

"Development of Application Software for Watershed Data Management, Visualization, and Analysis"

Nick Zinck, University of Massachusetts

March 2018

Contents

1	Introduction	4
1.1	Objective	4
1.2	Scope of Work	5
2	Background	7
2.1	Watershed Overview	7
2.1.1	Quabbin and Wachusett Reservoir	7
2.1.2	Water Quality Sampling Plan	8
2.1.3	Previous Watershed Studies	10
2.2	Watershed Data Management System.	12
2.2.1	Data Storage (needs work)	13
2.2.2	Application Development Frameworks	13
3	Application Development Process	15
3.1	Database Development	15
3.2	Code Development	15
3.2.1	WAVE	16

3.2.2	WIT	16
3.3	Application Deployment	17
4	WIT Features	19
4.1	Raw Data File Lookup	19
4.2	Data Processing	19
4.3	Quality Control	19
4.4	Data Importation	19
5	WAVE Features	20
5.1	Data Query and Export	20
5.2	Data Visualization	20
5.2.1	Time Series Scatter Plot	20
5.2.2	Correlation Scatter Plot	22
5.2.3	Distribution	22
5.2.4	Heatmap (Interpolated Color Profile Plot)	22
5.2.5	Profile Line Plot	23
5.2.6	Phytoplankton	23
5.3	Statistics	23
5.3.1	Spatial and Temporal Statistics	23
5.3.2	Pearson Correlation Matrix	23
5.4	Geospatial Data Mapping	24
5.5	MetaData	24
5.6	Data Import	24
6	Discussion and Reccomendations	25
6.1	Pros and Cons of Application	26
6.2	Future Work	26
6.2.1	Meteorological and Hydrological Data	26

6.2.2	Forestry	26
6.2.3	Reports	26
7	Appendix	27
7.1	WAVE Developer Manual	27
7.2	WIT Developer Manual	27
8	Extras and Trash Bin (Not a real Section)	27
	References	29

1 Introduction

A comprehensive watershed monitoring program includes the collection of water quality, meteorological, and hydrological data. Large amounts of data can be strenuous to manage if proper systems are not put in place for proper data management. Poor data management can result in poor data quality as well as underutilization of data due to the timely process of querying, visualizing, and analyzing poorly-managed data. Spreadsheet software is likely not an effective way to store large datasets and the visualization and analysis tools that spreadsheet software offer are often limited. Database software is a great solution to store and organize large datasets, yet database software often lack data visualization and analysis tools. Commonly, large datasets stored within a database are paired with an outside application specialized for data querying, visualization, and analysis. An application program was developed for watershed data visualization and analysis by means of a free and open source application development framework. A second application program was developed to facilitate importing watershed data into a database. This project is a product of the collaboration between UMass and the Massachusetts Department of Conservation and Recreation (DCR).

1.1 Objective

This project works to facilitate the DCR's data entry, querying, visualization, and analysis process through an R-based application creation tool called Shiny. An application can automate many of the tedious day to day processes of managing a watershed as well as allow for powerful and expeditious visualization and analysis. Developing an application, rather than using propriety software, allows for full customization by the developer to better target specific needs of the agency responsible for watershed management and protection. A well fitting application can greatly increase the timeliness and ability to explain data and generate insights which can direct decision making for these agencies. Increased water quality data insight can also greatly influence an agency's sampling plan to better represent the

watershed and focus on certain areas of high interest. The overall objective of this project is to maximize the efficiency of the DCR's ability to manage, visualize, and analyze data to inform decision making. The application developed in this project can be used as an example for other watershed management agencies.

1.2 Scope of Work

Two applications were created to facilitate with watershed data management at the DCR. These applications will be used across the three watersheds: Quabbin, Ware River, and Wachusett, which are under the management of two separate DCR offices, the Quabbin and Wachusett Offices. The applications are designed to meet the needs of both offices which require necessary variations in the application features to address the variation between the two offices needs including differing sampling plans. Although some variations are inevitable, efforts have been made to make the data management of the two offices more congruent. A large piece of this congruency has been in the form of database alterations to make the organization and formatting as similar as possible between the two offices. This effort has also included moving data that exists outside of a database into a common database where data observations all share a similar tidy data regimen. Naming conventions were also examined and modified to simplify the application creation and decrease potential for mistakes.

Watershed data Importer Tool (WIT) is the smaller of the two applications which is designed an interface for facilitated raw data import of watershed field and laboratory data. Raw data from a number of predetermined sources can be imported through a simple user interface. WIT transforms raw data into the desired database format, provides quality control checks, and imports the data to the right location in the database. WIT will help ensure in the future that all data remains stored in databases as this tool makes database storage timely and efficient.

Watershed data Analysis and Visualization Environment (WAVE) is the larger of the two

applications and as its name suggests it is designed to facilitate data querying, visualization, and analysis of water quality. The application opens in a web browser and allows the user to query data by user selection inputs including locations, parameters, and dates. More advanced filters for data querying are also offered including filtering data based on meteorological events and excluding flagged data. The corresponding queried dataset will be presented in an interactive table which the user will be able to output as a csv file. Discrepancies between watersheds and datasets exist (and are desired) which are reflected in variations of query selections and filters in WAVE.

