.5 | 28.5 | 29.5 |
| 1 | 24 | 24 | 29 | 29 |

| (c) Predict α-cut | $\tilde{x}_1$ Left Bound | $\tilde{x}_1$ Right Bound | $\tilde{x}_2$ Left Bound | $\tilde{x}_2$ Right Bound |
|---|---|---|---|---|
| 0 | 10 | 12 | 7 | 9 |
| 0.25 | 10.25 | 11.75 | 7.25 | 8.75 |
| 0.50 | 10.5 | 11.5 | 7.5 | 8.5 |
| 0.75 | 10.75 | 11.25 | 7.75 | 8.25 |
| 1 | 11 | 11 | 8 | 8 |

A simple example is provided to illustrate how the algorithm works.  With Eqn. 4.3 in mind, consider a network with 2 input nodes, $\tilde{x}_1$ and $\tilde{x}_2$, and using 2 training patterns (training weights $\tilde{w}_{11}, \tilde{w}_{12}, \tilde{w}_{21}$, and $\tilde{w}_{22}$, and corresponding pattern weights $\tilde{y}_1$ and $\tilde{y}_2$). Table 4.4 lists the training data in (a) and (b) and prediction inputs (c) after being passed through the discretization function.  For each α-cut, the left and right bound of the interval are given.  Recall that the weights connecting the *Pattern Unit* nodes to node B are equal to 1 and that a crisp number may be represented in interval form as [1,1]. Running the algorithm one α-cut at a time, the distance between the input nodes and the training weights is calculated producing 4 possible outcomes for each input node and weight combination.  For example, consider $\alpha = 0.5$, then $\tilde{x}_1 = [10.5,11.5]$ and $\tilde{w}_{11} =$

[17.5,18.5]. The possible squared distances for these two fuzzy numbers are:

$$([10.5,11.5] - [17.5,18.5])^2 = (10.5 - 17.5)^2 \text{ or } (10.5 - 18.5)^2$$
$$\text{or } (11.5 - 17.5)^2 \text{ or } (11.5 - 18.5)^2 \qquad (4.6)$$
$$= 49 \text{ or } 64 \text{ or } 36 \text{ or } 49.$$

The remaining distance calculations are listed in Table 4.5. Next, all possible summations must be performed at each *Pattern Unit* node. Given 2 weights connected to each node each with 4 possible outcomes, there are 16 possible summations. The 16 possible values for $I_1$ are the square root of all summation combinations of $(\tilde{x}_1 - \tilde{w}_{11})^2$ and $(\tilde{x}_2 - \tilde{w}_{21})^2$, and, similarly for $I_2$, the square root of all summation combinations of $(\tilde{x}_1 - \tilde{w}_{12})^2$ and $(\tilde{x}_2 - \tilde{w}_{22})^2$ (Table 4.6). These values are then passed through the activation function (Eqn. 4.2, Table 4.6). The output values from each *Pattern Unit* node connected to node A are multiplied by a pattern weight, $\tilde{y}_j$, producing 32 possible outcomes from each *Pattern Unit* node. So, if the first possible outcome from $I_1$, $f(I_1)$ = 6.89 x $10^{-4}$ (Table 4.6), is multiplied by $\tilde{y}_1$ = [22,26] (Table 4.4 (b)), two possible outcomes are produced (Table 4.7, associated with Node A). Similar calculations for the remaining possible *Pattern Unit* output ($f(I_j)$) results are listed in Table 4.7. All combinations are then summed at node A.

**Table 4.5: Results of taking the distance between the fuzzy input variables and the fuzzy training weights for the 0.5 $\alpha$-cut in the example in Section 4.2.3.**

| $(\tilde{x}_1 - \tilde{w}_{11})^2$ | $(\tilde{x}_1 - \tilde{w}_{12})^2$ | $(\tilde{x}_2 - \tilde{w}_{21})^2$ | $(\tilde{x}_2 - \tilde{w}_{22})^2$ |
|---|---|---|---|
| 49 | 25 | 4 | 25 |
| 64 | 36 | 1 | 36 |
| 36 | 16 | 9 | 16 |
| 49 | 25 | 4 | 25 |

75

**Table 4.6: Pattern Unit results for example in Section 4.2.3.**

| $I_1$ | $f(I_1)$ $\sigma^2 = 0.5$ | $I_2$ | $f(I_2)$ $\sigma^2 = 0.5$ |
|---|---|---|---|
| 7.28 | 6.89E-04 | 7.07 | 8.49E-04 |
| 7.28 | 6.89E-04 | 7.07 | 8.49E-04 |
| 7.07 | 8.49E-04 | 7.81 | 4.06E-04 |
| 7.62 | 4.93E-04 | 6.40 | 1.66E-03 |
| 7.28 | 6.89E-04 | 7.07 | 8.49E-04 |
| 7.28 | 6.89E-04 | 7.07 | 8.49E-04 |
| 7.07 | 8.49E-04 | 7.81 | 4.06E-04 |
| 7.62 | 4.93E-04 | 6.40 | 1.66E-03 |
| 8.25 | 2.62E-04 | 7.81 | 4.06E-04 |
| 8.25 | 2.62E-04 | 7.81 | 4.06E-04 |
| 8.06 | 3.15E-04 | 8.49 | 2.06E-04 |
| 8.54 | 1.95E-04 | 7.21 | 7.38E-04 |
| 6.32 | 1.79E-03 | 6.40 | 1.66E-03 |
| 6.32 | 1.79E-03 | 6.40 | 1.66E-03 |
| 6.08 | 2.28E-03 | 7.21 | 7.38E-04 |
| 6.71 | 1.22E-03 | 5.66 | 3.49E-03 |

For example, the possible outcomes for node A start with:

$$1.58\text{E-}02 + 2.38\text{E-}02 = 3.96\text{E-}02 \text{ or}$$

$$1.58\text{E-}02 + 2.55\text{E-}02 = 4.13\text{E-}02 \text{ or}$$

$$1.58\text{E-}02 + 2.38\text{E-}02 = 3.96\text{E-}02 \text{ or}$$

…

and end with

$$3.05\text{E-}02 + 1.05\text{E-}01 = 1.36\text{E-}01.$$

In this example, the 2 *Pattern Unit* nodes connected to node A, result in 1,024 possible outcomes.

Concurrently, similar calculations are being conducted for the *Summation Unit* node B and these divide the possibilities in node A; hence, there are 1,024 divisions. The minimum and maximum of these possible outcomes are selected as the result for the

**Table 4.7: Results of multiplication of output from Pattern Units by corresponding pattern weights for example in Section 4.2.3.**

| Associated with Node A | | Associated with Node B | |
|---|---|---|---|
| Possible Outcomes from $f(I_1) \times \tilde{y}_1$ | Possible Outcomes from $f(I_2) \times \tilde{y}_2$ | Possible Outcomes from $f(I_1) \times 1$ | Possible Outcomes from $f(I_2) \times 1$ |
| 1.58E-02 | 2.38E-02 | 6.89E-04 | 8.49E-04 |
| 1.72E-02 | 2.55E-02 | 6.89E-04 | 8.49E-04 |
| 1.58E-02 | 2.38E-02 | 6.89E-04 | 8.49E-04 |
| 1.72E-02 | 2.55E-02 | 6.89E-04 | 8.49E-04 |
| 1.95E-02 | 1.14E-02 | 8.49E-04 | 4.06E-04 |
| 2.12E-02 | 1.22E-02 | 8.49E-04 | 4.06E-04 |
| 1.13E-02 | 4.64E-02 | 4.93E-04 | 1.66E-03 |
| 1.23E-02 | 4.97E-02 | 4.93E-04 | 1.66E-03 |
| 1.58E-02 | 2.38E-02 | 6.89E-04 | 8.49E-04 |
| 1.72E-02 | 2.55E-02 | 6.89E-04 | 8.49E-04 |
| 1.58E-02 | 2.38E-02 | 6.89E-04 | 8.49E-04 |
| 1.72E-02 | 2.55E-02 | 6.89E-04 | 8.49E-04 |
| 1.95E-02 | 1.14E-02 | 8.49E-04 | 4.06E-04 |
| 2.12E-02 | 1.22E-02 | 8.49E-04 | 4.06E-04 |
| 1.13E-02 | 4.64E-02 | 4.93E-04 | 1.66E-03 |
| 1.23E-02 | 4.97E-02 | 4.93E-04 | 1.66E-03 |
| 6.03E-03 | 1.14E-02 | 2.62E-04 | 4.06E-04 |
| 6.56E-03 | 1.22E-02 | 2.62E-04 | 4.06E-04 |
| 6.03E-03 | 1.14E-02 | 2.62E-04 | 4.06E-04 |
| 6.56E-03 | 1.22E-02 | 2.62E-04 | 4.06E-04 |
| 7.25E-03 | 5.78E-03 | 3.15E-04 | 2.06E-04 |
| 7.88E-03 | 6.19E-03 | 3.15E-04 | 2.06E-04 |
| 4.48E-03 | 2.07E-02 | 1.95E-04 | 7.38E-04 |
| 4.87E-03 | 2.22E-02 | 1.95E-04 | 7.38E-04 |
| 4.12E-02 | 4.64E-02 | 1.79E-03 | 1.66E-03 |
| 4.48E-02 | 4.97E-02 | 1.79E-03 | 1.66E-03 |
| 4.12E-02 | 4.64E-02 | 1.79E-03 | 1.66E-03 |
| 4.48E-02 | 4.97E-02 | 1.79E-03 | 1.66E-03 |
| 5.25E-02 | 2.07E-02 | 2.28E-03 | 7.38E-04 |
| 5.70E-02 | 2.22E-02 | 2.28E-03 | 7.38E-04 |
| 2.81E-02 | 9.78E-02 | 1.22E-03 | 3.49E-03 |
| 3.05E-02 | 1.05E-01 | 1.22E-03 | 3.49E-03 |

particular $\alpha$-cut being analyzed. Here, the minimum and maximum values of the divisions for the $\alpha = 0.5$ are [23.41, 29.74]. These steps are repeated for each user-

defined α-cut. Once all α-cuts have been evaluated, the membership function can be constructed from these results (Figure 4.4).

Equations 4.1, 4.2, and 4.3 may be rewritten to accommodate the Vertex Method and fuzzy numbers. Taking a step back from Eqn. 4.1 and considering the distance metric itself yields:

$$\tilde{D}_{ijk} = (\tilde{w}_{ij} - \tilde{x}_i)^2, \qquad \text{for } k = 1, 2, 3, 4. \tag{4.7}$$

Then the *Pattern Unit* nodes may be constructed by summing all possible combinations of $\tilde{D}_{ijk}$ entering node $j$ and taking the square root. This results in $m = 2^n$ possible $I_j^m$,



**Figure 4.4: Example final membership function output from the fuzzy GRNN.**

78

where $n$ is the number of fuzzy numbers associated with the node, (i.e. in the above example, each $I_j$ node had 4 fuzzy numbers attached to it (2 input nodes and 2 weights), hence the $m = 16$ possible values for each $I_j$).

Nodes A and B use the output from $I_j^m$ to calculate the possible output values:

$$\tilde{y}_p(\tilde{X}) = \frac{\tilde{A}_p}{\tilde{B}_p}, \quad \text{where } p = (2^{n+1})^j. \tag{4.8}$$

The values of $\tilde{A}_p$ and $\tilde{B}_p$ are calculated by taking the product of the pattern weights and the $f(I_j^m)$ and then summing all possible combinations attached to $\tilde{A}_p$ and $\tilde{B}_p$.

## 4.3 Example Application:  Predicting RHA score

This section introduces an example application to test the fuzzy GRNN algorithm and briefly discuss the results.  Here, fuzzy RGA scores are used as inputs to predict RHA scores.  Instead of using two input nodes, like in the previous example, now only one input node, the total RGA score, is used.  Three training patterns (Table 4.8) were selected from the 20 reaches on the Lewis Creek that had both RGA and RHA scores.  One reach from each of the habitat conditions (*fair*, *good*, and *reference*) present in the Lewis Creek have been selected.  An imprecision of ±4 points was added to the total

**Table 4.8: Subset of Lewis Creek reaches used for demonstrating the fuzzy GRNN training, only center values of fuzzy number are shown.**

| Reach | RGA Score (Vertex of Fuzzy Number), ±4 | RHA Score (Vertex of Fuzzy Number), ±10 | Habitat Condition |
|---|---|---|---|
| M21A | 46 | 125 | fair |
| M20B | 38 | 143 | good |
| M07 | 71 | 175 | reference |

RGA score (training weights) and an imprecision of ±10 was added to the associated total

RHA score obtained from the VTANR database

([https://anrnode.anr.state.vt.us/SGA/default.aspx](https://anrnode.anr.state.vt.us/SGA/default.aspx)) to represent what an expert might

experience in the field.  As a result, the fuzzy numbers created for this example are all

symmetrical triangular membership functions.  The same assumptions hold for the

prediction data set; six reaches were randomly selected from the Middlebury River (Table

4.9), two from each of the habitat conditions present (*fair*, *good*, and *reference*).

The fuzzy GRNN predictions (Figure 4.5 (asterisks)) are plotted against the expected

values (solid triangles).  The three Middlebury River predictions that best match the

expected values are M13A, T3.03, and T3.10.  Comparing the RGA and RHA scores for

these three reaches to the Lewis Creek training patterns (Table 4.8) shows M13A best

matches Lewis Creek M20B; T3.03 and T3.10 matches Lewis Creek M07.  As expected,

when the RGA training weight is similar to a prediction input, but the associated pattern

weight is not similar to the real RHA score, then the fit is not as good (e.g. Figure 4.5,

Reach M13B).

**Table 4.9:  Middlebury River prediction data set, only center values of fuzzy number are shown.**

| Reach | RGA Score (Vertex of Fuzzy Number), ±4 | RHA Score (Vertex of Fuzzy Number), ±10 | Habitat Condition |
|-------|----------------------------------------|-----------------------------------------|-------------------|
| **M01** | 57 | 131 | good |
| **M04** | 39 | 127 | fair |
| **M13A** | 38 | 138 | good |
| **M13B** | 46 | 110 | fair |
| **T3.03** | 70 | 170 | reference |
| **T3.10** | 72 | 176 | reference |

**Figure 4.5: Predictions for six reaches in the Middlebury River using the fuzzy GRNN.**

This example is a first step in demonstrating and testing the applicability of the fuzzy GRNN algorithm. Ideally, one would use a training set larger than 3 patterns, but this example was limited due to the computational demand associated with the number of calculations necessary to predict a single $\alpha$-cut. Doubling the number of training patterns in this example (from 3 to 6) increases the number of calculations in the division step alone from $2^9$ to $2^{18}$. Future work includes a larger application using expert defined membership functions and operating the code on a faster computer system (e.g. IBM Bluemoon cluster machine located at the Vermont Advanced Computing Center).

