

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

Article type : Technical Paper

Statistical and Hybrid Methods Implemented in a Web Application for Predicting Reservoir Inflows During Flood Events

Tingting Zhao, Barbara Minsker, Fernando Salas, David Maidment, Vesselin Diev, Jacob Spoelstra, and Prashant Dhingra

Graduate Student (**Zhao**), Department of Civil and Environmental Engineering, University of Illinois at Urbana-Champaign, 4129 Newmark Civil Engineering Laboratory, Urbana, IL 61801; Professor (**Minsker**), Department of Civil and Environmental Engineering, Southern Methodist University, Dallas, TX 75275; Associate Scientist (**Salas**), UCAR, NOAA/NWS National Water Center, Tuscaloosa, AL 35401; Professor (**Maidment**), Center for Research in Water Resources, University of Texas at Austin, Austin, TX 78712; and Senior Data Scientist (**Diev**), Microsoft Azure Data Science Team, Director of Data Science (**Spoelstra**), Azure Machine Learning Platform, and Principle Project Manager (**Dhingra**), Data Science and Machine Learning, Microsoft, Redmond, WA 98052 (E-Mail/Zhao: tzhao6@illinois.edu).

Abstract: Reservoir management is a critical component of flood management, and information on reservoir inflows is particularly essential for reservoir managers to make real-time decisions given that flood conditions change rapidly. This study’s objective is to build real-time data-driven services that enable managers to rapidly estimate reservoir inflows from available data and models. We have tested the services using a case study of the Texas flooding events in the Lower Colorado River Basin in November 2014 and May 2015, which involved a sudden switch from drought to flooding. We have constructed two prediction models: a statistical model for flow prediction and a hybrid statistical and physics-based model that estimates errors in the flow. **This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/1752-1688.12575-16-0042](https://doi.org/10.1111/1752-1688.12575-16-0042)**

1 predictions from a physics-based model. The study demonstrates the statistical flow prediction
2 model can be automated and provides acceptably accurate short-term forecasts. However, for
3 longer-term prediction (2 hours or more), the hybrid model fits the observations more closely
4 than the purely statistical or physics-based prediction models alone. Both the flow and hybrid
5 prediction models have been published as Web services through Microsoft's Azure Machine
6 Learning (AzureML) service and are accessible through a browser-based Web application,
7 enabling ease of use by both technical and non-technical personnel.

8
9 **(Key Terms:** flooding; data-driven model services; AzureML; reservoir inflow.)

10 INTRODUCTION

11 In this paper we demonstrate a new framework for real-time flood management through data-
12 driven services to rapidly estimate reservoir inflows from available data and models. Physics-
13 based models are widely used in reservoir management: For example, the National Weather
14 Service (NWS) river forecast centers use physics-based models for daily forecasts. These models
15 often require extensive manual effort for calibration that can make real-time updates difficult.
16 Data-driven models, such as statistical or machine learning models, use historical data to rapidly
17 learn a functional map between concurrent input and output variables. Large and growing
18 volumes and varieties of data can be retrieved to derive these types of models using data services
19 from sensors, satellites, and other data sources. Data-driven models can be coupled with physics-
20 based models by fitting a data-driven model to the residual error from the physics-based model,
21 thereby reducing any persistent bias in the physics-based model (Singh & Woolhiser, 2002). This
22 paper explores these alternative approaches for real-time flood management and implements the
23 resulting models as real-time services using the AzureML service. More background for each of
24 these components is given below.

25
26 Traditional hydrologic models have evolved from lumped conceptual models to physics-based
27 distributed models where approximations of the partial differential equation or empirical
28 equations are applied (Abbott et al., 1986b). Models of the physical processes employ
29 mathematical functions that simulate hydrologic processes and usually involve complex
30 nonlinear processes with high spatial variability at the basin scale (Singh & Woolhiser, 2011).

1 Data sources for physically-based models can be complex and limited, and calibration can be
2 difficult and time consuming.

3 Data-driven modelling is an alternative approach that allows rapid construction of complex
4 models to estimate outcomes based on past experiences and historical events. Data-driven
5 models analyze relationships between concurrent input and output time series (Solomatine and
6 Ostfeld, 2008) and can be applied either alone (using a purely statistical model) or in conjunction
7 with physics-based models, which use mathematical equations derived from the physical
8 processes (creating a hybrid model). Machine learning methods are a type of statistical approach
9 that can be fit rapidly and automatically to represent highly complex relationships. The popular
10 data-driven machine learning methods used in river systems include artificial neural networks
11 (ANN), fuzzy rule-based systems, and support vector machines (SVM), among others.

12 Hybrid models often use machine learning approaches to fit error models to the residual errors in
13 a physics-based forecasting model. Gragne et al. (2015) implemented a filter updating
14 procedures to update error forecast to improve reservoir inflow forecasts. Gragne et al. (2015)
15 proposed an error model to improve hourly reservoir inflow forecasts over one day ahead.

16

17 Many applications of ANN focus on rainfall-runoff models (e.g., Sharma et al., 2000, Abrahart et
18 al., 2007, de Vos and Rientjes, 2007, Nourani et al., 2009). Rainfall is a common input feature
19 for data-driven models of river systems. Many reservoir inflow prediction studies also rely
20 mainly on ANN and rainfall data. Coulibaly et al. (2000) first used an ANN to forecast daily
21 reservoir inflow and a multi-layer feed-forward neural network (FNN) with an early stopped
22 training approach (STA) to improve prediction accuracy. EI-Shafie et al (2007) used historical
23 reservoir inflow and ANN to predict monthly reservoir inflows. Bae et al. (2007) implemented
24 Adaptive Neuro-Fuzzy Inference System (ANFIS) to predict monthly dam inflow using past
25 observed data and future weather forecasting information. Zhang et al. (2009) implemented
26 multilayer perceptron artificial neural networks (MLP-ANNs) using observed precipitation and
27 forecasted precipitation from QPF to predict daily reservoir inflow. Sharma and Chowdhury
28 (2011) reviewed static and dynamic ensemble methods in probabilistic reservoir system
29 forecasting models to reduce structural errors. Jothiprakash and Magar (2012) predicted daily
30 and hourly intermittent rainfall and reservoir inflow using ANN, an adaptive neuro-fuzzy
31 inference system (ANFIS), and linear genetic programming (GLP). Valipour (2013) compared

1 autoregressive moving average (ARMA) and autoregressive integrated moving average (ARIMA)
2 using increasing number of parameters with static and dynamic artificial neural networks. With
3 historical time series data as input, they demonstrated that static and dynamic autoregressive
4 ANNs perform best in forecasting monthly reservoir inflow. Kumar et al. (2015) developed an
5 ensemble model based on neural networks, wavelet analysis, and bootstrap data sampling to
6 generate a range of forecast instead of point predictions for reservoir inflow.

7
8 These previous studies have focused on non-linear regression models and the predictive
9 performance is good when compared with other statistical models. Although previous studies
10 have focused on predicting reservoir inflow from rainfall and historical reservoir inflow data,
11 they have not incorporated soil moisture as an input feature. Toukourou et al. (2010) showed that
12 rainfall and soil moisture data are the major relevant variables for reservoir inflow.

13
14 As described above, Artificial Neural Network (ANN) is a commonly used data-driven approach
15 in hydrology (see also Bowden et al., 2012, Abrahart et al., 2012, Maier et al., 2010), but the
16 convergence speed is low and training can require significant time that may be a barrier when
17 near-real-time model updating is required (e.g., Jain et al., 1999, Maier & Dandy, 2000). This
18 study uses boosted regression trees (BRT) as function approximators. Our early tests showed that
19 BRT has advantages in faster training and higher accuracy than ANN for this application. Others
20 have recently shown that BRT is effective as an ensemble machine learning approach for
21 hydrology. Erdal and Karakurt (2013) have applied BRT as an ensemble learning method, which
22 performed well in predicting a monthly streamflow forecast. Snelder et al. (2009) have used BRT
23 to map the flow regime class by predicting the likelihood of the class of gauge stations based on
24 watershed characteristics. BRT has the advantages of regression trees (which are based on
25 decision trees and built on a process of recursive partition) and boosting methods (creating
26 ensembles of multiple models that combine fast but weak learners to create a strong learner). The
27 approach combines multiple simple trees into an additive regression model to improve predictive
28 performance (Elith et al. 2008).

