Dissertations, Master's Theses and Master's Reports - Open

Dissertations, Master's Theses and Master's Reports

2010

# Use of hydroclimatic forecasts for improved water management in central Texas

Wenge Wei
*Michigan Technological University*

### Recommended Citation

Wei, Wenge, "Use of hydroclimatic forecasts for improved water management in central Texas ", Dissertation, Michigan Technological University, 2010.
https://digitalcommons.mtu.edu/etds/280

# USE OF HYDROCLIMATIC FORECASTS FOR IMPROVED WATER MANAGEMENT IN CENTRAL TEXAS

By

Wenge Wei

A DISSERTATION

Submitted in partial fulfillment of the requirements

for the degree of

DOCTOR OF PHILOSOPHY

Civil Engineering

MICHIGAN TECHNOLOGICAL UNIVERSITY

2010

This dissertation, "Use of Hydroclimatic Forecasts for Improved Water Management in Central Texas", is hereby approved in partial fulfillment of the requirements for the degree of DOCTOR OF PHILOSOPHY in the field of Civil Engineering.

Department of Civil and Environmental Engineering

Signatures:

Dissertation Advisor   _____

David W. Watkins

Department Chair   _____

William M. Bulleit

Date   _____

# Table of Contents

# List of Figures

# List of Tables

# Acknowledgements

I wish to express my gratitude to my advisor, Dave Watkins for his expertise and understanding and most of all, immense patience, throughout the course of this dissertation. This study would never have been possible without his support and encouragement.

I would like to thank my committee members, Brian Barkdoll, Alex Mayer, and Jianping Dong for taking interest in my works as well as giving me their time and suggestions.

I am very grateful to my friends for cheering me on during tough time. Their support and encouragement have been invaluable throughout my time at Tech.

Finally, I owe gratitude to my family, especially my wife who supported me throughout my education. I would never have made it this far without her great inspiration.

# Abstract

Accurate seasonal to interannual streamflow forecasts based on climate information are critical for optimal management and operation of water resources systems. Considering most water supply systems are multipurpose, operating these systems to meet increasing demand under the growing stresses of climate variability and climate change, population and economic growth, and environmental concerns could be very challenging. This study was to investigate improvement in water resources systems management through the use of seasonal climate forecasts. Hydrological persistence (streamflow and precipitation) and large-scale recurrent oceanic-atmospheric patterns such as the El Niño/Southern Oscillation (ENSO), Pacific Decadal Oscillation (PDO), North Atlantic Oscillation (NAO), the Atlantic Multidecadal Oscillation (AMO), the Pacific North American (PNA), and customized sea surface temperature (SST) indices were investigated for their potential to improve streamflow forecast accuracy and increase forecast lead-time in a river basin in central Texas.

First, an ordinal polytomous logistic regression approach is proposed as a means of incorporating multiple predictor variables into a probabilistic forecast model. Forecast performance is assessed through a cross-validation procedure, using distributions-oriented metrics, and implications for decision making are discussed. Results indicate that, of the predictors evaluated, only hydrologic persistence and Pacific Ocean sea surface temperature patterns associated with ENSO and PDO provide forecasts which are statistically better than climatology.

Secondly, a class of data mining techniques, known as tree-structured models, is investigated to address the nonlinear dynamics of climate teleconnections and screen promising probabilistic streamflow forecast models for river-reservoir systems. Results show that the tree-structured models can effectively capture the nonlinear features hidden in the data. Skill scores of probabilistic forecasts generated by both classification trees and logistic regression trees indicate that seasonal inflows throughout the system can be

predicted with sufficient accuracy to improve water management, especially in the winter and spring seasons in central Texas.

Lastly, a simplified two-stage stochastic economic-optimization model was proposed to investigate improvement in water use efficiency and the potential value of using seasonal forecasts, under the assumption of optimal decision making under uncertainty. Model results demonstrate that incorporating the probabilistic inflow forecasts into the optimization model can provide a significant improvement in seasonal water contract benefits over climatology, with lower average deficits (increased reliability) for a given average contract amount, or improved mean contract benefits for a given level of reliability compared to climatology. The results also illustrate the trade-off between the expected contract amount and reliability, i.e., larger contracts can be signed at greater risk.

# 1. Introduction

## 1.1 Motivation and Objectives

With seemingly ever-increasing demands for water, including urban water supply, recreation, hydroelectric power and environmental flow demands, it is becoming more important for water resources management to be as efficient as possible throughout the U.S., and indeed the world, to ensure reliable water supplies and ecosystem protection. The looming uncertainty about future supplies due to climate change, potentially increasing the frequency and severity of droughts and floods, presents another daunting challenge to water resources engineers and managers. To counter these challenges, it is imperative to develop optimal water resources management systems by utilizing the hydrologic and meteorological information readily available, including climate forecasts at monthly, seasonal and even longer-lead times.

Traditionally, water resources management and planning have been based on critical period hydrology, in which water supply operation decisions are made with the explicit goal of preparedness for the drought of record (Hall and Dracup, 1970). These and other heuristic methods primarily depend on past experience, observations of current conditions, and professional judgment (Lee, 1999). Decisions based on these methods may be problematic due to a lack of explicit consideration of risk and neglect of the effects of climate change and variability of water supplies, and thus tend to lead to reduced efficiency for multiple-objective water resources systems.

The use of climate forecasts could improve water resources management and planning, especially in light of changing conditions in which new information can help to mitigate adverse impacts. As significant progress has been made in understanding "teleconnections" between large-scale atmospheric circulation patterns and regional climate anomalies, streamflow forecasts have improved for a range of lead times. In particular, climatic predictors such as recurrent teleconnection patterns (e.g., El Niño/Southern Oscillation, Pacific Decadal Oscillation, North Atlantic Oscillation) can

provide sufficient lead-time and accuracy for long-lead streamflow forecasts in many regions (Piechota et al., 1997; Hamlet and Lettenmaier, 1999; Anderson et al., 2001; Gutierrez and Dracup, 2001). However, climate forecasts are not without limitations. In some cases, the skill of the climate forecast is not great enough for operational use due to limited understanding of climate processes and prediction capabilities and variability in climate signals. This may be particularly true for geographically small basins (i.e., requiring downscaling of global or regional climate models) (Hamlet et al., 2002). The complexity involved in using forecasts and the lack of extensive records and forecasts for verification indicate a need for developing new tools and management strategies. Additionally, water managers may be hesitant to apply new information and methods that could expose them and other system stakeholders to greater risk. These are some possible reasons why seasonal climate forecasts are not being used to the fullest extent possible.

The overall goal of this study is to investigate the potential benefits of seasonal climate forecasts through improved on water resources management. The study is to build on previous work in which a decision support model using stream flow ensembles was developed for the Lower Colorado River Authority (LCRA) in Austin, Texas (Watkins et al., 2000; Kracman et al., 2006). Whereas the previously used stream flow ensembles were based on climatology, this study will seek to add predictive skill to the model (or other appropriate decision models) by conditioning the ensemble forecasts on observable climate indicators such as the El Niño-Southern Oscillation (ENSO), Pacific Decadal Oscillation (PDO), the North Atlantic Oscillation (NAO), the Atlantic Multidecadal Oscillation (AMO), and the Pacific North American Pattern (PNA). Forecasts based on indicators of hydrologic persistence (e.g., soil moisture) will also be investigated, and the value of climate information and persistence-based forecasts will be estimated through retrospective comparison of management decisions with conditional and unconditional ensemble forecasts.

The project goal will be accomplished by completing the following five main tasks in a case study of the Highland Lakes system in central Texas:

(1) Analysis of potential hydroclimatic predictors for the case study region in central Texas.

(2) Derivation of maximal skill forecasts based on identified predictor variables.

(3) Generation of stream flow ensembles consistent with the skill and uncertainty of the forecasts.

(4) Development or modification of stochastic optimization models for water management decisions.

(5) Application of the model(s) with and without seasonal forecast information to evaluate the potential benefits of forecasts.

It is expected that this research will be of significant value to water managers at the Lower Colorado River Authority in Austin, Texas, as well provide a general framework that may be applied elsewhere. Furthermore, increased skill in seasonal forecasts may be incorporated in other products, such as the ensemble streamflow forecasts issued by the West Gulf River Forecast Center, which if adapted to neighboring basins, may prove useful to other water managers in Texas.

## 1.2   Case Study Background

**Lower Colorado River and Highland Lakes**

The Lower Colorado River Authority (LCRA) operates the Highland Lakes system in Central Texas, a series of six lakes on the watershed of the Lower Colorado River. As a water conservation and reclamation district, the LCRA provides water supply and flood control to a 33-county area, including the City of Austin and several rice irrigation districts along the Texas Gulf Coast (see Figure 1.1). In addition, the LCRA produces wholesale power for a 53-county service area and provides water resources for lake recreation activities and in-stream flow maintenance. To meet rapidly growing demands, reservoir inflow forecasts with lead times of up to 6 to 12 months would be very beneficial. However, seasonal and long-term forecasts are not used by the LCRA for a number of reasons, including high seasonal and annual variability of stream flow and the

absence of easily measured hydrologic indicators such as snowpack. Until recently, the LCRA has not had much experience with probabilistic planning methods, although they are now using a Monte Carlo simulation model to predict expected ranges of lake levels (Ron Anderson, LCRA, personal communication, 2009).



**Figure 1.1.** Lower Colorado River Authority District (Source: LCRA 2009)

The Colorado River of Texas runs from Southwest New Mexico, across Texas to the Matagorda Bay on the Gulf of Mexico. The river's watershed covers nearly 40,000 square miles, and the river flows a distance of approximately 600 miles from its headwaters to its mouth. The water quality ranks high and is important to abundant aquatic biota and wildlife such as migratory birds. The watershed includes a variety of land types, from the Central Texas Hill Country to the flat Coastal Plain. Land use is also varied, with urbanized areas such as the City of Austin, as well as smaller residential and agricultural regions, wetlands, and community parks. Perennial rivers exhibit a large

range of flows, which subject the region to frequent droughts and flooding (Kracman, 2002).

The Lower Colorado River starts in central Texas, and was legally distinguished from the upper portion of the river through legislation that appointed jurisdiction over this part of the river to the Lower Colorado River Authority (LCRA). There are a series of six reservoirs, known as the Highland Lakes, on the Lower Colorado (Figure 1.2). Development of the Highland Lakes system occurred between the years 1939 to 1951 with the construction of dams, mostly for the purpose of flood and drought mitigation, although additional uses, especially recreation, have become important over the years. The Owen H. Ivie Reservoir, built in 1990, marks the upstream boundary of the Lower Colorado River and releases flows upon which the Lower Colorado River flows depend. Downstream of O.H. Ivie are the confluence of the Colorado with the Pecan Bayou and the San Saba River, a major tributary. It then flows into Lake Buchanan, which was formed by Buchanan Dam, completed in 1947 with a capacity of about 918,000 acre-feet. Lake Buchanan is one of the two lakes with capacity for water storage. Immediately below Lake Buchanan is the much smaller Inks Reservoir (17,500 acre-feet, built in 1948 for hydropower purposes), the confluence with the Llano River, Lake Lyndon B. Johnson (138,500 acre-feet, built in1950), and then Lake Marble Falls (8,760 acre-feet, built in 1951 for hydropower purposes).

**Figure 1.2.** Lower Colorado River Basin and Highland Lakes (LCRA, n.d.)

Next, the river reaches Lake Travis, created by the construction of Mansfield Dam in 1941. Another important tributary, the Pedernales River, flows into this lake. The Mansfield Dam is considered the only flood control structure, by design, for the Lower Colorado and can hold 748,502 acre-ft above the conservation pool (Kracman 2002). The total capacity of Lake Travis is 1,170,752 acre-feet, and together, Lakes Buchanan and Travis hold approximately 2 million acre-feet of conservation storage. This upper part of the Lower Colorado River Basin, from San Saba County to Lake Austin, is considered the Texas Hill Country, below which slopes decrease and the river broadens significantly. Finally, the sixth reservoir is reached, Lake Austin, with a capacity of 21,000 acre-feet. The Tom Miller Dam, which formed Lake Austin in 1940, is operated by the LCRA but is owned by the City of Austin.

Downstream of the Tom Miller Dam, and not part of the Highland Lakes system, is Town Lake, which is controlled by the City of Austin. Within the Austin city limits, or below Austin, the Lower Colorado River is met by Barton Creek, Onion Creek, and

6

numerous other small tributaries before finally reaching Matagorda Bay. Outflows to the Gulf of Mexico average 2,600 cubic feet per second. Along this downstream reach of the river, four main rice irrigation districts withdraw water: Lakeside, Garwood, Pierce Ranch, and Gulf Coast. These irrigation districts are shown in Figure 1.3.



**Figure 1.3.** LCRA irrigation service areas (Source: LCRA)

The LCRA's initial goals were to moderate droughts and floods. With the construction of the Highland Lakes system these goals were better realized, and a stable water supply and a source of hydroelectric power encouraged growth in the region. In recent years, a major portion of water releases (about 50%) has gone to the irrigation districts (Figure 1.4). Hydropower has become a secondary concern since the development of fossil fuel energy plants, though deregulation of the energy market has changed the value of this power source. Overall the usage of the water has changed, as have community needs and goals. Though hydropower importance has diminished, and flood control has remained a primary purpose, residential and municipal water supply,

recreation, and environmental uses have become more important on the Lower Colorado River.



112 billion gallons released for irrigation

54%

25 billion gallons released for environmental purposes

12%

27%

7%

55 billion gallons released for municipal and industrial uses

14 billion gallons released to keep Lake Travis from overfilling

**Figure 1.4.** Distribution of water releases from the Highland Lakes (LCRA, n.d.)

**LCRA Water Management Plan**

The LCRA is a public agency that was established in 1934 by Texas legislature for the purpose of "conservation and reclamation" (LCRA, 1999). Today, the LCRA's water management plan includes consideration of private rights holders, recreational, environmental, and hydroelectric interests, as well as two main types of customers, those with firm contracts (municipal and industrial) and those who sign yearly interruptible contracts (agricultural). Firm water is diverted from storage under a long-term contract or resolution issued by the LCRA Board to high-priority users such as the City of Austin and is a guaranteed water right during repetition of drought of record. Interruptible water contracts are issued on a shorter time scale (typically one year or less) with the condition that supplies may be interrupted or curtailed in the event that firm supplies become endangered. In allocating interruptible water, priority is given to irrigation operations downstream of Austin. If the availability of interruptible water exceeds these irrigation needs, annual contracts can then be made with other entities within the Lower Colorado

basin. Currently, the LCRA uses beginning-of-year (January 1) storage levels to determine the amount of water available to meet firm and interruptible water demands in the coming year (Martin, 1991).

**Firm Contracts**

The combined firm yield of Lakes Buchanan and Travis has been established to be 536,312 acre-feet per year, considered to be the maximum demand that could be met during a recurrence of the drought of record. There are six types of firm demand customers. The O.H. Ivie Reservoir, upstream of the Lower Colorado, has rights to store up to 90,546 acre-feet per year. The city of Austin is supplemented by firm supply up to 148,300 acre-feet per year. Municipalities and industries are guaranteed a total of 95,789 acre-feet per year. Two other important interests are cooling water for LCRA's hydroelectric plants (63,851 acre-feet per year) and for the South Texas Project power plants (5,680 acre-feet per year). Finally, environmental needs including instream flow, bays, and estuaries, receive 12,860 acre-feet per year. An additional 50,000 acre-feet per year is reserved for future growth, and in preparation for potential depletion or pollution of ground water supplies. Continued regional growth and development will likely increase the load on these firm contracts until it meets the full legal amount. Better management today will help ensure that these future demands will be efficiently met (LCRA 1999).

**Rice Farming**

The four main irrigation districts--Lakeside, Garwood, Pierce Ranch, and Gulf Coast--are mainly concerned with supplying water for regional rice farms. These farms play an important role in regional economy, being part of a $300 million industry (Texas State Historical Association, 2002). Irrigation needs are met partially by groundwater, but 70% comes from surface water, mostly releases from the Highland Lakes (LCRA, 1999).

Rice typically takes 120 to 180 days to mature, and requires fields flooded to depths of 4 to 6 inches. In this region, rice farmers often use two growing seasons, the main one

being March through July, and a secondary growth season during July through December. Future water shortages may prevent second season harvests, even though irrigation water demand for these crops occurs only from August to October.

**Interruptible Contracts**

Interruptible contracts are primarily related to the irrigation districts' demand. Currently the LCRA operates on a rule curve to determine how to make the interruptible contracts each year. The rule curve is reapplied each month to check on status compared to historical levels and to detect possible problems. Analysis concerning the projected water availability is made in October, and firm contract holders submit their projected needs for the year. Then, January 1 reservoir levels are projected and the minimum upstream inflow for the coming year is added to these levels. The difference between projected levels and firm demands is considered to be available for the interruptible contracts. The final contracts for interruptible water are signed in November. Based on the minimum of April, May, and June maximum storage levels, contracts are updated in preparation for the second rice growth season. (As of this writing, LCRA water managers are investigating the ability to base first-season contract decisions on projected water levels in March or April.)

**Power Generation**

On the Lower Colorado River, the six dams together have the capacity for 270 megawatts of hydroelectric power (LCRA, 1999). Firm water is used for cooling at the fossil fuel plants, and other releases are used to generate hydroelectric power. Although hydropower is not given the status of a priority water demand, the LCRA uses brief high-volume releases to maximize daily power generation (Kracman, 2002).

**Recreation and Tourism**

The Highland Lakes are a popular recreation spot for fishers, boaters, birders, and others. Recreation and tourism have been recognized by the LCRA for their importance to the local economy and are included in the water management plan as part of the

LCRA's "public interest responsibilities" (LCRA, 1999). Such demand is considered during distribution of interruptible water. There is a tradeoff between releasing interruptible water to the irrigation districts, which have senior rights, and maintaining high lake levels for the economically important tourism industry. After the irrigation districts' contracts have been met, further sales are limited based on lake levels. The LCRA supports local tourism businesses by encouraging visitors to the area. They have built 25 parks on LCRA-owned land, which receive over one million visitors and bring in over $90 million each year. The economic importance of such tourism, along with political considerations, is also a factor in the maintenance of equitable lake levels.

**Environmental Concerns**

The LCRA also pays attention to environmental interests on the river, from both quality and quantity perspectives. They track the quality of the water with frequent monitoring and yearly assessments. In order to support habitat for waterfowl, fish, shrimp, aquatic plants, and other biological organisms, minimum daily flows must be met. The LCRA has target flows of 1.03 million acre-feet per year to maintain the streamflow, bays, and estuaries. These are met with firm water supply under drought conditions, but during normal conditions, only interruptible water is provided for instream flow maintenance.

**Drought Management Plan**

The LCRA's conceptual lake management policy for year 2010 projected demands calls for curtailment of interruptible supplies to begin when combined storage levels drop below 1,400,000 acre-ft, or about 70% of the maximum water supply storage (LCRA 2003). Aggressive curtailment begins at a January 1 storage level of 1,150,000 acre-ft (about 58% of maximum), and no interruptible water use will be sanctioned on January 1 if levels are below 325,000 acre-ft (about 16% of maximum). Additionally, interruptible water use will be stopped at any time during the year if combined storage levels drop below 200,000 acre-ft (10% of maximum). Conversely, in years of high storage levels, additional interruptible water supplies may be available for sale if combined storage

levels are greater than 1,865,000 acre-ft (about 94% of maximum).  Figure 1.5 illustrates a hypothetical "rule curve" that corresponds to the published conceptual lake management policy.



**Figure 1.5.** Hypothetical rule curve corresponding to LCRA's conceptual lake management policy. (1 AF = 1,233.5 m$^3$.)

# 2. Literature Review

## 2.1 Ocean-Atmosphere-Streamflow Teleconnections

There is an increasing awareness that climate variability is not necessarily randomly distributed in space and time. Instead, some climate anomalies appear to present certain patterns which may be useful for hydrologic forecasting (Piechota, et al. 2006). The term "teleconnections" refers to large and persistent ocean-atmospheric anomaly patterns (e.g., El Niño/Southern Oscillation, North Atlantic Oscillation), and apparent causal effects on regional climate conditions in adjacent or remote regions. Recent studies have shown that oceanic-atmospheric variability occurs on interannual, decadal and interdecadal timescales and has an impact on the climate of regions around the world. The results and information can be utilized to improve on long lead-time forecasts of water availability. This study investigates the influences of these wide-scale teleconnection patterns on streamflow in central Texas. These teleconnection patterns include the El Niño/Southern Oscillation (ENSO), the Pacific Decadal Oscillation (PDO) and the North Atlantic Oscillation (NAO).

**El Niño/Southern Oscillation**

ENSO is a contraction of names of two phenomena that were recognized to be different expressions of the same process: "El Niño" refers to anomalous strong warming of the surface waters of the eastern equatorial Pacific Ocean, while "Southern Oscillation" refers to concurrent changes in surface barometric pressure in the tropical Pacific (Ropelewski and Halpert 1987, Philander 1990, Piechota, et al., 2006). The ENSO phenomenon spans the equatorial Pacific and is associated with droughts in Australia, New Zealand, and Southern Africa, and simultaneously flooding in North America, Peru, and Ecuador (Ropelewski and Halpert 1987). The warm phase of ENSO is referred to as El Niño and the cool phase is referred to as La Niña (Philander, 1990). The ENSO generally experiences a two- to seven-year periodicity. The atmospheric mechanisms associated with ENSO are understood as follows:

Under normal conditions, the trade winds blow towards the west and push warm surface water to the western Pacific, so that the sea surface level and temperatures are about higher at Indonesia than off the coast of Ecuador. The sea surface temperature is higher in the west due to an upwelling of cold water off the coast of South America. In normal conditions, rainfall occurs in rising air over the warm water in the western Pacific, and the eastern Pacific is relatively dry, as illustrated in Figure 2.1.



**Figure 2.1.** "Neutral" ENSO conditions in the equatorial Pacific Ocean. (http://www.pmel.noaa.gov/tao/proj_over/diagrams/index.html. Accessed on Dec. 1, 2010).

During El Niño, the trade winds relax, causing the thermocline to drop in the eastern Pacific, and rise in the west, as illustrated in Figure 2.2. The eastward displacement of the warmest water results in large changes in the global atmospheric circulation, which in changes cthe limate in regions distant from the tropical Pacific through the movement of atmospheric wave-trains (e.g., Houseago et al., 1998).

**El Niño Conditions**

**Figure 2.2.** El Niño conditions in the equatorial Pacific Ocean. (http://www.pmel.noaa.gov/tao/proj_over/diagrams/index.html. Accessed on Dec. 1, 2010).

La Niña is characterized by unusually cold ocean temperatures in the Equatorial Pacific, compared to El Niño. During La Niña, the eastern Pacific is cooler than usual, and the cool water extends farther westward than usual and causes the depression of the thermocline in western Pacific. This leads to drier than normal conditions in the eastern Pacific, as illustrated in Figure 2.3.

A commonly used index to quantify the intensity of ENSO events is the Southern Oscillation Index (SOI), which compares the atmospheric pressure in Tahiti to that of Darwin, Australia, expressed as a standardized anomaly from normal pressure. Strong positive values are associated with La Niña events and strong negative values are associated with El Niño events. Other indicators of ENSO activity include equatorial Pacific sea surface temperature indices, e.g., NINO12, NINO3, and the Multivariate ENSO index (MEI), which integrates variations of oceanic and atmospheric variables (Piechota, 1999; Philander, 1990).