## 4.4 References

Besaw, L.E., Rizzo, D.M., Kline, M., Underwood, K.L., Doris, J.J., Morrissey L.A., and Pelletier, K., 2009. Stream classification using hierarchical artificial neural networks: A fluvial hazard management tool. Journal of Hydrology, 373(1-2): 34-43.

Clark, J. S., Rizzo, D. M., Watzin, M. C., and Hession, W. C., 2008. Spatial distribution and geomorphic condition of fish habitat in streams: An analysis using hydraulic modeling and geostatistics. River Research and Applications, 24: 885-899.

Dong, W. M., and Shah, H.C., 1987. Vertex method for computing functions of fuzzy Variables. Fuzzy Sets and Systems, 24(1): 65-78.

Dong, W. M., Shah, H. C., and Wong, F.S., 1985. Fuzzy computations in risk and decision analysis. Civil Engineering and Environmental Systems, 2(4): 201-208.

Doris, J.J., 2006. Master's Thesis: Application of Counterpropagation Networks to Problems in Civil and Environmental Engineering, University of Vermont, Burlington, VT, 139 pp.

Klir, G. and Yuan, B., 1995. Fuzzy Sets and Fuzzy Logic: Theory and Applications. Prentice Hall, 592 pp.

Kline, M., Alexander, C., Pytlik, S., Jaquith, S., and Pomeroy, S., 2007. Vermont Stream Geomorphic Assessment Protocol Hand-books. Vermont Agency of Natural Resources, Waterbury, Vermont. http://www.vtwaterquality.org/rivers/htm.

MATLAB version 7.10.0 (R2010a). Natick, Massachusetts: The MathWorks Inc., 2010.

Ross, T.J., 2004. Fuzzy logic with engineering applications. Second edition. Chichester, England: John Wiley & Sons, Chichester, 650 pp.

Schiff, R., Kline, M., and Clark, J., 2008. The Reach Habitat Assessment Protocol. Prepared by Milone and MacBroom, Inc. for the Vermont Agency of Natural Resources, Waterbury, Vermont.

Specht, D.F., 1991. A general regression neural network. IEEE Transactions on Neural Networks, 2(6): 568-576.

Sullivan, S.M.P, Watzin, M. C., and Hession, W.C., 2004. Understanding stream geomorphic state in relation to ecological integrity: evidence using habitat assessments and macroinvertebrates. Environmental Management, 34(5), 669-683.

Zadeh, L.A., 1975.  The concept of a linguistic variable and its application to approximate reasoning I, II, III.  Information Sciences, 8: 199-251, 301-357; 9: 43-80.

# CHAPTER 5
# ASSESSING LINKAGES BETWEEN STREAM GEOMORPHIC CONDITION AND HABITAT HEALTH USING A GENERALIZED REGRESSION NEURAL NETWORK

## 5.1 Abstract

Using physical geomorphic and habitat assessments to assist watershed management decisions regarding the biological health of a stream could help reduce cost and time to identify stream reaches that are most in need of management help. However, the complex linkages between the physical geomorphic and habitat conditions, and the biological health of stream reaches are not fully understood. In this study, a generalized regression neural network (GRNN) is used to explore these nonlinear relationships using Vermont streams as a model system. The GRNN was first used to examine correlations between Vermont Agency of Natural Resources (VTANR) River Management Program's legacy rapid habitat assessment (LRHA) scores from rapid geomorphic assessments (RGA) and channel evolution stage data. The GRNN, trained with 50% of the data set, was able to correctly predict 69.9% of the remaining 50% of the (testing) data set supporting its use as a tool to further explore relationships involving these variables. Fish and macroinvertebrate biological health assessment data, collected independently by the Biomonitoring and Aquatic Studied Section, were then investigated as input data (in combination with RGA and channel evolution stage) to predict LRHA. In another analysis, the biological health was used as the *output* of the GRNN. The prediction rates were better for fish than macroinvertebrate data in both cases; however, when the GRNN

was used to predict the biological health, the accuracy of prediction was significantly less than when the GRNN was used to predict LRHA. For the fish data, the prediction dropped from a 95.7% match (when predicting LRHA) to 48% (when predicting health) and for the macroinvertebrate data the drop was from 82.1% to 23.2%. A preliminary study was conducted using VTANR's "new" RHA protocol scores, which began in 2008. There was no clear improvement in the prediction rates involving biological health data; however, the datasets, to date, are not large enough to be truly representative, and further study is warranted. Ideally, a study involving both the physical and biological assessments conducted concurrently could provide a better understanding of the mechanisms and complex relationships among them.

## 5.2 Introduction

Identifying streams with high environmental risk and fluvial hazard is essential for a proactive adaptive watershed management approach. Such efforts require environmental managers to gather and assess various forms of information - quantitative, qualitative and subjective. The Vermont Agency of Natural Resources (VTANR) River Management Program (RMP) has developed and adopted protocols for physical stream geomorphic (Kline *et al*., 2007) and habitat assessments (Schiff *et al*., 2008) throughout the state of Vermont. Since physical stream processes form the habitat, habitat assessments study physical ecological parameters needed to understand the relationship between fluvial processes and aquatic communities (VTANR, 2008). From a management viewpoint, these geomorphic and habitat assessments, taken together, may be used to identify problem areas and the steps necessary for mitigation (Kline, 2007).

Separate from the VTANR River Management Program's habitat assessments, the Vermont Biomonitoring and Aquatic Studies Section (BASS) is responsible for monitoring the biological communities in streams. Ideally, in cases where the biological findings are unexpected, the hope is that the physical geomorphic and habitat reach assessments may be used to help understand the findings. In this work, a least-squares regression artificial neural network originally developed by Specht (1991), known as the generalized regression neural network (GRNN), is used to explore the nonlinear interactions between the physical geomorphic and habitat conditions, and the biological metrics collected at the reach-scale to assist watershed managers in making informed decisions. The GRNN, in particular, is an appropriate tool since: (1) the algorithm approximates complex, nonlinear relationships, (2) the method is data-driven thus allowing for continual updates and refinements as understanding/condition of fluvial geomorphology evolves, (3) large quantities of data can easily be passed through the algorithm, (4) its least-squares regression methodology is familiar, and (5) unlike more well-known regression methods, there is no need to know the best-fit polynomial (e.g. linear, quadratic, cubic) prior to data analysis, enabling a truly adaptive management approach.

## 5.3 Background

Over the past two centuries, human impacts (e.g. deforestation, channel straightening, urbanization) have greatly altered streams in Vermont from their original state (Vermont River Management Program, 2009). The VTANR protocols used to classify stream stability (Rapid Geomorphic Assessment – RGA), were developed from a combination of

classification systems by Rosgen (1994, 1996), Montgomery and Buffington (1997), Schumm (1977), Schumm *et al*. (1984) and Simon and Hupp (1986). Stream habitat health (i.e. the ability of the stream to sustain life) protocols, originally a modified version of the U.S. Environmental Protection Agency's Rapid Bioassessment Protocols, have been in use since 2002. Kline and Cahoon (2010) note that data from geomorphic and habitat assessments spanning a six-year period indicate almost three-quarters of Vermont's streams have lost connection with their historical floodplains. These induced changes likely reduce the abundance and diversity of the natural biota (Allan, 2004).

Several studies have demonstrated a relationship between stream geomorphic condition, physical habitat and biological health (Chessman *et al.*, 2006; Sullivan *et al*., 2004, 2006; Sullivan and Watzin, 2008). However, the complex linkages are not well understood or easily studied and include many factors such as variation in fish, macroinvertebrate, and bird species present, metrics used, and/or spatial and temporal measurement scales (Clark *et al*., 2008; Chessman *et al.*, 2006).

### 5.3.1 Generalized Regression Neural Network (GRNN)

ANNs, in general, are used in pattern classification, pattern completion, function approximation, prediction, optimization, and system control applications among others (Wasserman, 1993). Although more than 95% of ANNs used in environmental engineering applications have used either a feed-forward back-propagation network or a radial basis function neural network (Govindaraju and Ramachandra, 2000), here a GRNN is used to explore linkages between geomorphic conditions, physical habitat, and biological health for the reasons stated in Section 5.2.

The GRNN has extensive applications in the water resources and hydrological fields. Aksoy and Dahamsheh (2009) use a GRNN for forecasting monthly precipitation. Several studies have had success predicting leaf wetness (Chtioui *et al*., 1999a; Chtioui *et al*., 1999b) and evapotranspiration (Kim and Kim, 2008; Kisi, 2008a). Cigizoglu and Alp (2004) found the GRNN to be successful in predicting rainfall runoff and, unlike the radial basis function and multiple linear regression, did not produce negative flow estimations. Several studies found the GRNN outperformed the feed-forward back-propagation network when forecasting intermittent stream (Cigizoglu, 2005a) or monthly stream flow (Cigizoglu, 2005b; Kisi, 2008b). Firat (2008) explored its use in daily stream flow forecasting, while Ng *et al.* (2009) estimated missing observations in extreme daily stream flow records. Turan and Yurdusev (2009) predicted stream flow from measured upstream flow records, while Besaw *et al.* (2009a) used a recurrent GRNN to predict flow in ungauged streams. The GRNN has also been used to estimate daily mean sea level heights (Sertel *et al.*, 2008), to predict water quality as a function of rainfall, surface discharge and nutrient concentration (Kim and Kim, 2007) and to model river sediment transport (Cigizoglu and Alp, 2006; Cobaner *et al*., 2009; Kisi *et al.*, 2008). Wang *et al.* (2009) used the GRNN to model event-based suspended sediment concentration in rivers due to tropical storms given turbidity, water discharge, and suspended sediment concentrations collected in a weir during storm events over a one-year time frame.

## 5.4  Stream Assessment Data

### 5.4.1 Vermont Stream Geomorphic and Habitat Assessments

The VTANR developed a three-phase system to perform stream geomorphic assessments.  Each successive phase is more detailed and improves the assessor's certainty about the condition of the reach.  The first phase, remote sensing, uses data obtained from topographic maps, aerial photos, previous studies, and from very limited field studies.  This type of reach assessment is considered provisional, enabling large watersheds (100-150 square miles) to be assessed in a few months.  Using Phase 1 assessments, ~35% or 8,279 of Vermont's ~23,000 stream miles have been assessed to date (Kline and Cahoon, 2010).

The Phase 2, or the rapid field assessment phase, includes the RGA and reach habitat assessment (RHA, habitat assessments prior to 2008 are denoted in this work as legacy rapid habitat assessments – LRHA) where field data are collected at the stream reach or sub-reach scale.  A one-mile reach requires 1 to 2 days to assess; and to date, 6% or 1,371 stream miles (~2,500 stream reaches) have been assessed at the Phase 2 level (Kline and Cahoon, 2010). The geomorphic condition, physical habitat condition, adjustment processes, reach sensitivity, and channel evolution stage are determined from quantitative and qualitative field evaluation of erosion and depositional processes, changes in geometry, and riparian land use/land cover.  Phase 2 assessments identify "at risk" reaches and allow reaches to be flagged for protection, restoration, or further Phase 3 assessment.

Phase 3, the survey-level field assessment phase, requires detailed field measurements at the sub-reach scale that allow for stream types and adjustment processes to be further documented and confirmed. Quantitaive measurements of channel dimension, pattern, profile, and sediments are measured during this level of assessment. Phase 3 assessments require 3 to 4 days on average to survey a sub-reach of two meander wavelengths.

Data used in this study was obtained from VT Department of Environmental Conservation (DEC) and is available at https://anrnode.anr.state.vt.us/SGA/default.aspx. All Phase 2 assessments, quality assured by the River Management Program as of August 2009, that had RGA, LRHA, and channel evolution stage data were selected resulting in 1292 reaches (Figure 5.1).

### 5.4.1.1 VTANR Rapid Geomorphic Assessment (RGA)

The assessed stream reach condition is based on its perceived departure from reference condition. Reference condition for each reach is inferred based on watershed zone, confinement, and valley slope (from Phase 1), as well as, entrenchment, width/depth ratio, sinuosity, channel slope, substrate d50, and bed form collected during the Phase 2 assessment (Kline *et al*., 2007). Quantification of the adjustment processes involves assigning a score between 0 (poor) and 20 (reference) for each of the four adjustment processes (degradation, aggradation, widening and planform change) resulting in a summed total RGA score ranging from 0 to 80. The overall score is used to classify the stream reach as poor, fair, good, or reference condition.

**Figure 5.1: Map of the state of Vermont showing the Phase 2 reach locations used in this study. Note: only 1006 of the 1292 reaches used here are plotted since the remaining reaches were not part of the GIS database at the time this map was created.**

### 5.4.1.2 VTANR habitat assessment

Stream habitat assessments examine the physical processes that are key in determining aquatic habitat and hence the biota that inhabit it. These data complement biological data and may indicate problems with the biotic health in the reach where the biological data alone cannot explain the cause (Schiff *et al*., 2008).

Vermont's legacy rapid habitat assessments (LRHAs) are slightly modified versions of the EPA's Rapid Bioassessment Protocols (Barbour *et al*., 1999). The LRHAs comprise ten parameters that explore physical properties of the channel bed, bank, and riparian vegetation (Table 5.1). Each parameter is scored between 0 (*poor*) and 20 (*excellent*) and then summed to obtain a total score (no greater than 200) categorizing the reach as *poor*, *fair*, *good*, or *reference*. The LRHAs, implemented through 2007, were replaced in 2008 with new reach habitat assessment (RHA) protocols.

The new RHA was developed to allow for more specific assessment of the various stream types found in Vermont and more precise evaluation of the key ecological attributes and requirements for aquatic life. For example, while the LRHA categorized a stream as either low or high gradient, the new RHA allows the assessor to select a form from 1 of 5 possible stream habitat types: cascade, step-pool, plane bed, riffle-pool, or dune-ripple. The RHA uses only eight parameters (Table 5.1); although like the LRHA, each component is scored between 0 to 20 and the total score is used again to categorize the stream reach into *poor*, *fair*, *good*, or *reference*.