29
30 The data-driven models in this study were developed using AzureML Studio Predictive
31 Analytics, a Cloud-hosted user-friendly software toolkit that allows graphical construction of

1 data analysis steps (“workflows”) such as data requests, fitting data-driven models, and data
2 visualization (AzureML team Microsoft, 2015). The data-driven models built in AzureML
3 Studio can be published as Web services on the Azure Cloud, providing scalability and high
4 software availability and reliability, as well as easy integration into modern software systems.

5
6 This study’s purpose is to investigate the feasibility and accuracy of real-time data-driven
7 services to estimate reservoir inflows from available data. The Texas flooding events in the
8 Lower Colorado River Basin in November 2014 and May 2015, which involved a sudden switch
9 from drought to flooding, are used as a case study. The Lower Colorado River Authority
10 (LCRA), which is responsible for reservoir management in this basin, uses the physics-based
11 Hydrologic Engineering Center’s Hydrologic Modeling System (HEC-HMS) in the Corps Water
12 Management System (CWMS). HEC-HMS predicts reservoir inflows from real-time data,
13 including precipitation, reservoir information, and other hydro-meteorological data.

14
15 Currently LCRA uses a HEC-HMS rainfall-runoff model to predict reservoir inflows that does
16 not consider soil moisture as an input dataset. The observed streamflow and soil moisture data
17 are used only to calibrate reservoir inflows manually. Soil moisture may be an important factor
18 for predicting reservoir inflows (Kang et al., 2015) and a data-driven approach would allow
19 LCRA reservoir managers to automatically update the reservoir inflows as these conditions
20 change . In this study, we explore a workflow approach that allows the model set-up process to
21 be completed only once by a technical analyst and then executed by technical or non-technical
22 users through a Web browser. A workflow is a collection of tasks that build an automated
23 pathway for heterogeneous modeling steps.

24
25 The performance of data-driven modeling approaches, including both statistical and hybrid
26 (coupling statistical and physics-based) models is also assessed using boosted regression tree
27 modules from AzureML to predict reservoir inflows from real-time and historical precipitation
28 and soil moisture data. The models can be connected with other data services to obtain the input
29 data. The system is implemented as Web services on AzureML, which do not require any
30 software installation and can be rapidly updated as new data are obtained. The data-driven

1 services allow users and water managers to automatically fit model parameters, compute data-
2 driven models, and retrieve reservoir inflow information through a Web browser.

3 METHODOLOGY

4 Figure 1 shows the general data-driven framework developed in this study to support reservoir
5 management. The framework consists of two main components: 1) algorithms and tools from
6 Azure Predictive Analytics toolkit; and 2) Web application. Azure Predictive Analytics
7 (predictive analytics is a commercial term for machine learning) is a machine learning platform
8 that allows rapid training of statistical models to describe the relationships between inputs
9 (“features”) and outputs (“targets”), with execution on remote servers (in the “Cloud”). This first
10 component comprises data preparation, data preprocessing, and model development. The input
11 datasets, which include feature datasets and target values, are first uploaded into AzureML
12 Studio.

13
14 For this study, a wavelet analysis filter method is applied for data preprocessing to reduce data
15 noise, since noise or errors in the measured datasets may mask important features in the data.
16 Boosted Regression Tree modules in AzureML are then employed to statistically model the
17 reservoir inflows using data-driven models. These model execution steps have been constructed
18 as workflows in AzureML, and flow prediction models and hybrid prediction models have been
19 implemented as modules in a workflow to predict reservoir inflow. AzureML has significant
20 advantages in publishing the constructed workflows as Web services. A Web application, which
21 is Web browser-based software for executing the built models, has been built that enables users
22 to execute the data-driven model using Web services to predict reservoir inflow (named *flowin* in
23 this study).

24
25 Data-driven models use historical data to learn a functional map between input and output
26 variables that can be used to predict future output variables. Given input datasets that include
27 input features and output target values from historical data, a mapping can be built to predict
28 future outputs from known future input features (Mitchell, 1997). For instance, $y=f(x)$ is a
29 mapping (training model) between input variables x and output variable y . Once the future input
30 variables \hat{x} are available, the future outputs \hat{y} can be predicted using the training model. In this
31 study, we develop two types of data-driven models. The first type is a purely data-driven

1 statistical prediction model that is used to directly predict reservoir inflows from soil moisture,
2 precipitation, upstream reservoir outflow, and historical reservoir inflow. The second type of
3 model is a hybrid prediction model, which corrects the results of physics-based models that
4 predict reservoir inflows from weather, runoff, and streamflow predictions. The hybrid
5 prediction model applies the available input features to predict differences between the physics-
6 based model-predicted results and the observed data.

7 **[Insert Fig 1 here]**

8

9 *Data Preprocessing Using Wavelet Analysis*

10 Wavelet analysis is used to filter the reservoir inflow data into trend and noise parts. This step is
11 necessary because there are no direct measures of reservoir inflows. Reservoir inflows are
12 derived from reservoir storage and *flowout* (the flow out from the reservoir), which are subject to
13 fluctuations that may be caused by wave action when winds are high during storms or by
14 measurement errors of the sensors at the gauging stations (Tao, 1998). We use wavelet functions
15 to decompose the original data into high-pass filter (details) and low-pass filter (trend)
16 components (Valens, 1999, Polikar, 2001, Okkan, 2012).

17

18 Maximal Overlap Discrete Wavelet Transform (MODWT) is a linear filtering operation that
19 produces time-dependent wavelets and scaling coefficients (Cornish & Percival, 2005). It
20 performs better than other methods such as discrete wavelet transform (DWT) in fitting all
21 sample sizes since DWT requires sample size to be a multiple of 2^J where J is the decomposition
22 level (Cornish & Percival, 2005). In addition, MODWT is independent of the starting point of
23 the time series, which means that MODWT is not affected by circular shifting of the input time
24 series (Percival and Walden, 2000).

25

26 The wavelet coefficient generated by a high-pass filter is defined as

$$27 \quad W_{j,t} = \sum_{l=0}^{L_j-1} \tilde{h}_{j,l} X_{t-l} \quad (1)$$

28 The wavelet coefficient generated by a low-pass filter is defined as

$$29 \quad V_{j,t} = \sum_{l=0}^{L_j-1} \tilde{g}_{j,l} X_{t-l} \quad (2)$$

1 where j is the level of decomposition, L is the width of the $j=1$ base filter, $\{\tilde{h}_{j,l}\}$ and $\{\tilde{g}_{j,l}\}$ are
2 wavelet and scaling filters respectively.

3

4 The decomposition process is shown in Figure 2. Take the decomposition level = 3 as an
5 example. In each level, the original dataset X is decomposed as trend V and residual error W . In
6 the first level, X is decomposed as V_1 and W_1 . The level 2 decomposition is based on V_1 , which
7 is the trend component from the last level. W_1 is discarded. The decomposition continues until
8 the defined decomposition level is reached. The level of filtering selected for the particular case
9 study (in this case, level 2) is then selected based on best professional judgment of the reservoir
10 operators.

11

[Insert Fig 2 here]

12

13 *Prediction Modeling Using Boosted Regression Tree (BRT)*

14 Data-driven prediction models are computed using a boosted regression tree model, which is an
15 ensemble model that integrates multiple single regression trees. Regression tree models use
16 recursive binary splits to predict the target variable (Elith et al., 2008). Figure 3 demonstrates a
17 simple regression tree example. A tree model is built by splitting the input datasets into subsets
18 based on each selected input feature (such as x_1, x_2, x_3, x_4, x_5). The best partition (e.g., $x_1 < V_1$
19 and $x_1 \geq V_1$) is computed from each derived subset (called recursive partitioning) to maximize
20 improvement in the model prediction. This process continues until no further splitting improves
21 the predictions. Boosting is an adaptive method of combining simple models into a single strong
22 learner to improve model performance. Pseudo code for BRT has been included in the appendix
23 of the paper. Key features are the ability to fit complex nonlinear models and high accuracy
24 (Elith et al., 2008, Caruana & Niculescu-Mizil, 2006).

25

[Insert Fig 3 here]

26

27 *Performance Metrics*

28 We use five performance metrics to evaluate the developed models for predicting current and
29 future reservoir inflows.