15

**Figure 2.3.** La Niña conditions in the equatorial Pacific Ocean. (http://www.pmel.noaa.gov/tao/proj_over/diagrams/index.html. Accessed on Dec. 1, 2010).

**Pacific Decadal Oscillation (PDO)**

The Pacific Decadal Oscillation, or PDO, is often described as a long-duration pattern of Pacific climate variability, similar to El Niño (Zhang et al. 1997). Specifically, it is defined as the standardized difference between sea surface temperatures (SSTs) in the north-central Pacific and Gulf of Alaska (Mantua et al. 1997). As with ENSO, the phases of the PDO are classified as being either *warm* or *cool*, as defined by ocean surface temperatures in the northeast and tropical Pacific Ocean. Specifically, a PDO index value is defined as the leading principal component of North Pacific monthly sea surface temperature variability (poleward of 20N), with warm (cool) phase conditions corresponding to positive (negative) index values (Mantua et al. 1997). Although the PDO is similar to ENSO, two main characteristics distinguish the PDO from ENSO. First, PDO phases last much longer (typically 20 to 30 years for a single warm or cool phase) than ENSO events (6 to 18 months for a single phase) (Mantua et al. 1997, Minobe 1997). Second, the temperature patterns of the PDO are most visible in the North

16

Pacific/North American sector, while ENSO patterns exist in the tropics. Several studies find evidence for just two full PDO cycles in the past century (e.g. Mantua et al. 1997, Minobe 1997): cool PDO regimes prevailed from 1890-1924 and again from 1947-1976, while warm PDO regimes occurred from 1925-1946 and from 1977 through the mid-1990's. Recent changes in Pacific climate suggest a switch to cool PDO conditions in 1998. Figure 2.4 illustrates ocean and atmospheric patterns corresponding to PDO warm and cool phases.



**Figure 2.4.** Typical wintertime Sea Surface Temperature (colors), Sea Level Pressure (contours) and surface wind stress (arrows) anomalies and anomaly patterns during warm and cool phases of PDO (http://jisao.washington.edu/pdo/. Accessed on Dec. 1, 2010).

**North Atlantic Oscillation**

The NAO is the dominant mode of surface level pressure (SLP) variability in the North Atlantic region. The NAO index, defined as the SLP difference between the subtropical high pressure system located in the tropical Atlantic near the Azores and the subpolar low pressure system located near Iceland (Rogers, 1984), describes the magnitude of a north-south atmospheric pressure gradient across the North Atlantic

Ocean (Hurrell, 1995). Like ENSO, there is an atmospheric pressure oscillation between Iceland and the Azores, such that if the atmospheric pressure near Iceland is low, then the atmospheric pressure near the Azores is usually high, and vice-versa. While low pressure in the north and high pressure in the south characterizes average conditions, deviations from this mean can result in significant shifts in Northern Hemisphere climate. The corresponding index varies from year to year, but exhibits a tendency to remain in one phase for intervals lasting several years. A positive NAO (illustrated in Figure 2.5) is associated with strong westerlies, and a negative NAO is linked with a reorganization of the Jet Stream and associated changes in regional temperatures, storm tracking, and heat and moisture transport (Kushnir, 1999; Hurrell et al., 2001).



**Figure 2.5.** Positive phase NAO effects on sea level pressure for January, April, July and October. Values shown are the correlation (x100) of sea level pressure with the NAO index in the month indicated (http://www.cpc.ncep.noaa.gov/data/teledoc/nao.shtml. Accessed on Dec. 1, 2010).

**Atlantic Multidecadal Oscillation (AMO)**

The Atlantic Multi-decadal Oscillation (AMO) is a mode of sea surface temperature (SST) variability in the North Atlantic Ocean exhibiting a period of 60-80 years (Kerr, 2000; Knight et al., 2005). Warm AMO phases occurred from 1860 to 1880 and 1930 to 1990, while cool phases occurred from 1905 to 1925 and 1970 to 1990. Recent studies suggest that the AMO returned to a warm phase in 1995 (McCabe et al., 2004). The AMO index consists of detrended SST anomalies for the Atlantic Ocean region. Tootle et al. (2006) found that AMO affects continental U.S. streamflow variability--the middle Atlantic and central U.S. streamflow are influenced by the cold phase of AMO, while the upper Mississippi River basin, peninsular Florida, and Northwest U.S. streamflow are affected by the warm phase of AMO. Other studies have related the AMO to drought in the U.S. and demonstrated the potential coupling of AMO and PDO with ENSO (McCabe et al 2004; Hidalgo et al., 2004; Tootle et al., 2005).

**Pacific North American Pattern (PNA)**

The PNA teleconnection pattern is one of the most prominent modes of low-frequency climate variability, especially during the Northern Hemisphere winter (Horel and Wallace, 1981). It appears as anomalies in the geopotential height fields, and is usually depicted as the 500 and 700 mb levels. The PNA teleconnection pattern has two phases. The positive phase, illustrated in Figure 2.6, consists of higher than normal geopotential heights over the western U.S. and below normal geopotential heights over the eastern U.S., and the negative phase involves below normal geopotential heights over the western U.S. and above normal geopotential heights over the eastern U.S. The PNA is closely related to the upper-level flow patterns and surface temperature and precipitation conditions in North America (Yin, 1994a). During the positive phase, above normal temperatures are expected in the western United States, while the southeastern United States may experience drought conditions due to an upper-level pressure ridge. Also, the eastern and southeastern United States may experience cooler than average conditions due to intrusions of polar air masses, along with enhanced cyclonic activity. During the

negative phase, the western United States tends to be cool and wet, while the eastern United States tends to be warm and dry (Yin, 1994a).

The PNA has been found to be strongly influenced by ENSO, with the positive phase of the PNA pattern associated with Pacific warm episodes (El Niño), and the negative phase associated with Pacific cold episodes (La Niña). Researchers have demonstrated that the PNA pattern is important to understanding the low-frequency variability of the mean tropospheric flow over North America, and therefore it is very useful in explaining temperature patterns and precipitation patterns over North America (Yin, 1994b).



**Figure 2.6.** Positive phases of PNA patterns for January, April, July, and October. The plotted value at each grid point represents the temporal correlation (x100) between the monthly standardized height anomalies at that point and the PNA index for the specified month. (http://www.cpc.noaa.gov/data/teledoc/pna_map.shtml. Accessed on Dec. 1, 2010).

**Influences of Teleconnections on Streamflow in Central Texas**

At least two previous studies have identified teleconnections for central Texas. Piechota and Dracup (1996) found strong correlation between the Southern Oscillation Index (SOI) and the Palmer Drought Severity Index (PDSI), indicating the potential of improved climate forecasts for the region, up to a year in advance. However, a strong relationship between SOI and streamflow was not found. One possible reason for this is that PDSI is a mathematical function of temperature and precipitation and provides a general indication of drought, whereas streamflow tends to integrate climatic processes over interseasonal time scales, and this seasonal averaging may limit forecast accuracy. For instance, streamflow is a function of both surface runoff and groundwater discharge, and groundwater recharge and discharge processes often exhibit lag times markedly longer than those of rainfall-runoff processes. Furthermore, groundwater basins seldom align directly with surface watersheds, which may confound statistical analyses of climate and streamflow variables measured at specific gage locations.

In another study, Rajagopalan et al. (2000) found correlation between summer PDSI and winter Pacific Ocean sea surface temperature anomalies (Niño-3 index). However, they also found epochal variations in this correlation, with the period of 1963-1995 showing weaker teleconnections than the period 1895-1962. Without a means of predicting these epochal shifts in teleconnections, such variation tends to confound statistical forecasting methods based on the entire historical record. It is widely hypothesized that interdecadal North Pacific variability modulates ENSO-precipitation teleconnections (e.g., Gershunov and Barnett, 1998), but Rajagopalan et al. (2000) were not able to conclude that either NAO or PDO has any effect on ENSO-precipitation teleconnections in central Texas.

Tootle et al. (2005) completed a study of the influence of interdecadal, decadal, and interannual oceanic-atmospheric influences on streamflow in the United States. Unimpaired streamflow was identified for 639 stations for the period 1951–2002, and the phases (cold/negative or warm/positive) of ENSO, PDO, NAO, and the Atlantic

Multidecadal Oscillation (AMO) were identified for the year prior to the streamflow year. Statistical significance testing of streamflow indicated no spatially coherent teleconnections for central Texas. Although this study focused on a specific annual period (October-September) and forecast lead time (one year), it provides an indication that skillful long-lead forecasts may not be available in this region

## 2.2 Nonparametric Statistical Methods in Water Resources

Traditionally, statistical methods are based on rigid assumptions about the form of dependence between variables or the underlying joint or marginal probability density functions, and include assumptions of homogeneity and stationarity. In practice, hydroclimatic data or time series often show "unusual" features in the underlying dependence structure, such as asymmetry (or a large positive or negative coefficient of skewness) or multimodality, which are difficult to represent or model using analytical probability density functions (Sharma, et al., 1997; Lall and Sharma, 1996). For example, in many parametric models, streamflow data (monthly or seasonally) is assumed to be normally (Gaussian) distributed (Salas, 1985). However, streamflow data usually exhibits non-Guassian features that vary from month to month and are skewed towards the low flows, with an extended tail in the high flows (Prairie et al. 2006). To partially address these drawbacks, data are often transformed to a Gaussian distribution using a log or power transformation before fitting a parametric model to the transformed data, and the statistics generated from the model are then back-transformed into the original space. However, this process does not guarantee preservation of the original statistics (Sharma et al., 1997; Salas, 1985; Bras and Iturbe, 1985; Benjamin, 1970).

Nonparametric methods strive to approximate a target function locally, i.e., using data from a "small" neighborhood of the point of estimate (Lall, 1995). They impose only weak assumptions, such as continuity of the target function and its differentiability to some order in the neighborhood, rather than *a priori* assumption of the global form of the entire target function, as do parametric methods (e.g., linear regression or fitting probability density functions). As a result, outliers do not exert undue influence on the overall fit, any arbitrary underlying functional form may be captured, and local features present in the data may be reproduced. The trade-off for these features of nonparametric methods is increased computational requirements. However, with increasing computational power readily available, nonparametric techniques offer an attractive and efficient alternative to traditional parametric approaches.

Nonparametric methods, such as kernel density estimation and *K*-nearest-neighbor (K-NN) bootstrap methods, have been successfully applied to a variety of hydrologic problems. Kernel-based methods have been applied to rainfall modeling (Lall et al., 1996; Harrold et al., 2003); flood frequency (Lall et al., 1993; Moon and Lall 1994); streamflow simulation (Sharma et al., 1997; Tarboton et al., 1998); groundwater applications (Adamoski and Feluch 1991); and streamflow forecasting (Smith 1991). *K*-nearest-neighbor methods have been used in streamflow simulation (Lall and Sharma 1996; Prairie 2002) and multivariate stochastic daily weather generation (Rajagopalan and Lall 1999; Yates et al.2003). More recently, a modified *K*-nearest-neighbor method has been developed to overcome the drawback of the *K*-nearest-neighbor method, which is that values not seen in the historical records cannot be simulated (Prairie 2002). Granz et al. (2006) applied this modified K-NN approach to Truckee-Carson River basin streamflow forecasting; Prairie et al. [2006] applied the approach for stochastic streamflow simulation at the Lees Ferry gauge on the Colorado River; and Singhrattna et al. [2005] employed the approach to develop summer rainfall forecasts in Thailand. Their results showed that the modified K-NN approach had better performance in terms of capturing the features (especially nonlinearity) present in the data in comparison with both a parametric periodic autoregressive and a nonparametric index sequential method. One of the drawbacks of the mollified K-NN approach is that the number of neighbors used to bootstrap the residuals will be small and consequently will limit variety in the ensembles when the sample size is small (Prairie et al., 2006).

A powerful new class of non-parametric approaches known as data mining, also referred to as knowledge discovery, has attracted a great deal of attention in the information industry. Data mining involves extracting "hidden" information from large amounts of data (Hand et al. 2001). The data mining process consists of data selection, data cleaning, data transformation and reduction, developing a data mining model, interpretation and evaluation of model results, and knowledge presentation (Han et al., 2006). In general, data mining models can be classified into two categories: descriptive and predictive. Data mining algorithms, which are the mechanisms for creating data

mining models, include a wide array non-parametric methods, including nearest neighbor methods, genetic algorithms, decision trees, cluster analysis, and artificial neural networks. Each of these is described briefly as follows:

- Genetic algorithms are optimization techniques based on the process of natural evolution (Ting, 2005). A genetic algorithm has been shown to be successful in simple reservoir rule generation (Wardlaw et al., 1999). To apply genetic algorithms for streamflow forecasting, parameters in the forecast model (i.e., coefficients and exponents applied to predictor variables) could be selected using a genetic algorithm to maximize the predictive skill (or minimize prediction error) of the model.

- Decision trees represent decisions in a flowchart-like tree structure through a series of "if-then-else" constructs. The basic principle of using decision trees in data mining is to partition datasets to maximize the purity (homogeneity) of a response variable within each partition. Decision tree methods include classification and regression trees and chi square automatic interaction detection (Hand et al. 2001). Relatively little research has been done using decision trees for streamflow forecasting.

- Cluster analysis divides a dataset into mutually exclusive groups such that the members of each group will be similar (or related) to one other and different from (or unrelated to) the members in other groups. The greater the similarity (or homogeneity) within a group, and the greater the difference between groups, the "better" or more distinct the clustering (Han and Kamber, 2006). Various cluster analysis algorithms may be applied to streamflow forecasting in a manner analogous to decision trees.

- Artificial neural networks (ANN) are non-linear predictive models that learn through training and resemble biological neural networks in structure (Han and Kamber, 2006). Recently, numerous ANN-based models have been employed in water resource management because of its power and flexibility (Coulibaly et al., 2000). Notably, the applicability of ANNs in hydrology has been extensively

evaluated by the American Society of Civil Engineers Task Committee on the Application of ANNs in Hydrology (ASCE, 2000), as well as by Dawson and Wilby (2001). These studies reported that ANN can be an efficient and promising alternative to traditional (more physically based) hydrologic models. A disadvantage of ANN is their "black box" nature, which makes it impossible to interpret relations between the individual predictors and response variable.

## 2.3 Stochastic Optimization in Water Resources

In recent years, sustainable development and environmental conservation have become increasingly important due to population growth, climate change, and increased awareness of environmental problems. In many places it is difficult and controversial to construct new large-scale water storage projects as was done in the past. How to improve the operational effectiveness and efficiency of existing reservoir systems for maximizing benefits is a crucial issue. Most reservoirs are designed as multi-purpose systems operating for water supply, flood control, irrigation, hydropower generation, navigation, recreation and environmental and ecological concerns. More often than not, there is conflict and competition among these diverse purposes, particularly during drought conditions. This is one reason why multi-objective reservoir systems often perform more poorly than anticipated (WCD, 2000). In addition, the inherent uncertainties associated with future hydrologic conditions, including possible impacts of climate change, may increase the difficulty of reservoir operation.

Optimization techniques have been applied in water resources planning and management for several decades. A common and long used method of determining "optimal" reservoir operation policies is to use deterministic optimization (DO) models with the historical flow record or a sequence of synthetic data. This approach is intuitive, computationally tractable and evidently sound for systems in which supply typically exceeds demand significantly and the primary planning goal is to avoid some catastrophic failure. As a result, this method has been widely used in practice (e.g., Grygier and

Stedinger 1985, Martin 1991, Karamouz et al. 1992, Kirshen 1992). Several different optimization techniques are used to implement DO models, including linear programming (LP), network flow programming (NFP), dynamic programming (DP), and genetic algorithms. A drawback of all DO models, however, is that they select values of decision variables (e.g. reservoir releases, storage levels) with perfect knowledge of the future and ignore uncertainty. This has been known to lead to solutions that are suboptimal, or even infeasible (Dantzig 1955, Beale 1955). Techniques such as sensitivity analysis and parametric programming can be used to estimate the risks of sub-optimality or infeasibility, but these techniques do not provide a means of reducing or controlling the risks (Watkins and McKinney, 1997). Another drawback is the resulting optimal operational policies inferred from the DO approach are unique to the assumed hydrologic time series unless the period-of–analysis is extremely long (Lund and Fereira 1996). Although multiple regression analysis could be applied to the optimization results for developing operating rules, this method may result also in poor correlations that invalidate the operating rules (Labadie 2004).

Since reservoir operation planning and management is inherently stochastic given the uncertain nature of reservoir inflows, a large number of studies have used stochastic optimization methods. Stochastic optimization methods are designed to operate directly on probabilistic descriptions of random streamflow processes (as well as other random variables) rather than deterministic hydrologic sequences (Labadie 2004). This means that optimization is performed without the presumption of perfect foresight of future events. A wide variety of stochastic optimization methods, such as chance-constrained programming, stochastic linear programming, and stochastic dynamic programming have been applied to water resources problems.

Chance-constrained programming (CCP) is one approach which explicitly accounts for uncertainty in hydrologic inputs. CCP replaces the deterministic constraints involving uncertain parameters with probabilistic constraints, which are then transformed to their deterministic equivalent form using the distributions (means and variances) of the random variables. Release policies have been derived from linear decision rules (ReVelle

et al 1969, Houck and Datta 1981), which permit simple formulation of the chance-constrained problem. However, CCP can be overly pessimistic and conservative when more than one chance-constraint exists in the model (Loucks et al. 1981, Hogan et al. 1981), leading to operational rules that exceed the prescribed reliability levels (Labadie 2004). Furthermore, a number of studies have demonstrated that the operating policies derived from CCP models do not perform as well as some simple alternatives (Loucks and Dorfman 1975, Stedinger et al. 1983, Stedinger 1984).

One technique which overcomes the limitations of CCP is stochastic dynamic programming (SDP). SDP models strive to overcome the problem of dimensionality due to multiple decision stages. Based on estimated Markovian conditional probabilities of inflows, SDP uses a recursive relationship in each time stage to determine the policy which maximizes the expected benefit for each state of the system. SDP formulations can easily incorporate nonlinear and discrete features of water resources problems, and techniques have been developed for incorporating chance constraints (e.g., Askew 1974, Sniedocih 1979). SDP can also use predicted inflow instead of the previous flow as a hydrologic state variable (Stedinger et al. 1984). However, SDP models require discretization of state variables, which can lead to the "curse of dimensionality" if there are more than two or three state variables (Yeh 1985, Pereira and Pinto 1985). Principles from SDP were used to develop sampling stochastic dynamic programming (SSDP), which can capture the complex temporal and spatial structure of the streamflow process by using a large number of sample streamflow sequences (scenarios) (Labadie 2004; Kelman et al., 1990). SSDP can computationally outperform more traditional SDP methods; however, it does not alleviate the dimensionality problem associated with multiple state variables (Labadie 2004).

Stochastic linear programming (SLP) (also called stochastic programming with recourse or two-stage linear programming) is typically used for problems with multiple state variables (e.g., multiple reservoirs) In this method, only the first stage decisions are actually implemented, since future decisions are not known with certainty. Following implementation of the first stage decisions, the problem is reformulated starting with the

next stage and solved over the remainder of the operational horizon (Labadie 2004). To apply SLP with recourse, a number of scenarios corresponding to sequences of realizations of random variables at each stage are required. For multi-stage models, these scenarios can be represented by scenario trees, as illustrated in Figure 2.7 (Watkins et al. 2000; Kracman et al. 2006). The primary advantage of scenario-based stochastic programming over the other approaches is the flexibility it offers in modeling the decision process and defining scenarios, particularly if the number of state variables is high, e.g., more than a few (Watkins et al, 2000). One disadvantage, however, is that a large number of possible scenarios can result in a very large-scale linear programming problem requiring special solution algorithms. This can be overcome through decomposition methods such as L-shaped algorithm (Bender 1962, Van Slyke and Wets 1969), which can allow the large-scale problem to be decomposed by scenario and/or decision period (Birge 1985, Gassman 1990).



**Figure 2.7.** Scenario tree for a multistage reservoir optimization model (Watkins et al., 2000)

Multistage stochastic optimization models using linear programming have been developed for the LCRA (Watkins et al., 2000; Kracman et al., 2002). The water supply planning models were based on a simplified representation of the Highland Lakes System as shown in Figure 2.8. The general formulation of these economic optimization models was as follows:

$$\text{Max} \quad Z = c_1 x + \sum_s p_s c_2 y_s \tag{2.1}$$

Subject to

$$A_x + B y_s = b_s \qquad \forall s \tag{2.2}$$

$$y_s = y_{s'} \qquad \text{for } s \equiv s' \tag{2.3}$$

$$X, y_s \geq 0 \qquad \forall s \tag{2.4}$$

where $x$ is the first–stage water contract decision supported by model; $y_s$ are the subsequent stage contract and release decisions, and the resulting reservoir storage levels corresponding to scenario $s$; $p_s$ is the probability of scenario s; and $A$, $B$, $c$, and $b_s$ are the model parameters and data, some of which vary across scenarios. Equation (2.1) is the objective function for the model, including components representing the expected benefits from run-of-river and interruptible water diversions made to the irrigation districts, penalties for water demand deficits (municipal and irrigation), recreation benefits, and hydropower generation benefits. Equations (2.2) are constraints that represent reservoir mass balances and water demand requirements. Equations (2.3) are the non-anticipativity constraints which ensure that decisions are the same for scenarios that identical up to the point in time that the decisions are made.

The multiple-stage linear programming model represented by Eqs (2.1)–(2.4) requires stochastic inputs in the form of a scenario tree, with levels in the tree corresponding to decisions stages. A scenario in the model is a sequence of monthly available tributary inflows which are representative of flows which could occur in the future (Watkins et al., 2000). This stochastic optimization model, however, is not run in real time, but rather was formulated as a planning model to help "enhance the understanding of water

planning in the LCRA service area and provide a rational method of developing and comparing robust and reliable reservoir operation alternatives for the LCRA in the face of uncertainty" (McKinney et al., 2002). In contrast, the forecast models developed in this work may be applicable for real-time decision making.



**Figure 2.8.** Schematic of the Highland Lakes System (Kracman 2002)

Only a few studies have utilized forecasts within the context of stochastic optimization models. Faber and Stedinger (2001) used National Weather Service Ensemble Streamflow Prediction (ESP) forecasts in SSDP models. Prior to that, Kelman et al. (1990) discussed how inflow forecasts could be used in a multi-stage stochastic linear programming model for hydropower system operations. Jacobs et al. (1995) describe a multi-stage stochastic optimization model for scheduling hydroelectric power generation under uncertainty, where the scenario tree includes short- to medium-term streamflow forecasts. More recently, a number of studies by Kim and colleagues have assessed the value of seasonal forecasts using SSDP (Kim et al., 2007; Eum and Kim, 2010).