Since the RHA protocol was first implemented in 2008, LRHA data were used to show proof of concept due to the availability of data. The histogram of LRHA scores for

**Table 5.1:  Parameters that comprise the Vermont RGAs, LRHAs, and RHAs.**

| | Parameters (20 points each) | Condition (Based on total assessment score) | | | |
|---|---|---|---|---|---|
| | | Poor | Fair | Good | Reference |
| **RGA** | 1. Degradation<br>2. Aggradation<br>3. Widening<br>4. Planform Change | 0 - 27 | 28 -51 | 52 - 67 | 68 - 80 |
| **LRHA** | 1. Epifaunal Substrate/ Available Cover<br>2. Embeddedness or Pool Substrate<br>3. Velocity/Depth Patterns or Pool Variability<br>4. Sediment Deposition<br>5. Channel Flow Status<br>6. Channel Alteration<br>7. Frequency of Riffles/Steps or Channel Sinuosity<br>8. Bank Stability (score each bank)<br>9. Bank Vegetative Protection (score each bank)<br>10. Riparian Vegetative Zone Width   (score each side of channel | 0 - 68 | 69 - 128 | 129 - 168 | 169 - 200 |
| **RHA** | 1. Woody Debris Cover<br>2. Bed Substrate Cover<br>3. Scour and Depositional Features<br>4. Channel Morphology<br>5. Hydrologic Characteristics<br>6. Connectivity<br>7. River Banks<br>8. Riparian  Area | 0 - 55 | 56 - 103 | 104 - 135 | 136 - 160 |

the 1292 reaches used in this study is normally distributed (Figure 5.2, $p < 0.0579$ with a Shapiro-Wilkes W test of $W = 0.9976$), with most of the reach scores falling into *fair* or *good* habitat condition.

## 5.4.2  Biological Assessments

The biological health of Vermont streams and rivers is determined by protocols set forth by the Vermont Department of Environmental Conservation (VTDEC), within the

**Figure 5.2: Histogram of Legacy Rapid Habitat Assessment scores for the 1292 reaches used in this study.**

VTANR. Metric assessments of fish and macroinvertebrate assemblages are used to classify streams based on their departure from reference. To define reference streams for the current biomonitoring protocol, VTDEC Biologists from the Biomonitoring and Aquatic Studies Section selected macroinvertebrate and fish sites that appeared minimally impacted by human activity using data in the VTDEC biological database.

### 5.4.2.1 Macroinvertebrate Health

Combining professional judgment and statistical analyses at the reference sites, four categories for macroinvertebrate communities were identified: Small High Gradient Streams, Medium High Gradient Streams, Warm Water Moderate Gradient Streams and Rivers, and Slow Winders (BASS, 2004). Since few sites fall into the latter category, biocriteria evaluations do not exist at this time for Slow Winders. Currently, eight

metrics are used to assess reaches for macroinvertebrate health (Table 5.2). Table 5.2 also provides the metric thresholds for the three stream types and three macroinvertebrate community categories (Class A1: minimal impacts from human activity, Class B1: minor changes from reference, and Class B2, B3, and A2: moderate change from reference). Biomonitoring and Aquatic Studies Section experts assign rankings such as Excellent, Very Good, and Good to the above categories (Class A1, B1, and B2, B3, and A2, respectively) to capture the stream macroinvertebrate health. If the metrics do not satisfy one of these three criteria, the reach is categorized as "Fair" if there is greater than

**Table 5.2: Threshold values for macroinvertebrate assemblages in Vermont wadeable streams. Adapted from BASS (2004).**

| Class Criteria<br><br>Metric* | Small High Gradient Streams | | | Medium High Gradient Streams | | | Warm Water Moderate Gradient Streams and Rivers | | |
|---|---|---|---|---|---|---|---|---|---|
| | Excellent<br><br>A1 | Very Good<br><br>B1 | Good<br>A2,<br>B2,B3 | Excellent<br><br>A1 | Very Good<br><br>B1 | Good<br>A2,<br>B2,B3 | Excellent<br><br>A1 | Very Good<br><br>B1 | Good<br>A2,<br>B2,B3 |
| Richness | >35 | >31 | >27 | >43 | >39 | >30 | >40 | >35 | >30 |
| Ephemeroptera, Plecoptera, Trichoptera - EPT Index | >21 | >19 | >16 | >24 | >22 | >18 | >21 | >19 | >16 |
| Percent Model Affinity of Orders - PMA-O | >65 | >55 | >45 | >65 | >55 | >45 | >65 | >55 | >45 |
| Hilsenhoff Biotic Index - BI | <3.00 | <3.50 | <4.50 | <3.50 | <4.00 | <5.00 | <4.25 | <4.75 | <5.40 |
| % Oligochaeta | <2 | <5 | <12 | <2 | <5 | <12 | <2 | <5 | <12 |
| EPT/EPT+ Chironomidae | >0.65 | >0.55 | >0.40 | >0.65 | >0.55 | >0.40 | >0.65 | >0.55 | >0.40 |
| Pinkham-Pearson Coefficient of Similarity – Functional Groups - PPCS-FG | >0.50 | >0.45 | >0.40 | >0.50 | >0.45 | >0.40 | >0.50 | >0.45 | >0.40 |
| Density | >500 | >400 | >300 | >500 | >400 | >300 | >500 | >400 | >300 |

\* Metric details can be found at http://www.vtwaterquality.org/bass/docs/bs_wadeablestream1a.pdf.

moderate change from reference, or "Poor" if there is extreme change. Reach condition metric values falling on the threshold are hyphenated (e.g. Excellent-Very Good, Very Good-Good, Good-Fair, and Fair-Poor).

### 5.4.2.2 Fish Health

Fish community health is currently assessed using two Vermont calibrated Indices of Biotic Integrity (IBI) (BASS, 2004). The mixed-water IBI (MW IBI) designation is applied to any stream containing five or more native fish species and is comprised of nine metrics ranging from a total score of 9 to 45 (Table 5.3). The second index, the Coldwater IBI (CW IBI), applies to smaller coldwater streams that contain two to four native species and has six metrics (Table 5.3).

Biological and geomorphic data were not collected at the exact physical location and often not in the same year. Variation in physical location was accommodated by including biological survey data from locations within 200 m of the 1292 locations with Phase 2 assessments in this analysis resulting in 46 reaches for fish data and 133 for macroinvertebrate data. To retain sufficient sample sizes, no data were excluded due to differences in the time of the biological and geomorphic assessments. When biological assessments were performed over multiple years *at the same reach location*, the most recent assessment was used.

**Table 5.3: Fish MWIBI and CWIBI score thresholds for associated Water Quality Classes and Water Management Types. Adapted from BASS (2004).**

| Class Criteria[*]/ Metric | | Excellent | Very Good | Good | Fair | Poor |
|---|---|---|---|---|---|---|
| **Mixed Water Index of Biotic Integrity** | Total number of native species | >41 | >37 | >33 | >25 | <25 |
| | Number of intolerant species | | | | | |
| | Number of benthic insectivore species | | | | | |
| | Percent as white suckers and creek chub | | | | | |
| | Percent as generalist feeders | | | | | |
| | Percent of insectivores | | | | | |
| | Percent as top carnivores | | | | | |
| | Percent with DELT anomalies | | | | | |
| | Abundance | | | | | |
| **Cold Water Index of Biotic Integrity** | Number of intolerant species | >42 | >36 | >33 | >26 | <26 |
| | Percent coldwater species | | | | | |
| | Percent generalist feeders | | | | | |
| | Percent top carnivores | | | | | |
| | Brook trout density | | | | | |
| | Brook trout length class number | | | | | |

[*] Excellent corresponds to Class A1, Very Good to Class B1, and Good to Class B2,B3, and A2.

## 5.5  Methodology

### 5.5.1 Generalized Regression Neural Network

The generalized regression neural network (GRNN) introduced by Donald Specht (1991) is a parallel, one-pass network that does not require training like the more popular feed-forward backpropagation networks (i.e., the training data are used to set the network weights). The GRNN is distinguished from traditional least-squared regression, in that the algorithm does not require *a priori* knowledge of the best-fit polynomial. Figure 5.3

shows the structure of the GRNN algorithm as applied to the prediction of the LRHA scores. The network consists of four nodal layers. The *Input Layer* simply passes the *n* user-defined input variables, $X = \{x_{i=1}, x_{i=2}, ..., x_{i=n}\}$, (equivalent to the independent variables associated with traditional regression techniques) to the weights of the second network layer. The training weights, $w_{ij}$, connect the *Input Layer* to the next layer, the *Pattern Units* layer (e.g. $w_{12}$ connects input node $x_{i=1}$ with pattern unit node $I_{j=2}$, Figure 5.3). Each $j^{th}$ training pattern weight, $w_{ij}$, contains a value (e.g., degradation, aggradation,



**Figure 5.3:   GRNN structure showing the components of the RGA and channel evolution stage as inputs used to predict the total Legacy RHA score.**

widening, planform change, or channel evolution stage) for which there is a corresponding output (LRHA score). These weights are set by the training data and do not update as in other artificial neural network (ANN) algorithms. The corresponding training output (LRHA score) is stored in the pattern weights, $y_j$, associated with node A of the *Summation Units* layer. The *Pattern Units* layer has one node, $I$, for each of the $j$ training patterns and calculates a distance metric (e.g., the Euclidean distance) between all sets of training weights and the current input pattern (Eqn. 5.1):

$$I_j = \sqrt{\sum_{i=1}^{n}(w_{ij} - x_i)^2},$$  (5.1)

where $x_i$ refers to the $i^{\text{th}}$ input parameter, $w_{ij}$ are the $i^{\text{th}}$ input variable associated with the $j^{\text{th}}$ training pattern. The resulting, Euclidean distance, $I_j$ is passed through an exponential activation function (Eqn. 5.2):

$$f(I_j) = \exp\left(\frac{-I_j}{2\sigma^2}\right),$$  (5.2)

where $\sigma$ is a smoothing parameter explained in greater detail below.

The third layer, *Summation Units*, calculates the dot product of the output of the *Pattern Units* (Eqn. 5.2) and, for node A, the corresponding $y_j$ training weights. The pattern weights associated with node B are set equal to 1. Therefore, node B calculates the dot product between the output from the *Pattern Units* and the weights set equal to 1. The final output is the result of dividing the nodes in the *Summation Units* (Eqn. 5.3):

$$\hat{y}(X) = \frac{\sum_j y_j \bullet f(I_j)}{\sum_j 1 \bullet f(I_j)} = \frac{A}{B} \quad .$$  (5.3)

Note that σ, in the $f(I_j)$ term (Eqn. 5.2), is used to optimize the GRNN output and is the only parameter that can be changed. As σ approaches zero, the predicted network output, $\hat{y}$, tends to overfit the training data. When σ is large, $\hat{y}$ is smoothed and assumes the value of the sample mean. For further details the reader is referred to Specht (1991). The GRNN algorithm described in this paper was coded in MATLAB 7.10.0 (R2010a).

## 5.6 Results

Building on previous work by Besaw *et al.* (2009b), the nonlinear relationships between RGA and LRHA were explored using the GRNN. A scatter plot of the 1292 expert-assigned RGA and LRHA scores is shown in Figure 5.4 ($r^2 = 0.414$, $p < 0.05$). The majority of the *poor* habitat ranked reaches aligns with either *poor* or *fair* RGA scores (one exception is a reach with *poor* LRHA and *good* RGA). The *fair* ranked reach habitats overlap all four categories of RGA scores; however, only one (on the dividing line between *good* and *fair* LRHA) falls in the *reference* RGA category. Similarly, the *good* LRHA scores span the entire range of RGA scores with the majority assessed in the *good* and *fair* categories. The LRHA *reference* reaches coincide mostly with the RGA *reference* and *good* reaches; however, one reach is categorized with a *reference* LRHA and a *fair* RGA.

A summary of the GRNN trials conducted in this study to link geomorphology (RGA) to habitat (LRHA and RHA) is provided in Table 5.4. Figure 5.5 (a) shows the comparison of the GRNN predicted LRHA (trial LRHA1, Table 5.4) against the expert assigned LRHA. Fifty percent of reaches from each LRHA category were selected

**Figure 5.4: Correlation between RGA and LRHA scores. The vertical lines mark divisions between categories of poor (0-27), fair (28-51), good (52-67), and reference (68-80) for RGA scores. The dashed horizontal lines show the category endpoints for LRHA scores, poor (0-68), fair (69-128), good (129-168), and reference (169-200).**

randomly to construct the training set. The remaining 50% are used for testing/prediction. Figure 5.5 (b) displays the categorical (total LRHA score post-processed into categories) GRNN predictions (LRHA1) against the categorical expert-assigned LRHA score. The GRNN was able to correctly predict 69.9% (195 misclassified out of 647) of the data in the prediction set compared to a 66.8% match (215 misclassified out of 647) using traditional multiple linear regression. The boxes highlighted along the diagonal show correctly classified predictions. Thirteen stream reaches categorized as *poor* by VTANR experts were categorized as *fair* by the GRNN

**Table 5.4:** **Summary of GRNN trials including inputs, outputs and outcome predicted correctly.**

| Type of Data | Trial ID | GRNN Inputs [*] | GRNN Output | Correctly Classified /Total | % Match |
|---|---|---|---|---|---|
| **Original 1292 Reaches** | LRHA1 | Deg., Agg., Wid, PC, Channel Evolution | Total LRHA | 452/647 | 69.9 |
| | LRHA2 | Deg., Agg., Wid, PC | Total LRHA | 445/647 | 68.6 |
| **Fish Subset of LRHA Data (46 reaches)** | LRHA3 | Deg., Agg., Wid, PC, Channel Evolution | Total LRHA | 22/23 | 95.7 |
| | LRHA4 | Deg., Agg., Wid, PC, Channel Evolution, Fish Health | Total LRHA | 22/23 | 95.7 |
| | FISH1 | Deg., Agg., Wid, PC, Channel Evolution, Total LRHA | Fish Health | 12/25 | 48.0 |
| | FISH2 | Deg., Agg., Wid, PC, Channel Evolution (NO LRHA) | Fish Health | 10/25 | 40.0 |
| **Macro-invertebrate Subset of LRHA Data (133 reaches)** | LRHA5 | Deg., Agg., Wid, PC, Channel Evolution | Total LRHA | 56/67 | 83.6 |
| | LRHA6 | Deg., Agg., Wid, PC, Channel Evolution, Macroinvertebrate Health | Total LRHA | 55/67 | 82.1 |
| | MAC1 | Deg., Agg., Wid, PC, Channel Evolution, Total LRHA | Macro-invertebrate Health | 16/69 | 23.2 |
| | MAC2 | Deg., Agg., Wid, PC, Channel Evolution (NO LRHA) | Macro-invertebrate Health | 15/69 | 21.7 |
| **Fish Subset of New RHA Data (13 reaches)** | RHA1 | Deg., Agg., Wid, PC, Channel Evolution | Total RHA | 5/7 | 71.4 |
| | RHA2 | Deg., Agg., Wid, PC, Channel Evolution, Fish Health | Total RHA | 5/7 | 71.4 |
| | FISH3 | Deg., Agg., Wid, PC, Channel Evolution, Total RHA | Fish Health | 3/6 | 50 |
| | FISH4 | Deg., Agg., Wid, PC, Channel Evolution (NO RHA) | Fish Health | 4/6 | 66.7 |
| **Macro-invertebrate Subset of New RHA Data (36 reaches)** | RHA3 | Deg., Agg., Wid, PC, Channel Evolution | Total RHA | 16/19 | 84.2 |
| | RHA4 | Deg., Agg., Wid, PC, Channel Evolution, Macroinvertebrate Health | Total RHA | 16/19 | 84.2 |
| | MAC3 | Deg., Agg., Wid, PC, Channel Evolution, Total RHA | Macro-invertebrate Health | 7/19 | 26.3 |
| | MAC4 | Deg., Agg., Wid, PC, Channel Evolution (NO RHA) | Macro-invertebrate Health | 4/19 | 21.1 |