WAVE consists of numerous visualization and analysis tools to give more insight on a queried dataset. Visualization tools for tributary and reservoir data include time-series plots, correlation plots, and distribution charts (histograms, density curves, and box-and-whisker plots) to visualize trends and patterns on selected water quality parameters. Geospatial data visualization and analysis allows one to spatially view data statistics on an interactive map. Heatmap and profile line plot tools are also available for reservoir profile data. Statistics can be quickly generated with WAVE including min, max, average for user selected temporal and spatial groupings. More advanced statistical analysis includes Mann-kendall statistics and pearson correlation matrixes. WAVE also makes information related to the sampling history of a specific site or a specific parameters easily accesible.

Both WAVE and WIT are shared openly through Github and are ran locally on a computer with minimal setup, although hosting the application online is also a possibility. The application is organized in a modular manner which eases future updates to the code as well as minimizes code repitition. Future additions to WAVE can be added as a seperate and independent module. A developer manual was also created to help future developers of WAVE as well as user setup guides for WAVE and WIT. As the needs of the DCR changes, both applications can dynamically change with it.

2 Background

2.1 Watershed Overview

2.1.1 Quabbin and Wachusett Reservoir

The metropolitan area of Boston, Massachusetts, receives its drinking water from the Massachusetts Water Resource Association (MWRA) water supply system. The sources of the water supply are the Quabbin Reservoir and the Wachusett Reservoir which are managed in partnership with the Department of Conservation and Recreation (DCR). Both the reservoirs combined supply about 200 million gallons per day to consumers [MWRA " "]. The Quabbin and Wachusett Reservoirs are protected and over 85% of the watershed lands that surround the reservoirs are covered in forest and wetlands. The water supply system is rather unique in that the Quabbin Aqueduct, a 24.6 mile long tunnel, connects the Quabbin Reservoir to the Wachusett Reservoir. The MWRA transfers water from the Quabbin reservoir intermittently to the Wachusett Reservoir to maintain the water level and water quality of the Wachusett Reservoir [DCR, 2007] which most always makes up the majority of the total inflow to the Wachusett Reservoir. The MWRA can also divert water from the Ware River, located between the two reservoirs, to either the Quabbin Reservoir (preferable) or the Wachusett Reservoir through the Quabbin Aqueduct, but these diversions are rare. The Ware River usually has poorer water quality which makes an input location with a long hydraulic retention time desirable, thus the Ware River is most always diverted to a location in the Quabbin Reservoir that is hydraulically far from the Quabbin Aqueduct. Transfers generally occurs from June through November and can last for weeks at a time to meet higher water demands, maintain the water level, and mitigate water quality concerns in the Wachusett Reservoir (DCR, 2007).

The complex nature of the system allows for decision making in reservoir management that can alter the water quality that ends up in the Boston Supply System. Although all water enters from one location in the Wachusett Reservoir, it is essential to monitor the water quality in the whole watershed to best understand the reservoir processes. A greater

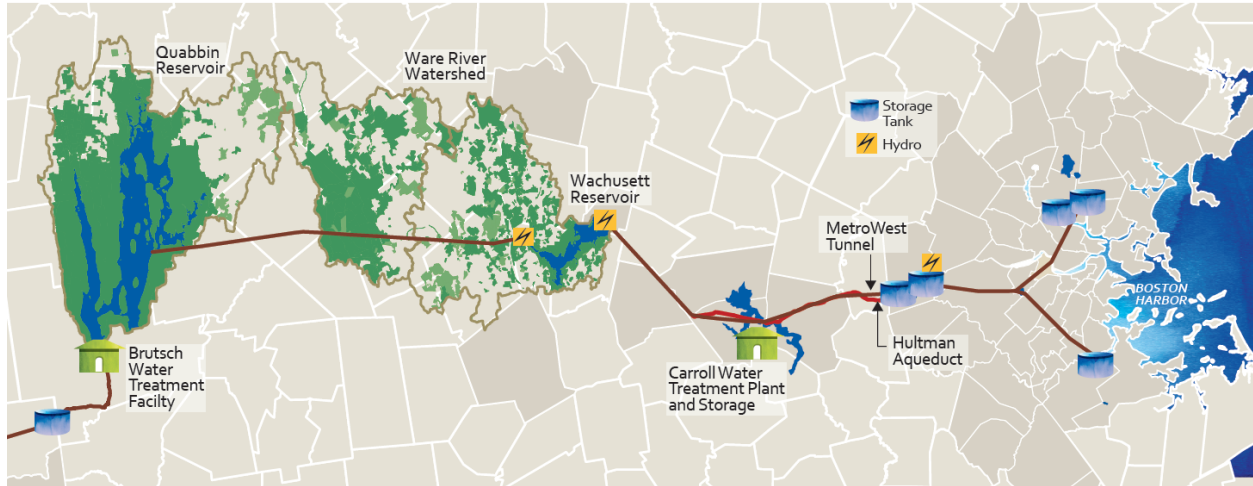


Figure 2.1: MWRA Water Supply System

understanding of reservoir processes can allow for more informed actions when water quality issues occur. Tributaries and reservoirs are routinely sampled by the DCR through specified sampling plans for various water quality parameters. Stream gauges and meteorological instruments also are used to collect tributary discharge and watershed meteorological data, respectively.