# 3. Seasonal Forecasts Using Logistic Regression[1]

## 3.1 Introduction

Reliable streamflow forecasts with lead times of even one season can have a significant effect on the performance of reservoir operation policies and operation efficiency (Karamouz *et al.* 2004; Sun *et al.* 2006). In recent years, much effort has been made to develop mid- to long-term (seasonal to annual) hydroclimatic and streamflow forecasting models for water management in the United States. For example, the NOAA Climate Prediction Center issues seasonal forecasts of temperature, precipitation, and soil moisture, as well as a drought outlook, for the entire U.S. In the western U.S., the USDA Natural Resources Conservation Service provides streamflow forecasts in the first half of the year based on observed snowpack conditions. At many stream gage locations throughout the U.S., the National Weather Service provides probabilistic seasonal flow forecasts through a procedure known as Ensemble Streamflow Prediction, or ESP (Day 1985; Smith *et al.* 1992). Traditionally, each of the meteorology traces has been assumed to represent an equally likely scenario for the future; more recently, methods have been developed to condition the probabilities of the historical meteorological traces based on seasonal climate forecasts (e.g., Croley 2000; Duan *et al.* 2006). On a global scale, the NOAA/Columbia University International Research Institute for Climate and Society (IRI) is one institution that issues seasonal forecasts of temperature and precipitation.

Significant progress has been made in understanding the influence of large-scale ocean-atmospheric patterns, such as El Niño–Southern Oscillation (ENSO), Pacific Decadal Oscillation (PDO), the North Atlantic Oscillation (NAO), and Atlantic Multidecadal Oscillation (AMO), on regional climate anomalies around world. Numerous studies have shown that statistical models incorporating large-scale ocean-atmospheric patterns can improve the ability to forecast streamflow with long lead times (e.g., Hamlet

---

[1] This chapter is constituted by the article by Wei and Watkins (2010) "Probabilistic Streamflow Forecasts Based on Hydrologic Persistence and Large-Scale Climate Signals in Central Texas," to be published in the *Journal of Hydroinformatics*. It has been reprinted from the *Journal of Hydroinformatics,* with permission from the copyright holders, IWA Publishing (Appendix A).

& Lettenmaier 1999; Sharma 2000; Tootle *et al.* 2006).  At least three previous studies have identified teleconnections for Central Texas. Piechota and Dracup (1996) found strong correlation between the Southern Oscillation Index (SOI) and the Palmer Drought Severity Index (PDSI), indicating the potential of improved hydroclimatic forecasts, with up to one year in lead time, for the region.  However, a strong relationship between SOI and streamflow was not found. One possible reason for this is that PDSI is a mathematical function of temperature and precipitation and provides a general indication of drought, whereas streamflow tends to integrate climatic processes over interseasonal time scales, and this seasonal averaging may limit forecast accuracy. For instance, streamflow is a function of both surface runoff and groundwater discharge, and groundwater recharge and discharge processes often exhibit lag times markedly longer than those of rainfall-runoff processes (Alley 1985).  Furthermore, groundwater basins seldom align with surface watersheds, which may confound statistical analyses of climate and streamflow variables measured at specific gage locations.

In another study, Rajagopalan *et al.* (2000) found correlation between summer PDSI and winter Pacific Ocean Sea Surface Temperature (SST) anomalies.  However, they also found epochal variations in this correlation, with the period of 1963-1995 showing weaker teleconnections than the period 1895-1962. Of course, without a means of predicting these epochal shifts in teleconnections, such variation tends to confound statistical forecasting methods based on the entire historical record.  It is widely hypothesized that interdecadal North Pacific variability modulates ENSO-precipitation teleconnections (e.g., Gershunov & Barnett 1998), but Rajagopalan *et al.* (2000) were not able to conclude that either NAO or PDO has any effect on ENSO-precipitation teleconnections in Central Texas.

Finally, Tootle *et al.* (2005) completed a study of the influence of interdecadal, decadal, and interannual oceanic-atmospheric influences on streamflow in the United States. Unimpaired streamflow was identified for 639 stations for the period 1951–2002, and the phases (cold/negative or warm/positive) of ENSO, PDO, NAO, and AMO were identified for the year prior to the streamflow year (i.e., long lead time). Statistical

significance testing of streamflow, based on the interdecadal, decadal, and interannual oceanic-atmospheric phase (warm/positive or cold/negative), indicated no spatially coherent teleconnections for Central Texas. Although this study focused on a specific annual period (October-September), and a particular forecast lead time, it provides an indication that long-lead climate forecasts may not be useful to water managers in this region.

Streamflow forecasts may be either deterministic or probabilistic, but probabilistic methods are often preferable for water management because they can provide more information about uncertainty. Categorical streamflow forecasts are common, providing the probabilities of flow being in each of a number of categories (e.g., low, medium, high). Probabilities of each category could be generated directly or indirectly. Piechota *et al.* (1998) proposed linear discriminant analysis to produce the probabilities of each category of streamflow directly. This method involves nonparametric kernel density estimation of the probability density function for each flow category. Regonda *et al.* (2006) employed logistic regression to directly predict the probability of streamflow above a given threshold. They applied this approach to categorical forecasts of the spring (April-June) streamflow at six locations in the Gunnison River Basin. However, this approach treats the response variable as binary, i.e., equal to 1 if the streamflow value exceeds a given threshold and zero otherwise. A drawback of this approach for multiple categories is that the logistic regression needs to be repeated to obtain the probability corresponding to each category threshold, and the sum of probabilities is not guaranteed to equal 1.

In this paper, a statistical method called polytomous logistic regression for ordinal response (Kutner *et al.* 2004) is proposed to generate probabilistic forecasts with seasonal lead times for the Highland Lakes system in Central Texas. In the method, the response variable (streamflow) has multiple discrete outcomes rather than binary. Further, the response categories (e.g., below normal, normal, above normal) could be considered as ordered, thus allowing a parsimonious and easily interpreted logistic model, called a proportional odds model, that may be employed to generate a probabilistic (multi-

category) forecast using a single model. A number of distributions oriented metrics, such as the Brier Skill Score and the Ranked Probability Skill Score, may be used to assess model performance (Wilks 1995). This is demonstrated for the Highland Lakes system in Texas for a number of potential predictor variables, including streamflow autocorrelation (hydrologic persistence), sea surface temperatures, and other large-scale climate signals.

## 3.2 Case Study Data

The Lower Colorado River Authority (LCRA) operates the Highland Lakes system in Central Texas, a series of six lakes on the Lower Colorado River. As a water conservation and reclamation district, the LCRA provides water supply and flood control to a 33-county area, including the City of Austin and several rice irrigation districts along the Texas Gulf Coast (see Figure 3.1). In addition, the LCRA produces wholesale power for a 53-county service area and provides water resources for lake recreation activities and in-stream flow maintenance. To meet rapidly growing demands, reliable reservoir inflow forecasts with seasonal lead times would potentially be very beneficial; however, hydrologic forecasts are not used by the LCRA for a number of reasons, including high seasonal and annual variability of stream flow, the absence of easily measured hydrologic indicators such as snowpack, and a lack of experience with probabilistic planning methods (Watkins & O'Connell 2005).

**Figure 3.1 |** Lower Colorado River Authority District

To explore the patterns of streamflow and the influence of teleconnections in Central Texas, monthly streamflow data are acquired from two sources: 1) aggregate inflows to the Highland Lakes (upstream), based on USGS gage measurements and adjustments made by LCRA staff to account for runoff from ungaged areas; and 2) unregulated tributary flows to the Colorado River downstream of the Highland Lakes, as determined by the Texas Water Availability Model (WAM) (Wurbs 2005). The reservoir inflow data spans a total of 57 years, from 1950 to 2006, and the naturalized downstream flow data spans 59 years, from 1940 to 1998. For most of the analyses, the flow data are normalized through a two-step process—first a logarithmic transformation, then conversion to a standardized anomaly by subtraction of the mean (of the log values) and division by the standard deviation (of the log values). While this transforms the data so that the statistical assumption of normality is more valid, it should be noted that the correlation coefficients are generally inflated by the log-transform, and thus an effort is made to illustrate the results in terms of the raw flow data.

Figures 3.2 and 3.3 show the monthly and seasonal autocorrelations of these two time series. Monthly autocorrelations of (upstream) reservoir inflows range from a high of nearly 0.8 for February and March flows to a low of essentially zero for July and August flows. Seasonal correlation coefficients also peak in the winter season, with a value of 0.66. (The correlation coefficient that is significant at the $p = 0.05$ level is 0.20.) It may be surprising that the seasonal correlation between OND and JFM flows is higher than the average of monthly correlation coefficients during this period. One reason for this may be that averaging over a three-month period reduces the "noise" that results from individual storm events which have a significant effect on monthly flow totals.

The tributary flows to the Colorado River downstream of the Highland Lakes, estimated as the WAM naturalized flows at Mansfield Dam minus the naturalized flows at Bay City, have monthly autocorrelation coefficients that reach a maximum of 0.8 for February and March and have a minimum of near 0.2 for September and October. The average monthly autocorrelation for the downstream data is about 0.27 higher than the upstream data. Seasonal autocorrelations peak in the spring season with a value of 0.82, which is lagged by one season in comparison with the upstream data. All seasonal autocorrelation coefficients for the downstream data are significant at $p = 0.10$ or less, and the average autocorrelation coefficient is about 0.23 higher than the upstream data.

**Figure 3.2 |** Monthly autocorrelations for aggregate inflows to the Highland Lakes (Upstream) and the Texas Water Availability Model (Downstream) data. Correlation coefficients are computed using raw flow data.



**Figure 3.3 |** Seasonal autocorrelations for aggregate inflows to the Highland Lakes (Upstream) and the Texas Water Availability Model (Downstream) data. Correlation coefficients are computed using raw flow data.

Based on these autocorrelation coefficients, seasonal streamflow forecasts for certain times of the year may be based solely on hydrologic persistence. We investigate the predictive skill of these forecasts, as well as the potential for large scale ocean-atmosphere interactions to provide additional forecast skill. The oceanic-atmospheric phenomena investigated as potential predictor variables for streamflow in Central Texas are the El Niño-Southern Oscillation (ENSO), the Pacific Decadal Oscillation (PDO), the North Atlantic Oscillation (NAO), the Atlantic Multidecadal Oscillation (AMO), and the Pacific North American (PNA). The ENSO and NAO generally have a two- to seven-year periodicity (Philander, 1990), while PDO and AMO exhibit long-term periodicity (about 25 to 60 years) (Mantua et al. 1997; Kerr 2000; Gray et al. 2004).

Various indices were selected to quantify the magnitude of these ocean-atmospheric oscillations. The Niño 3.4 index, which characterizes the tropical Pacific Ocean sea surface temperature (SST) anomalies between latitudes 5S and 5N and longitudes 170W and 120W, was selected as an indicator of ENSO, and monthly index data were obtained from the National Weather Service (NWS) Climate Prediction Center (CPC) (http://www.cpc.ncep.noaa.gov/data/indices). The PNA index values, a measure of atmospheric pressure anomalies at four locations in the northern hemisphere (Horel & Wallace 1981) were also obtained from the CPC. The PDO index values were obtained from the University of Washington (http://jisao.atmos.washington.edu/pdo). Finally, NAO index values were obtained from the National Center for Atmospheric Research (http://www.cgd.ucar.edu/cas/jhurrell/indices.html), and the AMO index values (Kerr 2000) were obtained from the National Oceanic and Atmospheric Administration (NOAA) Climate Diagnostics Center (CDC) (http://www.cdc.noaa.gov/Climateindices). In all cases, monthly index values for the period 1940-2006 were used for the analysis.

In addition, SST data was analyzed directly for correlations with streamflow. The data used was the extended reconstructed sea surface temperature (ERSST) analysis (Smith et al. 2008), obtained from the National Climatic Data Center (NCDC) through the KNMI Climate Explorer, an on-line data analysis tool (Oldenborgh & Burgers 2005). Six correlation patterns with high statistical significance ($p < 0.01$) were identified and

referenced as ERSST1 through ERSST6, corresponding to the following seasonal streamflows: 1) winter reservoir inflows, 2) spring reservoir inflows, 3) winter downstream flows, 4) spring downstream flows, 5) summer downstream flows, 6) fall downstream flows.  For each SST pattern, a normalized index was computed based on average seasonal temperatures over a 4-degree by 4-degree area, similar to the procedure of Block and Rajagopalan (2007).

## 3.3   Statistical Forecast Model

Multiple logistic regression is most frequently used to model the relationship between a binary response variable and a set of predictor variables, which may be either numerical or categorical. Let $p = Pr(Y = 1)$ denote the probability of success. The ratio of $p/(1-p)$ is called odds, and the function $\log(p/(1-p))$ is called logit($p$), which is in fact the logarithm of the ratio of probability of success to the probability of failure. A multiple logistic regression model can be expressed as follows:

$$\text{logit}(p) = \log(p/(1-p)) = \beta_0 + \beta_1 x_1 + \cdots \beta_k x_k$$

(3.1)

where the parameters $\beta_j$ are usually estimated using maximum likelihood theory (Menard 1995).

In the binary logistic case, logit($p$) compares the probability of a category-1 response (success) to the probability of a category-2 response (failure). If the response variable has more than two levels, logistic regression can still be employed by means of a polytomous or multicategory logistic regression model (Kutner *et al.* 2004). For a response variable with $J$ categories, it is necessary to develop $J$ - 1 logistic regression models. One category will be chosen as the baseline or reference category, and then all other categories will be compared to it. The choice of reference category is arbitrary. Frequently the last category is chosen.

Using category $J$ to denote the reference category, only $J$ - 1 logits need to be developed. For a nominal response, the $j$th logit expression for the $i$th observation is given as:

$$\log\left[\frac{p_{ij}}{p_{iJ}}\right] = X_i^T \beta_j \qquad \text{for } j = 1, 2, \ldots J - 1 \qquad (3.2)$$

where $\beta_j = \begin{bmatrix} \beta_{0j} & \beta_{1j} & \cdots & \beta_{kj} \end{bmatrix}^T$ and $X_i = \begin{bmatrix} 1 & X_{i1} & \cdots & X_{ik} \end{bmatrix}^T$.

(Note that vectors $\beta_j$ are different for each category $j$.).

Given the $J$ - 1 logit expressions, it is possible (algebra not shown) to obtain the $J$-1 direct expressions for the category probabilities in terms of the $J$-1 linear predictors, $X^T \beta_{jJ}$. The resulting expressions are

$$p_{ij} = \frac{\exp(X_i^T \beta_j)}{1 + \sum_{k=1}^{J-1} \exp(X_i^T \beta_k)} \qquad \text{for } j = 1, 2, \ldots\ldots J\text{-}1 \qquad (3.3)$$

The estimates of the $J$ -1 parameter vectors $\beta_1, \beta_2, \ldots\ldots\beta_{J-1}$ can be obtained simultaneously using maximum likelihood estimation. The sum of probabilities of each category for $i$th observation is equal to 1. For example, for 3 response categories, we use category $J = 3$ as the baseline category, and there are two comparisons to this reference category. Let $p_{ij}$ denote the probability that category $j$ is selected for the $i$th response, and then the logit for the two comparisons are:

$$\log_e \frac{p_{i1}}{p_{i3}} = X_i' \beta_1, \quad \log_e \frac{p_{i2}}{p_{i3}} = X_i' \beta_2 \qquad (3.4)$$

and we constrain $p_{i1} + p_{i2} + p_{i3} = 1$. Then we can obtain the probabilities of each category for $i$th observation by solving above 3 algebraic equations as below:

$$p_{i1} = \frac{\exp(X_i^T \beta_1)}{1 + \exp(X_i^T \beta_1) + \exp(X_i^T \beta_2)} \qquad (3.5)$$

$$p_{i2} = \frac{\exp(X_i^T \beta_2)}{1 + \exp(X_i^T \beta_1) + \exp(X_i^T \beta_2)} \qquad (3.6)$$

$$p_{i3} = \frac{1}{1 + \exp(X_i^T \beta_1) + \exp(X_i^T \beta_2)} \qquad (3.7)$$

If multiple response categories are treated as ordered, the logistic regression model could be reduced to $J$ - 1 cumulative logits as follows:

$$\log\left[\frac{p(Y_i \le j)}{1 - p(Y_i \le j)}\right] = \alpha_j + X_i^T \beta \qquad \text{for } j = 1, 2, \ldots J - 1 \qquad (3.8)$$

The difference between the ordinal response logits and the nominal response logits is that each of the $J$ - 1 parameter vectors $\beta_j$ is unique for the nominal case; for ordinal response, the slope coefficient vector $\beta$ is identical for each of the $J$ - 1 cumulative logits, and only the intercepts $\alpha_i$ differ. Finally, the cumulative probabilities $p(Y_i \le j)$ for the ordinal logistic regression model are given as follows:

$$p(Y_i \le j) = \frac{\exp(\alpha_i + X_i^T \beta)}{1 + \exp(\alpha_i + X_i^T \beta)} \qquad \text{for } j = 1, 2, \ldots J - 1 \qquad (3.9)$$

The goal of this study is to develop a framework for developing categorical forecasts of streamflows. Since these categories may be treated as ordered, thus the ordinal polytomous logistic regression model, which is also called the proportional odds model, can be used to produce tercile probability forecasts (below normal, normal, and above normal categories). This may be more effective, yielding a more parsimonious model with easily interpreted results. The software package VGAM developed in R (http://www.r-project.org) was used to derive the ordinal polytomous logistic regression model. This software employs the maximum likelihood method to estimate the model parameters (Yee 2010).

## 3.4 Predictors Selection

In this study, a total of seven potential predictor variables are examined for each of the two forecast locations and four seasonal forecast periods. Therefore, 128 ($2^7$) alternative models can be constructed with each predictor either included or excluded from each the 8 forecast models.  Automatic search procedures are employed to screen the most promising models according to a specified criterion without requiring the fitting of all of possible regression models. For logistic regression modeling, two commonly used criteria are Akaike's Information Criterion (AIC) and Bayesian Information Criterion (BIC), which are defined, respectively, as follows:

$$AIC = -2Ln\,(L(\boldsymbol{b})) + 2P \tag{3.10}$$

$$BIC = -2Ln\,(L(\boldsymbol{b})) + PLn(n) \tag{3.11}$$

where $\boldsymbol{b}$ denotes the vector of estimated parameters of the logistic regression model (using maximum likelihood method), $L(b)$ is the log-likelihood function, $P$ is the number of estimated parameters, and n is the total number of observations. Promising models will yield relatively small values for these criteria.

In this study, a forward stepwise search procedure is used to select the best logistic regression model (Seber *et al.* 2003). Essentially, this search method develops a sequence of regression models, at each step adding or deleting a predictor variable according to a decision rule. For logistic regression, the decision rule is based on the likelihood-ratio test and its significance ($p$-values), which are obtained from a chi-square distribution with the associated degree of freedom. In the forward stepwise procedure, a predictor variable will be added to the model at each step only if the chi-square statistic is greater than a critical value or if the corresponding $p$-value is less than a predetermined level (usually 0.05). Additionally, a predictor variable in the model will be deleted when its $p$-value associated with its test statistic exceeds a predetermined level.  The procedure will terminate until no further predictor variables can be added with resulting $p$-values less than a predetermined

level, i.e., there are no predictors considered sufficiently helpful to enter the regression model.

The forward stepwise method is applied to select predictor variables from a set of potential predictors, which include streamflows and large-scale climate signals observed in the seasonal period prior to the forecast. For convenience, the seasons are defined as winter (January-March), spring (April–June), summer (July–September), and fall (October–December). Semi-annual streamflow (January-June) is also predicted based on observations from the previous fall. A summary of the logistic regression models selected using the forward stepwise method for reservoir inflows is presented in Table 3.1. Climate indices and streamflow for the season prior to the predicted seasonal streamflow are designated by (-1). The results show that streamflow persistence is a statistically significant predictor ($p = 0.05$) for winter, spring, fall and Jan.-June streamflow forecasts. Amongst the large-scale climate signals, either the ERSST1 pattern (shown in Figure 3.4) or PDO is a significant predictor in the logistic regression model for winter streamflow forecasts, but not both, likely due to high colinearity between these predictors. For spring streamflow forecasts, however, both the ERSST2 (shown in Figure 3.5) and PNA are significant predictors to be retained in the model. The other large-scale signals (ENSO, NAO, and AMO) are not statistically significant in the logistic regression model. Also shown is the relative improvement of forecast models with the selected predictor(s) in terms of model selection criterion AIC and BIC over forecasts based on seasonal climatology (i.e., forecasts equal to the median historical value). The greatest improvement is observed for winter streamflow forecasts; whereas the stepwise selection procedure indicates there are no significant predictor variables for summer streamflow forecasts.

**Table 3.1 |** Best logistic regression models from forward stepwise regression for seasonal streamflow forecasts of aggregate inflows to the Highland Lakes (Upstream). Values in brackets are p-values for variable entering the model. Percent improvement is relative to climatology (median historical value).

| Model | Selected Predictors (p-value) | | Model Fitting Criteria | | | | Improvement (%) | |
|---|---|---|---|---|---|---|---|---|
| | | | Regression Model | | Climatology | | | |
| | | | AIC | BIC | AIC | BIC | AIC | BIC |
| Winter-A | Fall(-1)  (<0.001)<br>PDO(-1)  (0.017) | | 94.93 | 107.08 | 127.01 | 131.06 | 25.3 | 18.3 |
| Winter-B | Fall(-1)  (<0.001)<br>ERSST1(-1)  (0.010) | | 92.73 | 104.88 | 127.01 | 131.06 | 27.0 | 20.0 |
| Spring | Winter(-1)  (0.003)<br>ERSST2(-1)  (<0.001)<br>PNA(-1)  (0.002) | | 105.69 | 117.95 | 129.24 | 133.33 | 18.2 | 11.5 |
| Fall | Summer(-1)  (0.004) | | 122.41 | 130.58 | 129.24 | 133.33 | 5.3 | 2.1 |
| Jan-Jun | Fall(-1)  (0.004) | | 119.83 | 127.93 | 127.01 | 131.06 | 5.6 | 2.4 |

A summary of logistic regression models selected for seasonal streamflow forecasts of downstream flows is presented in Table 3.2. Results show that streamflow persistence is a significant predictor for all seasons. Amongst the large-scale climate signals, SST indices (based on the correlation patterns shown in Figures 3.6-3.8) have the most significant impact in the logistic regression models. The other large-scale signals do not have a statistically significant impact on the unregulated flows downstream of the Highland Lakes; however, it should be noted that each of the SST patterns used to develop the climate indices can be related to one of the named oscillations.

**Table 3.2 |** Best logistic regression models from forward stepwise regression for seasonal streamflow forecasts of unregulated tributary flows, from the Texas Water Availability Model (Downstream). Values in brackets are p-values for variable entering the model. Percent improvement is relative to climatology (median historical value).

| Model | Selected Predictors (p-value) | Model Fitting Criteria | | | | Improvement (%) | |
|---|---|---|---|---|---|---|---|
| | | Regression Model | | Climatology | | | |
| | | AIC | BIC | AIC | BIC | AIC | BIC |
| Winter | Fall(-1)　　　(<0.001)<br>ERSST3(-1)　(0.013) | 88.47 | 99.70 | 109.47 | 113.21 | 19.2 | 11.9 |
| Spring | Winter(-1)　　(<0.001)<br>ERSST4(-1)　(<0.001) | 83.82 | 91.39 | 111.62 | 115.41 | 24.9 | 20.8 |
| Summer | Spring(-1)　　(<0.001)<br>ERSST5(-1)　(0.007) | 107.39 | 119.86 | 133.60 | 137.76 | 19.6 | 13.0 |
| Fall | Summer(-1)　(0.002)<br>ERSST6(-1)　(0.005) | 100.97 | 108.54 | 111.62 | 115.41 | 9.5 | 6.0 |



**Figure 3.4 |** Correlation map of winter aggregate inflow to the Highland Lakes with fall sea surface temperatures. Circled regions indicate strong positive and negative correlations used to derive the ERSST1 index.