30 a. Mean Absolute Error (MAE)

$$31 \quad MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (3)$$

1 where \hat{y}_i is the prediction and y_i is the true value. MAE averages all of the errors in the
2 model. When MAE is closer to zero, the model fits better.

3 b. Root Mean Squared Error (RMSE)

$$4 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (4)$$

5 where \hat{y}_i is the prediction and y_i is the true value. RMSE is a measurement of the average
6 of the squares of the errors. RMSE=0 means a perfect fit of the model.

7 c. Relative Absolute Error (RAE)

$$8 \quad RAE = \frac{\sum_{i=1}^n |\hat{y}_i - y_i|}{\sum_{i=1}^n |y_i - \bar{y}|} \quad (5)$$

9 where \hat{y}_i is the prediction, y_i is the true value and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. RAE measures the
10 percentage of error over the true value. RAE = 0 if there is a perfect fit.

11 d. Relative Squared Error (RSE)

$$12 \quad RSE = \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (6)$$

13 where \hat{y}_i is the prediction, y_i is the true value, and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is the mean true value.

14 e. Coefficient of Determination (R^2)

$$15 \quad R^2 = \left(\frac{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \right)^2 \quad (7)$$

16 where \hat{y}_i is the prediction, y_i is the true value, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$. R^2
17 measures how close the data are to the fitted regression line. An R^2 of 1 indicates a
18 perfect fit of the regression line, and an R^2 of 0 indicates that the line does not fit the data
19 at all.

20
21 *Web Application*

22 AzureML is a Cloud service for machine learning experiments. The workflows are constructed
23 as directed acyclic graphs (DAGs) in a Web-based graphical user interface that enables module
24 operations on datasets (AzureML team Microsoft, 2015). AzureML includes machine learning
25 libraries from open source languages such as R and Python, in addition to libraries of statistical
26 methods and other data processing operations. In addition, Azure ML allows connections to
27 other infrastructure such as database servers to handle large amounts of data.

1

2 Machine learning models can be manipulated as data workflows by joining modules in AzureML
3 Studio as shown in Figure 4. Such data workflows, including data preprocessing, model building,
4 and results visualization, are more natural and intuitive than scripts. Non-technical users can
5 implement and update the data-driven approach without requiring advanced machine learning
6 skills or computing expertise (AzureML team Microsoft, 2015). After the complete workflow is
7 built in AzureML Studio, it can be published as a Web service and shared with other users as a
8 Web application.

9

10 A Web application builds the connection between client and server to enable Cloud-based Web
11 services to execute through a simple Web interface. For instance, a modeling Web application
12 can be built as an automated modeling system (workflow) that includes data access, model
13 execution and output visualization. Such a system can be published as Web services. A custom
14 Web User Interface (UI) is then built to allow non-technical users to access the Web services and
15 view the output directly through the Web browser.

16

17 In AzureML, a python Application Programming Interface (API) is provided to easily access
18 AzureML Web services. A custom UI allows users to download input data and execute the
19 prediction models, with the results made available through the UI. Reservoir managers who are
20 not familiar with machine learning and data-driven approaches and are interested in machine
21 learning approaches can use the Web application to predict reservoir inflow and compare or
22 incorporate results from physics-based models. The Web services provide a rapid approach for
23 reservoir managers to understand near-term impacts of current conditions on reservoir inflow and
24 provides a proof of concept for a real-time Cloud-based system for reservoir management.

25

[Insert Fig 4 here]

26

27

CASE STUDY

28 Lake Travis is in Travis County, located upstream of Lake Austin. Mansfield Dam, operated by
29 LCRA, creates Lake Travis, which serves to contain floodwaters and helps to manage flooding
30 downstream. The floodgate release is operated by LCRA under the direction of the U.S. Army
31 Corps of Engineers. The amount of release depends on weather and flood conditions, such as the

1 water level of the reservoir and downstream flow. Understanding the predicted reservoir inflow
2 during flooding events helps reservoir managers operate the dam more effectively based on such
3 information and their operating experience (Mateo et al., 2014).

4 **[Insert Fig 5 here]**

6 *Datasets*

7 The case study focuses on Texas flooding events in the Lower Colorado River Basin in
8 November 2014 and May 2015, using the input and output data given in Figure 6. Precipitation
9 and soil moisture input data were collected from 31 grid points in Lake Travis Basin in the
10 upstream of Mansfield Dam, as shown in Figure 5.b. The precipitation becomes direct runoff and
11 the soil moisture affects surface runoff by reducing infiltration, which physically affects
12 reservoir inflow. Other input features are the flow out of the upstream reservoir Starcke Dam
13 (*flowout*) and the previous *flowin* to Mansfield Dam, as shown in Figure 5.a.

14 **[Insert Fig 6 here]**

15
16 The precipitation data (in kg/m^2) were downloaded from Phase 2 of the North American Land
17 Data Assimilation System (NLDAS-2). NLDAS-2 forcing data are derived from: (1) Doppler
18 radar data, which are used in national weather forecasts (<http://radar.weather.gov/>), (2) CPC
19 MORPHing (CMORPH) Technique, which produces global precipitation data at a high spatial
20 and temporal resolution

21 http://www.cpc.ncep.noaa.gov/products/janowiak/cmorph_description.html), and (3) HPD
22 (Hourly Precipitation Datasets) data (http://www.srh.noaa.gov/ridge2/RFC_Precip/). The data
23 are in $1/8^{\text{th}}$ degree grid spacing (Rui & Mocko, 2013). The soil moisture data, in units of kg/m^2 ,
24 relied on the Noah land surface model (Noah soil moisture 0-100 cm). Data from both models
25 can be downloaded via Web application by providing spatial coordinates and specific time
26 periods.

27
28 The reservoir hourly data were collected by LCRA from November 1, 2014, 00:00, to December
29 3, 2014, 23:00, and from May 1, 2015, 00:00, to June 4, 2015, 23:00, which were the recent time
30 periods with severe flooding in the Lower Colorado River Basin. These data were retrieved from
31 the LCRA database for this study. The two flooding datasets were concatenated together. From

1 the available datasets, the first 85% (from Nov 1, 2014 00:00 to May 26, 2015 15:00) were
2 considered as the training dataset to train the model. The remaining 15% (from May 26th 2015
3 16:00 to June 4th 2015 23:00) were used for testing to evaluate the model predictions. To ensure
4 that the validation and training datasets were interchangeable, 80% of the training dataset was
5 designated as training and 20% as validation. The purpose of such splits is to keep the model
6 fitting completely separate from the validation so that the model is not overfit to this particular
7 dataset.

9 *Model Implementation*

10 **Wavelet analysis to filter data noise.**

11 Wavelet analysis is intended to smooth the fluctuations in the reservoir inflow data and keep the
12 trend. The decomposition level (Figure 2) is a key element to choose in wavelet analysis.
13 Nourani et al. (2008) estimated the optimum decomposition level for DWT using the following
14 equation:

$$15 \quad L = \text{int}[\log_{10}(N)] \quad (8)$$

16 where L is the decomposition level and N is the number of data values.

17
18 In this study, the number of time series data is 1656. Based on Equation 8, the decomposition
19 level $L = \text{int}[\log(1656)] = 3$. To select the best decomposition level, Figure 7 shows *flowin* after
20 each level. At level 1, the dataset still has significant fluctuations and the noise removal is
21 insufficient. At level 3, the dataset is smooth but the peak flow is significantly truncated. LCRA
22 staff advised that Figure 7.b, with level 2 noise removal, represents the best data filtering: the
23 dataset is smooth and the peak is not excessively truncated. Figure 8 shows the original reservoir
24 inflow versus the filtered reservoir inflow.

25 **[Insert Fig 7 here]**

26 **[Insert Fig 8 here]**

27 **Correlation.**

28 To assess appropriate time lags for inclusion in the model, cross correlation was performed and
29 the results are shown in Figure 9. The figure presents the respective correlations between soil
30 moisture and reservoir inflow, precipitation and reservoir inflow, and *flowout* from the upstream
31 reservoir and the downstream reservoir inflow.