[*]Deg. = Degradation; Agg. = Aggradation; Wid. = Widening; PC = Planform Change

**(a)**

| | (b) Expert Assigned RHA Category (69.9% match) | | | |
|---|---|---|---|---|
| | **poor** | **fair** | **good** | **reference** |
| **GRNN Predicted Category**   **poor** | 0 | 2 | 1 | 0 |
| **fair** | 13 | 303 | 98 | 1 |
| **good** | 1 | 58 | 134 | 14 |
| **reference** | 0 | 0 | 7 | 15 |

**Figure 5.5: (a) Results of GRNN predicted LRHA using degradation, aggradation, widening, planform change, and channel evolution stage as inputs to the algorithm (trial LRHA1, Table 5.4) plotted against the expert assigned total RHA score. (b) Frequency of predictions after output is categorized.**

and one reach was estimated as *good*. In addition, only 15 of the *reference* stream reaches were correctly classified; while 14 were predicted as *good* and 1 as *fair*.

Figure 5.6 (a) shows some correlation between the VTANR Biomonitoring and Aquatic Studies Section assigned fish health (plotted along horizontal axes) and the River Management Program LRHA ($r^2 = 0.053$, $p > 0.05$), but less correlation with RGA scores ($r^2 = 0.0002$, $p > 0.05$ - Figure 5.6 (b)). Figures 5.6 (c) and (d) show no obvious trend for

## Fish Health



a)

b)

## Macroinvertebrate Health



c)

d)

**Figure 5.6: Plot showing biological health versus RHA and RGA.  Results for fish at 46 VT stream reaches are shown in (a) and (b). Results for macroinvertebrates at 133 VT stream reaches are shown in (c) and (d).**

macroinvertebrate health plotted against RHA ($r^2 = 0.0004$, p > 0.05) and RGA ($r^2 = 0.0026$, p > 0.05).

In selecting the fish training data, only one *poor* LRHA reach and 2 reference reaches existed in the data set. As a result, this single *poor* reach and one reference reach, were placed into the training set, then 50% of the data in each of the other LRHA categories were randomly selected for training; while the remainder were held back for testing/prediction. The macroinvertebrate data set had one *poor* and one *reference* LRHA reach. Both the *poor* and *reference* reaches were included in the training set along with 50% of each of the other conditions; the remaining reaches were used for prediction.

When considering the fish data set and its relationship to LRHA prediction, the GRNN was able to correctly predict 22 of the 23 reaches (a 95.7% match, Table 5.4, trial LRHA3). Adding the fish health assessment data as a sixth input (Table 5.4, trial LRHA4) did not impact the results. The one misclassified reach was a *fair* reach that the GRNN predicted as *good* (Table 5.5 (a)). For the macroinvertebrate data, when only the geomorphic data was used as inputs to predict the LRHA (Table 5.4, Trial LRHA5), the GRNN classified 56 out of 67 correctly (or 83.6% match). Interestingly, when the macroinvertebrate health assessment data was added as an input (Table 5.4, trial LRHA6), the GRNN correctly classified one less reach (55 out of 67, Table 5.5 (b)).

The trials that are, perhaps, more interesting from a management standpoint, are FISH1, FISH2, MAC1, and MAC2 (Table 5.4) where the GRNN is used to predict biological health for fish and macroinvertebrates, respectively. This is because rapid assessment tools have the potential to identify reaches in need of more detailed fish or macroinvertebrate field assessments. The prediction capabilities of the GRNN are much

105

**Table 5.5: Results of GRNN prediction using (a) fish biological health and (b) macroinvertebrate health as the sixth input parameter.**

|  |  | (a) Expert Assigned RHA Category (95.7% match) | | | |
|---|---|---|---|---|---|
|  |  | poor | fair | good | reference |
| GRNN Predicted Category (Using Fish Data) | poor | 0 | 0 | 0 | 0 |
|  | fair | 0 | 13 | 0 | 0 |
|  | good | 0 | 1 | 8 | 0 |
|  | reference | 0 | 0 | 0 | 1 |

|  |  | (b) Expert Assigned RHA Category (82.1% match) | | | |
|---|---|---|---|---|---|
|  |  | poor | fair | good | reference |
| GRNN Predicted Category (Using Macroinvertebrate Data) | poor | 0 | 0 | 0 | 0 |
|  | fair | 0 | 39 | 6 | 0 |
|  | good | 0 | 5 | 16 | 1 |
|  | reference | 0 | 0 | 0 | 0 |

lower when predicting biological health than when predicting LRHA scores (Table 5.4: 95.7% match for LRHA3 versus 40% match for FISH2 and 83.6% match for LRHA5 versus 21.7% match for MAC2). Including the total LRHA score as an input improved the biological health predictions slightly (Table 5.4: 48% match for FISH1 versus 40% match for FISH2 and 23.2% match for MAC1 versus 21.7% match for MAC2). The prediction rate for fish health is higher (FISH1) than that for macroinvertebrate health (MAC1) as the rate decreases from 48.0% to 23.2%, respectively.

Although it was not a goal of this paper to explore the relationship(s) between the new RHA protocols, RGA, and biological health, the new RHA data available as of July 21, 2010 were used in a preliminary analysis. Since the new RHA protocols were designed to better assess the key ecological parameters that affect habitat in Vermont streams, the hope is that a better correlation will exist between the physical and biological conditions. Eight trials (Table 5.4) show the results of various GRNN predictions.

Note that there are significantly fewer reaches with both biological and physical geomorphic and habitat assessments (n = 13 in the fish subset and n = 36 in the macroinvertebrate subset). Neither subset provided full representation of the possible conditions that exist in Vermont. All RHA scores contained reaches ranked either *fair* or *good*; there were no *poor* or *reference* reaches. Also, the fish health conditions had one Poor, no Fair, 3 Good, 8 Very Good, and one Reference reach. The macroinvertebrate health conditions only represented the Good-Fair, Good, Very Good-Good, Very Good, Excellent-Very Good, and Excellent categories; no Poor, Fair-Poor, or Fair reaches yet exist in the dataset. Training and prediction sets were created in a similar manner as other trials.

Predicting the RHA score using the fish data set (n = 13) produced a 71.4% percent match (Table 5.4, trial RHA1). Adding the fish health assessment for the reaches as an input (Table 5.4, trial RHA2) produced the same results. These are not as strong as predictions using the LRHA data (trials LRHA3 and LRHA4); however, again the addition of the fish health as an input did not improve the prediction rate. When the fish health was predicted, compared to the LRHA dataset (FISH1 and FISH2), the prediction rates are slightly better (FISH3 and FISH4). Interestingly, the GRNN was able to predict one more reach health condition correctly when the RHA score was removed as an input parameter (FISH3 had a match of 50% and FISH4 had a match of 66.7%).

Predictions of the new RHA using macroinvertebrate health (Table 5.4, trials RHA3 and RHA4) were similar to predictions obtained for LRHA (trials LRHA5 and LRHA6). As in the fish case, there was no change in prediction when biological health was added as an input (RHA3 and RHA4). When the GRNN was used to predict the

macroinvertebrate health using the new RHA scores as input parameters (trial MAC3), the prediction rate improved slightly from using the LRHA (MAC1); however, the rates are still much lower than the fish predictions.

## 5.7 Discussion

The results of trial LRHA1 (Figure 5.5 (b)) show that the GRNN was unable to predict *poor* stream reaches. One possible explanation is that of the 14 *poor* LRHA reaches, only 3 of the data patterns are associated with a *poor* RGA score. Therefore, if a prediction input pattern (degradation, aggradation, widening, planform change, and channel evolution) in the LRHA1 trial is similar to a training reach with *fair* RGA and RHA condition, the GRNN output will be *fair*. Another possibility is that since the LRHA is more subjective than the RGA, there is information (in the expert's neural networks) that is currently not being used in the GRNN (e.g. water quality information). Also, the optimal boundaries for the habitat categories were originally selected prior to data collection. Now that VTANR has a large and growing data set, the category boundaries could be optimized. Besaw *et al*. (2009b) showed that VTANR current stream sensitivity classification may need to be adjusted based on analysis of RGA and stream inherent vulnerability.

The lack of strong linear correlations in Figures 5.6 (a) and (b) are not unexpected as the complexities between the physical geomorphic and habitat conditions, and biological health are not completely understood and are compounded by scale incompatibilities, species present, and metrics used (Clark *et al*., 2008; Chessman *et al.*, 2006). Adding fish health as an additional input to the GRNN (trial LRHA4), did not improve the

prediction of the LRHA scores (95.7% for both LRHA3 and LRHA4). Also, in the case of the macroinvertebrates, adding the health as an input (trial LRHA6) actually decreases the prediction rate of the LRHA score (from 83.6% for trial LRHA5 to 82.1% for trial LRHA6). It's possible that this is the result of a smaller sample size than the original data set ($n$ = 46 for fish and 133 for macroinvertebrates) and therefore, not truly representative of the relationships between the physical and biological conditions. It was suggested that weighting embeddedness more heavily for the macroinvertebrate health trial MAC1 might improve the prediction results. This was tested by adding the embeddedness score of the LRHA as an additional input, however, after weighting embeddedness up to six times (making up 50% of the other inputs), the GRNN was only able to correctly predict one more reach than when embeddedness was not included.

In addition, experts with different backgrounds collected the physical (RGA, LRHA, and RHA) and biological health assessments used in this study at separate times (in some cases spanning several years) and at different spatial scales. While it is important for these assessments to be conducted independently to prevent biased results, temporal gaps of several years can result in a loss of information (relationships) that may have existed. Geomorphic reach assessments are conducted with the intent of capturing the best representation of the reach as a whole. Biological assessments tend to be more specific to certain locations within a reach based on sampling preferences. This sampling scale incompatibility may hinder the discovery of linkages between the physical and biological assessments (Clark *et al*., 2008). The fact that the GRNN was able to predict fish health better than macroinvertebrate health (Table 5.4: FISH1, FISH2, FISH3, and FISH4 versus MAC1, MAC2, MAC3, and MAC4) may demonstrate that the fish assessments

are conducted on a scale more similar to that of habitat assessments (LRHA and RHA) than the macroinvertebrate assessments. Another possibility noted by King and Baker (2010) is that some community metrics used to determine the fish and macroinvertebrate biological health (richness, Index of Biotic Integrity used in this study) may allow for a loss of important information. They show the community metrics may be insensitive to changes in individual taxa or populations. Knowing which taxa respond to stressors in the environment can assist in understanding the mechanisms behind the changing habitat and assist managers in making appropriate remediation decisions.

Although no drastic improvement in the prediction rates occurred when the new RHA data were used versus the legacy RHA data in the cases using biological health, given the small sample size and lack of data spanning all categories, no definitive conclusions can be drawn about whether the new RHA captures habitat health better than the LRHA. The results do, however, stimulate curiosity for further study.

## 5.8 Conclusions

The idea that physical habitat conditions would influence the biological health of a stream seems obvious; however, understanding this relationship proves to be a challenging task. The results in this work show that drawing clear linkages in such systems is not obvious. The GRNN, however, does appear to be useful in exploring these complex relationships in Vermont stream reaches. The algorithm is a generalized regression algorithm and as such will provide comparable predictions to traditional generalized regression provided the function the data best fit is known; however, a key advantage of the GRNN is that one does not need to know the order of the best-fit

polynomial *a priori*. For this study, the GRNN was, therefore, easier to implement. The algorithm also allows for continual update and refinement as more data becomes available.

One possible conclusion that can be drawn is that since the fish data have better prediction rates than the macroinvertebrates in almost all the cases studied here, the LRHA and RHA are better at indicating habitat conditions for fish. Ideally, however, a more detailed study with additional physical and biological assessments conducted in tandem may help resolve the complex temporal, spatial, and assessment metric issues.

## 5.9 Acknowledgements

## 5.10 References

Aksoy, H., and Dahamsheh, A., 2009. Artificial neural network models for forecasting monthly precipitation in Jordan. Stochastic Environmental Research and Risk Assessment, 23(7): 917–931.

Allan, J. D., 2004. Landscapes and riverscapes: The influence of land use on stream ecosystems. Annual Review of Ecology, Evolution and Systematics, 35: 257–284.

Barbour, M. T., Gerritsen, J., Snyder, B. D., and Stribling, J. B., 1999. Rapid Bioassessment Protocols for Use in Streams and Wadeable Rivers: Periphyton, Benthic Macroinvertebrates and Fish, Second Edition. EPA 841-B-99-002. U.S. Environmental Protection Agency; Office of Water; Washington, D.C.

BASS (Biomonitoring and Aquatic Studies Section). 2004. Biocriteria for fish and macroinvertebrate assemblages in Vermont wadeable streams and rivers: Executive Summary. Vermont Agency of Natural Resources, Waterbury, VT.

Besaw, L. E., Rizzo, D. M., Bierman, P. R., and Hackett, W., 2009a. Advances in ungauged streamflow prediction using artificial neural networks. Journal of Hydrology, 386(1-4): 27-37.