2.1.2 Water Quality Sampling Plan

The task of water quality sampling of the three watersheds in the water supply system are divided between the Quabbin and Wachusett office of the DCR Water Supply Protection Division. Each office creates their own sampling plan for the reservoirs and tributaries in their watersheds each are responsible for which consists of a schedule of routine sampling and storm sampling at various sites. The Quabbin Office is responsible for the sampling of the tributaries in the Quabbin and Ware River watersheds and sampling in the Quabbin Reservoir itself. The Wachusett Office is responsible for the sampling of the tributaries in the Wachusett watershed and the Wachusett Reservoir itself. Each Office has sites that they sample consistently every year which are considered core sites. Other Site locations are temporarily selected for various reasons including to target areas of concerns and to target

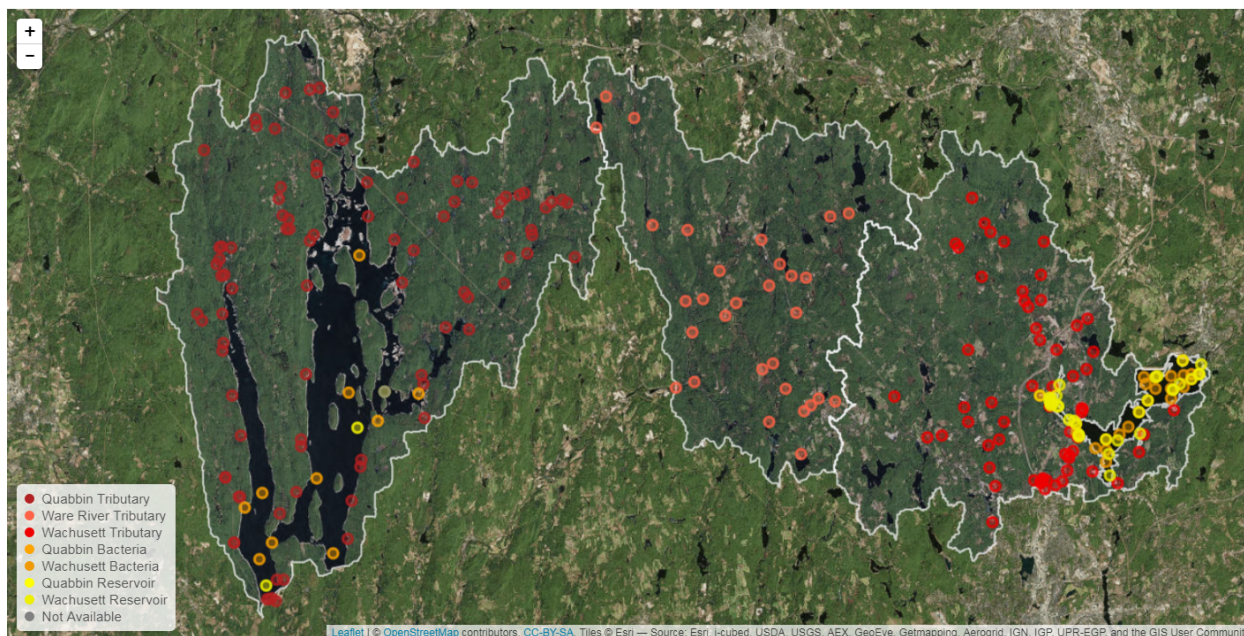


Figure 2.2: Watershed Sampling Locations

areas that are usually underrepresented in the sampling plan.

The tributary Water quality parameters sampled in both offices are turbidity, water temperature, E. coli, mean UV254, ammonia, nitrate, nitrite, specific conductance, total kjeldahl nitrogen, total phosphorus. Most of these parameters are sampled on a biweekly basis, although () are sampled (). The sampling frequencies of water quality parameters in each watershed is displayed in image 3. The Sampling frequency denoted in the graphic are for the years 2013 to 2017.

Reservoir water quality samples are collected by boat during the time of the year with no ice freeze. Various collection methods are used which are surface grab samples, samples at three depths, and profile sampling. Surface grab samples are collected at the Wachusett Reservoir for E. Coli and fecal coliform. () are collected at reservoir locations, at three depths: surface, mid, and deep. The parameters collected at both reservoirs include . Profile Data is collected in both reservoirs at about a meter depth. Bacteria Samples are collected in the Wachusett at the surface only which include (). Phytoplankton data is also collected at both Reservoirs. Do I include Bacteria? brief? The sampling frequencies of water quality

parameters sampled in within each reservoir are displayed in image 2. The Sampling frequency denoted in the graphic are for the years 2013 to 2017.