1 [Insert Fig 9 here]

2
3 Figure 9.a shows that the correlation between soil moisture and reservoir inflow reaches the
4 highest point at lag=0, indicating that the soil moisture at time t is correlated most strongly with
5 the reservoir inflow at time t . Figure 9.b demonstrates that the precipitation at time $t-1$ hour
6 affects the reservoir inflow most, as the precipitation in the past hour usually has the largest
7 influence on the reservoir inflow. The *flowout* of the upstream reservoir (Lake Marble Falls at
8 Starcke Dam) at time $t-2$ hours is correlated most strongly with the reservoir inflow, consistent
9 with LCRA's assessment that flow typically requires two hours to travel from the upstream
10 reservoir to the downstream reservoir inflow at Mansfield Dam.

11 12 **A flow prediction model to predict reservoir inflow.**

13 To develop the BRT model, different combinations of feature inputs were tested to identify
14 which combinations of variables are most predictive. The variables used in the best performing
15 models were those that decreased errors the most. Although the cross-correlation results
16 identified the lags corresponding to the strongest correlation, experimentation with different
17 combinations of time lags is still needed to assure the best performance. Seven experiments were
18 conducted:

- 19 1) soil moisture at time t and precipitation at time $t-1$ at all 31 grid points, *flowout* from
20 upstream reservoir at time $t-2$, and reservoir inflow at time $t-1$;
- 21 2) soil moisture at time t and precipitation at time $t-1$ at the grid point that is closest to the
22 reservoir, *flowout* from upstream reservoir at time $t-2$, and reservoir inflow at time $t-1$;
- 23 3) soil moisture at time $t-1$ and precipitation at time $t-1$ at the grid point that is closest to the
24 reservoir, *flowout* from upstream reservoir at time $t-2$, and reservoir inflow at time $t-1$;
- 25 4) soil moisture at time $t-2$ and precipitation at time $t-1$ at the grid point that is closest to the
26 reservoir, *flowout* from upstream reservoir at time $t-2$, and reservoir inflow at time $t-1$;
- 27 5) soil moisture at time $t-3$ and precipitation at time $t-1$ at the grid point that is closest to the
28 reservoir, *flowout* from upstream reservoir at time $t-2$, and reservoir inflow at time $t-1$;
- 29 6) soil moisture at time t , $t-1$, and $t-2$ and precipitation at time $t-1$ at the grid point that is
30 closest to the reservoir, *flowout* from upstream reservoir at time $t-2$, and reservoir inflow
31 at time $t-1$; and

1 7) soil moisture at time t, t-1, and t-2 and precipitation at time t, t-1, and t-2 at the grid point
2 that is closest to the reservoir, *flowout* from upstream reservoir at time t-2, and reservoir
3 inflow at time t-1 and t-2.

4
5 Since early tests indicated that the precipitation and soil moisture at the closest point to the
6 reservoir were more predictive of the reservoir inflow, most experiments were conducted using
7 data from the closest point to the reservoir.

8
9 AzureML facilitates ease of implementation of these alternative models using graphical
10 workflows, shown in Figure 10, for data manipulation, regression models, training models, score
11 models and other machine learning-related modules. The boosted regression tree module in
12 AzureML was used with the following settings: maximum number of leaves per trees = 10,
13 minimum number of samples per leaf node = 10, and learning rate = 0.1. The sweep parameter
14 module in AzureML was used to select the number of trees constructed. Users provided a range
15 of values for the number of trees ([5, 10, 15, 20, 30, 40, 50, 60, 70, 80] in this case) and the
16 module builds training models for each value and selects the best (20 in this case). The criteria to
17 choose the best number of trees was based on the MAE of the validation dataset.

18 Figure 11 shows the structure of one example regression tree in the BRT. The algorithm takes
19 the entire data set as an input and then splits the dataset at the value of one feature variable
20 (“node”) that maximizes the “separation” of the dataset. Separation is measured by the variance
21 reduction, shown in Equation (9), which measures the total reduction in variance of the output
22 variable due to the split of the node (Timofeev, 2004, Breiman, 1984). Selecting the feature
23 variable with the largest variance reduction minimizes the model error at each split (Timofeev,
24 2004). The dataset is then split into two parts based on this value (in Figure 11, the value for the
25 first split is $\text{flowin_lag} \leq 36965$ at the root node). Similar splitting continues, as shown in Figure
26 11, until the stopping criterion (the maximum number of leaves per tree=10 in this case) is met.
27 Finally, a prediction value of the output is obtained at each leaf node of the tree. For instance, if
28 the $\text{flowin_lag} > 36965$ and $\text{smLoc24} > 357.4$, the prediction value is 4469.

29
30
$$I_V(N) = \frac{1}{|S|^2} \sum_{i \in S} \frac{1}{2} (y_i - \hat{c}_0)^2 - \left(\frac{1}{|S_t|^2} \sum_{i \in S_t} \frac{1}{2} (y_i - \hat{c}_1)^2 + \frac{1}{|S_f|^2} \sum_{i \in S_f} \frac{1}{2} (y_i - \hat{c}_2)^2 \right) \quad (9)$$

1 where S is the original dataset, S_t and S_f are the split datasets, and \hat{c}_0 , \hat{c}_1 , \hat{c}_2 represent the
2 estimate of the average of output label in the respective dataset.

3 **[Insert Fig 10 here]**

4 **[Insert Fig 11 here]**

5
6 **A hybrid prediction model to predict residual error between observed reservoir
7 inflow and the predicted inflow from physics-based model.**

8 Figure 12 shows the plot of the residual errors, which were calculated as the filtered observed
9 reservoir inflow minus the predicted reservoir inflow from the HEC-HMS model. HEC-HMS is a
10 lumped parameter watershed model that simulates watershed response to precipitation and
11 predicts flows throughout the watershed, including reservoir inflows (Hydrologic Engineering
12 Center, 2011). Based on the flow information, LCRA staff simulate reservoir operation using the
13 HEC Reservoir System Simulation (HEC-ResSim) in CWMS, assess the impacts of the
14 operations using HEC Flood Impact Analysis (HEC-FIA), and make decisions for reservoir
15 management (e.g., determine reservoir releases to meet reservoir and downstream operational
16 goals). The same input features as the above flow prediction model were applied here. The seven
17 experiments described above were repeated for the hybrid model, with the best-performing
18 experiment selected.

19 **[Insert Fig 12 here]**

20
21 **RESULTS**

22 *Physics-based Model Performance*

23 Figure 13 shows the predicted reservoir inflow from the physics-based model HEC-HMS in
24 CWMS and Table 1 shows the performance metrics for the physics-based model. The results
25 show that the physics-based model fits the general trend of the reservoir inflows but a residual
26 error remains that can be fit with the hybrid model.

27
28 **[Insert Fig 13 here]**

29 **[Insert Table 1 here]**

30
31 *Data-Driven Flow Prediction Model*

1 Table 2 shows the performance of the data-driven flow prediction model for the seven
2 experiments. Experiment #4 (soil moisture at time $t-2$ at the reservoir-located grid point,
3 precipitation at time $t-1$ at the reservoir-located grid point, *flowout* at time $t-2$, and *flowin* at time
4 $t-1$) demonstrates the best performance metrics. We can see that the flow prediction, shown in
5 Figure 14, is close to the real reservoir inflow, with the prediction capturing both the general
6 trend of the reservoir inflow and closely matching the peak values.

7
8 A comparison of experiment #1 and experiment #2 shows that the closest soil moisture estimate
9 (experiment #2) is more effective than all 31 available estimates in the area (experiment #1),
10 indicating that some input features are not improving predictions of reservoir inflow.

11 Experiments #2 through #5 demonstrate that a time lag of 2 hours for soil moisture input
12 (experiment #4) is the best option, despite the correlation results showing a time lag of zero
13 having maximum correlation. Experiment #7 has similar performance to that of experiment #6,
14 possibly because the additional input variables in experiment #7 (precipitation at time $t-2$ and
15 reservoir inflow at time $t-2$) provide trivial information to improve the prediction performance.

16
17 We also conduct experiments to predict reservoir inflow 1 to 9 hours ahead using the same input
18 variables in Tables 2 and 3. Figure 15 shows the RMSE of future predictions from the data-
19 driven flow prediction model. After 1 hour, the RMSE increases sharply, then fluctuates,
20 indicating that while the flow prediction model can be used to predict reservoir inflow one hour
21 ahead, later performance drops off significantly.