**Figure 3.5 |** Correlation map of spring aggregate inflow to the Highland Lakes with winter sea surface temperatures. Circled regions indicate strong negative correlations used to derive the ERSST2 index



**Figure 3.6 |** Correlation map of winter downstream flows with fall sea surface temperatures. Circled regions indicate strong positive and negative correlations used to derive the ERSST3 index.

**Figure 3.7 |** Correlation map of spring downstream flows with winter sea surface temperatures. Circled regions indicate strong positive and negative correlations used to derive the ERSST4 index.



**Figure 3.8 |** Correlation map of summer downstream flows with spring sea surface temperatures. Circled regions indicate strong positive and negative correlations used to derive the ERSST5 index.

## 3.5 Forecast Verification

Forecast verification aims to evaluate the agreement between forecasts and observations (Katz & Murphy 1997; Stephenson 2003). The Brier skill score (BSS) and the ranked probability skill score (RPSS) are common statistics used to measure the improvement in the accuracy of multicategory probability forecasts over a naïve forecasting method such as climatology. The Brier skill score (BSS) (Dogget 1998) is defined as

$$BSS = 1 - BS / BSC \tag{3.12}$$

where $BS$ is the Brier score, and is defined as

$$BS = (1/n)\Sigma(f_i\text{-}l(\text{obs}_i))^2 \tag{3.13}$$

where $f_i$ is the forecast probability of event $i$ occurring; $l(\text{obs}_i)$ is an indicator variable (1 if event in category $i$ occurs, else 0); $n$ is the number of events; and $BSC$ is the climatologically expected value of $BS$, equal to $BSC = P_i * (1 - P_i)$, where $P_i$ is the climatological probability of the event. The Brier score is often applied to events that exceed a given threshold, but it can also be applied to categorical events. In this study, we consider three tercile categories, i.e., below normal, normal and above normal. Accordingly, $BSC = (1/3)*(1\text{-}1/3) = 0.222$ for each category. A perfect forecast has a value of BSS of 1; positive values between zero and one indicate forecast performance better than climatology, and negative values indicate forecast performance worse than climatology.

The ranked probability score ($RPS$) evaluates the sum of the squared differences in the cumulative probability space, so that

$$RPS = \frac{1}{K-1}\sum_{m=1}^{K}\left[\left(\sum_{k=1}^{m}f_k\right) - \left(\sum_{k=1}^{m}o_K\right)\right]^2 \tag{3.14}$$

where $K$ is the number of forecast categories (below normal, normal and above normal), $f_k$ is the forecast probability for the $k^{th}$ point, and $o_k$ equals zero or one to indicate whether or not the observed flow is in the $k^{th}$ category. The use of *RPS* results in higher penalties for forecasts farther away from actual outcomes, rather than scoring based on only hit and miss. The *RPS* can assume a number between zero and one, with a perfect forecast scoring zero. The ranked probability skill score (*RPSS*) then measures the relative improvement of using a forecast over climatology alone, and is given by

$$RPSS = \frac{\overline{RPS} - \overline{RPS}_{c\lim ato\log y}}{o - \overline{RPS}_{c\lim ato\log y}} = 1 - \frac{\overline{RPS}}{\overline{RPS}_{c\lim ato\log y}} \tag{3.15}$$

The probabilities of streamflow in each category in the climatology forecast (i.e., prior probabilities) are the same and equal 1/3 due to the definition of the flow regimes. Thus, the *RPS* values of the three categories (below normal, normal and above normal) in the climatology forecast are 0.278, 0.111 and 0.278, respectively. A perfect *RPSS* is 1, and negative scores indicate that forecasts performed worse than climatology.

In this study, forecast performance is evaluated using a leave-one-out cross-validation method, in which one observed streamflow value is held out and the remaining data are used to generate a prediction. This process is repeated for each value in the data set, and the cross-validated forecasts are then evaluated using BSS and RPSS.


## 3.6  Results and Discussion

The ordinal polytomous logistic regression models with the minimum AIC and BIC values, indicated in Tables 3.1 and 3.2, were applied to generate tercile probability forecasts for flows in the Highland Lakes system. The Brier skill score (BSS) and the ranked probability skill score (RPSS) were used to evaluate the performance of cross-validated (leave-one-out) forecasts. Summaries of the results for both of the data sets (upstream and downstream of the reservoirs) are presented in Tables 3.3 and 3.4.

Results show that (upstream) reservoir inflows for the winter season can be predicted with significant skill with one season lead time based on either persistence and ERSST1 or persistence and PDO. The BSS and RPSS values indicate that winter streamflow forecasts have an average improvement in skill of 25.6% and 22.5% improvement, respectively, over climatology. Table 3.3 also shows the BSS and RPSS for winter streamflow forecasts based only on persistence (18.7% and 12.7%, respectively), which indicates that including ERSST1 or PDO can provide a significant improvement over a forecast based on streamflow persistence only. For spring streamflow forecasts, BSS and RPSS values indicate that forecast skill can be improved by about 8-13% over climatology based on persistence, ERSST2 and PNA together. Spring streamflow forecasts based only on persistence have no skill. Hydrologic persistence also shows no skill as a predictor for fall streamflows, and the other large-scale climate indices (ENSO, PDO, NAO, AMO and PNA) are not useful predicators for any season.

**Table 3.3** | Brier Skill Score (BSS) and Ranked Probability Skill Score (RPSS) for cross-validated seasonal forecasts of inflows to the Highland Lakes reservoir system. "*" denotes values for a forecast model based on hydrologic persistence only.

| Forecast Model | Forecast Skill Score | |
|---|---|---|
| | BSS (%) | RPSS (%) |
| Winter-A | 22.9 | 19.2 |
| Winter-B | 28.2 | 25.8 |
| Winter * | 18.7 | 12.7 |
| Spring | 13.8 | 8.4 |
| Spring * | 0.2 | -2.0 |
| Fall | -3.2 | -4.7 |
| Jan-Jun | 2.7 | 1.4 |

Results in Table 3.4 show that streamflow persistence (autocorreation) is a useful predictor for downstream unregulated flows for winter, spring, and summer seasons. In

addition, the derived SST indices ERSST3, ERSST4, and ERSST5 provide additional forecast skill for these three seasons, respectively. The forecast skill score indicate between 12 and 33% improvement over climatology-based forecasts for these seasons. For fall forecasts, however, skill scores are very close to zero, indicating no improvement over climatology, despite the correlations found in the regression analysis using all data together (i.e., not holding data out as in cross-validation).

**Table 3.4** | Brier Skill Score (BSS) and Ranked Probability Skill Score (RPSS) for cross-validated seasonal forecasts of downstream unregulated flows. "*" denotes values for a forecast model based on hydrologic persistence only.

| Forecast Model | Forecast Skill Score | |
|:---:|:---:|:---:|
| | BSS (%) | RPSS (%) |
| Winter | 32.3 | 14.1 |
| Winter * | 27.6 | 8.5 |
| Spring | 36.1 | 16.7 |
| Spring * | 33.4 | 12.6 |
| Summer | 16.2 | 14.8 |
| Summer * | 9.2 | 8.6 |
| Fall | -4.8 | -3.9 |

These results for both upstream and downstream flows forecasts are generally similar to the findings of Rajagopalan *et al.* (2000) and Tootle *et al.* (2005), but with some important differences. Both of these previous studies concluded that weak relationships exist between climate indices and streamflow in Central Texas. In this study, only one identified oceanic-atmospheric mode, PDO, was found to provide significant improvement in winter streamflow forecasts (a second, PNA, provided small improvement in spring forecasts). However, derived SST indices, each of which have a spatial correlation pattern similar to ENSO or PDO, were found to provide significant

improvements in forecast skill when included in the regression models. Some differences are also attributed to different lead times—this study focuses on a seasonal lead time while the investigation by Tootle *et al.* (2005) was based on an annual lead time.

Comparisons of observations and cross-validated forecasts are shown in Figures 3.9 and 3.10 for winter reservoir inflows and spring downstream flows, respectively. The ordinal polytomous logistic regression models provide tercile probability forecasts for flows. So, the following method is used to obtain the streamflow volume forecast corresponding to the categorical probabilities.

1) First we calculate the empirical cumulative probabilities based on climatology forecasts. This assumes that each flow observation is equally likely, and the probability of each category is equal to 1/3.

2) Given the forecasted tercile probability of each flow, we adjust the empirical cumulative probabilities. The probability of each category is then adjusted to the forecasted probability instead of 1/3. This is similar to the approach of Croley (2000).

3) To obtain the streamflow volumes corresponding to the 33% and 66% non-exceedance probabilities, we simply interpolate quantiles of a log-normal distribution to according to the adjusted cumulative probabilities.

Results show many years in which the climate-based forecasts provide a significant improvement over climatology. For instance, the forecasts accurately predict high reservoir inflows in 10 of the 12 years in which winter inflows exceeded 442,400 acre-ft ($Ln(442,400) = 13.0$). Perhaps more importantly, the forecasts accurately predict low reservoir inflows ($Ln(\text{Streamflow}) < 11.0$) in 8 of 11 cases. Downstream flow forecasts for the spring season have even better skill, with accurate predictions of high flows ($Ln(\text{Streamflow}) > 13.5$) in 9 of 11 cases, and accurate predictions of low flows ($Ln(\text{Streamflow}) < 11.5$) in all 11 cases.

**Figure 3.9 |** Plot of observations and cross-validated forecasts for winter reservoir inflows to the Highland Lakes. Units on the y-axis are the natural logarithm of flow in acre-feet per month (1 acre-ft equals approximately 1,234 m$^3$).



**Figure 3.10 |** Plot of observations and cross-validated forecasts for spring streamflows downstream of the Highland Lakes. Units on the y-axis are the natural logarithm of flow in acre-feet per month (1 acre-ft equals approximately 1,234 m$^3$).

The results of this cross-validation exercise indicate that hydrologic persistence (streamflow autocorrelation) can provide skillful forecasts of reservoir inflows during the winter and spring seasons and downstream flows during the winter, spring, and summer seasons. One reason for this is that streamflow in the winter months is closely related to soil moisture, which tends to be higher during fall and winter, with persistence from fall through winter and early spring in Central Texas. This is not the case for late spring and summer in the upper watershed, when soils dry and high runoff mainly results from convective storm events. Persistence in soil moisture extends further into the spring, and sometimes early summer, in the more humid portion of the watershed downstream of the reservoirs, which may explain why downstream flows are somewhat predictable during the summer season.

To further evaluate the forecast models, we performed the forecasts by randomly holding out two observations from each of the three (climatology based) categories, thus dropping approximately 10-12% at random, then making the forecasting on the dropped observation. Repeated above procedure100 times. The BSS and RPSS skills for both upstream and downstream flows during winter and spring were shown as box plots in Figure 3.11 and 3.12. The median skill scores were almost the same as the leave-one-out cross-validation forecasts, although there is considerable variability in the skill scores due to resampling.

**Figure 3.11 |** Box plot of BSS (left) and RPSS (right) of upstream flows forecasts for winter and spring. Forecasts are based on dropping 10-12% of the observations randomly. The boxes correspond to the interquartile ranges (IQR), the horizontal line in each box is the median, and whiskers extend to the 1.5 IQR of lower quartile and upper quartile. Individual symbols o represent "mild" outliers.



**Figure 3.12 |** Box plot of BSS (left) and RPSS (right) of downstream flows forecasts for winter and spring. Forecasts are based on dropping 10-12% of the observations randomly. The boxes correspond to the interquartile ranges (IQR), the horizontal line in each box is the median, and whiskers extend to the 1.5 IQR of lower quartile and upper quartile. Individual symbols o represent "mild" outliers.

It is more difficult to explain how large-scale climate signals affect streamflow in Central Texas.  Either the derived SST index ERSST1 or the PDO index was found to be useful in predicting winter reservoir inflow, and there is strong correlation between the two indices ($r = 0.88$), indicating that ERSST1 is essentially a surrogate for PDO.  There is also strong correlation between ERSST2 and PDO ($r = 0.62$).  Warm-phase PDO winters correspond to blocking high pressure over the northeastern Pacific Ocean, which shifts the jet stream northward, leading to warmer and drier than average conditions, and thus lower soil moisture, in Central Texas.  Conversely, cool-phase PDO winters tend to be cooler and wetter than average, with higher soil moisture.  There is strong persistence in the PDO index from October to March, and so the fall PDO index (October to December) is a good indicator of soil moisture through the winter season. This persistence may extend to early spring, which may explain why ERSST2 is a significant indicator for spring season.  PNA has a similar but weaker effect on soil moisture; nonetheless it is a statistically significant indicator for spring conditions.

In contrast to the upstream flows, none of the climate indices considered were selected by the stepwise method for inclusion in the logistic regression models for downstream flow forecasts.  However, each of the identified SST patterns (Figures 3.6-3.8), from which new indices were derived, can be related to typical PDO or ENSO SST patterns in the Pacific and Atlantic Oceans.  (PDO and Nino3.4 indices were likely not selected due to colinearity with these patterns.)  Figure 3.6 shows the temporal correlation map between fall SSTs and winter downstream flows.  While ENSO is most well known to affect SSTs in the equatorial Pacific, it also strongly affects the climate of northeastern Brazil (e.g., Souza Filho *et al.* 2003).  The SST pattern shown in Figure 3.7, used to forecast spring downstream flows, exhibits a classic ENSO pattern in the western Pacific.  Similarly, the SST pattern exhibited in Figure 3.8, used to forecast summer downstream flows, is typical of PDO patterns in the northern Pacific.

## 3.7 Conclusion

This work aimed to develop seasonal streamflow forecast models for the Highland Lakes system in Central Texas. Hydrologic persistence (streamflow autocorrelation), five large-scale climate indices (Nino3.4, PDO, NAO, AMO, and PNA), and six derived SST indices were screened for inclusion in an ordinal polytomous logistic regression model. Results indicate that hydrologic persistence is a useful predictor of seasonal streamflows both upstream and downstream of the Highland Lakes reservoir system during the winter and spring. Summer downstream flow forecasts based on persistence also exhibit significant skill. In addition, winter reservoir inflow forecasts may be significantly improved by including either a derived SST index or the PDO index, and spring reservoir inflow forecasts may be improved by including a derived SST index and PNA. Similarly, including derived SST indices, related to ENSO and PDO SST patterns, improves downstream flow forecasts during the winter, spring and summer.

The methods presented here are completely transferable to other regions where significant hydrologic persistence and/or teleconnections between seasonal streamflow and large-scale climate anomalies exist. Stepwise linear regression with selection of predictor variables based on information criteria proved an effective method of screening a large number of potential predictors. Ordinal polytomous logistic regression proved an effective and parsimonious method for producing the probabilistic (categorical) streamflow forecasts. Both of these methods assume linearity, however, while relationships between streamflow and climate anomalies are likely to be nonlinear and include multivariate interactions due to the complexity of ocean-atmosphere dynamics (Araghinejad *et al.* 2006). Linear regression and logistic regression models are difficult to interpret if nonlinearity and/or interactions are present. To this end, nonlinear statistical methods such as data mining (machine learning) may be considered in the future work to improve the predictive skill of seasonal forecasts.

# 4. Seasonal Forecasts Using Data Mining[2]

## 4.1 Introduction

Data mining is the automated analysis of (often large) data sets to classify the data and uncover relationships that are both understandable and useful to the data owner (Hand et al. 2001). As an automated technique, data mining involves the integration of multiple disciplines such as database technology, statistics, machine learning, data visualization, and information science. In general, data mining tasks can be classified into two categories: descriptive and predictive. Data mining algorithms typically search databases for trends, patterns, and relationships that describe data (e.g., knowledge discovery), such as those that can be represented as regression models, rules, clusters, graphs, tree structures or recurrent patterns in time series. The patterns generated from a data mining system should be novel, easily understandable, and potentially useful for prediction. To effectively extract information from large amounts of data, it is necessary for data mining algorithms to be efficient and scalable. The mining process will be ineffective if the samples are not a good representation of the large body of data. Therefore, another important issue is the verification and validation of patterns on new or test data. The capability of handling noise, exceptional cases, or incomplete data objects is also required.

In recent years, a number of data mining algorithms have been developed to infer models or patterns from large datasets in many different fields of application, including marketing, surveillance, fraud detection and scientific discovery (Han and Kamber, 2006; Hand et al., 2001). Methods used in hydrology include cluster analysis, nearest-neighborhood methods, tree models, and artificial neural networks. Clustering uses iterative techniques to identify relationships and group data into clusters that contain similar characteristics, which can then be used to generate predictions (Han and Kamber, 2006). Nearest-neighborhood methods try to classify or predict the new objects based on

---

[2] This chapter is constituted by the article by Wei and Watkins (2010) "Data Mining Methods for Hydroclimatic Forecasting," currently under review by the journal *Advances in Water Resources*.

nearest neighbors in the training dataset (Hand et al., 2001), with Euclidean distance or Mahalanobis distance typically used to define nearest or closest neighbors. (For observational data with $n$ dimensions, or variables, the Euclidean distance D between two points $X=(x_1, x_2, ....x_n)$ and $Y=(y_1, y_2, .....y_n)$ is defined as: $D = [\sum(x_i - y_i)]^{1/2}$; the Mahalanobis distance accounts for correlation among variables and is scale-invariant.) This algorithm is sensitive to the local structure of the data but can perform poorly in problems with many variables (McLachlan, 1992). Tree models or decision tree models are also known as classification and regression trees or induction trees (Bessler et al., 2003). The basic principle of tree-based models is to partition datasets to maximize the purity (homogeneity) of a response variable within each partition. This method can explain and/or predict a response that is either categorical (classification) or continuous (piecewise regression).

Artificial neural networks (ANN) is a related method based on the operation of biological neural networks. Recently, ANN has attracted a great deal of attention for hydrologic forecast modeling (Maier and Dandy, 1996; Shamseldin, 1997; Clair and Ehrman, 1998; Coulibaly et al., 2001; Giustolisi and Laucelli, 2005) because of its power and flexibility. Notably, the applicability of ANNs in hydrology has been extensively evaluated by the American Society of Civil Engineers Task Committee on the Application of ANNs in Hydrology (ASCE, 2000), as well as by Dawson and Wilby (2001). These studies reported that ANN can be an efficient and promising alternative to traditional (more physically based) hydrologic models. A disadvantage of ANN is their "black box" nature, which makes it impossible to interpret relations between the individual predictors and response variable. In this respect, tree models are considered more comprehensible for decision makers (Tu, 1996).

In this study, we investigate the potential applicability of tree models for long-lead time streamflow forecasting. The methods of classification trees (CT) and logistic regression trees (LRT) are used to examine a set of potential streamflow predictors, including large-scale climate indices (teleconnections) and hydrological persistence, and screen the most promising predictive models accordingly. Data mining is particularly

attractive because the methods can effectively address the nonlinear dynamics of oceanic-atmospheric interactions with regional climate (Araghinejad et al, 2006)—nonlinearities which make traditional modeling approaches such as multiple and multivariate linear regression models statistically invalid (Piechota et al., 1998). Although there have been a number of studies using artificial neural networks (ANNs) to forecast streamflow (Coulibaly, et al., 2001), and Bessler et al. (2003) report a study using induction trees to screen multi-reservoir control rules, to the authors' knowledge tree methods have not been applied to long-lead streamflow forecasting.

With increasing water demands in many watersheds, increasing environmental awareness and conflicts over water resources, and growing concern about the hydrologic impacts of climate change, long-lead streamflow forecasts may play a critical role in water resources planning and management. Especially, the looming uncertainty about future supplies due to climate change, presents a daunting challenge to water resources engineers and managers. Many researchers have been investigating the relationship between hydrological variables, particularly streamflows, and the large-scale climate indices, such as El Niño–Southern Oscillation (ENSO), Pacific Decadal Oscillation (PDO), Pacific North American (PNA), the North Atlantic Oscillation (NAO), and Atlantic Multidecadal Oscillation (AMO). Recent studies have shown that incorporation of large-scale ocean-atmospheric patterns can improve the ability to forecast streamflow at seasonal to annual lead times in particular regions (Hamlet and Lettenmaier, 1999; Sharma, 2000; Piechota et al., 2001; Tootle, et al, 2006). Studies also indicate that the relationship between the large-scale climate indices and streamflow is usually nonlinear due to the complex dynamics of the ocean-atmosphere interaction with regional climates (Araghinejad et al, 2006).

## 4.2 Methodology

The basic goal of data analysis using tree-structured algorithms is to determine a set of if-then logical conditions (or "splits") that permit accurate predictions or classification of observational data. In this study, we use two tree-structured data mining techniques—classification trees (CT) and logistic regression trees (LRT) to develop seasonal streamflow prediction models. The details of the two methods are described below.

**Classification and Regression Trees**

The methodology of classification and regression trees (CRT) was developed by Breiman et al. (1984). This method is intended to explain and predict a dependent (response) variable, using a set of independent (predictor) variables, also referred to as explanatory variables, through a binary partitioning procedure. Both the response and the explanatory variables can be either categorical or numerical. Typically, a classification tree is used when the response variable is categorical, whereas a regression tree is used when the response variable is continuous. In this study, although the response variable (streamflow) is continuous, we use classification trees to estimate the probability that the observed value will be within one of three categories (high, medium, and low).

Classification trees have much in common with the traditional methods of discriminant analysis (Breiman et al., 1984), but the flexibility of classification trees makes them an attractive analysis option. Discriminant analysis determines the class of an observation based on a set of linear functions of the predictors, known as discriminant functions. The maximum number of discriminant functions will be equal to the degrees of freedom or the number of predictor variables in the analysis. The recursive approach to constructing classification trees does not face this limitation. Additionally, classification trees can be computed for categorical predictors, continuous predictors, or any mix of the two types of predictors, while discriminant analysis requires that predictor variables are continuous or at least measured on an interval scale (Breiman et al., 1984; Lim et al., 1997). Similarly, regression trees parallel analysis of variance (ANOVA) techniques. In

the ANOVA model, interaction is represented by cross-products between predictors, while in the regression tree model, it is represented by branches from the same node which have different splitting predictors lower in the tree.

In contrast to other methods of analyzing classification and regression problems, such as generalized linear/nonlinear models, interpreting results summarized in a tree is often very straightforward, and there are no implicit assumptions about underlying relationships between the response variable and predictor variables (such as the variables being linearly related or normally distributed). Thus, tree-structured methods are well suited for data mining tasks with little *a priori* knowledge about the data being analyzed, and they are powerful for screening variables, summarizing large multivariate datasets, constructing and evaluating predictive models, and assessing the adequacy of alternative linear models (Ripley, 1996).

CRT analysis consists of three basic steps: (1) construction of the maximal-tree, (2) pruning of the tree, and (3) selection of the optimal tree. CRT builds trees by recursively splitting the data into mutually exclusive subgroups. Each such step may give rise to new branches, called nodes. The goal of this process is to maximize homogeneity (purity) of the values of the dependent variable in each subgroup or node, i.e. minimize the variability (impurity) of the response variable in each node (Ripley, 1996). To this end, the CRT algorithm searches through all possible splits for all variables included in the analysis. The best split then is chosen by evaluation of impurity of the nodes resulting from all possible splits. For numerical explanatory variables, a split value is selected to generate two groups (nodes) at each node. For categorical explanatory variables, a split is made by relating one or more levels of the variable to a specific node. If the splitting procedure is repeated until no further split can perform, the resulting tree thus is called the maximal tree, and the terminal nodes are referred to as leaves.