In addition to sampling plans for water quality, stream gauges are installed to estimate the flow of major tributaries to the reservoirs, both by the United States Geological Survey (USGS) and by the DCR. The Quinapoxet and the Stillwater River are gauged by the USGS which were installed in the 1990's. Some meteorological data is also collected at various locations in the watersheds and additional data can be attained from the NOAA / NCDC. Water quality data related to forestry practices is also collected to compare water quality parameters between managed and unmanaged forests.

2.1.3 Previous Watershed Studies

Many studies have been conducted under the partnership of University of Massachusetts and the DCR. Many of these studies have benefitted from the vast amount of water quality and hydrological data that is collected in the Quabbin, Ware River, and Wachusett watersheds for data analysis and model calibration and validation. The water quality impacts from extreme precipitation events have been studied through statistical forecasting of potential loads coupled with a reservoir hydraulic and water quality model (Jeznach HageMan). The fate of a contaminant spill in the Wachusett Reservoirs have also been modeled (Lily 2013, 2011 Devonis) (Sojkowski) (more). Other water quality modeling topics have included the effects of climate change on the Wachusett reservoir (Lily) and modeling the fate of Natural Organic Matter and fecal pollution (). Having a proper sampling plan can allow for more informed research studies. Water Quality (profile) data is essential for the calibration and validation of hydraulic and water quality models of a Reservoir. Well maintained data further facilitates this process.

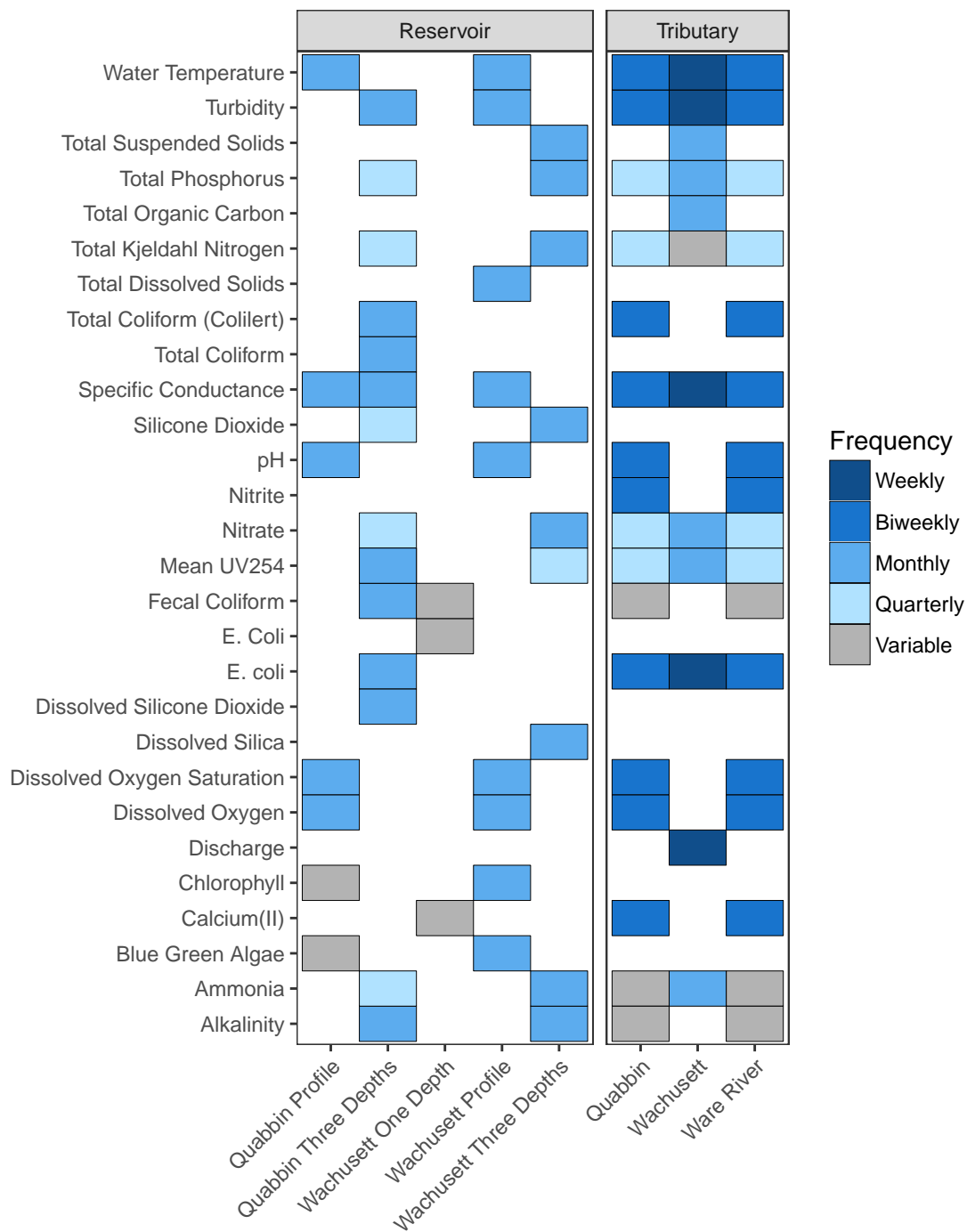


Figure 2.3: Reservoir and Tributary Sampling Frequencies (Double Check)

2.2 Watershed Data Management System.

A simple depiction of a watershed data management system can be divided into three parts:

1. data collection, 2. data storage, and 3. data utilization. Data collection includes the act measuring data observations in the field or a lab as well as importing data measurements into a database. Data storage is the location in which the data lives and is ever changing as more data is collected. Data utilization includes the querying, visualization, analysis, or any other process from which conclusions can be drawn from the data to benefit decision making.