Besaw, L. E., Rizzo, D. M., Kline, M., Underwood, K. L., Doris, J. J., Morrissey L. A., and Pelletier, K., 2009b. Stream classification using hierarchical artificial neural networks: A fluvial hazard management tool. Journal of Hydrology, 373(1-2): 34-43.

Chessman, B. C., Fryirs, K. A., and Brierley, G. J., 2006. Linking geomorphic character, behaviour and condition to fluvial biodiversity: implications for river management. Aquatic Conservation: Marine and Freshwater Ecosystems, 16: 267–288.

Chtioui, Y., Franci, L., Panigrahi, S., 1999a. Moisture prediction from simple micrometeorological data. Phytopathology, 89(3): 668-672.

Chtioui, Y., Panigrahi, S., Franci, L., 1999b. A generalized regression neural network and its application for leaf wetness prediction to forecast plant disease. Chemometrics and Intelligent Laboratory Systems, 48 (1): 47-58.

Cigizoglu, H. K., 2005a. Application of generalized regression neural networks to intermittent flow forecasting and estimation. Journal of Hydrologic Engineering, 10(4): 336-341.

Cigizoglu, H. K., 2005b. Generalized regression neural network in monthly flow forecasting. Civil Engineering and Environmental Systems, 22(2): 71-84.

Cigizoglu, H. K. and Alp, M., 2004. Rainfall-runoff modelling using three neural network methods. Lecture Notes in Artificial Intelligence, 3070: 166-171.

Cigizoglu, H. K. and Alp, M., 2006. Generalized regression neural network in modelling river sediment yield. Advances in Engineering Software, 37(2): 63-68.

Clark, J. S., Rizzo, D. M., Watzin, M. C., and Hession, W. C., 2008. Spatial distribution and geomorphic condition of fish habitat in streams: An analysis using hydraulic modeling and geostatistics. River Research and Applications, 24: 885-899.

Cobaner, M., Unal, B., Kisi, O., 2009. Suspended sediment concentration estimation by an adaptive neuro-fuzzy and neural network approaches using hydro-meteorological data. Journal of Hydrology, 367 (1-2): 52-61.

Firat, M., 2008. Comparison of artificial intelligence techniques for river flow forecasting. Hydrology and Earth System Sciences, 12(1): 123-139.

Govindaraju, R. S. and Ramachandra, R. A., 2000. Artificial Neural Networks in Hydrology. Kluwer Academic Publishers, Dordrecht, The Netherlands.

Kim, M. Y. and Kim, M. K., 2007. Dynamics of surface runoff and its influence on the water quality using competitive algorithms in artificial neural networks. Journal of

Environmental Science and Health Part A-Toxic/Hazardous Substances & Environmental Engineering, 42(8): 1057-1064.

Kim, S. and Kim, H. S., 2008. Neural networks and genetic algorithm approach for nonlinear evaporation and evapotranspiration modeling. Journal of Hydrology, 351(3-4): 299-317.

King, R. S., and Baker, M. E., 2010. Considerations for analyzing ecological community thresholds in response to anthropogenic environmental gradients. Journal of the North American Benthological Society, 29(3): 998–1008.

Kisi, O., 2008a. The potential of different ANN techniques in evapotranspiration modeling. Hydrological Processes, 22(14): 2449–2460.

Kisi, O., 2008b. River flow forecasting and estimation using different artificial neural network techniques. Hydrology Research, 39(1): 27-40.

Kisi, O., Yuksel, I. and Dogan, E., 2008. Modelling daily suspended sediment of rivers in Turkey using several data-driven techniques. Hydrological Sciences Journal-Journal Des Sciences Hydrologiques, 53(6): 1270-1285.

Kline, M., 2007. Draft Vermont Agency of Natural Resources River Corridor Planning Guide to Identify and Develop River Corridor Protection and Restoration Projects. Vermont River Management Program, Waterbury, VT.

Kline, M. and Cahoon, B., 2010. Protecting river corridors in Vermont. Journal of the American Water Resources Association (JAWRA), 46(2):227-236.

Kline, M., Alexander, C., Pytlik, S., Jaquith, S. and Pomeroy, S., 2007. Vermont Stream Geomorphic Assessment Protocol Handbooks. Vermont Agency of Natural Resources, Waterbury, Vermont. http://www.vtwaterquality.org/rivers/htm.

Montgomery, D. R. and Buffington, J. M., 1997. Channel-reach morphology in mountain drainage basins. Geological Society of America Bulletin, 109(5): 596-611.

Ng, W. W., Panu, U. S., and Lennox, W. C., 2009. Comparative studies in problems of missing extreme daily streamflow records. Journal of Hydrologic Engineering, 14(1): 91-100.

Rosgen, D.L., 1994. A classification of natural rivers. Catena, 22:169-199.

Rosgen, D. L. and Silvey, H. L., 1996. Applied River Morphology. Wildland Hydrology Books, Pagosa Springs, Colorado, USA. 325 pp.

Schiff, R., Kline, M., and Clark, J., 2008. The Reach Habitat Assessment Protocol. Prepared by Milone and MacBroom, Inc. for the Vermont Agency of Natural

Resources, Waterbury, Vermont.

Schumm, S. A., 1977. The Fluvial System. Wiley-Interscience, New York, N.Y., 338 pp.

Schumm, S. A., Harvey, M. D., and Watson, C. C., 1984. Incised Channels: Morphology, Dynamics and Control. Water Resources Publications, Littleton, Co. 200. pp.

Sertel, E., Cigizoglu, H. K. and Sanli, D. U., 2008. Estimating daily mean sea level heights using artificial neural networks. Journal of Coastal Research, 24(3): 727-734.

Simon, A. and Hupp, C. R., 1986. Channel evolution in modified Tennessee channels, Proceedings, Fourth Federal Interagency Sedimentation Conference, Las Vegas, 2(5): 71-82.

Specht, D.F., 1991. A general regression neural network. IEEE Transactions on Neural Networks, 2(6): 568-576.

Sullivan, S. M. P, and Watzin, M. C., 2008. Relating stream physical habitat condition and concordance of biotic productivity across multiple taxa. Canadian Journal of Fisheries and Aquatic Sciences, 65: 2667-2677.

Sullivan, S. M. P, Watzin, M. C., and Hession, W. C., 2004. Understanding stream geomorphic state in relation to ecological integrity: evidence using habitat assessments and macroinvertebrates. Environmental Management, 34(5), 669-683.

Sullivan, S. M. P., Watzin, M. C., and Hession, W. C., 2006. Influence of stream geomorphic condition on fish communities in Vermont, U.S.A. Freshwater Biology, 51: 1811–1826.

Turan, M., Yurdusev, M., 2009. River flow estimation from upstream flow records by artificial intelligence methods. Journal of Hydrology, 369: 71–77.

Wang, Y. M., Kerh, T., Traore, S., 2009. Neural networks approaches for modeling river suspended sediment concentration due to tropical storms. Global NEST Journal, 11(4): 457-466.

Wasserman, P. D., 1993. Advanced Methods in Neural Computing. Van Nostrand Reinhold, New York, N. Y.

Vermont River Management Program, 2009. Data Management System (DMS), Corridor Plans, and ArcGIS Tool: Vermont Stream Geomorphic and Habitat Assessment Data, Map Serve, and River Corridor Plans and Delineation Tools. Vermont Agency of Natural Resources, Waterbury, Vermont. http://www.vtwaterquality.org/rivers/htm.

# CHAPTER 6:
# CONCLUSIONS

In this work, fuzzy set theory and generalized regression neural networks are applied and modified as necessary to address groundwater and watershed management problems. Given the applications in this dissertation, the non-traditional analysis methods used (Dempster-Shafer Theory, fuzzy least-squares regression, and GRNNs) prove to perform as well or better than more traditional methods and warrant consideration for appropriate future applications.

In Chapter 2, uncertainty information from two permeability experts and three measurement techniques are combined using various combination rules under Dempster-Shafer theory. First, measurement uncertainty bounds associated with pump-test, drill-stem, or core data were obtained independently from experts. The uncertainty was applied to the data and combined using Dempster's rule of combination, Yager's rule, and the Hau-Kashyap method. The latter two methods were compared to the previously criticized Dempster's rule of combination. Since the conflict amongst the experts was realtively low, the three methods yield similar results, however it was clear how high levels of conflict could produce results that are not as meaningful.

In Chapter 3, fuzzy least-squares regression was used in place of ordinary linear least-squares regression in the Cooper-Jacob method. A modified version of the fuzzy least-squares regression was created to remove one of the fundamental problems with the existing methods: if crisp numbers were used in the algorithm, the results were not the same as ordinary least-squares regression. Our modified version corrected that issue.

The fuzzy least-squares regression was then used to calculate fuzzy slope and intercept values that were then used in the Cooper-Jacob equation. The Cooper-Jacob equation was solved using the Extension Principle to produce membership functions for storativity and transmissivity. The vertex values of the membership functions compared well to the results of the traditional analysis technique (i.e. using ordinary linear least-squares regression). Using the modified fuzzy least-squares regression to solve for transmissivity and storativity allows for incorporation of uncertainty that is typically not used and, therefore, a better understanding of the heterogeneous subsurface results.

In Chapter 4, the GRNN algorithm was modified to allow for the use of fuzzy numbers as input and training data. The Vertex Method was used to alter the equations in the algorithm to approximate the Extension Principle. The motivation behind the development of the fuzzy GRNN was to capture imprecision in experts assessments of stream reach geomorphic and physical habitat condition in Vermont, while linkages between the two are explored. The fuzzy GRNN algorithm was tested using a small subset of Vermont stream reach physical geomorphic and habitat data. The results are promising in capturing expert imprecision and the ability to better define stream reach habitat condition; however, due to the computational demand of the algorithm, a larger application needs to be conducted on a more powerful computing system to test this theory further.

In Chapter 5, a GRNN was used to explore linkages between physical habitat and geomorphic conditions and biological health using fish and macroinvertebrate assessments throughout the state of Vermont. Initially, a study of 1292 reaches with geomorphic assessments was used to predict habitat conditions (based on legacy habitat

assessments).  The algorithm provides comparable predictions to generalized regression; however, a key advantage of the GRNN is that one does not need to know the order of the best-fit polynomial *a priori*.  For this study, the GRNN was, therefore, easier to implement.  The algorithm also allows for easy manipulation of data as more becomes available and, as more is learned, input parameters can be quickly added or removed and new results obtained.  The results of the GRNN trials support that drawing clear linkages between the systems is not obvious.  The GRNN, however, appears to be viable tool to explore these complex relationships in Vermont stream reaches.  Further study with larger data sets and use of the new habitat protocols are needed to further understand the complex relationships between the physical and biological health conditions.

# REFERENCES

Agarwal, H., Renaud, J.E., Preston, E.L., and Padmanabhan, D., 2004, Uncertainty quantification using evidence theory in multidisciplinary design optimization: Reliability Engineering and System Safety, v. 85, p. 281-294.

Aksoy, H., and Dahamsheh, A., 2009.  Artificial neural network models for forecasting monthly precipitation in Jordan. Stochastic Environmental Research and Risk Assessment, 23(7): 917–931.

Allan, J. D., 2004.  Landscapes and riverscapes: The influence of land use on stream ecosystems.  Annual Review of Ecology, Evolution and Systematics, 35: 257–284.

Asefa, T., Wanakule, N. and Adams, A., 2007. Field-scale application of three types of neural networks to predict ground-water levels. Journal of the American Water Resources Association, 43(5): 1245-1256.

Bailly, K., and Milgram, M., 2009. Head pan angle estimation by a nonlinear regression on selected features. Proceedings of the 2009 16th IEEE International Conference on Image Processing (ICIP 2009)**:** 3589-92.

Barbour, M. T., Gerritsen, J., Snyder, B. D., and Stribling, J. B., 1999. Rapid Bioassessment Protocols for Use in Streams and Wadeable Rivers: Periphyton, Benthic Macroinvertebrates and Fish, Second Edition. EPA 841-B-99-002. U.S. Environmental Protection Agency; Office of Water; Washington, D.C.

Bardossy, A., Bogardi, I. and Duckstein, L., 1990. Fuzzy regression in hydrology. Water Resources Research, 26(7): 1497-1508.

BASS (Biomonitoring and Aquatic Studies Section). 2004.  Biocriteria for fish and macroinvertebrate assemblages in Vermont wadeable streams and rivers: Executive Summary. Vermont Agency of Natural Resources, Waterbury, VT.

Belitz, K. and Bredehoeft, J. D., 1988, Hydrodynamics of Denver Basin - Explanation of subnormal fluid pressures:  American Association of Petroleum Geologists Bulletin, v. 72, no. 11, p. 1334-1359.

Besaw, L. E., Rizzo, D. M., Bierman, P. R., and Hackett, W., 2009a. Advances in ungauged streamflow prediction using artificial neural networks. Journal of Hydrology, 386(1-4): 27-37.

Besaw, L. E., Rizzo, D. M., Kline, M., Underwood, K. L., Doris, J. J., Morrissey L. A., and Pelletier, K., 2009b. Stream classification using hierarchical artificial neural networks: A fluvial hazard management tool.  Journal of Hydrology, 373(1-2): 34-43.

Bi, Y., Bell, D., Wang, H., Guo, G., and Guan, J., 2007. Combining multiple classifiers using Dempster's rule for text categorization. Applied Artificial Intelligence, 21(3): 211-239.

Binaghi, E., Luzi, L., Madella, P., Pergalani, F., and Rampini, A., 1998, Slope instability Zonation: a comparison between certainty factor and fuzzy Dempster-Shafer approaches: Natural Hazards, v. 17, p. 77-97.

Bowden, G.J., Nixon, J.B., Dandy, G.C., Maier, H.R., Holmes, M., 2006. Forecasting chlorine residuals in a water distribution system using a general regression neural network. Mathematical and Computer Modeling, 44(5-6): 469-484.

Carranza, E.J.M., and Hale, M., 2003. Evidential belief functions for data-driven geologically constrained mapping of gold potential, Baguio district, Philippines. Ore Geology Reviews, 22(1-2): 117-132.

Cayuela, L., Golicher, J.D., Salas Rey, J., and Rey Benayas, J.M., 2006, Classification of a complex landscape using Dempster-Shafer theory of evidence: International Journal of Remote Sensing, v. 27, no. 10, p. 1951-1971.

Chang, Y.H.O., 2001. Hybrid fuzzy least-squares regression analysis and its reliability measures. Fuzzy Sets and Systems, 119(2): 225-246.

Chang, Y.H.O. and Ayyub, B.M., 2001. Fuzzy regression methods - a comparative assessment. Fuzzy Sets and Systems, 119(2): 187-203.