Maximal trees usually turn out to be very complex and fit the training set perfectly. In modeling, this is called overfitting (Heyden et al. 2002). Such trees may be difficult to interpret, and their ability to predict new observations is generally poor because they tend

to extract all information from the training set, even the random variation, or noise, in the data. The selection of a more parsimonious tree is then necessary for predictive purposes. The tree pruning is performed based on a best compromise between complexity and accuracy. For classification trees, a cost–complexity measure may be used to determine the best one. The cost-complexity measure $R_\alpha$ is defined as a linear combination of the cost (estimated prediction error) of the tree and its complexity (Caelli, et al., 2005):

$$R_\alpha(T) = R(T) + \alpha|\bar{T}| \qquad\qquad\qquad (4.1)$$

where $R(T)$ is the resubstitution estimated error, which for a classification tree is given by the misclassification error; $|\bar{T}|$ represents the tree complexity, which is the size of the sub-tree (number of terminal nodes); and $\alpha$ is the complexity parameter. During the pruning procedure $\alpha$ takes values between 0 and 1, starting at 0 for the maximal tree and increasing to generate the optimal tree. The cost-complexity measure is thus analogous to information criteria such as the Akaike Information Criterion (AIC) (Ripley, 1996), except that it is adaptive according to the stage of the modeling process.

The procedure of tree pruning will generate a sequence of smaller trees, with the optimal tree selected from the sequence of subtrees by evaluating the predictive error of the trees. The predictive error is often estimated using a cross-validation method, in which samples are randomly drawn from the data set to test the tree grown with the rest of the data. The optimal tree may be selected as the one with the minimal cross-validation error (most accurate tree). In practice, the optimal tree is generally obtained by selecting the simplest tree with a predictive error comparable to the predictive error of the most accurate tree (Put et al., 2003).

In this study, the software CART® developed by Salford Systems was used for data mining analysis using classification trees (Breiman et al., 1984; Steinberg et al., 1997).

**Logistic Regression Trees**

Logistic regression is a statistical method used to model the probability of occurrence of an event, represented as a binary-valued response, in terms of explanatory or predictor

variables that may be either numerical or categorical (Kutner 2004). Let $p = \Pr(Y = 1)$ denote the probability of an event occurring, or a "success." In statistics, the ratio $p/(1\text{-}p)$ is called the odds, and the function $\log(p/(1\text{-}p))$ is called logit($p$), which is in fact modeling the logarithm of the ratio of probability of success to the probability of failure. In linear logistic regression, logit($p$) can be expressed as a function of one (simple linear logistic regression) or more predictor variables $x_i$ (multiple linear logistic regression) as follows:

$$\text{logit}(p) = \log(p/(1-p)) = \beta_0 + \beta_1 x_1 + \cdots \beta_k x_k$$

(4.2)

The unknown parameters $\beta_j$ are usually estimated using maximum likelihood theory. Although multiple linear logistic regression models can provide accurate estimates of the probability of an event affected by many variables, the resulting model may be complex and difficult to interpret, especially if the number of predictor variables is large, and if collinearity, nonlinearity, or interaction exists among the predictor variables. On the other hand, an overly simple model may have little predictive power (Chan et al 2004).

To overcome these problems, a logistic regression tree method known as Logistic Tree with Unbiased Selection (LOTUS) was developed by Chan and Loh (2004), which can retain both the graphical interpretability of simple models and also the predictive accuracy of more complex models. LOTUS is an algorithm for automatic construction of logistic regression trees, based on the underlying principle of fitting a piecewise (simple or multiple) linear logistic regression model by recursively partitioning the data and fitting a different linear logistic regression to the data in each partition. LOTUS uses a trend-adjusted chi-square test to control bias in variable selection at the intermediate nodes (Cochran, 1954; Armitage, 1955). This can distinguish nonlinear from linear effects and ensure the integrity of inferences drawn from the tree structure.

Once the initial binary tree is grown, analogous to the maximal classification tree, it is pruned back by minimizing a cross-validation estimate of the predicted deviance per degree of freedom (similar to the cost-complexity measure in Eq. 4.1), instead of simply the sum of square residuals, which would tend to lead to overfitting. Deviance is a

standard measure of variation for generalized linear models and also the impurity measure for tree-based models (McCullagh and Nelder, 1989) and the degrees of freedom is defined as the number of fitted observations minus the number of estimated parameters, including the intercept terms. For logistic regression, the deviance is defined as:

$$D = -2 \sum_{i=1}^{n} [y_i \log (\hat{p}_i / y_i) + (1 - y_i) log\{(1 - \hat{p}_i)/(1 - y_i\})] \tag{4.3a}$$

$$\text{or} \quad D = -2 \sum_{i=1}^{n} [y_i \log (\hat{p}_i) + (1 - y_i) \log (1 - \hat{p}_i)] \tag{4.3b}$$

where $\hat{p}_i$ is the estimated probability for the $i$th observation, and $y_i$ is the $i$th binary response. The total impurity for a tree is the sum of the deviances in all the partitions. LOTUS allows the choice of one of three roles for each quantitative predictor variable: f-variables for fitting only, acting as a regressor; s-variables for splitting only, serving as split selection; and n-variables for both splitting and fitting. These features allow nonlinearity of the data to be modeled without requiring variable transformations. Furthermore, by fitting linear logistic regressions for each node, the tree model is visualizable and hence more comprehensible than standard multiple linear logistic regression.

## 4.3 Case Study Data

The Lower Colorado River Authority (LCRA) is a water conservation and reclamation district established by the State of Texas, USA. It supplies electricity, manages water supplies and floods in the lower Colorado River basin, supports water and wastewater utilities, provides public parks for water-based recreation, and promotes community and economic development in 58 counties in Central Texas (see Figure 4.1). To meet rapidly growing water demands through more efficient operation of the Highland Lakes reservoirs, seasonal river flow forecasts would be very beneficial.

**Figure 4.1.** Lower Colorado River Authority District in Central Texas (provided by Ron Anderson, LCRA).

To explore the patterns of streamflow and the influence of ocean-atmosphere teleconnections in Central Texas, monthly streamflow data are acquired from two sources: 1) aggregate Highland Lakes inflows (upstream), based on USGS gage measurements and adjustments made by LCRA staff to account for inflows from ungaged areas; and 2) unregulated flows downstream of the Highland Lakes, as determined by the Texas Water Availability Model (WAM) (Wurbs, 2008). The reservoir inflow data spans a total of 57 years, from 1950 to 2006, and the naturalized downstream flow data spans 59 years, from 1950 to 1998. For most of the analyses, the raw flow data are normalized through a two-step process—first a logarithmic transformation, then conversion to a standardized anomaly by subtraction of the mean (of the log values) and division of the standard deviation (of the log values). While this transforms the data so that the statistical assumption of normality is valid, it should be noted that the correlation coefficients are then inflated (due to the log-transform), and thus an effort is made to illustrate the results in terms of the raw flow data.

67

Monthly autocorrelations of (upstream) reservoir inflows range from a high of nearly 0.8 for February and March flows to a low of essentially zero for July and August flows. Seasonal correlation coefficients also peak in the winter season, with a value of 0.66. (The correlation coefficient that is significant at the $p = 0.05$ level is 0.20.) It may be surprising that the seasonal correlation between OND and JFM flows is higher than the average of monthly correlation coefficients during this period. One reason for this may be that averaging over a three-month period reduces the "noise" that results from individual storm events which have a significant effect on monthly flow totals.

The tributary flows to the Colorado River downstream of the Highland Lakes, estimated as the WAM naturalized flows at Mansfield Dam minus the naturalized flows at Bay City, have monthly autocorrelation coefficients that reach a maximum of 0.8 for February and March and a minimum of about 0.2 for September and October. The average monthly autocorrelation for the downstream data is about 0.27 higher than the upstream data. Seasonal autocorrelations peak in the spring season with a value of 0.82, which is lagged by one season in comparison with the upstream data. All seasonal autocorrelation coefficients for the downstream data are significant at $p = 0.10$ or less, and the average autocorrelation coefficient is about 0.23 higher than the upstream data, most likely due to the humid conditions along the Texas coast compared to semi-arid central and western Texas.

Based on these autocorrelation coefficients, seasonal streamflow forecasts for certain times of the year may be based solely on hydrologic persistence. The predictive skills of these forecasts are investigated, as well as the potential for large scale ocean-atmosphere interactions to provide additional forecast skill. The oceanic-atmospheric phenomena investigated as potential predictor variables for streamflow are the El Niño-Southern Oscillation (ENSO), the Pacific Decadal Oscillation (PDO), the North Atlantic Oscillation (NAO), the Atlantic Multidecadal Oscillation (AMO), and the Pacific North American (PNA). The ENSO and NAO generally have a two- to seven-year periodicity (Philander, 1990), while PDO and AMO exhibit long-term periodicity of about 25 to 60 years (Mantua *et al.* 1997; Kerr 2000; Gray *et al.* 2004). Various indices were selected to

quantify the magnitude of these ocean-atmospheric oscillations, as summarized in Table 4.1 (Wei and Watkins, 2010). In all cases, seasonal (three-month average) index values for the period 1950-2006 were used for the analysis.

In addition, SST data was analyzed directly for correlations with streamflow. The data used was the extended reconstructed sea surface temperature (ERSST) analysis (Smith et al., 2008), obtained from the National Climatic Data Center (NCDC) through the KNMI Climate Explorer, an on-line data analysis tool (Oldenborgh & Burgers, 2005). Six correlation patterns with high statistical significance ($p < .01$) were identified and referenced as SST1 through SST6. For each SST pattern, a normalized index was computed based on average seasonal temperatures over a 4-degree by 4-degree area, similar to the procedure of Block and Rajagopalan (2007). For details, see Wei and Watkins (2010).

**Table 4.1**. Predictor variables indentified for streamflow in Central Texas, based on sea surface temperatures (SST) or sea level pressures (SLP). Data last accessed on July 17, 2010.

| Climate Index | Climate variable | Periodicity (years) | On-line Data Source |
|---|---|---|---|
| Niño 3.4 | SST | 2-7 | http://www.cpc.ncep.noaa.gov/data/indices |
| PNA | SLP | 0.25-10 | http://www.cpc.ncep.noaa.gov/data/indices |
| PDO | SST | 25-60 | http://jisao.atmos.washington.edu/pdo |
| NAO | SLP | 2-7 | http://www.cgd.ucar.edu/cas/jhurrell/indices.html |
| AMO | SST | 25-60 | http://www.cdc.noaa.gov/Climateindices |

## 4.4  Forecast Model Development

Both CT and LRT models were used to screen significant predictor variables from the set of potential predictors, observed in the season prior to the forecast. For convenience, the seasons are defined as winter (January-March), spring (April–June), summer (July–September), and fall (October–December). Climate indices and streamflow for the season prior to the predicted seasonal streamflow are designated by (-1).

First, the CT approach is applied to both the reservoir inflow and downstream flow data sets.  In both cases, the response variable, i.e. seasonal streamflow, is treated as a categorical variable with three levels, or terciles (above normal, normal, and below normal), and all potential predictor variables are continuous. Accordingly, in a classification tree, the forecast probabilities for each category are given by the empirical relative frequencies of the classes in the terminal nodes of the optimal tree. For example, if there are 30 cases at a certain terminal node, 15 of which are in the below normal category, 10 of which are in the normal category, and 5 of which are in the above normal category, then this terminal node will correspond to the following probabilistic (tercile probability) forecast: 50% chance of below normal, 33.3% chance of normal, and 16.7% chance of above normal.

A summary of the optimal classification trees for downstream flows, selected based on cross-validation, is presented in Table 4.2. The results show all trees have three or more levels, indicating nonlinearity in the relationships between streamflow and the predictor variables for all four seasons. The more levels a tree has, the more complicated the nonlinearity of the relationship is. Generally, the variable used to split data in the first level of the tree is more important than the variables used to split data in lower levels. Thus, for winter downstream flow forecasts, SST3 appears to be the most important predictor, though the other predictors (SST6 and fall streamflow) are also statistically significant ($p = 0.05$).

A plot of the optimal classification tree for winter downstream flow forecasts, with four levels and four terminal nodes, is shown in Figure 4.2(a). For spring downstream flow forecasts, a classification tree with five levels indicates more complicated nonlinearity. In this case, streamflow persistence is likely the most important predictor, while both SST3 and AMO are also significantly related to spring downstream flows. For summer downstream flow forecasts, persistence is not important, but SST5 and AMO are, with SST5 being the more important predictor. For fall downstream flows forecasts, streamflow persistence is again important, and SST3 is also found significant.

A summary of the optimal classification trees for reservoir inflows (upstream) based on cross-validation is presented in Table 4.3. Results show that all classification trees have two or three levels, again indicating nonlinearity in the relationships between streamflow and predictor variables, but in this case the relationships are less complicated than those for downstream flows. Results also show that fewer predictors are significant for reservoir inflow forecasts. Streamflow persistence is a statistically significant predictor ($p = 0.05$) for winter, spring and fall reservoir inflow forecasts. In addition, PNA is also found to be a significant predictor for spring inflow forecasts, and both PNA and NAO are significant predictors for summer reservoir inflows. A plot of the optimal classification tree for spring reservoir inflow forecasts is shown in Figure 4.2(b).

**Table 4.2**. The optimal classification trees based on cross-validation for seasonal streamflow forecasts of unregulated tributary flows, from the Texas Water Availability Model (Downstream).

| Classification trees | Winter(JFM) | Spring (AMJ) | Summer (JAS) | Fall (OND) |
|---|---|---|---|---|
| Tree levels | 4 | 5 | 3 | 3 |
| Terminal nodes | 4 | 5 | 3 | 3 |
| Split variables and values | SST3(-1)=-0.974<br>SST6(-1) = 0.217<br>Fall(-1) = -0.177 | Winter(-1)=-0.974<br>SST3(-1)=-0.120<br>AMO(-1) = -0.082<br>Winter(-1) = 0.782 | SST5(-1)=0.037<br>AMO(-1)= -0.073 | summer(-1)= -0.543<br>SST3(-1)= -0.482 |

**Table 4.3**. The optimal classification trees based on cross-validation for seasonal streamflow forecasts of aggregate inflows to the Highland Lakes (Upstream).

| Classification trees | Winter(JFM) | Spring (AMJ) | Summer (JAS) | Fall (OND) |
|---|---|---|---|---|
| Tree levels | 2 | 3 | 3 | 2 |
| Terminal nodes | 2 | 3 | 3 | 2 |
| Split variables and values | Fall(-1) = 0.238 | Winter(-1) = 0.255<br>PNA(-1) = 0.030 | PNA(-1) = -0.102<br>NAO(-1) = -0.195 | summer(-1) = 0.17 |

**Figure 4.2.** Optimal classification trees for (a) winter streamflows downstream of the Highland Lakes, and (b) spring inflows to the Highland Lakes. Intermediate and terminal nodes are represented by circles and squares, respectively. The number inside an intermediate node is the splitting value, and splitting variable is given beneath it. If a case is equal to or less than the splitting value, it goes to the left branch; otherwise the right branch. The number inside a terminal node indicates the dominant category level, i.e., above normal (3), normal (2), and below normal (1).

Next, the LRT approach is applied to both the reservoir inflows (upstream) and downstream flow data sets. In logistic regression trees, the response variable, i.e. seasonal streamflow, is considered binary (either a threshold flow is exceeded or it is not), and the potential predictor variables are continuous and used both for splitting selection during tree construction and for fitting the linear logistic models at each terminal node. To develop tercile probability forecasts using LRT, the following steps are taken:

1. The 33.3 percentile is first chosen as a threshold, and the response variable $Y$ is equal to 1 if seasonal streamflow (Q) is equal to or less than the threshold, and 0 otherwise.
2. A logistic regression tree is generated to obtain the probability for the below-normal category: $p(BN) = p(Q <= 33.3$ percentile).
3. Steps 1 and 2 are repeated for the 66.6 percentile, and the probability for normal and above-normal categories, $p(N)$ and $p(AN)$, are derived as follows:

    $p(N) = p(Q <= 66.6$ percentile$) - p(BN)$

    $p(AN) = 1 - p(Q <= 66.6$ percentile$)$ or $p(AN) = 1 - p(BN) - p(N)$

    Following this procedure, the sum of probabilities for each category is guaranteed to equal 1.

Summaries of the logistic regression trees generated using LOTUS for both downstream flows and (upstream) reservoir inflows are presented in Tables 4 and 5, respectively. The logistic regression trees for both downstream flows and upstream flows are much simpler in structure than the corresponding classification trees. This is likely due to the logistic regression essentially being a nonlinear regression in terms of response functions (Kutner, et al. 2004), capturing the nonlinear features between the response variable and independent variables without the need for more complex splitting in the tree. Results show that logistic regression trees with two levels are generated for winter and summer downstream flow forecasts, while only one-node trees are grown for spring and fall downstream flows and all seasonal reservoir inflows. One-node trees indicate

there is a global relationship between streamflows and predictors, and a single logistic regression function can be used to capture the relationship.

Results in Table 4.4 show that streamflow persistence is a statistically significant predictor (*p*= 0.05) as a fitting variable for all seasons of downstream flows. SST3 through SST6 are also statistically significant predictors as fitting variables for corresponding seasonal downstream flows. Forecast models for spring and fall streamflow can be represented by a single multiple logistic regression function, respectively. Nonlinear features present in winter and summer streamflow data are accounted for by partitioning the data into two parts using splitting variable SST3 and Nino3.4, respectively, and fitting a different multiple logistic regression for each partition. For example, the logistic regression tree for winter streamflow forecasts, shown in Figure 4.3, has two terminal nodes, designated as Nodes 1 and 2. The corresponding fitting multiple logistic regression functions are given as follows:

Node 1:    *Logit(p) = -1.282 – 1.704\*Fall(-1) + 1.432\*SST3(-1)*        (4.4)

Node 2*:    Logit(p) = -1.061 + 1.308\*Fall(-1) – 6.084\*SST3(-1)*        (4.5)

The estimated probabilities are then derived from these logit equations as follows:

Node 1: $p_i = \left[ 1 + \exp(-1.282 - 1.704 * Fall(-1) + 1.432 * SST3(-1) \right]^{-1}$        (4.6)

Node 2: $p_i = \left[ 1 + \exp(-1.061 + 1.308 * Fall(-1) - 6.084 * SST3(-1) \right]^{-1}$        (4.7)

Results in Table 4.5 show that streamflow persistence is a statistically significant predictor (p = 0.05) for each season of reservoir inflows (upstream) except for summer. This is consistent with results from the classification tree analysis. In addition, both SST1 and SST2 are significant predictors for winter streamflow forecasts, and SST2 is also a significant predictor for spring streamflow forecasts. PNA is found to be a significant predictor for fall streamflow forecasts. Each of the global relationships present in the winter, spring and fall can be modeled using a single multiple logistic regression

function. The logistic regression tree analysis indicates there are no significant predictor variables for summer streamflow forecasts.

**Table 4.4**. The best logistic regression tree based on the cross-validation method for seasonal streamflow forecasts of unregulated tributary flows downstream of reservoirs.

| Logistic regression trees | Winter(JFM) | Spring (AMJ) | Summer (JAS) | Fall (OND) |
|---|---|---|---|---|
| Tree level | 2 | 1 | 2 | 1 |
| Terminal nodes | 2 | 1 | 2 | 1 |
| Splitting variables | SST3(-1) = -0.163 | None | Nino34(-1) =0.17 | None |
| Fitting variables | Fall(-1)<br><br>SST3(-1) | Winter(-1)<br><br>SST4(-1) | Spring(-1)<br><br>SST5(-1) | Summer(-1)<br><br>SST6(-1) |

**Table 4.5.** The best logistic regression tree based on the cross-validation method for seasonal streamflow forecasts of aggregate inflows to the Highland Lakes reservoirs (Upstream).

| Logistic Regression trees | Winter(JFM) | Spring (AMJ) | Summer (JAS) | Fall (OND) |
|---|---|---|---|---|
| Tree level | 1 | 1 | 1 | 1 |
| Terminal nodes | 1 | 1 | 1 | 1 |
| Splitting variables | None | None | None | None |
| Fitting variables | Fall(-1)<br><br>SST1(-1), SST2 (-1) | Winter(-1)<br><br>SST2(-1) | None | Summer(-1)<br><br>PNA(-1) |

**Figure 4.3.** Plot of stepwise logistic regression tree for winter streamflows downstream of the Highland Lakes. Intermediate and terminal nodes are represented by circles and squares. The splitting variable is SST3, and the splitting value is -0.163. If a case is equal to or less than the splitting value, it goes to the left branch; otherwise to the right branch. The ratio of cases with Y =1 to the node sample size is given beneath each terminal node. Total sample size is 48.

## 4.5  Forecast Verification

Forecast verification aims to evaluate the agreement between forecasts and observations (Katz and Murphy, 1997; Stephenson, 2003). In this study, the Brier skill score (BSS) and the ranked probability skill score (RPSS) are used to measure the improvement in the accuracy of multicategory probability forecasts over a naïve forecasting method such as climatology.

The Brier skill score (BSS) (Dogget, 1998) is defined as

$$BSS = 1 - BS \,/\, BSC \tag{4.8}$$

where $BS$ is the Brier score, and is defined as

$$BS = (1/n)\Sigma(f_i - l(\mathrm{obs}_i))^2 \tag{4.9}$$

where $f_i$ is the forecast probability of event $i$ occurring; $l(\text{obs}_i)$ is an indicator variable (1 if event in category $i$ occurs, else 0); $n$ is the number of events; and *BSC* is the climatologically expected value of *BS*, and is defined as $BSC = P_i * (1 - P_i)$, where $P_i$ is the climatological probability of the event. In this study, we consider three tercile categories, i.e., below normal, normal and above normal. Accordingly, $BSC = (1/3)*(1-1/3) = 0.222$ for each category. A perfect forecast has a value of BSS of 1; positive values between zero and one indicate forecast performance better than climatology, and negative values indicate forecast performance worse than climatology.