Proprietary software exists that incorporate all three parts of a watershed data management system. Aquarius, created by Aquatic Informatics, is a group of leading software services for water supply management and include AQUARIUS Samples and AQUARIUS Time-Series. AQUARIUS Samples is designed for watershed data management for a variety of field and laboratory data. amplex streamlines the production and management of environmental lab and field sample data, saving time while increasing the quality of final data [aquaticinformatics]. Aquatic Informatics also created Aquarius Time-Series software is the most powerful platform for managing water resources and is used by the USGS as well as many other agencies [aquaticinformatics]. This paper will focus on water quality data management and thus AQUARIUS Samples better resembles the watershed data management system that this paper describes. Proprietary software like AQUARIUS Samples attempt to supply versatile a product that can be widely used on any watershed. This can be a difficult task due to the variations in water quality data and variations of specific agency needs. Data visualization and analysis of the software is often limited and cannot be customized within the software. Proprietary software is also costly.

A personalized watershed management system is another approach for a watershed data management system which can be developed using an open source application development framework paired with a database. With the advancement of programming languages, it is getting easier to build a custom dashboard or application with a application development framework. Creating an application, rather than using a proprietary software allows for more

flexibility and customization. This is a less expensive solution than paying for proprietary software, though it can be time consuming creating an application, so cost and time are a trade off.

2.2.1 Data Storage (needs work)

There are various types of data collected in the watershed. Due to the vast amount of data from various sources the raw data can be in very different formats. Storing data in a central database in similar format is crucial to set up a platform for efficient data analysis. The application provides a connection to this data in the database with an easy-to-use user interface for a human user to interact with the data. Relational databases are a common way to store data which all use a common API SQL. MYSQL, Microsoft Access, are examples of a database. There are many open source and free sources of databases.

2.2.2 Application Development Frameworks

Applications can be created to assist in data import, querying, visualization, and analysis. Application frameworks can assist on application creation which are essentially a reusable, “semi-complete” template application that can be specialized to produce custom applications []. Selection of a proper application development framework can greatly decrease the amount of work must do and knowledge one must possess to create an application. Most application development frameworks can be considered either a front-end framework or a back-end framework which help create a front-end server and a back-end server, respectively. The front-end server and the back-end server communicate through a common API, usually JSON, which is considered the universal binary []. A front-end server is responsible for the construction and layout of the user interface, which is what the user sees and interacts with, commonly in their web browser. The user does not see the back-end server but this does the computational work. The separation of front-end server and back-end server can allow for increased customization as various front-end and back-end frameworks can be paired. Some

frameworks serve as both a front-end and back-end framework which usually have a benefit of simplicity in the development process. This paper focuses on this latter type of application framework that function as both a front-end and back-end due to the type of application desired. Only a minimilistic user interface is needed and most of the focus will be on the data science potential of the the programming language that the framework uses.

R and Python are common programming languages for data science and both have their own collection of data science application framework libraries. Python is a much more widely used, more developed language than R, though R specializes in statistics and data visualization. Shiny is a development framework library in R which allows a relatively unexperienced developer to build an application in the R language. Shiny is an open source application that consists of a front-end user interface component which can be opened with any common web browser as well as a back-end server component which utilizes the data science power of R for visualization and statistical analysis. Shiny is a development framework library in R created by RStudio and is very well documented with thorough tutorials. Application Frameworks libraries in Python focused on data science include Bokeh, Spyre, and Dash [source]. All of these libraries are free and open source and require little application development experience, although programming experience is required (or highly preferred). These application frameworks leverage JavaScript and HTML to render the user interface inputs and outputs in a way that little knowlege of JavaScript or HTML is neccesary. Each framework has its differences and some might work better than others for a particular application design. Familiarity of a partcular programming language can also influence the decision of which application development framework to choose.

3 Application Development Process

3.1 Database Development

Access database.

4 databases

Spreadsheets to databases. Now i can have opinions slightly.

A tidy data format is not only better for data science, it is also essential for proper database storage. This is the best way to link different data tables and relate data tables together via keys to allow for joins to be possible. Storing data in multiple data tables of related information is much more efficient than. Metadata is considered. An example of tidy data is a “Parameter” and “Result” column rather than a column for each parameter. This also allows for a single “Units” column which would not have been possible with the wide format.

In the future it may makes sense to move towards a different database software due to the the outdating of Microsoft Access or whether be the DCR terminating the Microsoft Access license. If so this will likley to a relatively smooth process in converting to anohter relational database since the all have a common API, SQL.