22 **[Insert Table 2 here]**

23 **[Insert Fig 14 here]**

24 **[Insert Fig 15 here]**

25 26 *Hybrid Prediction Model*

27 The hybrid prediction model is used to predict the residual error [residual(t)] between observed
28 *flowin* and predicted reservoir inflow from the physics-based model, shown in Figure 12. The
29 predicted reservoir inflow is then calculated using the predicted residual error plus the predicted
30 reservoir inflow from the physics-based model. Table 3 summarizes the performance of the
31 hybrid model for each of the seven experiments, using the same input variables as for the flow

1 prediction model. The best performance comes from experiment #2, followed by that of
2 experiment #6. Since the hybrid model is intended to rapidly enhance the physics-based model's
3 performance, it makes sense that the model including soil moisture has the best performance
4 since CWMS does not consider soil moisture as an input (Hydrologic Engineering Center, 2011).
5 Figure 16 shows the performance of the physics-based model, the hybrid prediction model, and
6 the observed *flowin* for Experiment #2. The hybrid prediction model improves upon the
7 performance of the physics-based model in terms of the peak value prediction, but does not
8 perform as well as the data-driven model in the short term (Figure 14).

9
10 Figure 17 shows the future prediction performance of the hybrid model. Within four hours, the
11 RMSE curve fluctuates under $170 \text{ m}^3/\text{s}$. However, after four hours, the model's performance
12 begins to drop off.

13 **[Insert Table 3 here]**

14 **[Insert Fig 16 here]**

15 **[Insert Fig 17 here]**

16 *Web Interface*

17 In AzureML, the built workflows were published as Web services using "Set Up Web Service"
18 function. The Uniform Resource Locator (URL) and Application Programming Interface (API)
19 Web Service keys were generated. The resulting data-driven services allow users and water
20 managers to automatically fit model parameters, compute data-driven models, and retrieve
21 reservoir inflow information through a Web browser. A Web application was built that enables
22 users to give input parameters and retrieve output (Figure 18). Figure 18.a shows the user
23 interface. The models can be executed in AzureML by filling the input parameter boxes and
24 selecting the "Compute" button; the result (the value of the predicted reservoir inflow) is
25 retrieved and shown in the Web interface. The input parameters include the "StartTime" and
26 "EndTime," which will automatically download precipitation and soil moisture from NLDAS2,
27 as well as *flowout* (which is the flow exiting the upstream reservoir) and *flowin_lag* (which is the
28 reservoir inflow in the previous time step).

1 Figure 18.b shows a prototype Web application that allows users to see the reservoir inflow
2 prediction based on the prediction models. Users can provide a prediction starting time and a
3 future prediction steps to examine how predictions compare with the measured data in the recent
4 past, which will provide a sense for potential errors in the predicted reservoir inflows. In the
5 future, when predicted soil moisture, precipitation and upstream reservoir *flowout* are available,
6 such data can be incorporated into the prediction model to improve performance.

7
8 Furthermore, the Web application can be extended to other river basins. For instance, for any
9 ungauged basin, users only need to upload the longitude and latitude of grid points affecting
10 reservoir inflow to Azure ML. These points are then used to automatically download
11 corresponding precipitation and soil moisture data from NLDAS2 using our workflow in
12 AzureML. Then users can predict reservoir inflows based on the start time, end time, *flowin_lag*,
13 and *flow_out*, as shown in Figure 18.a. Using this interface, the Web application provides an
14 easy way for reservoir operators to forecast reservoir inflows and explore multiple scenarios
15 without modeling or computational expertise.

16 **[Insert Fig 18 here]**

17 18 DISCUSSION AND CONCLUSIONS

19 In this study, we propose a data-driven framework for real-time reservoir inflow prediction using
20 a service-oriented approach that enables ease of access through a Web browser. Statistical and
21 hybrid models are developed to predict flow and residual errors from a physics-based model,
22 respectively. We created a workflow in Microsoft AzureML, a machine learning studio, for end-
23 to-end downloading of the data, executing the models, and visualizing the results. Azure ML
24 provides fast and easy implementation of the whole workflow as well as publishing of the
25 workflow as Web services. In addition, the input datasets and workflow can be updated when
26 new data are available. One of the workflows that predicts reservoir inflow has been published
27 at <https://gallery.cortanaintelligence.com/Experiment/Predict-Reservoir-Inflow-1>. Users who
28 wants to use AzureML to predict reservoir inflow can update the input data and the model will
29 be automatically updated without manual calibration or tuning of model parameters.

1 The framework was implemented and tested in the Lower Colorado River Basin. The results
2 show that the statistical flow prediction model is more accurate for short-term forecasts than the
3 hybrid prediction model, while the hybrid model performs better for longer-term prediction (2
4 hours or more), as it considers forecasts from a physics-based model.

5
6 The flow prediction model has a peak prediction value close to the actual value. Of the set of
7 experiments shown in Table 2, experiment #4 has the best performance. Using soil moisture at
8 time $t-2$ at the reservoir-located grid point, precipitation at time $t-1$ at the reservoir-located grid
9 point, $flowout$ at time $t-2$, and $flowin$ at time $t-1$ will lead to the best prediction of $flowin$ at time t .

10
11 From a physical process perspective, soil moisture affects surface runoff by reducing infiltration.
12 When flooding happens, infiltration has reached a saturated level. Therefore, high soil moisture
13 conditions are indicative of wet conditions that are well correlated with high reservoir inflows
14 and are thus useful for prediction.

15
16 The hybrid prediction model improves upon the performance of the physics-based model. Based
17 on the set of experiments shown in Table 3, experiment 2 gives the best performance. Using soil
18 moisture at time t at the reservoir-located grid point, precipitation at time $t-1$ at the reservoir-
19 located grid point, $flowout$ at time $t-2$, and $flowin$ at time $t-1$ will lead to the best prediction in
20 $flowin$ at time t . The hybrid model's short-term performance is worse than that of the $flowin$
21 prediction model. The hybrid model is affected by complex processes, as shown by the high
22 fluctuations in Figure 12, and available data to build the model are limited to just two flooding
23 events. With more flooding events available in the future, the incorporation of more data will
24 likely improve the model's performance.

25
26 In considering longer-term predictions, the hybrid prediction model is better than the data-driven
27 flow prediction model in terms of RMSE (Figure 19). The flow prediction model's RMSE is
28 lower than that of the hybrid prediction model one hour ahead. Later, the flow prediction
29 model's RMSE is higher than that of the hybrid prediction model, indicating that the flow
30 prediction model's performance declines after two hours. Because the hybrid prediction model's
31 performance remains reasonably high within the following five hours, in the future the Web

1 application could allow the user to create a combined prediction model that uses the data-driven
2 model for the first two hours and the hybrid prediction model for time steps further in the future.

3
4 Further research is needed to explore how these findings generalize to other locations and storms.
5 The models and tools developed in this work can be generalized to other reservoirs by updating
6 the input data in the workflow. The workflow can also be combined with other modeling services
7 requesting the Web service using URL and API keys, as mentioned previously.

8 **[Insert Fig 19 here]**

9
10 In the future, the hybrid prediction model for long-term prediction will need to be improved.
11 Currently the only available CWMS forecasts from the LCRA database were *nowcasts* (forecasts
12 for the current time period only). If longer-term CWMS predictions could be obtained, then the
13 hybrid model might perform better for longer-term forecasts.

14
15 In addition, the current Web application is a prototype and further user-centered design and
16 development is necessary before the system should be adopted for operational reservoir
17 management. Feedback from LCRA's testing and evaluation of the Web application can be used
18 to improve the interface and add more features as needed to support effective decision making.
19 Moreover, when more flooding data are available, the data-driven and hybrid models can readily
20 be updated and improved using the AzureML framework. Replacing historical data for soil
21 moisture, precipitation, and upstream reservoir *flowout* with model predictions might improve
22 reservoir inflow prediction in later time periods. For instance, the precipitation might be replaced
23 by the Quantitate Precipitation Forecast (QPF) or local LCRA rain gauge data. In the future,
24 other data preprocessing approaches such as partial information approach [Sharma & Mehroma,
25 2014, Sharama, et al., 2016] could also be implemented to automatically choose the best input
26 parameters for data-driven models to improve reservoir inflow forecast.