The ranked probability score (*RPS*) evaluates the sum of the squared differences in the cumulative probability space, so that

$$RPS = \frac{1}{K-1} \sum_{m=1}^{K} \left[ \left( \sum_{k=1}^{m} f_k \right) - \left( \sum_{k=1}^{m} o_K \right) \right]^2$$

(4.10)

where $K$ is the number of forecast categories (below normal, normal and above normal), $f_k$ is the forecast probability for the $k^{th}$ point, and $o_k$ equals zero or one to indicate whether or not the observed flow is in the $k^{th}$ category. The use of *RPS* results in higher penalties for forecasts farther away from actual outcomes, rather than scoring based on only hit or miss. The *RPS* can assume a number between zero and one, with a perfect forecast scoring zero. The *RPSS* then measures the relative improvement of using a forecast over climatology alone, and is given by

$$RPSS = \frac{\overline{RPS} - \overline{RPS}_{c\lim ato\log y}}{o - \overline{RPS}_{c\lim ato\log y}} = 1 - \frac{\overline{RPS}}{\overline{RPS}_{c\lim ato\log y}}$$

(4.11)

The probabilities of streamflow in each category in the climatology forecast (i.e., prior probabilities) are the same and equal 1/3 due to the definition of the flow regimes. Thus, the *RPS* values of the three categories (below normal, normal and above normal) in the climatology forecast are 0.278, 0.111 and 0.278, respectively. A perfect *RPSS* is 1, and negative scores indicate that forecasts performed worse than climatology

Using the BSS and RPSS as metrics for evaluation, the optimal classification tree and logistic regression tree models, indicated in Tables 4.2-4.5, are applied to forecast seasonal flows upstream and downstream of the Highland Lakes reservoir system in Central Texas. Summaries of the results are presented in Tables 4.6 and 4.7. Results in Table 4.6 show that downstream unregulated flows can be predicted with significant skill based on either the CT or LRT approach. The BSS and RPSS values indicate that forecasts using CT have an average improvement in skill over climatology of 39.4% and 25.4%, respectively, while using LRT the average corresponding improvement in skill is 43.5% and 33.2%. In particular, downstream flows for winter and spring seasons can be predicted very well, with skill score improvements of about 40-50% using either the CT or LRT model.

Results also indicate the tree-structured models can capture the nonlinear features of the downstream flows data. For example, SST3 appears as a splitting variable in the first level of both the classification and logistic trees for winter streamflow forecasts. This implies not only that SST3 is an important predictor of streamflows, but also nonlinear features present in winter streamflow data can be accounted for by partitioning data based on SST3 values. In the logistic regression tree, there is no further partitioning, but a multiple logistic regression model using streamflow persistence and SST3 as fitting variables is selected for each partition. In the classification tree, SST6 and streamflow persistence are used as splitting variables for further partitioning. The nonlinear relationships of winter streamflow with these predictors are illustrated in Figures 4.4 and 4.5. It is the capability of tree-structured models to capture such nonlinear features that make them attractive for developing forecast models.

Results in Table 4.7 show that reservoir inflows can also be predicted with significant skill for winter, spring and fall seasons using either CT or LRT models. The BSS and RPSS values indicate that reservoir inflows forecasts for these seasons using CT have an average improvement in skill over climatology of 14.1% and 10.4%, respectively, while using LRT the average corresponding improvement in skill is 21.6% and 20.6% improvement. For summer forecasts, however, skill scores obtained from CT are very

close to zero, indicating no improvement over climatology. There is no predictor selected by LRT for summer forecasts, so the skill score is zero.

Forecast skill scores indicate that LRT generally performs better than CT. This is likely because the probabilities derived from LRT are based on a regression function fitted for each terminal nodes of the tree, while the probabilities given by CT are based on simply the empirical relative frequencies of each category represented in each terminal node. Thus, the logistic regression tree can more accurately model variability within each classification (terminal node).

**Table 4.6.** Brier Skill Score (BSS) and Ranked Probability Skill Score (RPSS) using classification trees and logistic regression trees for seasonal forecasts of unregulated tributary flows downstream of the reservoirs. Skill scores are based on cross-validation forecasts.

| Forecast Models | Seasonal Forecast Skill Scores | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Winter | | Spring | | Summer | | Fall | |
| | BSS | RPSS | BSS | RPSS | BSS | RPSS | BSS | RPSS |
| Classification Tree | 40.1% | 39.3% | 49.9% | 42.1% | 34.7% | 7.0% | 32.8% | 13.1% |
| Logistic Regression Tree | 53.7% | 43.5% | 50.8% | 38.7% | 40.7% | 34.5% | 28.9% | 16.3% |

**Table 4.7.** Brier Skill Score (BSS) and Ranked Probability Skill Score (RPSS) using classification trees and logistic regression trees for seasonal streamflow forecasts of aggregate inflows to the Highland Lakes reservoirs (Upstream). Skill scores are based on cross-validation forecasts.

| Forecast Models | Seasonal Forecast Skill Scores | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Winter | | Spring | | Summer | | Fall | |
| | BSS | RPSS | BSS | RPSS | BSS | RPSS | BSS | RPSS |
| Classification Tree | 13.8% | 8.9% | 16.3% | 12.2% | 5.7% | 1.2% | 12.1% | 10.0% |
| Logistic Regression Tree | 29.3% | 27.1% | 24.2% | 23.1% | 0.0% | 0.0% | 11.4% | 11.6% |

**Figure 4.4.** Surface plot of winter streamflow as a function of SST3 and SST6 predictors, demonstrating nonlinearities in the relationship represented by the logistic regression tree. Units on the z-axis (winter streamflow) are the natural logarithm of flow in acre-feet per month (1 acre-ft equals approximately 1,234 m3).

**Figure 4.5.** Surface plot of winter streamflow as a function of SST3 and streamflow persistence predictors, demonstrating nonlinearities in the relationship represented by the classification tree. Units on the y- and z-axis (streamflow persistence and winter streamflow) are the natural logarithm of flow in acre-feet per month (1 acre-ft equals approximately 1,234 m3).

Comparisons of observations and cross-validated flow forecasts based on either the logistic regression tree or the classification tree models are shown in Figures 4.6 and 4.7 for winter reservoir inflows and spring downstream flows, respectively. Shown are many years in which the tree-based forecasts provide a significant improvement over climatology. For instance, as shown in Figure 4.6, the forecasts from LRT accurately predict high winter reservoir inflows in 8 of the 11 years in which inflows exceeded 442,400 acre-ft ($Ln(442,400) = 13.0$). Perhaps more importantly, the forecasts accurately predict very low winter reservoir inflows ($Ln$(Streamflow) < 11.0) in 6 of 9 cases. As

shown in Figure 4.7, downstream flow forecasts from CT for the spring season have even better skill, with accurate predictions of high flows (*Ln*(Streamflow) > 13.5) in 10 of 11 cases, and accurate predictions of low flows (*Ln*(Streamflow) < 11.5) in 8 of 11 cases.



**Figure 4.6.** Plot of observations and cross-validated forecasts for winter reservoir inflows to the Highland Lakes using the logistic regression tree model. Units on the y-axis are the natural logarithm of flow in acre-feet per month (1 acre-ft equals approximately 1,234 m3).

**Figure 4.7.** Plot of observations and cross-validated forecasts for spring streamflows downstream of the Highland Lakes using the classification tree model. Units on the y-axis are the natural logarithm of flow in acre-feet per month (1 acre-ft equals approximately 1,234 m3).

One additional observation from Figure 4.7 is that the tercile probability forecasts are "flat" in places, with many of the same values re-occurring. This is an artifact (and potential drawback) of the CT models developed in this study, in which the probabilities are based on the empirical relative frequencies in each category in the terminal nodes of the classification tree. In contrast to the LRT models, the CT models are not able to model variability within a terminal node. Future work may extend the CT approach to classification and regression trees (CRT), with a method to derive tercile probability forecasts based on analysis of regression residuals (errors).

## 4.6 Conclusion

For developing statistical hydroclimatic forecast models, tree-structured data mining techniques offer a flexible and attractive alternative to traditional modeling approaches such as multiple linear regression. Since these data mining methods require no *a priori* knowledge about the data being analyzed, they are powerful tools for screening large numbers of potential predictor variables. Tree-structured models also have the ability to deal effectively with multicollinearity, nonlinearity and/or interactions present in the data. In addition, the results represented by the binary trees are easy to understand and explain to decision makers as a set of if-then-else rules. These features may be valuable in searching for useful predictors or improving the reliability of existing statistical forecast models.

In this study, classification trees (CT) and logistic regression trees (LRT) are used to screen 12 predictor variables (including hydrologic persistence, large-scale climate indices, and derived sea surface temperature patterns) and identify seasonal streamflow prediction models for a reservoir system in Central Texas. Application of the tree-structured models to flows both upstream and downstream of the reservoirs resulted in significantly improved forecast skill for both locations and all seasons, except for summer flows upstream of the reservoirs. Comparing the CT and LRT approaches, classification trees are easier to understand, but logistic regression trees are more accurate due to their ability to model variability in each node of the tree.

The tree-structured data mining techniques presented here are completely transferable to other regions with significant hydrologic persistence and where teleconnections between seasonal streamflow and large-scale ocean-atmospheric patterns exist. Whenever data mining is used to identify predictor variables, as in this study, further research is needed to better understand the physical mechanisms behind these interactions. Future work may extend the classification tree approach to include linear regression models for each terminal node, and should also include development of a

statistical forecast-decision model to evaluate the benefits of the streamflow forecasts in the context of water resources decision making.

# 5. Use of Seasonal Forecasts in Water Management

## 5.1 Introduction

Accurate seasonal to interannual streamflow forecasts based on climate information are critical for optimal management and operation of water resources systems. Considering most water supply systems are multipurpose, operating these systems to meet increasing demand under the growing stresses of climate variability and climate change, population and economic growth, and environmental concerns could be very challenging.

In the last decade, significant improvement in the skill of seasonal climatic forecasts has been achieved based on the output of general circulation models or statistical models developed from historical data (Goddard et al., 2003). There is increasing evidence that the continental scale rainfall and streamflow patterns are modulated by the large-scale oceanic- atmospheric circulation patterns such *as El* Niño-Southern Oscillation (ENSO), the Pacific Decadal Oscillation (PDO), the Atlantic Multidecadal Oscillation (AMO), and the North Atlantic Oscillation (NAO) (Dettinger et al. 2000; Souza et al. 2003; Tootle et al. 2005).  Recent results show that incorporating the large-scale oceanic-atmospheric phenomena could improve the accuracy of seasonal to interannual streamflow forecasts (Clark et al. 2001; Hamlet et al. 2002).  Seasonal streamflow forecasts based on climate information rely on statistical and dynamic modeling approaches. The statistical modeling frequently employs the statistical relationship between the related climate indicators, historical rainfall or streamflow volume at a site to forecast streamflow (Souza et al., 2003; Devineni et al., 2008). With the dynamic modeling approach, seasonal streamflow forecasts can be obtained by applying climate predictions from a regional climate model (RCM) coupled with the General Circulation Model (GCM) outputs to a hydrological model (Seo et al., 2003). To address the uncertainties of streamflow forecasts arising from initial conditions and boundary conditions, model structure and modeled processes, significant attention has been given to multimodel forecasts techniques by combining different individual  models. The results show that multimodel

forecasts have considerable improvement in the overall predictability of seasonal streamflow forecasts and reducing the overall model error (Regonda et al., 2006; Devenieni et al., 2008).

A wide array of seasonal hydroclimatic forecast products is now available in the public domain. For example, the NOAA Climate Prediction Center issues seasonal forecasts of temperature, precipitation, and soil moisture, as well as a drought outlook, for the entire U.S. In the western U.S., the USDA Natural Resources Conservation Service forecasts streamflows in the first half of the year based on observed snowpack conditions. At many stream gage locations throughout the U.S., the National Weather Service provides probabilistic seasonal flow forecasts through a procedure known as Ensemble Streamflow Prediction, or ESP (Day, 1985; Smith et al., 1992). The ESP method uses conceptual or physically based hydrologic models to issue streamflow forecasts based on the current soil moisture, river, and reservoir conditions by assuming that past meteorological events will recur in the future with historical probabilities (Schaake and Larsen 1998). Recently, methods have been developed to condition the probabilities of the historical meteorological traces based on seasonal climate forecasts (e.g., Croley 2000). On a global scale, the NOAA/Columbia University International Research Institute for Climate and Society (IRI) issues seasonal forecasts of temperature and precipitation.

Recent studies have demonstrated that seasonal streamflow forecasts based on climate information can significantly improve management of water supply systems (Georgakakos, et al., 2007; Golembesky et al., 2009). Notably, Hamlet et al. (2002) estimated $161 million/year in potential benefits from use of long-lead streamflow forecasts to improve hydropower system operations in the Columbia River basin. Grantz et al. (2007) showed that incorporating seasonal streamflow forecasts based on climate information into a decision-making model for water management in the Truckee-Carson River Basin can offer skillful, longer lead-time forecasts of decision variables. Golembesky et al. (2009) showed that multimodel streamflow forecasts with season-ahead lead time could provide a more reliable way to develop water management

strategies such as invoking restrictions during below-normal flow years for the Falls Lake Reservoir in the Neuse River Basin, N.C.

Although verification of seasonal climatic forecasts and the corresponding seasonal streamflow forecasts often shows that they have significant skill, adoption by water management agencies appears to be slow. Some proposed reasons for this include a lack of understanding of probabilistic forecasts and the associated uncertainty (Pagano et al., 2001). Furthermore, there are no structured framework and policy instruments for water managers and reservoir operators to incorporate such information into water-resources decision-making (Pagano et al., 2002). In addition, lack of confidence in forecast accuracy also discourages water managers from utilizing such probabilistic forecasts in the current operating systems.

In this study, a seasonal streamflow forecast model is developed based on hydrological persistence and large-scale climate indicators. A simple water resources economic-optimization model is proposed to investigate the potential value of these forecasts for seasonal water contracts under water availability uncertainty. Some recommendations for adoption of this approach by the water management agency are provided.

## 5.2 LCRA Water Management

The Lower Colorado River Authority (LCRA) is a water conservation and reclamation district that operates a series of six lakes and dams on the watershed of the Lower Colorado River in Central Texas. The purposes of the LCRA are to supply low-cost electricity for Central Texas; manage water supplies and floods in the lower Colorado River basin, including the City of Austin and four rice irrigation districts along the Texas Gulf Coast; develop water and wastewater utilities; provide public recreation; and support community and economic development in 58 Texas counties (Figure 5.1). According to the LCRA Revised Water Management Plan (LCRA 2003), the LCRA supplies water to two general categories of water demands: firm and interruptible. Firm

demands include municipal and industrial, steam-electric power generation, some irrigation, and in-stream flow and estuarine flow maintenance. Currently, interruptible stored water is used almost entirely for agricultural irrigation, specifically rice irrigation, and environmental flow maintenance.



**Figure 5.1**. LCRA Water Service Area (Source: Ron Anderson, LCRA)

In year 2000, surface water demands within the lower Colorado River basin totaled approximately 675,800 acre-ft annually (1 acre-ft = 1,233.5 m³), including stored water and pass through of storable inflows from Lakes Buchanan and Travis to maintain in-stream flows and freshwater inflows to the bay and estuary in the lower Colorado River. About 56 percent of surface water diversions are used for rice irrigation in the four major irrigation operations located in Colorado, Wharton and Matagorda Counties in the Gulf Coastal Plain. The next largest demand for surface water is the City of Austin, which in year 2000 used approximately 163,800 acre-ft for municipal use and steam-electric power

generation under its own run-of-river rights and contracts for stored water from Lakes Buchanan and Travis (LCRA 2003).

The LCRA uses beginning-of-year (January 1) combined storage levels in the two lakes used for water supply, Lakes Buchanan and Travis, to determine the amount of water available to meet firm and interruptible water demands in the coming year. Firm water is that which is diverted from storage under a contract or resolution issued by the LCRA Board to high-priority users such as the City of Austin. Interruptible water contracts are issued on a shorter time scale (typically one year) with the condition that supplies may be interrupted or curtailed in the event that firm supplies become endangered. In allocating interruptible water, priority is given to irrigation operations downstream of Austin. If it is projected that the availability of interruptible water exceeds these irrigation needs, annual contracts can then be made with other entities within the Lower Colorado basin. Seasonal and long-term forecasts are not used formally by the LCRA for a number of reasons, including high seasonal and annual variability of stream flow and the absence of easily measured hydrologic predictors such as snowpack.

The LCRA's conceptual lake management policy for year 2010 projected demands calls for curtailment of interruptible supplies to begin when combined storage levels drop below 1,400,000 acre-ft, or about 70% of the maximum water supply storage, decreasing at a rate of approximately 31,200 acre-feet for each 100,000 acre-foot decrease in combined storage, (LCRA 2003). "Aggressive" curtailment begins at a January 1 storage level of 1,150,000 acre-ft (about 58% of maximum), decreasing at a rate of approximately 4,250 acre-feet for each 100,000 acre-foot decrease in combined storage. No interruptible water use will be sanctioned on January 1 if levels are below 325,000 acre-ft (about 16% of maximum). Additionally, interruptible water use will be stopped at any time during the year if combined storage levels drop below 200,000 acre-ft (10% of maximum). Conversely, in years of high storage levels, additional interruptible water supplies may be available for sale if combined storage levels are greater than 1,865,000 acre-ft (about 94% of maximum). Figure 5.2 illustrates a hypothetical "rule curve" that

corresponds to the published conceptual lake management policy (Wei and Watkins, 2006).



**Figure 5.2.** Hypothetical rule curve corresponding to LCRA's conceptual lake management policy. (1 AF = 1,233.5 m$^3$.)

## 5.3 Development of Seasonal Streamflow Forecast model

Recent studies have shown that incorporation of large-scale ocean-atmospheric patterns can improve the ability to forecast streamflows at seasonal to annual lead times in particular regions (Hamlet and Lettenmaier, 1999; Sharma, 2000; Piechota et al., 2001; Tootle, et al, 2006). Studies also indicate that the relationship between the large-scale climate indices and streamflow is usually nonlinear due to the complex dynamics of the ocean-atmosphere interaction with regional climates (Araghinejad et al, 2006). According to the LCRA Revised Water Management Plan, interruptible stored water may be contracted for sale with six-month ahead of each year based on January 1 combined storage levels. According to this decision horizon, a data mining technique known as logistic regression trees (LRT) is used to develop a seasonal streamflow forecast model

for January-June reservoir inflows. A number of potential predictors were evaluated for forecasting these 6-month reservoir inflows, including hydrological persistence (streamflow and precipitation), large-scale climate indices related to the El Niño/Southern Oscillation (ENSO), Pacific Decadal Oscillation (PDO), North Atlantic Oscillation (NAO), and customized sea surface temperature (SST) indices (for detail, see Wei and Watkins (2010). In all cases, monthly index values for the period 1950-2006 were used for the analysis. The SST data used was the extended reconstructed sea surface temperature (ERSST) analysis (Smith et al., 2008), obtained from the National Climatic Data Center (NCDC) through the KNMI Climate Explorer, an on-line data analysis tool (Oldenborgh and Burgers, 2005). The SST pattern having the highest correlation with January-June inflows (p < .01) was identified and referenced as ERSST8, shown in Figure 5.3. A normalized index for ERSST8 was computed based on average seasonal temperatures over a 4-degree by 4-degree area, similar to the procedure of Block and Rajagopalan (2007).



**Figure 5.3 |** Correlation map of Jan.-June reservoir inflows to the Highland Lakes with fall (Oct.-Dec.) sea surface temperatures. Circled regions indicate strong positive and negative correlations used to derive the ERSST8 index.

An algorithm for building accurate and comprehensible logistic regression trees known as Logistic Tree with Unbiased Selection (LOTUS), developed by Chan and Loh (2004), was employed to screen significant predictor variables from all potential predictors observed in fall season (October- December) prior to the forecast (for detail, see Wei and Watkins (2010). The underlying principle of LOTUS is to fit a piecewise (simple or multiple) linear logistic regression model by recursively partitioning the data and fitting a different linear logistic regression to the data in each partition. LOTUS uses a trend-adjusted chi-square test to control bias in variable selection at the intermediate nodes (Cochran, 1954; Armitage, 1955). This can distinguish nonlinear from linear effects and ensure the integrity of inferences drawn from the tree structure. The logistic regression tree obtained from LOTUS for January-June reservoir inflow forecasts is shown in Figure 5.4. Results indicate that both hydrological persistence (streamflow and precipitation) and ERSST8 are statistically significant predictors ($p= 0.05$) for January-June reservoir inflows forecasts. The tercile probabilities for January-June reservoir inflows then are predicted using the logistic regression tree based on these three predictors. The following steps are taken to derive tercile probability forecasts from LRT

1) The 33.3 percentile is first chosen as a threshold, and the response variable $Y$ is equal to 1 if seasonal streamflow (Q) is equal to or less than the threshold, and 0 otherwise.

2) A logistic regression tree is generated to obtain the probability for the below-normal category: $p(BN) = p(Q <= 33.3$ percentile).

Steps 1 and 2 are repeated for the 66.6 percentile, and the probability for normal and above-normal categories, $p(N)$ and $p(AN)$, are derived as follows:

$p(N) = p(Q <= 66.6$ percentile$) – p(BN)$

$p(AN) = 1- p(Q <= 66.6$ percentile$)$ or $p(AN) = 1- p(BN) - p(N)$

Following this procedure, the sum of probabilities for each category is guaranteed to equal 1.

For example, the estimated logistic regression functions for terminal node 1 and 2 are given as follows:

$$\hat{p}_1 = \left[1 + \exp(-0.873 - 0.784 * Streamflow + 0.904 * ERSST8\right]^{-1} \tag{5.1}$$

$$\hat{p}_2 = \left[1 + \exp(-0.811 + 0.973 * Precipitation)\right]^{-1} \tag{5.2}$$

Using the Brier skill score (BSS) and the ranked probability skill score (RPSS) as metrics, forecast performance is assessed through a leave-one-out cross-validation procedure. The BSS and RPSS values show that the January-June reservoir inflows forecasts have an improvement in skill of 16.7% and 13.6% improvement, respectively, over climatology-based forecasts (i.e., tercile probability forecasts equal to 1/3 for each category).



**Figure 5.4**. Plot of stepwise logistic regression tree for January-June reservoir inflows of the Highland Lakes. Intermediate and terminal nodes are represented by circles and squares, respectively. The splitting variable is streamflow, and the splitting value is -0.573. If a case is equal to or less than the splitting value, it goes to the left branch; otherwise to the right branch. The selected fitting variables for each terminal node are given below.

## 5.4  Decision Modeling

To illustrate the beneficial use of information provided by the forecast model developed above, we consider the LCRA's decision of whether or not to sell additional interruptible stored water, exclusive of priority allocation to the Gulf Coast irrigation districts, when January 1 combined storage levels are greater than 94% of maximum storage.  According to the LCRA Revised Water Management Plan, up to 13,000 acre-ft of stored water may be contracted for sale during the first six months of the year. Presumably, this is a "safe" allocation based on a repetition of the drought of record (DOR) beginning on January 1.  We consider whether or not larger contracts may be safely signed when reliable forecast information is available.

To formulate the decision model, we first consider that "stored water" is water that may be stored in the reservoirs after pass-through releases are made for senior water rights holders, in-stream flow maintenance, channel losses, etc. (LCRA 2003).  While the LCRA estimates historical values of storable inflows using a detailed daily simulation model, It is found that the following linear relationship provides a good approximation of this model on a seasonal (6-month) basis ($r^2 = 0.8487$ for a regression of the total January-June inflows and actual stored inflows volumes for the period January-June):

$$Q^{storable} = 0.708\ Q^{total} - 88209 \quad \text{acre-ft} \tag{5.3}$$

Thus, we can consider a certain volume of inflow during the first six months of each year to be passed through and not available for storage. Second, we also must consider that stored water is committed to firm uses (firm yield) and priority irrigation demands. Simulations of the DOR, though lasting for 11 years (historically, 1946-1956), indicate a critical drawdown period of approximately 6 years (1946-1952), at the end of which the combined storage in Lakes Travis and Buchanan would drop to a critically low "reserve" level under current demand levels.  During this drawdown period, the average aggregate storable inflow to the lakes in the first 6 months of each year was 232,905 acre-ft, and so we assume that this volume would be needed to meet firm demands during a repeat of the DOR Thus, our goal is to forecast when the total January-June inflows will exceed the

sum of these two values (88,209 plus 232,905 acre-ft) and to allow the "excess" inflow to be made available for sale.

$$Q^{available} = \text{MAX } (0,\ 0.708 \times Q^{total} - 88209 - 232{,}905) \quad \text{acre-ft} \qquad (5.4)$$

We propose a simplified two-stage stochastic economic-optimization model to investigate improvement in water use efficiency and the potential value of using seasonal forecasts, under the assumption of optimal decision making under uncertainty (Israel and Lund, 1995; Watkins and McKinney, 1997). The model uses ensemble streamflow forecasts as input and includes a risk aversion parameter, which makes it somewhat general for a range of practical applications. Real-world application, however, would likely require a more complex system representation and consideration of additional objectives and constraints.