3.2 Code Development

R is the primary language used to write the code, although there are also direct uses of CSS and one or two instances where direct javascript is used. Other languages are used indirectly if a function in R is written in another language.

R is a language and environment for statistical computing and graphics. It is a GNU project which is similar to the S language and environment which was developed at Bell Laboratories (formerly AT&T, now Lucent Technologies) by John Chambers and colleagues. R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, .) and graphical techniques, and is highly

extensible. One of R's strengths is the ease with which well-designed publication-quality plots can be produced, including mathematical symbols and formulae where needed. R is available as Free Software under the terms of the Free Software Foundation's GNU General Public License in source code form. It compiles and runs on UNIX, Windows and MacOS. (???)

Many users think of R as a statistics system. We prefer to think of it of an environment within which statistical techniques are implemented. R can be extended (easily) via packages. There are about eight packages supplied with the R distribution and many more are available through the CRAN family of Internet sites covering a very wide range of modern statistics. R has its own LaTeX-like documentation format, which is used to supply comprehensive documentation, both on-line in a number of formats and in hardcopy.

RStudio is the leading (IDE) for the R language. RStudio creates a number of packages that have driven R to be highly used. Among some of the more highly used packages in R are ggplot, dplyr, tidyr, RMarkdown, and Shiny. Many of RStudio's packages are contained under the " " called tidyverse which can savetime and be loaded all together.

Shiny is

Github

3.2.1 WAVE

Modules and Functions

Naming Conventions

Developer Manual

3.2.2 WIT

There are many different effective ways to go about data importation but importing data by hand is no longer one of them. Most always data that one collects is not in the same format as how the data is store. Data transformations must be made including but not limited to column name changes, adding or removing columns, data type, and spreading or

gathering columns. Spreadsheets also can be used to import data, though manually typing data and making format changes can be timely and prone to errors. Programming languages can be used as a tool to transform data to a particular format. R and Python are just two of the many programming languages that can be used for this process. A code script written with a programming languages can additional perform set quality control measures including checking for duplicate data or alerting a user when their is an usual data value. Even further, a user interface can be created for a user to be able to run this code without being familiar with the programming language or integrated development environment (IDE) is note required. A common user interface can be built to input data of all types.

3.3 Application Deployment

There were many possibilities of ways to launch hte application. A large driver of this decision is based on the answer to the question: who do we want to be able to have access to the app. Becuase this application is to be used primarily by the DCR internally, it was not neccesary to host the application on the worled wide web. Since there was a discrete small set of computers that need to have access to the application, it is feasible to install this application on each of these computers to run locally. This would not be possible if it was desired for an application with access from users on any computer connected to teh internet. If it was neccesary to do so we would have had to set uop our own server, pay for the hosting service taht Shiny directly offers, or pay for another outside cloud hosting sevice. This also felt like the safer bet for any data security purposes since not hosting the data on the web allows for the data to stay on the DCR's internal network.

There are two primary launching options to launch a shiny App locally. The first is by executing the code directly, or using RStudio Runapp button, which requires the application script to be stored locally. The second launching option is made possible by a built in function in Shiny called RunGitHub which will fetch the most updated code stored in a Github Repository. Github also happens to be, not coincidentaly, the leading code sharing,

storage

A desktop shortcut can be created to allow the user. figure __ shows the picture of the desktop icon used for WAVE which is a modified version of Department of Conservation logo. The desktop shortcut when clicked executes batch file to run the R script to install/load packages and the Rungithub function. A configuration file was created to allow for customized computer settings including personalized user settings. The configuration file includes information

Although a user does not need a copy of the code on the user's computer, using RunGitHub command in Shiny still requires R to be installed on a computer. To avoid this requirement of downloading R, Shiny can be packaged with a porta. Also, the computer will need a web browser installed. This is most definitely already available on a working computer, yet there can be discrepancies between how various web browser's interpret HTML and Javascript, so it is safer to use have a consistent browser. (Reference for packaged) .

The Application can also be packaged. Include link to this details. A portable Chrome and portable R

4 WIT Features

4.1 Raw Data File Lookup

4.2 Data Processing

4.3 Quality Control

4.4 Data Importation

5 WAVE Features

5.1 Data Query and Export

It is advantageous, if not essential, to have fast access to water quality and other watershed data. It is beneficial for any scientist or engineer to be able to access this data with ease. In practice, this is often not the case due to all data not being stored in a location known by all scientists and engineers. Lack of experience with certain technologies can also inhibit a person from being able to access timely data. Although spreadsheets are commonplace among most workplace settings, being able to query this data for the exact data that a person is looking for can be timely and troublesome. Queries in a relational database is a better approach to this common task.

The application queries the data based on the user selected input. The user is prompted to first select one or more Site Locations with a map that indicates which sites are selected. The user is then prompted to select the parameters and date range in which data is available for the selected sites. Additional filters can also be applied to the data which include selection of seasons, months, years, flags, storm event, and other eventually meteorological and hydrological conditions. Queried data can be exported to a csv file with the click of a button.