27
28 The findings clearly indicate promise for this type of approach and potential value in making
29 datasets and model forecasts more readily available in real time to support such analyses. In
30 addition to reservoir inflow forecasting, the framework can be extended to other water resources
31 applications with rich data sets using the AzureML framework.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

APPENDIX

Algorithm (Friedman 2001, Hastie & Friedman, 2008)

Input: training dataset $\{(x_i, y_i)\}_{i=1}^n$ where x_i represents input datasets ('features') and y_i represents output dataset ('targets'), number of iterations.

Algorithm:

1. Initialize model with a constant value

$$F_0(x) = \operatorname{argmin} \sum_{i=1}^n L(y_i, F(x_i))$$

2. For each iteration:

- a. Compute pseudo-residuals:

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right] \text{ for } i = 1, \dots, n$$

- b. Fit a decision tree learner $f_t(x)$ to pseudo-residuals using the training dataset.
- c. Add $f_t(x)$ to the model $F_t(x) = F_{t-1}(x) + \epsilon f_t(x)$, where ϵ is called step-size or shrinkage. In this study, it was set to 0.1 to prevent overfitting by not doing a full optimization in each step.

3. Output $F_t(x)$

ACKNOWLEDGMENTS

The authors are grateful to the 2015 National Flooding Interoperability Experiment (NFIE) Summer Institute for providing an opportunity to conduct this study. The authors also acknowledge Microsoft Research for funding support, the Microsoft Azure data science team for technical support, and LCRA staff David Walker and Chris Riley for providing data and expertise on the case study.

LITERATURE CITED

Abbott, M.B., J.C. Bathurst, J.A. Cunge, and P.E. O'connell, 1986. An Introduction to the European Hydrological System - Systeme Hydrologique Europeen, "SHE," 1: History and Philosophy of a Physically-Based, Distributed Modelling. Journal of Hydrology. DOI:10.1016/0022-1694(86)90115-0.

- 1 Abrahart, R.J. and L.M. See, 2007. Neural Network Modelling of Non-Linear Hydrological
2 Relationships. *Hydrology and Earth System Sciences* 11(5):1563-1579.
- 3 Abrahart, R.J., F. Anctil, P. Coulibaly, C.W. Dawson, N.J. Mount, L.M. See, A.Y. Shamseldin,
4 D.P. Solomatine, E. Toth, and R.L. Wilby, 2012. Two Decades of Anarchy? Emerging
5 Themes and Outstanding Challenges for Neural Network River Forecasting. *Progress in*
6 *Physical Geography* 36:480-513.
- 7 Bae, D.H., D.M. Jeong, and G. Kim, 2007. Monthly Dam Inflow Forecasts Using Weather
8 Forecasting Information and Neuro-Fuzzy Technique. *Hydrological Sciences Journal* 52(1):
9 99-113. DOI:10.1623/hysj.52.1.99.
- 10 Bowden, G.J. and H.R. Maier, 2012. Real-Time Deployment of Artificial Neural Network
11 Forecasting Models: Understanding the Range of Applicability. *Water Resources Research*
12 48(10). DOI:10.1029/2012WR011984.
- 13 Breiman, L., J. Friedman, C.J. Stone, and R.A. Olshen, 1984. Classification and Regression
14 Trees.
- 15 Caruana, R. and A. Niculescu-Mizil, 2006. An Empirical Comparison of Supervised Learning
16 Algorithms. In *Proceedings of the 23rd international conference on Machine learning, ACM,*
17 *New York, USA*, pp. 161-168
- 18 Cornish, C.R., C.S. Bretherton, and D.B. Percival, 2006. Maximal Overlap Wavelet Statistical
19 Analysis with Application to Atmospheric Turbulence. *Boundary-Layer Meteorology*
20 119:339-374.
- 21 Coulibaly, P., F. Anctil, and B. Bobee, 2000. Daily Reservoir Inflow Forecasting Using Artificial
22 Neural Networks with Stopped Training Approach. *Journal of Hydrology* 230:244-257.
- 23 De Vos, N.J. and T.H.M. Rientjes, 2007. Multi-Objective Performance Comparison of an
24 Artificial Neural Network and a Conceptual Rainfall-Runoff Model. *Hydrological Sciences*
25 *Journal* 52:397- 413.
- 26 El-Shafie, A., M.R. Taha, and A. Noureldin, 2006. A Neuro-Fuzzy Model for Inflow Forecasting
27 of the Nile River at Aswan High Dam. *Water Resources Management* 21:533-556.

- 1 Elith, J., J.R. Leathwick, and T. Hastie, 2008. A Working Guide to Boosted Regression Trees.
2 Journal of Animal Ecology 77:802-813.
- 3 Erdal, H.I. and O. Karakurt, 2013. Advancing Monthly Streamflow Prediction Accuracy of
4 CART Models Using Ensemble Learning Paradigms. Journal of Hydrology 477:119-128.
- 5 Friedman, J.H., 2001. Greedy Function Approximation: a Gradient Boosting Machine. Annals of
6 Statistics. doi:10.2307/2699986.
- 7 Gragne, A.S., K. Alfredsen, A. Sharma, and R. Mehrotra, 2015. Recursively Updating the Error
8 Forecasting Scheme of a Complementary Modelling Framework for Improved Reservoir
9 Inflow Forecasts. Journal of Hydrology. doi:10.1016/j.jhydrol.2015.05.039.
- 10 Gragne, A.S., A. Sharma, R. Mehrotra, and K. Alfredsen, 2015. Improving Real-Time Inflow
11 Forecasting Into Hydropower Reservoirs Through a Complementary Modelling Framework.
12 Hydrology and Earth System Sciences 19:3695–3714.
- 13 Hastie, T., R. Tibshirani, and J. Friedman, 2008. Boosting and Additive Trees. The Elements of
14 Statistical Learning, Springer Series in Statistics. Springer New York, New York, NY, pp. 1-
15 51.
- 16 Hydrologic Engineering Center, 2011. Accelerated Corps Water Management System (CWMS)
17 Deployment Campaign. [http://www.hec.usace.army.mil/publications/ProjectReports/PR-](http://www.hec.usace.army.mil/publications/ProjectReports/PR-79.pdf)
18 [79.pdf](http://www.hec.usace.army.mil/publications/ProjectReports/PR-79.pdf), accessed December 2015
- 19 Jain, S.K., A. Das, and D.K. Srivastava, 1999. Application of ANN for Reservoir Inflow
20 Prediction and Operation. Journal of Water Resources. DOI:10.1061/(ASCE)0733-
21 9496(1999)125:5(263).
- 22 Jothiprakash, V. and R.B. Magar, 2012. Multi-Time-Step Ahead Daily and Hourly Intermittent
23 Reservoir Inflow Prediction by Artificial Intelligent Techniques Using Lumped and
24 Distributed Data. Journal of Hydrology 450:293-307.
- 25 Kang, B., Y.H. Ku, and Y. Do Kim, 2015. A Case Study for ANN-Based Rainfall–Runoff Model
26 Considering Antecedent Soil Moisture Conditions in Imha Dam Watershed, Korea.
27 Environmental Earth Sciences:1-12.

- 1 Kumar, S., M.K. Tiwari, C. Chatterjee, and A. Mishra, 2015. Reservoir Inflow Forecasting Using
2 Ensemble Models Based on Neural Networks, Wavelet Analysis and Bootstrap Method.
3 Water Resources Management 29:4863-4883.
- 4 Maier, H.R. and G.C. Dandy, 2000. Neural Networks for the Prediction and Forecasting of
5 Water Resources Variables: a Review of Modelling Issues and Applications. Environmental
6 Modelling and Software 15:101-124.
- 7 Maier, H.R., A. Jain, G.C. Dandy, and K.P. Sudheer, 2010. Methods Used for the Development
8 of Neural Networks for the Prediction of Water Resource Variables in River Systems:
9 Current Status and Future Directions. Environmental Modelling and Software 25:891-909.
- 10 Mateo, C.M., N. Hanasaki, D. Komori, K. Tanaka, M. Kiguchi, A. Champathong, T.
11 Sukhappunnaphan, D. Yamazaki, and T. Oki, 2014. Assessing the Impacts of Reservoir
12 Operation to Floodplain Inundation by Combining Hydrological, Reservoir Management,
13 and Hydrodynamic Models. Water Resources Research 50:7245-7266.
- 14 Microsoft, A.M.T., 2015. AzureML: Anatomy of a Machine Learning Service.
- 15 Mitchell, T.M., 1997. Machine Learning. 1997. Burr Ridge.
- 16 Nourani, V., M. Komasi, and A. Mano, 2009. A Multivariate ANN-Wavelet Approach for
17 Rainfall-Runoff Modeling. Water Resources Management. DOI:10.1007/s11269-009-9414-
18 5.
- 19 Okkan, U., 2012. Wavelet Neural Network Model for Reservoir Inflow Prediction. Scientia
20 Iranica 19:1445-1455.
- 21 Percival, D.B. and A.T. Walden, 1993. Spectral Analysis for Physical Applications. Cambridge
22 University Press.
- 23 Polikar, R., 2001. The Wavelet Tutorial Part I. Fundamental Concepts and an Overview of the
24 Wavelet Theory.
- 25 Roncak, P., J. Bezak, and Z. Bajtek, 2012. Comparison of a Physically-Based and Data-Driven
26 Model in the Bela River Basin. International Multidisciplinary Scientific GeoConference:
27 SGEM: Surveying Geology & mining Ecology Management 3 (2012): 595.