The conceptual model applied here involves the sale of interruptible water contracts under water supply uncertainty, with the objective of maximizing the seller's expected net revenues from contract sales. In the event that a given contract amount cannot be provided due to lower than expected water availability, the seller is to pay a penalty. Considering a set of scenarios, denoted by $s$, to represent hydrologic uncertainty, the model is formulated as the following linear programming problem:

$$\text{Max} \quad Z = Contract - pen \sum_{s \in S} p_s \times Deficit_s \qquad (5.5)$$

Subject to

$$Contract - Deficit_s \leq Inflows \qquad \forall s$$

$$Deficit_s \geq 0 \qquad \forall s$$

where *Contract* is the contract amount, *pen* is a penalty coefficient selected by the decision maker, $p_s$ is the probability assigned to scenario $s$, $Deficit_s$ is the deficit occurring under scenario $s$, and $Inflow_s$ is the random inflow (water availability) realized under scenario $s$.

In this study, the inflow data set is taken from 56 years (from 1951 to 2006) of aggregate inflows to the Highland Lakes , based on USGS gage measurements upstream of the reservoirs and adjustments made by LCRA staff to account for runoff from ungaged areas. Assuming climatology, each scenario is assumed equally likely, $p_s$ = 1/56. Reservoir inflow forecasts provided by the above seasonal forecast model (based on hydrological persistence and large-scale ocean-atmospheric patterns) enter the optimization model through adjustment of these probabilities to be consistent with the derived tercile probability forecasts. This is similar to the approach of Croley (2000).

## 5.5  Results and Discussion

The water contract optimization model was applied using climatology-based forecasts (equal probabilities of low, medium, and high inflows in all years) and the hydroclimatic forecasts (with a leave-one-out cross validation approach).  Model results are shown in Tables 5.1 and 5.2 for these two cases, respectively, using a range of penalty coefficients to evaluate the tradeoff between maximizing contract amounts and minimizing the risk of deficits occurring. The values listed are the mean seasonal (first six-month of each year) interruptible water contracts and deficits over all 56 years, assuming that the optimal amount is contracted according to the inflow forecast.  If the actual inflow is less than the (predicted) optimal contract, a contract deficit results.  The reliability is computed as the fraction of years in which the actual inflow (water availability) is greater than or equal to the forecast contract. A plot of contract-deficit trade-off curves for climatology and seasonal forecasts is shown in Figure 5.5.

These results demonstrate that incorporating the probabilistic inflow forecasts into the optimization model can provide a significant improvement in seasonal water contract benefits over climatology, with lower average deficits (increased reliability) for a given average contract amount, or improved mean contract benefits for a given level of reliability compared to climatology. Comparing results with LCRA Revised Water Management Plan which calls up to 13,000 acre-ft of stored water for sale during the first

six months of the year when storage levels are high, an additional 16,000 acre-ft of water is available for seasonal contracts with a reliability of 93% (in any year) if reservoir inflow forecasts are used optimally. The results also illustrate the trade-off between the expected contract amount and reliability, i.e., larger contracts can be signed at greater risk.  Keep in mind that the "no risk" option under the assumed circumstances is no additional interruptible water contract, as specified by the LCRA Revised Water Management Plan whenever reservoir storage levels are below 94% of capacity.

Table 5.1.  Trade-off between mean contract amount and deficit (reliability) using climatology-based forecasts (1 AF = 1233.5 m$^3$).

| Mean Contract (AF) | Mean Deficit (AF) | Reliability (%) |
|---|---|---|
| 116537 | 59505 | 41 |
| 95485 | 47238 | 45 |
| 57025 | 26864 | 52 |
| 31220 | 14422 | 52 |
| 21640 | 9804 | 54 |
| 13644 | 6091 | 55 |

**Table 5.2.** Trade-off between mean contract amount and deficit (reliability) using seasonal streamflow forecasts (1 AF = 1233.5 m$^3$).

| Mean Contract (AF) | Mean Deficit (AF) | Reliability (%) |
|:---:|:---:|:---:|
| 137189 | 59382 | 50 |
| 91404 | 36138 | 64 |
| 49989 | 13519 | 77 |
| 29242 | 3239 | 93 |
| 23188 | 2097 | 95 |
| 17491 | 719 | 96 |
| 9437 | 0 | 100 |



**Figure 5.5.** Plot of contract-deficit tradeoff curves for both climatology and seasonal forecasts

Future research will continue to investigate other potential predictor variables for streamflow in Central Texas. In this study, a customized Pacific Ocean SST index was derived because the standard indices of the large-scale climate phenomena were not good predictors, as has been observed in other locations (e.g., Grantz et al., 2006). However, by no means was our search exhaustive. Using new predictors, more skillful streamflow forecast models may also be derived with longer lead times (for instance, one-year lead times would provide valuable information for annual water interruptible water contracts). Other future work will be to develop a more complex multi-reservoir systems optimization model with consideration of additional objectives and constraints and potential to prescribe reservoir releases based on downstream (unregulated) inflow forecasts. Furthermore, an insurance mechanism for limiting risk in the event of bad forecasts may be incorporated into revised operating policies that improve risk-based water resources management and planning.

# 6. Conclusions and Future Work

The overall goal of this study was to investigate improvement in water resources systems management through the use of seasonal climate forecasts. Hydrological persistence (streamflow and precipitation) and large-scale recurrent oceanic-atmospheric patterns such as the El Niño/Southern Oscillation (ENSO), Pacific Decadal Oscillation (PDO), North Atlantic Oscillation (NAO), the Atlantic Multidecadal Oscillation (AMO), the Pacific North American (PNA), and customized sea surface temperature (SST) indices were investigated for their potential to improve streamflow forecast accuracy and increase forecast lead-time in a river basin in Central Texas. The study allowed the following conclusions to be drawn for this specific case study location:

1. Hydrological persistence, alone, shows the greater potential for skillful forecasts than oceanic-atmospheric teleconnections alone. Monthly and seasonal streamflow autocorrelations are much stronger both upstream and downstream of the Highland Lakes reservoir system, especially during the winter and spring seasons. Furthermore, the average monthly and seasonal autocorrelations for downstream data are higher than the corresponding correlations for upstream data. Hydrological persistence does not provide for longer lead time (six months or more) forecasts, however.

2. Nonlinearity exists in hydrologic persistence and in the relationships between streamflow the large-scale ocean-atmospheric patterns examined for flows both upstream and downstream of the Highland Lakes reservoir system. This study used linear correlation analysis as a means of screening potential predictors, but by means of other algorithms such as data mining methods, the nonlinear features are revealed. Some climate indices, such as the derived sea surface temperature patterns SST1 through SST4, PNA, PDO, and AMO were identified and employed as important predicators for both upstream and downstream flows during certain seasons of the year.

3. Seasonal streamflow forecasts with considerable skill were achieved based on distributions-oriented metrics. After developing forecast models for flows both

102

upstream and downstream of the reservoirs, forecast performance was assessed through a leave-one-out cross-validation procedure using the Brier skill score (BSS) and the ranked probability skill score (RPSS) as metrics. The results show that seasonal streamflow forecasts based on either tree-structured models or polytomous logistic regression have significant skill compared to climatology-based forecasts. In particular, both upstream and downstream flow forecasts during winter and spring offer a great improvement over climatology. Results also show that forecast skills for downstream flow forecasts are higher than for upstream flow (reservoir inflow) forecasts.

4. Incorporating the probabilistic inflow forecasts into a simple water contract optimization model indicated significantly increased benefits. Using seasonal streamflow forecasts for January-June can provide a significant improvement in water contract benefits over climatology, with lower average deficits (increased reliability) for a given average contract amount, or improved mean contract benefits for a given level of reliability.. For example, an additional 16,000 acre-ft of water is available for seasonal contracts with a reliability of 93% (in any year) if reservoir inflow forecasts are used optimally. Comparing results with LCRA Revised Water Management Plan, water contracts could increase by 125%. The results also illustrated the trade-off between the expected contract amount and reliability, i.e., larger contracts can be signed at greater risk.

Several of the findings and conclusions from the Central Texas case study are relevant for other locations. First, non-traditional streamflow forecast models (i.e., based on tree-structured data mining techniques) were developed to deal effectively with multicollinearity, nonlinearity and/or interactions present in the data. The tree-structured data mining techniques, i.e., classification and regression trees (CRT) and logistic regression tree (LRT), were also useful in screening large numbers of potential predictor variables and establishing the corresponding forecast models. Compared with traditional modeling approaches such as multiple linear regression, tree-structured data mining techniques offer a flexible and attractive alternative without a *priori* knowledge about the

103

data being analyzed. In addition, the results represented by the binary trees are easy to understand and explain to decision makers as a set of if-then-else rules.

As have many other studies (e.g., Araghinejad et al., 2006), this study demonstrated that more useful information can be provided by probabilistic forecast models than by deterministic forecast models. Probabilistic forecasting methods are preferable for water management and planning because of uncertain initial conditions, limited data resources, and complexity and nonlinearity of hydrometeorological processes. In this study, new methods of obtaining categorical streamflow forecasts, such as tercile probability forecasts, were used for tree-structured models (LCT and CRT) and logistic regression models. In particular, through the polytomous logistic regression model, multi-category probabilities can be generated directly using a single model and the sum of probabilities is guaranteed to equal 1.

Future work in the case study region should include further investigation of other potential predictor variables and their relationships to streamflow. Predictors identified in this study can provide significant improvement in skill and reliability of seasonal streamflow forecasts. However, streamflow forecasts with longer lead times, for example, one-year lead times would provide valuable information for annual interruptible water contracts. In this study, customized Pacific Ocean SST indices were used to forecast streamflow at seasonal time scales. Other studies (e.g., Grantz et al., 2006) have demonstrated that the standard indices of the large-scale climate phenomena were not good predictors in certain locations. To increase the lead times, further investigation may focus on other customized indices of large-scale climate phenomena. For example, Grantz et al. (2006) found that the 500mb geopotential height index could improve the skill of streamflow forecasts at longer lead times than other predictors in the Truckee and Carson River basins in the Sierra Nevada Mountains. Other predictors such as antecedent soil moisture content also show potential for streamflow forecasting in Central Texas and should be incorporated into the prediction model (Watkins et al. 2006). The NCEP NARR soil moisture data set (Mesinger et al., 2005), however, only covers the shorter

time period from 1979 to1999. This is why soil moisture data was not used as a predictor in our study.

Data mining methods were used in this study to screen and identify the potential predictor variables because the relationships between streamflow and climate anomalies are likely to be nonlinear and include multivariate interactions due to the complexity of ocean-atmosphere dynamics. Further research is needed to better understand the physical mechanisms behind these interactions. One approach is to conduct simulation experiments using fully coupled land-oceanic-atmospheric models to investigate the relative importance of different physical mechanisms (e.g., Anyah et al., 2006).

A potential drawback of the classification tree (CT) models developed in this study is that the probabilities are based on the empirical relative frequencies in each category in the terminal nodes of the classification tree. In contrast to the logistic regression tree (LRT) models, the CT models are not able to model variability within a terminal node. Future work may extend the CT approach to classification and regression trees (CRT), with a method to derive tercile probability forecasts based on analysis of regression residuals (errors).

To further address the uncertainty in forecast models, multimodel (superensemble) techniques need to be investigated in the future. Recent studies have demonstrated that multimodel forecasts techniques that combine different individual models, can provide considerable improvement in the overall predictability of streamflows and reduce the overall model error (Devenieni et al., 2008).

In this study, only a simplified stochastic economic-optimization model was developed for seasonal water contracts in the Highland Lakes system in Central Texas. Future work should continue to develop a more complex multi-reservoir systems optimization model with consideration of additional objectives and constraints, with the ability to prescribe reservoir releases based on both reservoir inflow and downstream (unregulated) flow forecasts. Furthermore, an insurance mechanism for limiting risk in

the event of bad forecasts may be incorporated into revised operating policies that improve risk-based water resources management and planning.

# References

Adamoski, K., and W. Feluch (1991). "Application of nonparametric regression to groundwater level prediction." *Can. J. Civ. Eng.*, 18, 600–606.

Alley, W.M. (1985). "Water balance models in one-month-ahead streamflow forecasting," *Water Resources Research*, 21(4), 597-606.

Anderson , M.L., M.L. Kavvas, and M.D. Mierzwa (2001). " Probability/ensemble forecasting; a case study using hydrologic response distributions associated with El Niño/Southern Oscillation(ENSO)," *J. Hydrol*. 249, 134-147.

Anyah, R.O., F.H.M. Semazzi, and L. Xie (2006). "Simulated Physical Mechanisms Associated with Climate Variability over Lake Victoria Basin in East Africa," *Monthly Weather Review*, 134, 3588-3609.

Araghinejad, S., D.H. Brun  and M. Karamouz (2006). "Long-lead probabilistic forecasting of streamflow using oceanic-atmospheric and hydrological predictors," *Water Resour. Res.,* 42, W07411, doi:10.1029/2004WR003853

Armitage, P. (1955). "Tests for linear trends in proportions and frequencies," *Biometrics*, 11, 375–386.

Askew, A. J. (1974). "Chance-constrained dynamic programming and the optimization of water resource systems," *Water Resour. Res*., 101(1), 51-56.

Beale, E. M. L. (1955). "On minimizing a convex function subject to linear inequalities," *J. R. Stat. Soc*., B17, 173-184.

Bender, J. F. (1962). "Partitioning procedures for solving mixed variables programming problems," *Numerische Mathematik*, 4, 238-252.

Benjamin, J. R., and Cornell, C. A. (1970). *Probability, statistics, and decision for civil engineers*, McGraw-Hill, New York.

Bessler, T. F., D.A. Savic, and G.A. Walter (2003). "Water reservoir control with data mining," *J. Water Resour. Plng. & Mgmt*., 129(1), 26-34.

Birge, J.R. (1985). "Decomposition and portioning methods for multistage stochastic linear programs," *Oper. Res*., 33, 989-1007.

Block, P., and B. Rajagopalan (2007). "Interannual variability and ensemble forecast of the Upper Blue Nile Basin Kiremt season forecast," *J. Hydrometeorology*, 8, 327-342.

Bras, R. L., and I. Iturbe (1985). *Random functions and hydrology*, Addison-Wesley, Reading, Mass.

Breiman, L., J. H. Friedman, R. A. Olshen, and C. J. Stone (1984). *Classification and Regression Trees*, Wadsworth, Belmont, California.

Caelli, T. and T. S. Caetano (2005). "Graphical models for graph matching: Approximate models and optimal algorithms," *Pattern Recognition Letters*, 26(3), 339-346.

Chan Kin-Yee and Wei-Yin LOH, (2004). "LOTUS: An algorithm for building accurate and comprehensible Logistic regression trees," *Journal of Computational & Graphical Statistics*, 13(4), 826-852.

Chan, K-Y. and W-Y. Loh (2004). "LOTUS: An algorithm for building accurate and comprehensible logistic regression trees," *Journal of Computational & Graphical Statistics*, 13(4), 826-852.

Clair, T.A., and J.M. Ehrman (1998). "Using neural networks to assess the influence of changing seasonal climates in modifying discharge, dissolved organic carbon, and nitrogen export in eastern Canadian rivers," *Water Resour. Res.,* 34(3), 1031-1022

Clark, M.P., M.C. Serreze, and G.J. McCabe (2001). "Historical effects of El Nino and La Nino events on the seasonal evolution of the montane snowpack in the Columbia and Colorado River basins*," Water Resources Research*, 37(3), 741-757.

Cochran, W. G. (1954). "Some methods of strengthening the common $\chi^2$ tests," *Biometrics*, 10, 417–451.

Cong, S., J. Schaake, and E. Welles (2003). "Retrospective Verification of Ensemble Streamflow Prediction (ESP): A Case Study," Proceedings of the American Meteorological Society Conference, <http: //ams.confex.com/ams/ annual2003/techprogram/paper_54667.htm>

Coulibaly, P., F. Anctil, and B. Bobée (2001). "Multivariate reservoir inflow forecasting using temporal neural networks," *J.Hydrologic Engrg.,* ASCE, 6 (5), 367-376.

Croley, T. E. (2000). *Using Meteorology Probability Forecasts in Operational Hydrology*, ASCE Press, Reston, VA.

Croley, T. E., II, and D.H. Lee (1993). "Evaluation of Great Lakes net basin supply forecasts," *Water Resources Bulletin*, AWRA, 29(2), 267-282.

Dantzig, G.B. (1955). "Linear programming under uncertainty," *Mgmt. Sci.*, 1: 197-206.

Dawson, C.W. and R.L. Wilby (2001). "Hydrological modeling using artificial neural networks," *Progr. Phys. Geogr.* 25, 80-108

Day, G.N. (1985). "Extended streamflow forecasting using NWSRFS," *J. Water Resour. Plng. & Mgmt*, 111(2), 157-170.

Dettinger, M. D., and H. F. Diaz (2000). "Global characteristics of streamflow seasonality and variability," *J Hydrometeorol.*, 1(4), 289–310.

Devineni, D., Sankarasubramanian, A., and Ghosh, S. (2008). "Multi model ensembling of streamflow forecasts: Role of predictor state in developing optimal combinations," *Water Resour. Res.*, 44(9), 1–22.

Dogget, K. (1998). Glossary of Verification Terms, <http://www.sel.noaa.gov/ forecast_verificaton/verif_glossary.html>

Duan, Q., N.K. Ajami, X. Gao and S. Sorooshian (2006). "Multi-model ensemble hydrologic prediction using Bayesian model averaging," *Advances in Water Resources*, 30, 1371-1386.

Eum, H.-I., and Y-O. Kim (2010). "The value of updating ensemble streamflow prediction in reservoir operations," *Hydrological Processes*, 24(20): 2888–2899.

Faber, B. A., and J. R. Stedinger (2001). "Reservoir optimization using sampling SDP with ensemble streamflow prediction (ESP) forecasts," *J. Hydro.* 249(1-4), 113-133.

Fan, Y., and van den Dool, H. (2004). "Climate Prediction Center global monthly soil moisture data set at 0.5° resolution for 1948 to present," *Journal of Geophysical Research*, 109, D10102, doi:10.1029/2003JD004345.

Gassman, H. I. (1990). "MSLiP: A computer code for the multistage stochastic linear programming problem," *Mathematical Programming*, 47, 407-423.

Georgakakos, K. P., and N. E. Graham (2008). "Potential benefits of seasonal inflow prediction uncertainty for reservoir release decisions," *J. Appl. Meteorol. Climatol.,* 47, 1297 – 1321, doi:10.1175/2007JAMC1671.1.

Gershunov, A., and T.P. Barnet (1998). "ENSO influence on intraseasonal extreme rainfall and temperature frequencies in the contiguous United States: Observations and model results," *Journal of Climate*, 11(7), 1575-1586.

Gershunov A. and T. P. Barnet (1998). " ENSO Influence on Intraseasonal Extreme Rainfall and Temperature Frequencies in the Contiguous United States: Observations and Model Results," *Journal of Climate*, 11(7), 1575-1586.

Giustolisi, O. and and D. Laucelli (2005). "Improving generalization of artificial neural networks in rainfall-runoff modeling," *Hydrol. Sci. J.,* 50(3), 439-457.

Goddard, L., A. G. Barnston, and S. J. Mason (2003). "Evaluation of the IRI's ''net assessment'' seasonal climate forecasts: 1997 – 2001," *Bull. Am. Meteorol. Soc.*, 84, 1761 – 1781, doi:10.1175/BAMS-84-12-1761.

Golembesky, K., A. Sankarasubramanian, and N. Devineni (2009). "Improved drought management of Falls Lake Reservoir: Role of multimodel streamflow forecasts in setting up restrictions," *J. Water Resour. Plann. Manage.*, 135(3), 188–197, doi:10.1061/(ASCE)0733-9496(2009)135:3(188).

Govindaraju, S. R. (2001). "ASCE Task Committee on Application of Artificial Neural Networks in Hydrology: Artificial neural networks in hydrology, Parts I and II," *J. Hydrologic Engrg.,* ASCE, 5(2), 115-137.

Grantz, K., B. Rajagopalan, E. Zagona, and M. Clark (2007). "Water Management Applications of Climate-Based Hydrologic Forecasts: Case Study of the Truckee-Carson River Basin," *J. Water Resour. Plann. Manage.*, 133(4), 339-350.

Grantz, K., B. Rajagopalan, M. Clark, and E. Zagona (2006). "A technique for incorporating large-scale climate information in basin-scale ensemble streamflow forecasts," *Water Resour. Res.*, 41, W10410.

Gray, S.T., L.J. Graumlich, J.L. Betancourt and G.T. Pederson (2004). "A tree-ring based reconstruction of the Atlantic Multidecadal Oscillation sing 1567 A. D.," *Geophys, Res. Lett.*, *31*, L12205, doi:10.1029/2004GL019932.

Grygier, J., and J. Stedinger (1985). "Algorithms for optimizing hydropower system operation." *Water Resour. Res.*, 21(1), 1-10.

Gutierrez, F., and J .A. Dracup (2001). "An analysis of the feasibility of long-range streamflow forecasting for Columbia using El Niño-Southern Oscillation indicators," *J. Hydrol.*, 246, 181-196.

Hall, W. A., and J. A. Dracup (1970). *Water Resources Systems Engineering*. New York: McGraw Hill.

Hamlet, A. F., and Lettenmeier, D.P. (1999). "Columbia River streamflow forecasting based on ENSO and PDO climate signals," *J. Water Resour. Plng. & Mgmt.*, 125(6), 333-334.

Hamlet, A. F., Huppert, D., and Lettenmeier, D.P. (2002). "Economic Value of Long-lead Streamflow Forecasts for Columbia River Hydropower," *J. Water Resour. Plng. & Mgmt.*, 128(2), 91-101.

Hamlet, A.F. and Lettenmeier, D.P. (1999). "Columbia River streamflow forecasting based on ENSO and PDO climate signals," *J. Water Resour. Plng. & Mgmt.*, 125(6), 333-334.

Hamlet, A.F., D. Huppert, and D.P. Lettenmaier (2002). "Economic value of long-lead streamflow forecasts for Columbia River hydropower," *J. Water Resour. Plng. And Mgmt.*, ASCE, 128, 91-101.

Han, J., and M. Kamber (2006). *Data Mining: Concepts and Techniques*. Second Edition, Morgan Kaufmann.

Hand, D. J., H. Mannila, and P. Smyth (2001). *Principles of data Mining*. The MIT Press. Cambridge, Massachusetts.