5.2 Data Visualization

This application allows for water quality visualization of temporal trends and temporal statistical analysis.

5.2.1 Time Series Scatter Plot

A scatter plot is available to see a water quality parameters trends over a specified duration of time. The data is queried in a similar manner as discussed previously and a user is able to plot data from multiple Locations and one or two parameters. The user can choose to group or facet by Location or Flags. A Facet creates multiple plots, all with similar mapping

techniques, to allow for easy comparison across plots.

The user is given many options to customize how the plot displays the data. the user can choose a log-scale for the Y-axis as well as to start at zero. These options are not applicable for the x-axis because the X-axis is in date. The user can adjust point size and point color (if the user has not already specified that the group by color). The user can choose from many display themes that are offered in ggplot2. The user can add a horizontal line, vertical line, or floating text anywhere on the plot. The plot automatically creates logical axis labels and a descriptive plot title, which the user can override with custom text or choose to have no labels or title. The user can save a plot with in multiple formats to a specified plot width and plot height. Figures " " are examples of these saved plots. The user can choose to turn on the interactive plot features that are allowed by Plotly which allows the user to hover over a data point for info and also toggle the plot in various ways like zooming in and out.

Temporal Trend lines can be added to the plots to help visualize if there are any temporal trends, if it is not clear when just looking at data points. The trendlines will automatically group in a fashion identical to the points. If the data points on the scatter plots are grouped by Locations as indicated by colors, then the trendlines will also appear grouped by these same colors. Grouping by shape and faceting works similarly.

Three methods for trendlines are available in this application. The first method is a linear trendline which is the linear line that minimizes the residuals of the fit. The second method is the Loess method which is a . The third method is a Generalized Additive method which . The user can choose to show a confidence ribbon with choices of confidence intervals of 0.90, 0.95, and 0.99. If the user selects a confidence ribbon with 90% confidence, a shaded region will appear on the plot where the data is 95% likely to . This should be used with caution because the linear confidence interval does not take into account the seasonal variation, and assumes that all variation is " " .

The user can choose to add a secondary parameter to the plot by introducing a secondary y-axis. This feature allows the user to compare temporal trends of two water quality

parameters. A plot comparig of two parameters, on the x-axis and y-axis can usually display a clearer picture of the relationship between two parameters which will be discussed in the Correlation Section. A benefit of the two parameter temporal plot is to keep more of the temporal information and can allow one to visuallly see a more complex trend like a delayed response trend. (Is there a scientific word for this?)

5.2.2 Correlation Scatter Plot

The application has a scatter plot feature designed for a scatter plot between two water quality variables. Water Quality data of two parameters are paired based on location and day of sampling. The two water quality observations are thus converted into one data point on the plot with the x location determined by the value of the first water quality parameter and the y location determined by the value of the second water quality parameter.

Trendlines can be added to. Similiar to the time-series plots, the methods for the trendlines are linear, Loess, and Generative Additive.

This anaylsis will be extended for water quality parameters to be correlated with metereological data and hydrological data.

5.2.3 Distribution

The distribution of Results of a particular water quality parameter can be visualized in the application. Based on the user selected Locations, Parameter, and Date Range, the user can create a histogram, Density Plot, or box plots.

5.2.4 Heatmap (Interpolated Color Profile Plot)

The

5.2.5 Profile Line Plot

5.2.6 Phytoplankton

5.3 Statistics

5.3.1 Spatial and Temporal Statistics

Temporal Statistics can be computed with minimal effort in the application. Based on the user selected query of Locations, Parameters, and Date Range, the following statistics will be calculated for each parameter: number of samples, average result, minimum result, maximum result, 1st quartile (25 percentile), median, 3rd quartile (75 percentile), variance, standard deviation, geometric mean, and Mann-Kendall statistic. Any blank data (represented by NA in R) is ignored for these statistic calculations. The geometric mean is " ". The Mann-Kendall statistic

Before Statistical calculation, the data can be grouped by Location as well as various types of temporal schemes. The data can be grouped by year, season (independent of year), month (independent of year), season and year, and month and year.

5.3.2 Pearson Correlation Matrix

A pearson correlation matrix can be created in this application. Based on the user selected query information, a correlation matrix is generated to show the correlation of parameters across all of the parameters that the user has selected (the user must select more than one parameter). Positive correlation statistics, R values, are shown in red and negative R values are shown in blue.

Confidence intervals are calculated to determine the significance of the pearson correlation coefficient. This is crucial in case the user was to falsely interpret the correlation matrix as significant, when not.

“Find more out about significance. Add to App”

5.4 Geospatial Data Mapping

Home Tab

Spatial trend analysis is incorporated in many locations in the app but primarily lives on the map plot tab. Geospatial plots allows the user to easily compare a parameter statistic across all sites in a visual of plots on their choice of map. Spatial analysis also exists in any tab when multiple sites are chosen for analysis.