Chessman, B. C., Fryirs, K. A., and Brierley, G. J., 2006. Linking geomorphic character, behaviour and condition to fluvial biodiversity: implications for river management. Aquatic Conservation: Marine and Freshwater Ecosystems, 16: 267–288.

Chtioui, Y., Franci, L., Panigrahi, S., 1999a. Moisture prediction from simple micrometeorological data. Phytopathology, 89(3): 668-672.

Chtioui, Y., Panigrahi, S., Franci, L., 1999b. A generalized regression neural network and its application for leaf wetness prediction to forecast plant disease. Chemometrics and Intelligent Laboratory Systems, 48 (1): 47-58.

Cigizoglu, H. K., 2005a. Application of generalized regression neural networks to intermittent flow forecasting and estimation. Journal of Hydrologic Engineering, 10(4): 336-341.

Cigizoglu, H. K., 2005b. Generalized regression neural network in monthly flow forecasting. Civil Engineering and Environmental Systems, 22(2): 71-84.

119

Cigizoglu, H. K. and Alp, M., 2004. Rainfall-runoff modelling using three neural network methods. Lecture Notes in Artificial Intelligence, 3070: 166-171.

Cigizoglu, H. K. and Alp, M., 2006. Generalized regression neural network in modelling river sediment yield. Advances in Engineering Software, 37(2): 63-68.

Clark, J. S., Rizzo, D. M., Watzin, M. C., and Hession, W. C., 2008. Spatial distribution and geomorphic condition of fish habitat in streams: An analysis using hydraulic modeling and geostatistics. River Research and Applications, 24: 885-899.

Cobaner, M., Unal, B., Kisi, O., 2009. Suspended sediment concentration estimation by an adaptive neuro-fuzzy and neural network approaches using hydro-meteorological data. Journal of Hydrology, 367 (1-2): 52-61.

Cooper, H.H. and Jacob, C.E., 1946. A generalized graphical method for evaluating formation constants and summarizing well field history. American Geophysical Union Transactions, 27: 526-534.

Dempster, A. P., 1967, Upper and lower probabilities induced by a multivalued mapping: Annals of Mathematical Statistics, v. 38, p. 325-339.

Dong, W. M., and Shah, H.C., 1987. Vertex method for computing functions of fuzzy Variables. Fuzzy Sets and Systems, 24(1): 65-78.

Dong, W. M., Shah, H. C., and Wong, F.S., 1985. Fuzzy computations in risk and decision analysis. Civil Engineering and Environmental Systems, 2(4): 201-208.

Doris, J.J., 2006. Master's Thesis: Application of Counterpropagation Networks to Problems in Civil and Environmental Engineering, University of Vermont, Burlington, VT, 139 pp.

Du, C., Tang, D., Zhou, J., Wang, H., Shaviv, A., 2008. Prediction of nitrate release from polymer-coated fertilizers using an artificial neural network model. Biosystems Engineering, 99(4): 478-486.

Dubois, D., Martin-Clouaire, R. and Prade, H., 1988. Practical computing in fuzzy logic. In: M.M. Gupta and T. Yamakawa (Editors), Fuzzy Computing. North-Holland, Netherlands.

Dubois, D. and Prade, H., 1986, A set-theoretic view on belief functions: logicaloperations and approximations by fuzzy sets: International Journal of General Systems, v. 12, p. 193-226.

Dubois, D. and Prade, H., 1991. Random Sets and Fuzzy Interval-Analysis. Fuzzy Sets and Systems, 42(1): 87-101.

Dubois, D. and Prade, H., 1992, On the combination of evidence in various mathematical frameworks, *in* Flamm , J. and Luisi, T., eds., Reliability Data Collection and Analysis, Kluwer Academic Publishers, Brussels, p. 213-241.

Durdu, O.F., 2009.  Spatial predictions of surface water quality based on general regression neural network:  A case study of the Buyuk Menderes Catchment, Turkey.  Fresenius Environmental Bulletin, 18(9): 1603-1613.

Ferson, S. and Kreinovich, V., 2002, Representation, Propagation, and Aggregation of Uncertainty: Sandia National Laboratories, Technical Report, Albuquerque, New Mexico.

Ferson, S., Kreinovich, V., Ginzburg, L., Myers, D.S., and Sentz, K., 2002, Constructing Probability Boxes and Dempster-Shafer Structures: Sandia National Laboratories, Technical Report SAND2002-4015, Albuquerque, New Mexico, Available at: http://www.sandia.gov/epistemic/Reports/SAND2002-4015.pdf.

Firat, M., 2008.  Comparison of artificial intelligence techniques for river flow forecasting.  Hydrology and Earth System Sciences, 12(1): 123-139.

Firat, M., Yurdusev, M.A. and Turan, M.E., 2009. Evaluation of Artificial Neural Network Techniques for Municipal Water Consumption Modeling. Water Resources Management, 23(4): 617-632.

Gharavol, E.A., Khademi, M., and Akbarzadeh-T., M.-R., 2007. A new variable bit rate (VBR) video traffic model based on fuzzy systems implemented using generalized regression neural network (GRNN). 2006 IEEE International Conference on Fuzzy Systems**:** 2142-2148.

Goulermas, J.Y., Zeng, X.J., Liatsis, P., and Ralph, J.F., 2007. Generalized regression neural networks with multiple-bandwidth sharing and hybrid optimization.  IEEE Transactions on Systems Man and Cybernetics Part B-Cybernetics, 37(6): 1434-1445.

Govindaraju, R. S. and Ramachandra, R. A., 2000.  Artificial Neural Networks in Hydrology.  Kluwer Academic Publishers, Dordrecht, The Netherlands.

Hau, H.Y. and Kashyap, R.L., 1990, Belief combination and propagation in a lattice-structured inference network:  IEEE Trans. on Systems, Man, and Cybernetics, v. 20, no. 1, p. 45-57.

Husain, H., Khalid, M., and Yusof, R., 2004.  Automatic clustering of generalized regression neural network by similarity index based fuzzy c-means clustering. TENCON 2004.  2004 IEEE Region 10 Conference (IEEE Cat. No. 04CH37582), 2: 302-305.

Inagaki, T., 1991, Interdependence between Safety-Control Policy and Multiple-SensorSchemes Via Dempster-Shafer Theory: IEEE Transactions on Reliability, v. 40, no. 2, p. 182-188.

Joslyn, C. and Booker, J.M., 2004, Generalized Information Theory for Engineering Modeling and simulation, *in* Nikolaidid, E., Ghiocel, D., and Singhal, S., eds., Engineering Design Reliability Handbook, CRC Press, p. 9-1 - 9-40.

Joslyn, C. and Ferson, S., 2004, Approximate representations of random intervals for hybrid uncertainty quantification in engineering modeling, *in* Hanson, K.M., and Hemez, F.M., eds., Sensitivity Analysis of Model Output (SAMO04), LANL, Los Alamos, p. 453-469 http://library.lanl.gov/cgi-bin/getdoc?event=SAMO2004&document=samo04-83.pdf.

Kanevski, M.F., 1999. Spatial predictions of soil contamination using general regression neural networks. International Journal of Systems Research and Information Systems, (8)4: 241-256.

Kim, M., Choi, C.Y. and Gerba, C.R., 2008. Source tracking of microbial intrusion in water systems using artificial neural networks. Water Research, 42(4-5): 1308-1314.

Kim, M. Y. and Kim, M. K., 2007. Dynamics of surface runoff and its influence on the water quality using competitive algorithms in artificial neural networks. Journal of Environmental Science and Health Part A-Toxic/Hazardous Substances & Environmental Engineering, 42(8): 1057-1064.

Kim, S. and Kim, H. S., 2008. Neural networks and genetic algorithm approach for nonlinear evaporation and evapotranspiration modeling. Journal of Hydrology, 351(3-4): 299-317.

King, R. S., and Baker, M. E., 2010. Considerations for analyzing ecological community thresholds in response to anthropogenic environmental gradients. Journal of the North American Benthological Society, 29(3): 998–1008.

Kisi, O., 2008a. The potential of different ANN techniques in evapotranspiration modeling. Hydrological Processes, 22(14): 2449–2460.

Kisi, O., 2008b. River flow forecasting and estimation using different artificial neural network techniques. Hydrology Research, 39(1): 27-40.

Kisi, O., Yuksel, I. and Dogan, E., 2008. Modelling daily suspended sediment of rivers in Turkey using several data-driven techniques. Hydrological Sciences Journal-Journal Des Sciences Hydrologiques, 53(6): 1270-1285.

Kline, M., 2007.  Draft Vermont Agency of Natural Resources River Corridor Planning Guide to Identify and Develop River Corridor Protection and Restoration Projects. Vermont River Management Program, Waterbury, VT.

Kline, M. and Cahoon, B., 2010. Protecting river corridors in Vermont.  Journal of the American Water Resources Association (JAWRA), 46(2):227-236.

Kline, M., Alexander, C., Pytlik, S., Jaquith, S. and Pomeroy, S., 2007. Vermont Stream Geomorphic Assessment Protocol Handbooks. Vermont Agency of Natural Resources, Waterbury, Vermont. http://www.vtwaterquality.org/rivers/htm.

Klir, G.J., 2003, Uncertainty:  Encyclopedia of Information Systems, v. 4, p. 511-521.

Klir, G. and Yuan, B., 1995.  Fuzzy Sets and Fuzzy Logic: Theory and Applications. Prentice Hall, 592 pp.

Kriegler. E., and Held, H., 2005, Utilizing belief functions for the estimation of future climate change: International Journal of Approximate Reasoning, v. 39, p. 185-209.

Kumara, S.S.P., Chin, H.C., and Weerakoon, W.M.S.B., 2003.  Identification of accident causal factors and prediction of hazardousness of intersection approaches. Statistical Methods and Modeling and Safety Data, Analysis, and Evaluation-Safety and Human Performance, Transportation Research Record, 1840: 116-122.

Lee E.W.M., Lim, C.P., Yuen, R.K.K., and Lo, S.M., 2004.  A hybrid neural network model for noisy data regression. IEEE Transactions on Systems Man and Cybernetics Part B-Cybernetics, 34(2): 951-960.

Lee E.W.M., Lee, Y.Y., Lim, C.P., and Tang, C.Y., 2006. Application of a noisy data classification technique to determine the occurrence of flashover in compartment tires.  Advanced Engineering Informatics, 20(2): 213-222.

Li, J., and Fenli, L., 2008. Applying general regression neural network in digital image watermarking. 2008 Fourth International Conference on Natural Computation (ICNC), 5: 452-6.

Li, X., Smith, D.W., Prepas, E.E., 2008. Artificial neural network modelling of nitrogen in streams: with emphasis on accessible databases. Annual Conference of the Canadian Society for Civil Engineering 2008. "Partnership for Innovation." 2: 793-796.

Ligang, Z., Shuijun, Y. and Minggao, Y., 2008. Monitoring NOx emissions from coal-fired boilers using generalized regression neural network. 2008 2nd International Conference on Bioinformatics and Biomedical Engineering (ICBBE '08): 1916-1919.

Lu, W.B., and Caselton, B., 1997. Using Dempster-Shafer Theory to represent climate change uncertainties. Journal of Environmental Management, 49: 73-93.

Mathon, B.M., Ozbek, M.M., and Pinder, G.F., 2008. Transmissivity and storage coefficient estimation by coupling the Cooper–Jacob method and modified fuzzy least-squares regression. Journal of Hydrology, 353(3-4): 267-274.

Mathon, B.M., Ozbek, M.M., and Pinder, G.F., 2010. Demspter-Shafer Theory applied to uncertainty surrounding permeability. Mathematical Geosciences, 42(3): 293-307.

MATLAB version 7.10.0 (R2010a). Natick, Massachusetts: The MathWorks Inc., 2010.

Montgomery, D. R. and Buffington, J. M., 1997. Channel-reach morphology in mountain drainage basins. Geological Society of America Bulletin, 109(5): 596-611.

Ng, W. W., Panu, U. S., and Lennox, W. C., 2009. Comparative studies in problems of missing extreme daily streamflow records. Journal of Hydrologic Engineering, 14(1): 91-100.

Ozelkan, E.C. and Duckstein, L., 2001. Fuzzy conceptual rainfall-runoff models. Journal of Hydrology, 253(1-4): 41-68.

Ravanbod, H., 2005. Application of neuro-fuzzy techniques in oil pipeline ultrasonic nondestructive testing. NDT&E International, 38(8): 643-653.

Ricciardi, K.L., 2002, Optimal groundwater remediation design subject to uncertainty: Ph.D dissertation, University of Vermont, USA, p. 50-66.

Rosgen, D.L., 1994. A classification of natural rivers. Catena, 22:169-199.

Rosgen, D. L. and Silvey, H. L., 1996. Applied River Morphology. Wildland Hydrology Books, Pagosa Springs, Colorado, USA. 325 pp.

Ross, J., Ozbek, M., and Pinder, G.F., 2008, Kalman filter updating of possibilistic hydraulic conductivity: Journal of Hydrology, v. 354, p. 149-159.

Ross, T.J., 2004. Fuzzy logic with engineering applications. Second edition. Chichester, England: John Wiley & Sons, Chichester, 650 pp.

Savic, D.A. and Pedrycz, W., 1991. Evaluation of fuzzy linear-regression models. Fuzzy Sets and Systems, 39(1): 51-63.

Schiff, R., Kline, M., and Clark, J., 2008. The Reach Habitat Assessment Protocol. Prepared by Milone and MacBroom, Inc. for the Vermont Agency of Natural Resources, Waterbury, Vermont.

Schumm, S. A., 1977. The Fluvial System. Wiley-Interscience, New York, N.Y., 338 pp.

Schumm, S. A., Harvey, M. D., and Watson, C. C., 1984. Incised Channels: Morphology, Dynamics and Control. Water Resources Publications, Littleton, Co. 200. pp.

Seng, T.L., Khalid, M., Yusof, R., Omatu, S., 1998. Adaptive neuro-fuzzy control system by RBF and GRNN neural networks. Journal of Intelligent & Robotic Systems, 23(2-4): 267-289.

Sentz, K. and Ferson, S., 2002, Combination of Evidence in Dempster-Shafer Theory: Sandia National Laboratories, Technical Report SAND2002-0835, Albuquerque, New Mexico, Available at: http://www.sandia.gov/epistemic/Reports/SAND2002-0835.pdf.