Harrold, T. I., A. Sharma, and S. J. Sheather (2003). "A nonparametric model for stochastic generation of daily rainfall occurrence," *Water Resour. Res.*, 39(10), 1300.

Harrold, T. I., A. Sharma, and S. J. Sheather (2003). "A nonparametric model for stochastic generation of daily rainfall amounts," *Water Resour. Res.*, 39(12), 1343.

Heyden, Y. V, S. T. Popovici, and P. J. Schoenmakers (2002). "Evaluation of size-exclusion chromatography and size-exclusion electrochromatography calibration curves," *Journal of Chromatography A*, 957( 2),127-137.

Hogan, A. J., J. G. Morris, and H. E. Thompson (1981). "Decision problems under risk and chance constrained programming: Dilemmas in the transition," *Manage. Sci.*, 27(6), 698-716.

Horel, J. D., and J. M. Wallace (1981). "Planetary scale atmospheric phenomena associated with the Southern Oscillation," *Mon. Weather Review*, 109: 813-829.

Houck, M. H., and B. Datta (1981). "Performance evaluation of a stochastic optimization model for reservoir design and management with explicit reliability criteria," *Water Resour. Res.*, 17(4), 827-832.

Houseago, R.E., G.R. McGregor, J.C. King, and S.A. Harangozo (1998). "Climate anomaly wave-train patterns linking sothern low and high latitudes during South Pacific warm and cold events," *Int. J. Climatol.*, 18, 1181-1193.

Hurrell, J.W. (1995). Decadal Trends in the North Atlantic Oscillation Regional Temperatures and Precipitation," *Science*, 269: 676-679.

Hurrell, J.W., Y. Kushnir and M. Visbeck (2001). "Perspectives: Climate. The North Atlantic Oscillation," *Science*, 291, 603-604.

Ingram, J.J., M.D. Hudlow, and D.L. Fread (1995). "Hydrometeorological coupling for extended streamflow predictions," *Proceedings, American Meteorological Society (AMS) Conference on Hydrology*, Dallas, Texas, January 15-20, 186-191.

Israel, M. and J. R. Lund (1995). "Recent California water transfers: implications for water management," *Natural Resources Journal, 35(1), 1-32*.

Jacobs, J., G. Freeman, J. Grygier, D. Morton, G. Schultz, K. Staschus, and J. Stedinger (1995). "SOCRATES: A system for scheduling hydroelectric generation under uncertainty." *Annals of Operations Research*, 59(1): 99-133.

Johnson, A.R., and D. W. Wichern (2007). *Applied Multivariate Statistical Analysis*. Upper Saddle River, NJ. Pearson Prentice Hall.

Jolliffe, I.T. 2002. *Principal Component Analysis*. Springer, New York.

Karamouz, M., Houck, M., and J. Delleur (1992). "Optimization and simulation of multiple reservoir systems," J. Water Resour. Plan. Manage., 118(1), 71–81.

Karamouz, M.F. and B. Zahraie (2004) "Seasonal streamflow forecasting using snow budget and El Niño-Southern Oscillation climate signals: Application to the Salt River Basin in Arizona," *J. Hydrologic Engrg*. 9(6), 523-533.

Katz, R. W., and A. H. Murphy (1997). *Economic Value of Weather and Climate Forecasts*, Cambridge University Press, Cambridge, United Kingdom, Ch. 2, pp. 19-74.

Kelman, J., J. Stedinger, L. Cooper, E. Hsu, and S.-Q. Yuan (1990). "Sampling stochastic dynamic programming applied to reservoir operation," *Water Resour. Res.*, 26(3), 447–454.

Kerr, R.A. (2000). "A North Atlantic climate pacemaker for the centuries," *Science*, 228, pp. 1984-1986.

Kim Y-O., H.-I. Eum, E.G. Lee, and I.H. Ko (2007). "Optimizing operational policies of a Korean multi-reservoir system using sampling stochastic dynamic programming with ensemble streamflow prediction," *Journal of Water Resources Planning and Management*, 131(1): 4–14.

Kracman, D.R. (2002). "Stochastic Optimization of the Highland Lakes System in Texas." M.S. Thesis, Engineering, University of Texas at Austin, Austin, TX.

Kracman, D.R., D.C. McKinney, D.W. Watkins Jr., and L.S. Lasdon. (2006). "Stochastic Optimization of the Highland Lakes System in Texas," *Journal of Water Resources Planning and Management*, ASCE, 132(2): 62-70.

Kushnir, Y. (1999). Europe's Winter Prospects. *Nature*, 398: 289-291.

Kutner, H.M., J.C. Nachtsheim and J. Neter (2004). *Applied Linear Regression Models*. McGraw-Hill, New York.

Labadie, J.  (2004). "Optimal Operation of Multireservoir Systems: State-of-the-Art Review." *J. Water Resour. Plan. Manage*. 130(2), 93-111.

Lall, U. (1995). Recent advances in nonparametric function estimation: Hydrologic applications, *U.S. National Report to IUGG, 1991-1994, Reviews of Geophysics*, 33 Suppl. 1995, 1093-1102.

Lall, U., and A. Sharma (1996). "A nearest neighbor bootstrap for resampling hydrologic time series." *Water Resour. Res.*, 32(3), 679–693.

Lall, U., Y.-I. Moon, and K. Bosworth (1993). "Kernel flood frequency estimators: Bandwidth selection and kernel choice." *Water Resour.Res.* 29(4), 1003–1015.

Lee, D. H. (1999). "Institutional and technical barriers to risk-based water resources management: A case study." *J. Water Resour. Plng. Mgmt.*, ASCE, 125(4), 186-193.

Lim, T.-S., W.-Y. Loh, and Y.-S. Shih (1997). An empirical comparison of decision trees and other classification methods. Technical Report 979, Department of Statistics, University of Wisconsin, Madison.

Loucks, D., and P. Dorfman (1975). "An evaluation of some linear decision rules in chance-constrained models for reservoir planning and operation," *Water Resour. Res.*, 11(6), 777–782.

Loucks, D., J. Stedinger, and D. Haith (1981). *Water resource systems planning and analysis*, Prentice-Hall, Englewood Cliffs, N.J.

Lower Colorado River Authority (1999). Water Resource Plan for the Lower Colorado River Basin. Retrieved Janury 2002, from Austin, Texas. Lower Colorado River Authority, http:// www.lcra.org/water/wmp/.

Lower Colorado River Authority (2003). *LCRA Revised Water Management Plan*. Austin, TX.

Lund, J., and I. Ferreira (1996). "Operating rule optimization for Missouri River reservoir system," J. Water Resour. Plan. Manage., 122(4), 287–295.

Maier, H. R., and G. C. Dandy (1996). "Use of artificial neural networks for prediction of water quality parameters," *Water Resour. Res.,* 32(4), 1031-1022

Mantua, N. (1997). Relationships between a naturalized Columbia River flow record and large scale climate variations over the north Pacific. JISAO Climate Impacts Group, Year 2 Progress Rep., University of Washington. Seattle, WA.

Mantua, N. J., S. R. Hare, Y. Zhang, J. M. Wallace, and R. C. Francis (1997). "A pacific interdecadal climate oscillation with impacts on salmon production," *Bull. Am. Meteorol. Soc., 78,* 1069-1079.

Mantua, N.J., S.R. Hare, Y. Zhang, J.M. Wallace, and R.C. Francis (1997). "A Pacific decadal climate oscillation with impacts on salmon," *Bulletin of the American Meteorological Society*, 78: 1069-1079.

Martin, Q.W. (1991). "Drought management plan for Lower Colorado River in Texas." *J. Water Resour. Plng. & Mgmt.*, 117(6), 645-660.

McCabe, G. J., M. A. Palecki, and J. L. Betancourt (2004). "Pacific and Atlantic Ocean Influences on Multidecadal Drought Frequency in the United States." *Proc. Natl. Acad. Sci. U.S.A.*, 101, 4136-4141.

McCullagh, P. and J. A. Nelder (1989). *Generalized Linear Models*, 2nd ed., Chapman & Hall, London.

McLachlan, J.G. (1992). *Discriminant analysis and statistical pattern recognition*. John Wiley & Sons, Inc., Hoboken, New Jersey.

Menard, S. 1995 *Applied Logistic Regression Analysis*. Sage Publishing, Thousand Oaks, CA.

Mesinger, F., et al. (2005). "North American regional reanalysis." Submitted to the *Bulletin of the American Meteorological Society*.

Minobe, S. (1997). "A 50-70 year climatic oscillation over the North Pacific and North America," *Geophysical Research Letters*, 24: 683-686.

Moon, Y.-I., and U. Lall (1994). "Kernel function estimator for flood frequency analysis," *Water Resour. Res.*, 30(11): 3095–3103.

Oldenborgh, G.J. van and G. Burgers (2005). "Searching for decadal variations in ENSO precipitation teleconnections," *Geophys. Res. Lett.,* 32, 15, L15701, doi:10.1029/2005GL 023110.

Pagano, T. C., H. C. Hartmann, and S. Sorooshian (2002). "Factors affecting seasonal forecast use in Arizona water management: A case study of the 1997 – 98 El Nino*," Clim. Res.,* 21(3), 259 – 269, doi:10.3354/cr021259.

Pagano, T.C., H.C. Hartmann, and S. Sorooshian. (2001). "Using climate forecasts for water management: Arizona and the 1997-1998 El Niño," *Journal of the American Water Resources Association*, 37(5), 1139-1153.

Pereira, M. V. F., and L. M. V. G. Pinto (1985). "Stochastic optimization of a multireservoir hydroelectric system: A decomposition approach," *Water Resour. Res.*, 21(6), 779-792.

Philander, S. G. (1990). *El Niño, La Niña, and the Southern Oscillation*. Academic Press, San Diego, CA.

Piechota, T. C., F. H.S. Chiew, J. A. Dracup, and T. A. McMahon (2001). "Development of an exceedance probability streamflow forecast," *J, Hydrologic Engrg.,* 6(1), 20-28.

Piechota, T.C. and J.A. Dracup (1996). "Drought and regional hydrologic variation in the United States: Associations with the EI Niño-Southern Oscillation." *Water Resour. Res.,* 32 (5), 1359-1373.

Piechota, T.C. and J.A. Dracup (1999). Long-Range Streamflow Forecasting Using El Niño-Southern Oscillation Indicators. *Journal of Hydrologic Engineering*, 4(2): 144-151.

Piechota, T.C., F.H.S. Chiew, J.A. Dracup and T.A. McMahon (1998). "Seasonal streamflow forecasting in eastern Australia and the El Nino–Southern Oscillation," *Water. Resour. Res.*, 34, 3035– 3044.

Piechota, T.C., J.A. Dracup and R.G. Fovell (1997). "Western US steamflow and atmospheric circulation patterns during El Niño-Southern Oscillation," *J. Hydrol.* 201, 249-271,

Piechota, T.C., J.D. Garbrecht, and J.M. Schneider, eds. (2006). *Climate Variability and Climate Change*. Reston, VA: ASCE.

Prairie, J. R. (2002). Long-term Salinity Prediction with Uncertainty Analysis: Application for Colorado River above Glenwood Springs, M.S. Thesis, Colorado, University of Colorado at Boulder.

Prairie, J.R., B. Rajagopalan, T. Fulp, and E. Zagona (2006). "Modified K-NN Model for Stochastic Streamflow Simulation," *Journal of Environmental Engineering,* 11(4): 371-378.

Put, R., C. Perrin, F. Questier, D. Coomans, D. L. Massart and Y. Vander Heyden (2003). "Classification and regression tree analysis for molecular descriptor selection and

retention prediction in chromatographic quantitative structure–retention relationship studies," *J.Chromatogr*. A, 988(2), 261-276.

Rajagopalan, B. and U. Lall (1999). "A k-nearest-neighbor simulator for daily precipitation and other weather variables," *Water Resources Research*, 35(10): 3089-3101.

Rajagopalan, B., E. Cook, U. Lall and B. K. Ray (2000). "Spatiotemporal Variability of ENSO and SST Teleconnections to Summer Drought over the United States during the Twentieth Century," *Journal of Climate*,13(24): 4244-4255.

Regonda, S. K., B. Rajagopalan, and M. Clark (2006). "A new method to produce categorical streamflow forecasts," *Water Resour. Res*., 42, W09501, doi:10.1029/2006WR004984.

Regonda, S. K., B. Rajagopalan, M. Clark, and E. Zagona (2006). "A multi-model ensemble forecast framework: Application to spring seasonal flows in the Gunnison River Basin," *Water Resour. Res., 42*, W09404, doi:10.1029/2005WR004653.

ReVelle, C., E. Joeres, and W. Kirby (1969). "Linear decision rules in reservoir management and design 1: Development of the stochastic model," *Water Resour. Res*., **5**(4), 767–777.

Ripley, B. D. (1996). *Pattern Recognition and Neural Networks*. Cambridge University Press, Cambridge.

Rogers, J.C. (1984). "The Association between the North Atlantic Oscillation and the Southern Oscillation in the Northern Hemisphere," *Monthly Weather Review*, 112: 1999-2015.

Ropelewski, C. F., and M. S. Halpert, (1987). "Global and Regional Scale Precipitation Patterns Associated with the El Niño/Southern Oscillation," *Monthly Weather Review*,115: 1606-1626.

Salas, J. D. (1985). "Analysis and modeling of hydrologic time series," *Handbook of hydrology*, D. R. Maidment, ed., McGraw-Hill, New York, 19.1–19.72.

Salas, J. D., J. W. Delleur, V. Yevjevich, and W. L. Lane (1988) *Applied Modeling of Hydrologic Time Series*, Water Resources Publications, Littleton, CO.

Schaake, J., and L. Larson (1998). "Ensemble streamflow prediction (ESP): Progress and research needs," *Special Symp. on Hydrology*, American Meteorological Society, Boston, J19–J24, 1998.

Seber. G.A.F. and A.S. Lee (2003). *Linear Regression Analysis*, 2nd ed., John Wiley & Sons, New York, NY.

Seo, D. J., V. Koren, and N. Cajina (2003). "Real-time variational assimilation of hydrologic and hydrometeorological data into operational hydrologic forecasting," *J. Hydrometeorol., 4:* 627 -641.

Shamseldin, A. Y. (1997). "Application of neural network technique to rainfall-runoff modeling," *J. Hydro.,* 199, 272-294.

Sharma, A. (2000). "Seasonal to interannual rainfall probabilistic forecasts for improved water supply management: 1. A strategy for system predictor identification," *J. Hydrol.*, 239, 249-258.

Sharma, A., D. G. Tarboton, and U. Lall (1997). "Streamflow simulation: A nonparametric approach." *Water Resour. Res.*, 33(2): 291–308.

Singhrattna, N., B. Rajagopalan, M. Clark and K. Krishna Kumar (2005). "Forecasting Thailand Summer Monsoon Rainfall," *Int. J. Climatology*, 25: 649-664.

Smith, J. A. (1991). "Long-range streamflow forecasting using nonparametric regression." *Water Resour. Res.*, 27(1): 39–46.

Smith, J., G. N. Day, and M. D. Kane (1992). "Nonparametric framework for long range streamflow forecasting," *J. Water Resour. Plng. and Mgmt.*, ASCE, 118(1), 82-92.

Smith, T.M., R.W. Reynolds, T.C. Peterson and J. Lawrimore (2008). "Improvements to NOAA's Historical Merged Land-Ocean Surface Temperature Analysis (1880-2006)," *Journal of Climate*, 21, 2283-2296.

Sniedovich, M. (1979). "Reliability-constrained reservoir control problems, 1, Methodological issues," *Water Resour. Res.*, 15(6), 1574-1582.

Souza Filho, F.A. and U. Lall (2003). "Seasonal to interannual ensemble streamflow forecasts for Ceara, Brazil: Applications of a multivariate, semi-parametric algorithm," *Water Resources Research*, 39(11), 1307-1325.

Stedinger J., and Y. –O., Kim (2002). "Updating ensemble probabilities based on climate forecasts", *Conf. On Water Resources Planning and Management Session C2 (CD), ASCE*, Roanoke, VA,

Stedinger, J. (1984). "The performance of LDR models for preliminary design and reservoir operation," *Water Resour. Res*., 20(2), 215–224.

Stedinger, J., B. F. Sule, and D. Pei (1983). "Multiple reservoir system screening models," *Water Resour. Res*., 19(6), 1383-1393.

Stedinger, J., B. Sule, and D. Loucks (1984). "Stochastic dynamic programming models for reservoir operation optimization," *Water Resour. Res*., 20(11), 1499–1505.

Steinberg, D., and P. Colla (1997). *CART—Classification and Regression Trees*. Salford Syst., San Diego, CA.

Stephenson, D. B. (2003). "Glossary," in Jolliffe, I. T., and D. B. Stephenson (eds.) *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, Wiley and Sons Ltd., UK.

Sun, L., D.F. Moncunill, H. Li, A.D. Moura, F.A. Souza Filho and S.E. Zebiak (2006). "An operational dynamical downscaling prediction system for Nordeste Brazil and the 2002-04 Real-Time Forecast Evaluation," *Journal of Climate*, 19, 1990-2007.

Tarboton, D. G., A. Sharma, and U. Lall (1998). "Disaggregation procedures for stochastic hydrology based on nonparametric density estimation," *Water Resour. Res.*, 34(1): 107–119.

Ting, C.-K. (2005). "On the Mean Convergence Time of Multi-parent Genetic Algorithms Without Selection," *Advances in Artificial Life, Vol(3630)*: 403–412.

Thompson, D. W. J., and J. M. Wallace (1998). "The Arctic Oscillation signature in the wintertime geopotential height and temperature fields," *Geophys. Res. Lett*., 25(9): 1297-1300.

Tootle, G. A., T. C. Piechota and A. Singh (2005). "Coupled oceanic-atmospheric variability and U.S. streamflow," *Water Resour. Res*., 41, W12408.

Tootle, G.A. & T.C. Piechota (2006). "Relationships between Pacific and Atlantic Ocean sea surface temperatures and U.S. streamflow variability," *Water Resour. Res.,* 42, W07411, doi:10.1029/2005WR004184.

Tootle, G.A., T.C. Piechota and A. Singh (2005). "Coupled oceanic-atmospheric variability and U.S. streamflow," *Water Resour. Res*., 41, W12408.

Tu, J.V. (1996). "Advantages and disadvantages of using artificial neural networks versus logistic regression for predicting medical outcomes**,"** *J. Clin. Epidemiol*., 49, 1225-1231.

Van Slyke, R. M., and R. J-B. Wets (1969). "L-shaped linear programs with applications to optimal control and stochastic programming," SLAM J. of Appl. Math., 17, 638-663.

Vander Heyden, Y., S.T. Popovici and P.J. Schoenmakers (2002). "Evaluation of size-exclusion chromatography and size-exclusion electrochromatography calibration curves," *J. Chromatogr., A 957*.

Wardlaw, R., and M. Sharif (1999). "Evaluation of genetic algorithms for optimal reservoir system operation." *J. Water Resour. Plan. Manage.,* 125(1),25-33

Watkins, D.W. Jr. & S.M. O'Connell (2005). "Teleconnections and disconnections in Central Texas: A guide for water managers," in *Climate Variations, Climate Change and Water Resources Engineering*, eds. J. Garbrecht and T. Piechota, ASCE Press, Reston, VA, pp. 103-114.

Watkins, D.W. Jr., and D.C. McKinney. (1997). "Finding Robust Solutions to Water Resources Problems," *Journal of Water Resources Planning and Management*, ASCE, 123(1): 49-58.

Watkins, D.W. Jr., W. Wei, and D.K. Nykanen (2006). "Simple Forecast-Operations Model Using Hydrologic Persistence," Proceedings, 7th ASCE/EWRI Operations Management Workshop, Sacramento, CA, August.

Watkins, D.W., and W. Wei (2008). "The Value of Seasonal Climate Forecasts and Why Water Managers Don't Use Them," ASCE World Environmental and Water Resources Congress, Ahupua'a, Hawaii..

Watkins, D.W., Jr., McKinney, D.C., Lasdon, L.S., Nielsen, S.S., and Martin, Q.W. (2000). "A scenario-based stochastic programming model for water supplies from the highland lakes," *Intl. Trans. in Op. Res*. **7**, 211-230.

Wei W. and D. W. Watkins (2010). "Probabilistic Streamflow Forecasts Based on Hydrologic Persistence and Large-Scale Climate Signals in Central Texas," *J. Hydroinf.,* in press.

Wilks, D. (1995). *Statistical methods in atmospheric science: An Introduction*. Academic Press: San Diego, California, USA.

World Commission on Dams. (2000). Dams and development: A new framework for decision-making, Earthscan Publications Ltd., London and Sterling, Va.

Wurbs, R. (2005). "Texas water availability modeling system," *J. Water Resour. Plng. and Mgmt*., 131(4), 270-279.

Wurbs, R.A. (2008). Fundamentals of Water Availability Modeling with WRAP, Technical Report 283, Texas Water Resources Institute, 4[th] Edition.

Yates, D.S., S. Gangopadhyay, B. Rajagopalan, and K. Strzepek (2003). "A technique for generating regional climate scenarios using a nearest neighbor bootstrap," *Water Resources Research*, 39, No. 7, 1199.

Yee, T.W. (2010). "The VGAM package for categorical data analysis," *J. Statistical Software*, 32(10), 1-34.

Yeh, W. (1985). "Reservoir management and operations models: A state-of-the-art review," *Water Resour. Res*., 21(12), 1797–1818.

Yin, Z.-Y. (1994a). "Moisture conditions in the southeastern USA and teleconnection patterns," *International Journal of Climatology*, 14: 947-967.

Yin, Z.-Y. (1994b). "Reconstruction of the winter Pacific-North American teleconnection pattern during 1895-1947 and its application in climatological studies," *Climate Research*, 4(2): 79-94.

Zhang, Y., J.M. Wallace, D.S. Battisti (1997). "ENSO-like interdecadal variability: 1900-93," *J. Climate*, 10: 1004-1020.

# Appendix A: Copyright Permission for Chapter 3

**IWA** | **Publishing**

Alliance House
12 Caxton Street
London SW1H 0QS
United Kingdom
Tel:  +44 (0)20 7654 5500
Fax:  +44 (0)20 7654 5555
Email: publications@iwap.co.uk
www.iwapublishing.com

Wenge Wei
PhD candidate
Dept. of Civil & Environmental Engineering
Michigan Technological University
Houghton
MI 49931

23 November 2010

Dear Wenge Wei,

In response to your request for copyright clearance to reproduce the following article:

Wenge Wei and David W. Watkins, Jr. (2010) "Probabilistic streamflow forecasts based on hydrologic persistence and large-scale climate signals in central Texas", *Journal of Hydroinformatics*, In Press, Uncorrected Proof © IWA Publishing 2010 |

in the thesis of Wenge Wei, to be published by the Michigan Technological University, Houghton; we are very happy to grant you permission to reproduce the material specified above without charge, provided that:

- the material to be used has appeared in our publication without credit or acknowledgement to another source;
- suitable acknowledgement to the source is given in accordance with standard editorial practice, e.g.,

"Reprinted from the *Journal of Hydroinformatics*, with permission from the copyright holders, IWA Publishing"

- reproduction of this material is confined to the purpose for which this permission is given.

I trust this permission will be satisfactory; if any point needs clarification or you have any further queries, please do not hesitate to contact us again.

Yours sincerely

Victoria Beddow
Publishing Assistant