- 1 Rui, H. and D. Mocko, 2013. Readme Document for North America Land Data Assimilation
2 System Phase 2 (NLDAS-2) Products. Greenbelt.
- 3 Sharma, A. and S. Chowdhury, 2011. Coping with Model Structural Uncertainty in Medium-
4 Term Hydro-Climatic Forecasting. Hydrology Research. doi:10.2166/nh.2011.104.
- 5 Sharma, A. and R. Mehrotra, 2014. An Information Theoretic Alternative to Model a Natural
6 System Using Observational Information Alone. Water Resources Research.
7 doi:10.1002/2013WR013845.
- 8 Sharma, A., R. Mehrotra, J. Li, and S. Jha, 2016. A Programming Tool for Nonparametric
9 System Prediction Using Partial Informational Correlation and Partial Weights.
10 Environmental Modelling and Software 83:271–275.
- 11 Sharma, A., K.C. Luk, I. Cordery, and U. Lall, 2000. Seasonal to Interannual Rainfall
12 Probabilistic Forecasts for Improved Water Supply Management: Part 2 - Predictor
13 Identification of Quarterly Rainfall Using Ocean-Atmosphere Information. Journal of
14 Hydrology. doi:10.1016/S0022-1694(00)00348-6.
- 15 Singh, V.P. and D.A. Woolhiser, 2002. Mathematical Modeling of Watershed Hydrology.
16 Journal of hydrologic engineering, 7(4), 270-292
- 17 Snelder, T.H., N. Lamouroux, J.R. Leathwick, and H. Pella, 2009. Predictive Mapping of the
18 Natural Flow Regimes of France. Journal of Hydrology, 373(1), 57-67.
- 19 Solomatine, D.P. and A. Ostfeld, 2008. Data-Driven Modelling: Some Past Experiences and
20 New Approaches. Journal of Hydroinformatics 10:3-20.
- 21 Tao, T., 1999. Local Inflow Calculator for Reservoirs. Canadian Water Resources Journal.
22 doi:10.4296/cwrj2401053.
- 23 Timofeev, R., 2004. Classification and Regression Trees (CART) Theory and Applications.
- 24 Toukourou, M., A. Johannet, G. Dreyfus, and P.-A. Ayrat, 2010. Rainfall-Runoff Modeling of
25 Flash Floods in the Absence of Rainfall Forecasts: the Case of “Cévenol Flash Floods.”
26 Applied Intelligence 35:178-189.
- 27 Valens, C., 1999. A Really Friendly Guide to Wavelets. DOWNLOAD DA INTERNET:
28 <https://www.cs.unm.edu/~williams/cs530/arfgtw.pdf>, accessed December 2015

1 Valipour, M., M.E. Banihabib, and S.M.R. Behbahani, 2013. Comparison of the ARMA,
 2 ARIMA, and the Autoregressive Artificial Neural Network Models in Forecasting the
 3 Monthly Inflow of Dez Dam Reservoir. *Journal of Hydrology* 476:433-441.

4 Zhang, J., C. Cheng, S. Liao, and X. Wu, 2009. Daily Reservoir Inflow Forecasting Combining
 5 QPF Into ANNs Model. *Hydrology and Earth System Sciences*: 121-150

6
7
8 TABLES

9 **TABLE 1.** Performance Metrics of Physics-based Model

Mean Absolute Error (m ³ /s)	Root Mean Squared Error (m ³ /s)	Relative Absolute Error	Relative Squared Error	Coefficient of Determination
67.677	130.541	0.381	0.282	0.718

10
11 **TABLE 2.** Performance Metrics for Flow Predicted Model

	Input Variables	Output Variable	Performance				
			Mean Absolute Error(m ³ /s)	Root Mean Squared Error(m ³ /s)	Relative Absolute Error	Relative Squared Error	Coefficient of Determination
1	SM(t) ₃₁ , Precip(t-1) ₃₁ , flowout(t-2), flowin_lag(t-1)	flowin(t)	28.883	60.032	0.167	0.061	0.939
2	SM(t) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	flowin(t)	28.600	66.545	0.165	0.075	0.925
3	SM(t-1) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	flowin(t)	26.873	48.705	0.155	0.040	0.960
4	SM(t-2) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	flowin(t)	23.984	46.723	0.139	0.037	0.963
5	SM(t-3) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	flowin(t)	27.269	50.404	0.157	0.043	0.957
6	SM(t) _{closest} , SM(t-1) _{closest} , SM(t-2) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	flowin(t)	24.607	50.121	0.142	0.042	0.958
7	SM(t) _{closest} , SM(t-1) _{closest} , SM(t-2) _{closest} , Precip(t) _{closest} , Precip(t-1) _{closest} , Precip(t-2) _{closest} ,	flowin(t)	27.666	52.953	0.160	0.048	0.953

	flowout(t-2), flowin_lag(t-1), flowin_lag(t-2)						
--	---	--	--	--	--	--	--

1

2

3 **TABLE 3.** Performance Metrics for Hybrid Prediction Model

	Input Variables	Output Variable	Performance				
			Mean Absolute Error (m ³ /s)	Root Mean Squared Error (m ³ /s)	Relative Absolute Error	Relative Squared Error	Coefficient of Determination
1	SM(t) ₃₁ , Precip(t-1) ₃₁ , flowout(t-2), flowin_lag(t-1)	residual(t)	81.836	167.636	0.472	0.475	0.525
2	SM(t) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	residual(t)	57.200	97.976	0.331	0.163	0.838
3	SM(t-1) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	residual(t)	71.925	121.196	0.415	0.249	0.751
4	SM(t-2) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	residual(t)	80.986	146.398	0.467	0.362	0.638
5	SM(t-3) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	residual(t)	69.943	138.186	0.404	0.322	0.678
6	SM(t) _{closest} , SM(t-1) _{closest} , SM(t-2) _{closest} , Precip(t-1) _{closest} , flowout(t-2), flowin_lag(t-1)	residual(t)	69.659	108.170	0.403	0.197	0.803
7	SM(t) _{closest} , SM(t-1) _{closest} , SM(t-2) _{closest} , Precip(t) _{closest} , Precip(t-1) _{closest} , Precip(t-2) _{closest} , flowout(t-2), flowin_lag(t-1), flowin_lag(t-2)	residual(t)	68.527	112.984	0.396	0.216	0.784