Sertel, E., Cigizoglu, H. K. and Sanli, D. U., 2008. Estimating daily mean sea level heights using artificial neural networks. Journal of Coastal Research, 24(3): 727-734.

Shafer, G., 1976, A mathematical theory of evidence: Princeton University Press, Princeton, New Jersey, 312 p.

Si, B.C. and Bodhinayake, W., 2005. Determining soil hydraulic properties from tension infiltrometer measurements: Fuzzy regression. Soil Science Society of America Journal, 69(6): 1922-1930.

Simon, A. and Hupp, C. R., 1986. Channel evolution in modified Tennessee channels, Proceedings, Fourth Federal Interagency Sedimentation Conference, Las Vegas, 2(5): 71-82.

Singh, T.N., Sinha, S., and Singh, V.K., 2007. Prediction of thermal conductivity of rock through physico-mechanical properties. Building and Environment, 42(1): 146-155.

Smarandache, F., 2004, An in-depth look at information fusion rules and the unification of fusion theories: arXiv electronic archives, available at: http://xxx.lanl.gov/ftp/cs/papers/0410/0410033.pdf.

Smarandache, F. and Dezert, J., eds., 2004, Applications and Advances of DSmT for Information Fusion: American Research Press, Rehoboth, NM: http://www.gallup.unm.edu/~smarandache/DSmT-book1.pdf.

Smets, P., 2005, Analyzing the combination of conflicting belief functions: available at: http://iridia.ulb.ac.be/%7epsmets/Combi_Confl.pdf.

Smets, P. and Kennes, R., 1994, The transferable belief model: Artificial Intelligence, v. 66, p. 191-234.

Specht, D.F., 1991.  A general regression neural network. IEEE Transactions on Neural Networks, 2(6): 568-576.

Sullivan, S. M. P, and Watzin, M. C., 2008.  Relating stream physical habitat condition and concordance of biotic productivity across multiple taxa.  Canadian Journal of Fisheries and Aquatic Sciences, 65: 2667-2677.

Sullivan, S. M. P, Watzin, M. C., and Hession, W. C., 2004.  Understanding stream geomorphic state in relation to ecological integrity: evidence using habitat assessments and macroinvertebrates.  Environmental Management, 34(5): 669-683.

Sullivan, S. M. P., Watzin, M. C., and Hession, W. C., 2006.  Influence of stream geomorphic condition on fish communities in Vermont, U.S.A.  Freshwater Biology, 51: 1811–1826.

Sun, G., Hoff, S.J., Zelle, B.C. and Nelson, M.A., 2008. Development and comparison of backpropagation and generalized regression neural network models to predict diurnal and seasonal gas and PM10 concentrations and emissions from swine buildings. Transactions of the Asabe, 51(2): 685-694.

Sun, W.X., Liang, S.L., Xu, G., Fang, H.L., and Dickinson, R., 2008. Mapping plant functional types from MODS data using multisource evidential reasoning.  Remote Sensing of Environment, 112(3): 1010-1024.

Tanaka, H., Uejima, S. and Asai, K., 1982. Linear-regression analysis with fuzzy model. Ieee Transactions on Systems Man and Cybernetics, 12(6): 903-907.

Tham, C.-W., Tian, S.-H., Ding, L., 2002.  Weather forecasting system based on satellite imageries: using neuro-fuzzy techniques. Advances in Soft Computing - AFSS 2002. 2002 AFSS International Conference on Fuzzy Systems. Proceedings (Lecture Notes in Artificial Intelligence), 2275**:** 267-73.

Turan, M., and Yurdusev, M., 2009.  River flow estimation from upstream flow records by artificial intelligence methods.  Journal of Hydrology, 369: 71–77.

Theis, C.V., 1940. The source of water derived from wells - essential factors controlling the response of an aquifer to development. American Society of Civil Engineers, 10: 277-280.

Uddameri, V., 2004. Relationships of longitudinal dispersivity and scale developed from fuzzy least-squares regressions. Environmental Geology, 45(8): 1172-1178.

Uddameri, V. and Honnungar, V., 2007. Interpreting sustainable yield of an aquifer using a fuzzy framework. Environmental Geology, 51(6): 911-919.

Ustaoglu, B., Cigizoglu, H. K., Karaca, M., 2008. Forecast of daily mean, maximum and minimum temperature time series by three artificial neural network methods. Meteorological Applications, 15: 431 – 445.

Vermont River Management Program, 2009. Data Management System (DMS), Corridor Plans, and ArcGIS Tool: Vermont Stream Geomorphic and Habitat Assessment Data, Map Serve, and River Corridor Plans and Delineation Tools. Vermont Agency of Natural Resources, Waterbury, Vermont. http://www.vtwaterquality.org/rivers/htm.

VTANR, 2008. The Vermont Agency of Natural Resources Reach Habitat Assessment (RHA). Vermont Agency of Natural Resources, Departments of Environmental Conservation and Fish and Wildlife, Waterbury, VT.

Wang, Y. M., Kerh, T., Traore, S., 2009. Neural networks approaches for modeling river suspended sediment concentration due to tropical storms. Global NEST Journal, 11(4): 457-466.

Wasserman, P. D., 1993. Advanced Methods in Neural Computing. Van Nostrand Reinhold, New York, N. Y.

Yager, R.R., 1987, On the Dempster-Shafer framework and new combination rules: Information Sciences, v. 41, p. 93-138.

Zadeh, L.A., 1975. The concept of a linguistic variable and its application to approximate reasoning I, II, III. Information Sciences, 8: 199-251, 301-357; 9: 43-80.

Zadeh, L. A., 1984, Review of Books: A Mathematical Theory of Evidence: The AI Magazine, v. 5, no. 3, p. 81-83.

Zadeh, L. A., 1986, A simple view of the Dempster-Shafer theory of evidence and its implication for the rule of combination: The AI Magazine. v. 7, p. 85-90.

Zhang, L., 1994, Representation, independence, and combination of evidence in the Dempster-Shafer theory, *in* Yager, R.R., Kacprzyk, J., and Fedrizzi, M., eds., Advances in the Dempster-Shafer Theory of Evidence: John Wiley & Sons, Inc., New York, p. 51-69.

Zhao, S., Zhang, J., Li, X., and Song, W., 2007. A generalized regression neural network based on fuzzy means clustering and its application in system identification. 2007 International Symposium on Information Technology Convergence - ISITC '07: 13-16.

# APPENDIX A: FUZZY GRNN CODE

```matlab
%This code is a GRNN that allows the inputs to be fuzzy
numbers, the output is a fuzzy number as well
%Bree Mathon
%5/6/10
%Last modified 10/20/10

%Functions called: Discretize() and DiscretizeY()
tic;
clear all
close all
clc


%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%Read in training and testing data

NumInPar=4;

[XTrain]=xlsread('ExampleXTrain090510');
[YTrain]=xlsread('ExampleYTrain');
[XPredict]=xlsread('ExampleXPredict');
[YRealScore]=xlsread('ExampleYPredict');

sigsq=0.5;
ScoreTotal=40;

[RowsTrain,ColsTrain]=size(XTrain);
[RowsPredict,ColsPredict]=size(XPredict);
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%Discretize the inputs
y=[0 1 0];
%alphaLt=[0,0.25,0.5,0.75,1];
alphaLt=[0,0.25,0.5,0.75,1];
%alphaRt=[1,0.75,0.5,0.25,0];
XTrainDis=Discretize(XTrain,y,alphaLt);
  %XTrainDis=XTrainDis(1,:);
XPredictDis=Discretize(XPredict,y,alphaLt);
  %XPredictDis=XPredictDis(1,:);
YTrainDis=DiscretizeY(YTrain,y,alphaLt);
  %YTrainDis=YTrainDis(1,:);
[RowsTrainDis,ColsTrainDis]=size(XTrainDis);
```

```matlab
[RowsPredictDis,ColsPredictDis]=size(XPredictDis);
[RowsTrainYDis,ColsTrainYDis]=size(YTrainDis);

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%cycle through the calculations one alpha cut at a time
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
for ii=1:length(alphaLt)
    %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
                %Assign the weights
    %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

    %Training weights
    W1=XTrainDis(ii,:);
    [RW1,CW1]=size(W1);
    %Pattern weights
    WA=YTrainDis(ii,:);
    [RWA,CWA]=size(WA);
    %Weights connected to Node B
    WB=ones(1,ColsTrainYDis);
    [RWB,CWB]=size(WB);

    %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
                %Calculate pattern units
    %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
    prod=RowsPredict*2;

  n=0; %used to keep track of spot in XPredictDis so we
        can pull out the current prediction set
    tt=0;
    for i=1:RowsPredict
        Predict=[XPredictDis(ii,i+n:i+n+1)];...
            %XPredictDis(ii,(i+n+prod):(i+n+prod+1))];...
            %XPredictDis(ii,(i+n+2*prod):(i+n+2*prod+1))...
            %XPredictDis(ii,(i+n+3*prod):(i+n+3*prod+1))];
%Calls in left and right bounds for prediction input one
%alpha cut at a time, currently written to accommodate four
%input variables/comment or uncomment to accommodate number
%of variables

    n=n+1;
        [RowPred,ColPred]=size(Predict);
%        for j=0:ColPred/2-1
%            leftPred(1,1)=1;
%            leftPred(1,j+1)=2*j+1;
```

```matlab
%           end
%           for j=1:ColPred/2
%               %rt(1,1)=2;
%               rtPred(1,j)=2*j;
%           end
        m=1;
        %%%%Create degradation, aggradation, widening, and
         planform change matrices%%%%%%%%%%%%
        for j=1:2*RowsTrain
            Deg(1,j)=W1(1,j);
        end
        v=1;
%       for j=2*RowsTrain+1:4*RowsTrain
%           Agg(1,v)=W1(1,j);
%           v=v+1;
%       end
        v=1;
%       for j=4*RowsTrain+1:6*RowsTrain
%           Wid(1,v)=W1(ii,j);
%           v=v+1;
%       end
        v=1;
%       for j=6*RowsTrain+1:8*RowsTrain
%           PC(1,v)=W1(ii,j);
%           v=v+1;
%       end

%%%%%%%  BEGIN TO CALCULATE PATTERN UNIT NODES  %%%%%%%%%
%%% Calculate the distance matrices for each parameter %%%
    %for j=1:RW1
        s=0;

        for k=1:RowsTrain
            m=1;
%               I(k,m)=abs(Predict(1,1));
%               TempI(j,k+s)=min(DistSumI(j,:));
%               TempI(j,k+s+1)=max(DistSumI(j,:));
        %%%%%%%%%%%%%%%%%% L1 Norm %%%%%%%%%%%%%%%%%%%%
        %%%%%%% Uncomment for one input %%%%%%%%%%%%%
%               DistDeg(k,1)=abs(Predict(1,1)-Deg(1,k+s));
%               DistDeg(k,2)=abs(Predict(1,2)-Deg(1,k+s+1));
%               DistDeg(k,3)=abs(Predict(1,1)-Deg(1,k+s+1));
%               DistDeg(k,4)=abs(Predict(1,2)-Deg(1,k+s));
        %%%%%%% Uncomment for two inputs %%%%%%%%%%%%%
```
131

```matlab
%               DistAgg(k,1)=abs(Predict(1,3)-Agg(1,k+s));
%               DistAgg(k,2)=abs(Predict(1,4)-Agg(1,k+s+1));
%               DistAgg(k,3)=abs(Predict(1,3)-Agg(1,k+s+1));
%               DistAgg(k,4)=abs(Predict(1,4)-Agg(1,k+s));
%       %%%%%% Uncomment for three inputs %%%%%%%%%
%               DistWid(k,1)=abs(Predict(1,5)-Wid(1,k+s));
%               DistWid(k,2)=abs(Predict(1,6)-Wid(1,k+s+1));
%               DistWid(k,3)=abs(Predict(1,5)-Wid(1,k+s+1));
%               DistWid(k,4)=abs(Predict(1,6)-Wid(1,k+s));
%       %%%%%% Uncomment for four inputs %%%%%%%%%%
%               DistPC(k,1)=abs(Predict(1,7)-PC(1,k+s));
%               DistPC(k,2)=abs(Predict(1,8)-PC(1,k+s+1));
%               DistPC(k,3)=abs(Predict(1,7)-PC(1,k+s+1));
%               DistPC(k,4)=abs(Predict(1,8)-PC(1,k+s));

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%%Sum the Dist in order to create node I, each row are the
%possible values for the node I (number of nodes I equals
%number of training patterns)
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%
%           t=1;
%            for u=1:4 %Deg
%               %for w=1:4 %Agg
%                   %for p=1:4  %Wid
%                       %for q=1:4   %PC
%                           I(k,t)=DistDeg(k,u);
%                               +DistAgg(k,w);
%                               +DistWid(j,p);
%                               +DistPC(j,q);
%                           t=t+1;
%                       %end
%                   %end
%               %end
%           end
%    %%%%%%%%%%%%%%%%%%% End L1 Norm %%%%%%%%%%%%%%%%%%%

     %%%%%%%%%%%%%%%%%%%%% L2 Norm %%%%%%%%%%%%%%%%%%%%%

     %%%%%%%%% Uncomment for one input %%%%%%%%%%%%%%%
            DistDeg(k,1)=(Predict(1,1)-Deg(1,k+s))^2;
            DistDeg(k,2)=(Predict(1,2)-Deg(1,k+s+1))^2;
            DistDeg(k,3)=(Predict(1,1)-Deg(1,k+s+1))^2;
            DistDeg(k,4)=(Predict(1,2)-Deg(1,k+s))^2;
```

```matlab
        %%%%%%%%% Uncomment for two inputs %%%%%%%%%%%%%%
%               DistAgg(k,1)=(Predict(1,3)-Agg(1,k+s))^2;
%               DistAgg(k,2)=(Predict(1,4)-Agg(1,k+s+1))^2;
%               DistAgg(k,3)=(Predict(1,3)-Agg(1,k+s+1))^2;
%               DistAgg(k,4)=(Predict(1,4)-Agg(1,k+s))^2;
        %%%%%%%%% Uncomment for three inputs %%%%%%%%%%%%
%               DistWid(k,1)=(Predict(1,5)-Wid(1,k+s))^2;
%               DistWid(k,2)=(Predict(1,6)-Wid(1,k+s+1))^2;
%               DistWid(k,3)=(Predict(1,5)-Wid(1,k+s+1))^2;
%               DistWid(k,4)=(Predict(1,6)-Wid(1,k+s))^2;
        %%%%%%%%% Uncomment for four inputs %%%%%%%%%%%%%
%               DistPC(k,1)=(Predict(1,7)-PC(1,k+s))^2;
%               DistPC(k,2)=(Predict(1,8)-PC(1,k+s+1))^2;
%               DistPC(k,3)=(Predict(1,7)-PC(1,k+s+1))^2;
%               DistPC(k,4)=(Predict(1,8)-PC(1,k+s))^2;

        %%%%%%%%%%%%% Sum & Take Square Root %%%%%%%%%%%%%
                t=1;
                for u=1:4 %Deg
                    %for w=1:4 %Agg
                        %for p=1:4  %Wid
                            %for q=1:4  %PC
                                I(k,t)=sqrt(DistDeg(k,u));
%                                       +DistAgg(k,w);
%                                       +DistWid(j,p);
%                                       +DistPC(j,q));
                                t=t+1;
                            %end
                        %end
                    %end
                end

        %%%%%%%%%%%%%%%%%%%%%%% End L2 Norm %%%%%%%%%%%%%%%%
                s=s+1;
            end

        fI=exp(-I/(2*sigsq)); %output from the pattern unit
        [RfI,CfI]=size(fI);

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
    %%%%%%%%%%%%%  CALCULATE SUMMATION UNITS %%%%%%%%%%%%%
    %Summation unit A
    p=1;
    %Calculate product of fI and the corresponding weights
```

```matlab
in WA
s=0;
for j=1:RfI
    t=1;
    WANode=WA(1,j+s:j+s+1);
    s=s+1;
     for jj=1:CfI
         for jjj=1:length(WANode)
             Product(j,t)=(WANode(1,jjj)*fI(j,jj));
             t=t+1;
         end
     end
end

[RProd,CProd]=size(Product);
%%% HAVE TO ADD/SUBTRACT LOOP FOR EACH TRAINING PATTERN %%%
t=1;
for u=1:CProd
    for w=1:CProd
        for p=1:CProd
            SumA(t,1)=Product(1,u)
                    +Product(2,w)
                    +Product(3,p);
            t=t+1;
        end
    end
end

%Summation unit B
%Calculate product of fI and the corresponding weights
in WB
s=0;
for j=1:RfI
    t=1;
    WBNode=WB(1,j+s:j+s+1);
    s=s+1;
     for jj=1:CfI
         for jjj=1:length(WBNode)
             ProductB(j,t)=(WBNode(1,jjj)*fI(j,jj));
             t=t+1;
         end
     end
end
```

```matlab
    [RProdB,CProdB]=size(ProductB);
   %% HAVE TO ADD/SUBTRACT LOOP FOR EACH TRAINING PATTERN %%
    t=1;
    for u=1:CProdB
        for w=1:CProdB
            for p=1:CProdB
                SumB(t,1)=ProductB(1,u)
                        +ProductB(2,w)
                        +ProductB(3,p);
                t=t+1;
            end
        end
    end

    %%%%%%%%%%%%%% CALCULATE OUTPUT NODE %%%%%%%%%%%%%%%%%%

    for j=1:length(SumA)
        DivAB(j,1)=SumA(j,1)/SumB(j,1);
    end
    % Take min/max of possible output values
    tt=0;
    Out(ii,i+tt)=min(DivAB);
    Out(ii,i+tt+1)=max(DivAB);
    tt=tt+1;

end
end

plot(Out(:,1:2),alphaLt,'b*');
xlim([0 200]);
xlabel('RHA Score');
ylabel('Membership Degree');

%%%%%%%%%%%%%%%%%%%%%%%%  END MAIN CODE   %%%%%%%%%%%%%%%%%%%%

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
function [ xout ] = Discretize(x,y,alphaLt)
%Discretize function: This function takes the input
triangular membership functions, assumes linear edges, and
discretizes the function based on user defined alpha cuts

[rowx,colx]=size(x);
ylt=y(1,1:2);
yrt=y(1,2:3);
```

```matlab
m=1;

for i=1:rowx
    xlt=x(i,1:2);
    xrt=x(i,2:3);
    plt=polyfit(xlt,ylt,1);
    prt=polyfit(xrt,yrt,1);
    xoutlt1(i,:)=(alphaLt-plt(1,2))/plt(1,1);
    xoutrt1(i,:)=(alphaLt-prt(1,2))/prt(1,1);
    xout1(:,m)=(xoutlt1(i,1:end));
    xout1(:,m+1)=(xoutrt1(i,1:end));
    m=m+2;
end
n=1;
for i=1:rowx
    xlt=x(i,4:5);
    xrt=x(i,5:6);
    plt=polyfit(xlt,ylt,1);
    prt=polyfit(xrt,yrt,1);
    xoutlt2(i,:)=(alphaLt-plt(1,2))/plt(1,1);
    xoutrt2(i,:)=(alphaLt-prt(1,2))/prt(1,1);
    xout2(:,n)=(xoutlt2(i,1:end));
    xout2(:,n+1)=(xoutrt2(i,1:end));
    n=n+2;
end
% p=1;
% for i=1:rowx
%     xlt=x(i,7:8);
%     xrt=x(i,8:9);
%     plt=polyfit(xlt,ylt,1);
%     prt=polyfit(xrt,yrt,1);
%     xoutlt3(i,:)=(alphaLt-plt(1,2))/plt(1,1);
%     xoutrt3(i,:)=(alphaLt-prt(1,2))/prt(1,1);
%     xout3(:,p)=(xoutlt3(i,1:end));
%     xout3(:,p+1)=(xoutrt3(i,1:end));
%     p=p+2;
% end
% q=1;
% for i=1:rowx
%     xlt=x(i,10:11);
%     xrt=x(i,11:12);
%     plt=polyfit(xlt,ylt,1);
%     prt=polyfit(xrt,yrt,1);
%     xoutlt4(i,:)=(alphaLt-plt(1,2))/plt(1,1);
```

```matlab
%       xoutrt4(i,:)=(alphaLt-prt(1,2))/prt(1,1);
%       xout4(:,q)=(xoutlt4(i,1:end));
%       xout4(:,q+1)=(xoutrt4(i,1:end));
%       q=q+2;
% end
xout=[xout1 xout2];% xout3 xout4];
end
%%%%%%%%%%%%%%%% END FUNCTION Discretize %%%%%%%%%%%%%%%%%%

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
function [ xout ] = DiscretizeY(x,y,alphaLt)
DiscretizeY function: This function takes the training
output triangular membership functions, assumes linear
edges, and discretizes the function based on user defined
alpha cuts
[rowx,colx]=size(x);
ylt=y(1,1:2);
yrt=y(1,2:3);
m=1;
for i=1:rowx

    xlt=x(i,1:2);
    xrt=x(i,2:3);
    plt=polyfit(xlt,ylt,1);
    prt=polyfit(xrt,yrt,1);
    xoutlt1(i,:)=(alphaLt-plt(1,2))/plt(1,1);
    xoutrt1(i,:)=(alphaLt-prt(1,2))/prt(1,1);
    xout1(:,m)=(xoutlt1(i,1:end));
    xout1(:,m+1)=(xoutrt1(i,1:end));
    m=m+2;
end
xout=xout1;
end
%%%%%%%%%%%%%%%% END FUNCTION DiscretizeY %%%%%%%%%%%%%%%%%%
```