5.5 MetaData

5.6 Data Import

The Watershed data Importer Tool (WIT) facilitates importing raw data from multiple sources into the database. Each Data Type has a formatting function script that is written in R to format the data. As more data sources are added or data sources are changed, these can be uploaded into the database. The user is prompted to select dataset type and the user will be shown a list of raw data files in the appropriate dataset type location on their computer. The user then selects a file from their computer and then press a button to format the data. A Warning message will be sent to the user if there was a problem with the data or if the data already exists in the database. After a successful formatting, the user will be able to see the formatted data in a table on the screen and an import button will appear that the user can press to import the data if they are satisfied with how the data looks. Once the data is imported, the raw data file is moved from the unprocessed folder on their computer to the processed folder.

6 Discussion and Recommendations

Through the insight that this application brings, future data collection needs can be better assessed which will direct changes to the current watershed monitoring program. Automation is a useful, yet underutilized tool that can save an organization much time day to day on decreasing the number of repeated tasks that come along with searching, displaying, and analyzing water quality data. A custom application with dashboard to allow a user to easily perform data science with the power of R, with little knowledge of R programming language is highly desired. This application should include but not limited to facilitated data entry, querying, visualization, and analysis of water quality data and other related watershed data. Save Tiem and dMoney, increase quality control

Visualization dashboards were created to view data in a specific customized fashion based on a user selected inputs. Statistical Analysis of parameters is useful to understand watershed and reservoir water quality in regards to temporal and spatial trends of water quality parameters and also correlations between water quality parameters.

Temporal Analysis and spatial analysis (do more (some) reading)

Questions should be asked the sampling plan to make informed decisions to update the sampling plan on a continuous basis. The reservoir and watershed is always changing, as well should the strategies to maintain adequate water quality. It is likely that as more information is discovered about the reservoir, more questions can be asked. Are there locations that we should be sampling more? Is the sampling frequent enough to see trends? Is the sampling more frequent than necessary to see trends? Are there parameters that should be added to the sampling plan? Are there any parameters that should be removed from the sampling plan? This Application will give insight to data collection needs

A Reservoirs water quality effects the type and extent of necessary treatment process of the water supply prior to distribution. Generally, less treatment is required for a water utility whose source is a remote, healthy reservoir than a reservoir which has been degraded by anthropogenic or other means. It is often infeasible to obtain water from completely

pristine areas, especially for water utilities who supply water to urban areas. Watershed management can help ensure the water quality of a reservoir and as well as predict, lessen, or prevent reservoir water quality degradation. from occurring. (Quantity as well.)

6.1 Pros and Cons of Application

Custombility

Upkeep

6.2 Future Work

6.2.1 Meteorological and Hydrological Data

6.2.2 Forestry

6.2.3 Reports

7 Appendix

7.1 WAVE Developer Manual

7.2 WIT Developer Manual

8 Extras and Trash Bin (Not a real Section)

Shiny is web application framework for R data projects. A user can create an application in the form of a website, an html document, or a dashboard. Shiny offers hosting services which cost some money in which a user can easily launch an application without the knowlege or hassle of hosting their own application as well as recieve customer support. Shiny is designed for people with people who have experience with R but do not have nay application developement experience. No web development skills are required. Shiny is well ducumented including many tutorial documents and videos as well as extra webinars and Github and Stackoverflow help.

Bokeh is a web application framework for Python Projects and is similiar to Shiny. Bokeh seems to offer more in depth interactive ability of plots and other graphics, yet Shiny does offer basic interactive graphic features as well. “Bokeh is a Python interactive visualization library for large datasets that natively uses the latest web technologies. Its goal is to provide elegant, concise construction of novel graphics in the style of Protovis/D3, while delivering high-performance interactivity over large data to thin clients.” - Bokeh. A benefit of Bokeh is that the visualizations can be connected to almost any web tool, widget, or framework, outside of Bokeh itself (Bokeh). Python frameworks like Django, Pyramid, or Flask have greater customobility than Shiny and can be used for things outside of data science.

Spyre is another web application framework for providing a simple user interface for Python data projects. Dash created by Plotly is another alternative to build dashboards using Python which utilizes plotly.js, a leading web chart library, without the use of Javascript.

Pyxley python package makes it easier to deploy Flask-powered dashboards using a collection of common JavaScript charting libraries. UI components are powered by PyxleyJS. Bowtie is also an interactive dashboard toolkit in python which can be used to create web applications for data science. D3 is a javascript library which combines powerful visualization and interaction techniques.

For the reason of simplicity, documentation, and previous familiarity in R, decided that Shiny is the best option.

Aquarius Time-Series and Aquarius Sample have a data import feature. Aquarius has the ability to import data from common more main stream instruments which basically has an internally built code to correctly transform this data into the necessary format for the database. A user can also import their create custom transformation " " to transform a particular dataset. Aquarius has the ability to perform queries and see basic plots of data including scatter plots and box plots.

References

“Shiny - Tutorial.” 2018. Accessed January 7. <https://shiny.rstudio.com/tutorial/>.

“Shiny - Widget Gallery.” 2018. Accessed January 7. <https://shiny.rstudio.com/gallery/widget-gallery.html>.