# APPENDIX B: GRNN CODE

```matlab
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%A generalized regression neural network (GRNN)
%Bree Mathon
%Originally created: 11/10/09
%Last modified: 10/1/10

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
% Input:
%   XTrain - training input, patterns are assigned to rows
%   YTrain - training output (pattern weights)
%   XPredict - prediction input variables
%   YRealScore - actual values of output in testing data
%              set (if using)
% Output:
%   ScaleOut - GRNN prediction rounded to the nearest
%              integer
%   Class - prediction data calssified, if using for RGA,
%           LRHA, RHA
%   YRealClass - testing output values classed like GRNN
%                output
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

clear all
close all
clc

%%%%%%%%%%%%%% READ IN TRAINING & PREDICTION DATA %%%%%%%%%
% The input data is read in as each row is a pattern and
each column is a parameter (variable)

[XTrain]=xlsread('XTrain');
[YTrain]=xlsread('YTrain');

[XPredict]=xlsread('XPredict');
[YRealScore]=xlsread('YReal');

[RowsTrain,ColsTrain]=size(XTrain);
[RowsPredict,ColsPredict]=size(XPredict);

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```matlab
%ScoreTotal=200;  %use when predicting LRHA condition
%ScoreTotal=160;  %use when predicting new RHA condition
ScoreTotal=1;     %use when predicting bio health class or
                  when using data
                  %that does not require classification of
                  GRNN output


%%%%%%%%%%%%%%%%%%%% SET SIGMA VALUE %%%%%%%%%%%%%%%%%%%%%

sigma=[0.05:0.05:1];  %Use to determine optimal sigma value
%sigma=0.5;  %NOTE: this is sigma not sigma squared, it
will be squared later in the activation equation

%%%%%%%%%%%%% NORMALIZE THE DATA IF NECESSARY %%%%%%%%%%%%

NormTrain=XTrain;
%    for i=1:RowsTrain
%        for j=1:ColsTrain
%            Sqrd(i,j)=XTrain(i,j)^2;
%            SumSq(i)=sum(Sqrd(i,:));
%            EuclidLength(i)=sqrt(SumSq(i));
%            NormTrain(i,j)=XTrain(i,j)/EuclidLength(i);
%        end
%
%    end

NormPredict=XPredict;
%    for k=1:RowsPredict
%        for m=1:ColsPredict
%            SqrdPr(k,m)=XPredict(k,m)^2;
%            SumSqPr(k)=sum(SqrdPr(k,:));
%            EuclidLengthPr (k)=sqrt(SumSqPr(k));
             NormPredict(k,m)=XPredict(k,m)/…
                             EuclidLengthPr(k);
%        end
%    end

NormPattern=YTrain;
% MinY=min(YTrain); MaxY=max(YTrain);
% NormPattern=(YTrain(:,1)-MinY)/(MaxY-MinY);

%%%%%%%%%%%%%%% ASSIGN THE WEIGHTS %%%%%%%%%%%%%%%%%%%%%%%%
```

```matlab
W1=NormTrain;
WA=NormPattern;
WB=ones(RowsTrain,1);

[RW1,CW1]=size(W1);
[RWA,CWA]=size(WA);
[RWB,CWB]=size(WB);

%%%%%%%%%%%%%% CALCULATE PATTERN UNITS %%%%%%%%%%%%%%%%%%

for s=1:length(sigma)
    for i=1:RowsPredict
        for j=1:RW1
            for k=1:ColsTrain
              Dist(j,k)=(NormPredict(i,k)-W1(j,k))^2;   %L2 norm
              %Dist(j,k)=abs(NormPredict(i,k)-W1(j,k));%L1 norm
            end
        end
        %I=sum(Dist,2);%should be a col vector with the same
        # of elements as the training set (if training has 20
        inputs (j) with 3 parameters each (k) then there
        would be 20 I elements for each prediction set - here
        it is being overwrittien for each prediction)
        I=sqrt(sum(Dist,2)); %L2 norm

fI=exp(-I/(2*(sigma(s))^2)); %output from the pattern unit

%%%%%%%%%%%%%% CALCULATE SUMMATION UNITS %%%%%%%%%%%%%%%%%
%Summation unit A
A= dot(fI,WA);

%Summation unit B
B=dot(fI,WB);

%%%%%%%%%%%%%% CALCUALTE OUTPUT NODE %%%%%%%%%%%%%%%%%%%%%

ScaleOut(i)=A/B;

%Re-scale the output if data was normalized
%Out(i)=(ScaleOut(i)*(MaxY-MinY))+MinY;
%end

ScaleOut=round(ScaleOut);
Scale=ScaleOut';
```

140

```matlab
Out=(Scale/ScoreTotal); % Scales data between 0 and 1 if
                        %          ScoreTotal > 1 so it can be
                        %          classed
YReal=(YRealScore/ScoreTotal);
end

%%%%%%%%%%%%%%%% CALCULATE RMSE %%%%%%%%%%%%%%%%%%%%%%%%%%%

for i=1:RowsPredict
        Diff(i)=(ScaleOut(i)-YRealScore(i))^2;
    end
    SumDiff=sum(Diff);
    RMS(s)=sqrt(SumDiff/RowsPredict);

%%%%%%% CLASS GRNN OUTPUT & TEST DATA (IF USING) %%%%%%%%%
% %Class
% for j=1:RowsPredict
%     if Out(j,1)>=0 && Out(j,1)<=0.34
%         Class(j,1)=1;
%     elseif Out(j,1)>0.34 && Out(j,1)<=0.64
%         Class(j,1)=2;
%     elseif Out(j,1)>0.64 && Out(j,1)<=0.84
%         Class(j,1)=3;
%     elseif Out(j,1)>0.84 && Out(j,1)<=1.0
%         Class(j,1)=4;
%     end
% end
% for j=1:RowsPredict
%     if YReal(j,1)>=0 && YReal(j,1)<=0.34
%         YRealClass(j,1)=1;
%     elseif YReal(j,1)>0.34 && YReal(j,1)<=0.64
%         YRealClass(j,1)=2;
%     elseif YReal(j,1)>0.64 && YReal(j,1)<=0.84
%         YRealClass(j,1)=3;
%     elseif YReal(j,1)>0.84 && YReal(j,1)<=1.0
%         YRealClass(j,1)=4;
%     end
% end

%%%%%%%%%% CALCULATE NUMBER CORRECTLY PREDICTED %%%%%%%%%%
%
%%%% USE WHEN CHECKING PREDICTIONS THAT WERE CLASSIFIED %%%
%
% correct = 0;
```

```matlab
% wrong(s)=0;
% for p = 1:RowsPredict
%      if Class(p,1) == YRealClass(p,1)
%          correct=correct+1;
%      else
%          wrong(s)=wrong(s)+1;
%      end
% end
% end
% percent_correct = (correct/RowsPredict)*100

%%%%% USE WHEN PREDICTING DATA THAT IS NOT CLASSIFIED %%%%

correct = 0;
wrong(s)=0;
for p = 1:RowsPredict
    if Scale(p,1) == YRealScore(p,1)
        correct=correct+1;
    else
        wrong(s)=wrong(s)+1;
    end
end
end
percent_correct = (correct/RowsPredict)*100

plot(sigma,RMS)
xlabel('sigma')
ylabel('RMS')

%%%%%%%%%%%%%%%%% END GRNN MAIN CODE %%%%%%%%%%%%%%%%%%%%
```