4

5

6

LIST OF FIGURES

7

FIGURE 1. Framework of the Data-Driven Services

8

FIGURE 2. Decomposition Based on Wavelet Analysis

9

FIGURE 3. Example of a Regression Tree

- 1 FIGURE 4. An Example of AzureML Graphical Workflow
- 2 FIGURE 5. Case Study Location and Data Points
- 3 FIGURE 6. Inputs to and Outputs from the Data-driven Models
- 4 FIGURE 7. *Reservoir Inflow* Graph after Each Decomposition Level
- 5 FIGURE 8. **Original** Observed *Flowin* vs Filtered Observed *Flowin* During 2014-2015 Flooding
- 6 Events
- 7 FIGURE 9. Correlation Plot between Input Features and Output Label
- 8 FIGURE 10. AzureML Workflow for Data-Driven Flow Prediction Model
- 9 FIGURE 11. Structure of an Example Regression Tree in BRT
- 10 FIGURE 12. Residual errors between Filtered Observed *Flowin* and *Flowin* from Physics-based
- 11 Models during 2014-2015 Flooding Events
- 12 FIGURE 13. **Filtered** Observed *Flowin* vs Physics-based *Flowin* during 2014-2015 Flooding
- 13 Events
- 14 FIGURE 14. **Filtered** Observed *Flowin* vs Predicted *Flowin* from May 26 2015 to Jun 5 2015
- 15 FIGURE 15. Flow Prediction Model Performance for Future Prediction
- 16 FIGURE 16. **Filtered** Observed *Flowin* vs Predicted *Flowin* from Hybrid Prediction Model vs
- 17 Physics-based Model *Flowin* from May 26 2015 to June 5 2015
- 18 FIGURE 17. Hybrid Predicted Model Performance for Future Prediction
- 19 FIGURE 18. Screenshot of Web Interface for the Data-driven Model Services
- 20 FIGURE 19. Prediction Performance of Flow Prediction Model and Hybrid Prediction Model

Input Feature Data

- Soil moisture and precipitation from NLDAS
- Other data from LCRA database
- Predicted reservoir inflow from physics-based model

Target Values

- Observed reservoir inflow
- Residual between observed reservoir inflow and predicted reservoir inflow from physics-based model

Data Preprocessing

Flow Prediction Model

Reservoir inflow as the dependent variable

Residual Prediction Model

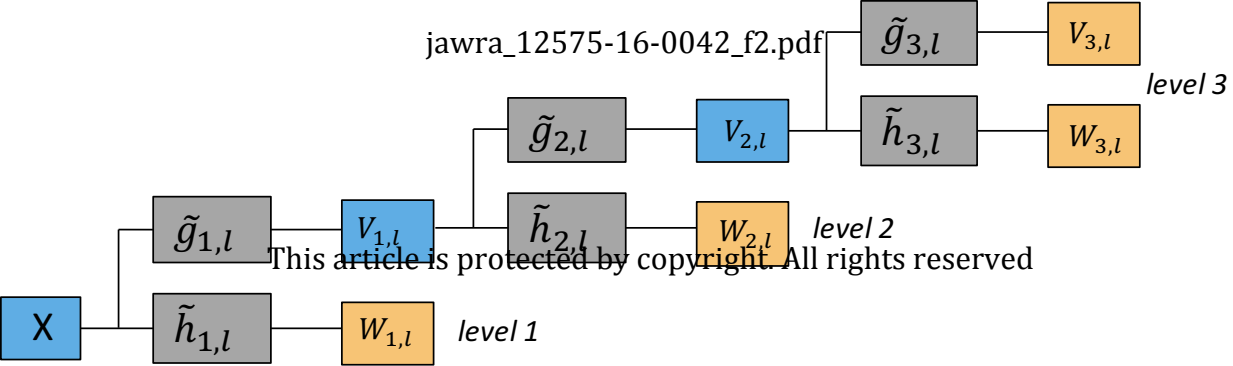
Residual between observed reservoir inflow and predicted reservoir inflow from physics-based model as the dependent variable

Web Application

This article is protected by copyright. All rights reserved.

prediction

Decision Support for Water Management



Root Node

jawra_12575-16-0042_f3.pdf

$x_1 < V_1$

Node 1

$x_1 \geq V_1$

Node 2

$x_2 < V_2$

Node 3

$x_2 \geq V_2$

Node 4

$x_3 < V_3$

Node 5

$x_3 \geq V_3$

Node 6

This article is protected by copyright. All rights reserved

$x_4 < V_4$

Node 7

$x_4 \geq V_4$

Node 8

$x_5 < V_5$

Node 9

$x_5 \geq V_5$

Node 10

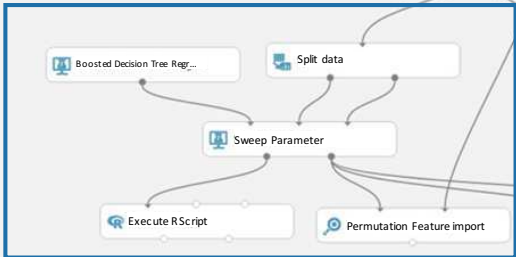
Temporary Dataset
jawra_12575-16-0042_f4.pdf

input data
data clean

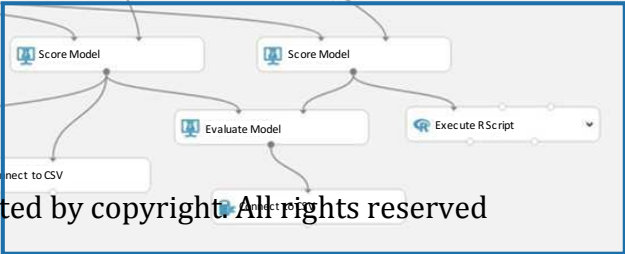
Clean Missing Data

Split Data

split data into training and testing



choose model
structure and build
training model



scoring and evaluation

This article is protected by copyright. All rights reserved

Inputs

- Soil moisture from NLDAS
- Precipitation from NLDAS
- Flow out of the upstream reservoir
- Previous flow in to Mansfield Dam

jawra_12575-16-0042_f6.pdf

Outputs

- Reservoir inflow
- Residual between observed reservoir inflow and predicted reservoir inflow from physics-based model

This article is protected by copyright. All rights reserved.

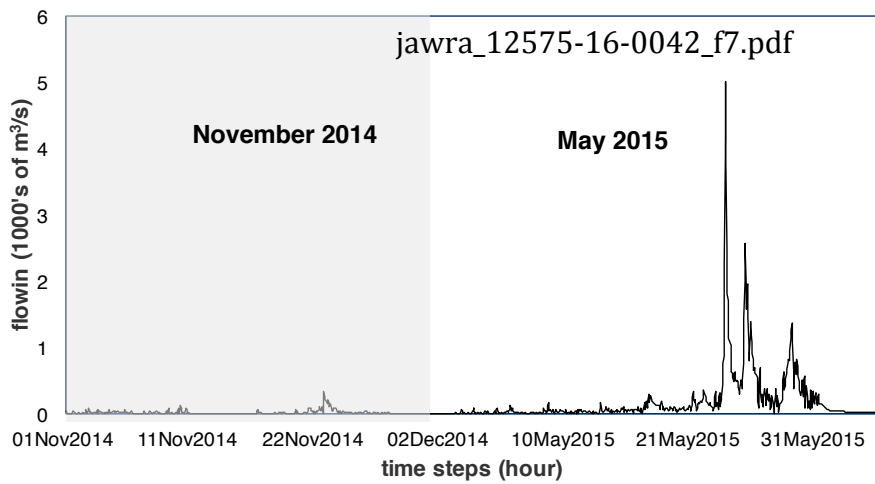


Figure 7.a. Reservoir Inflow at Level 1

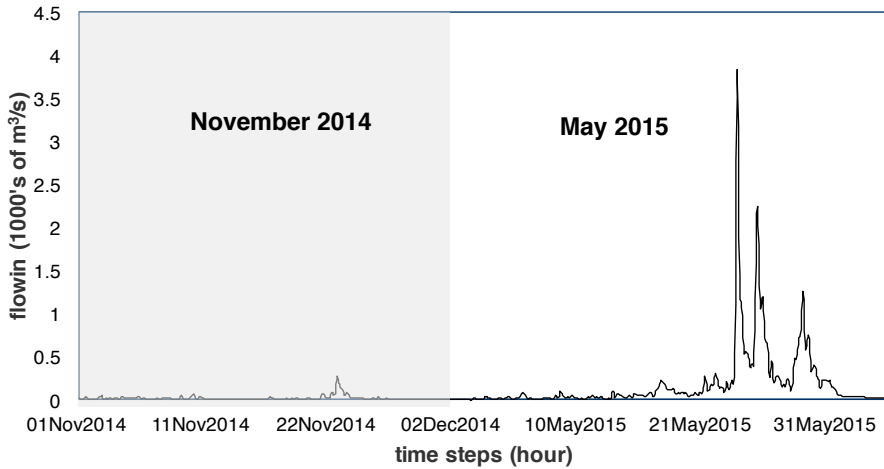


Figure 7.b. Reservoir Inflow at Level 2

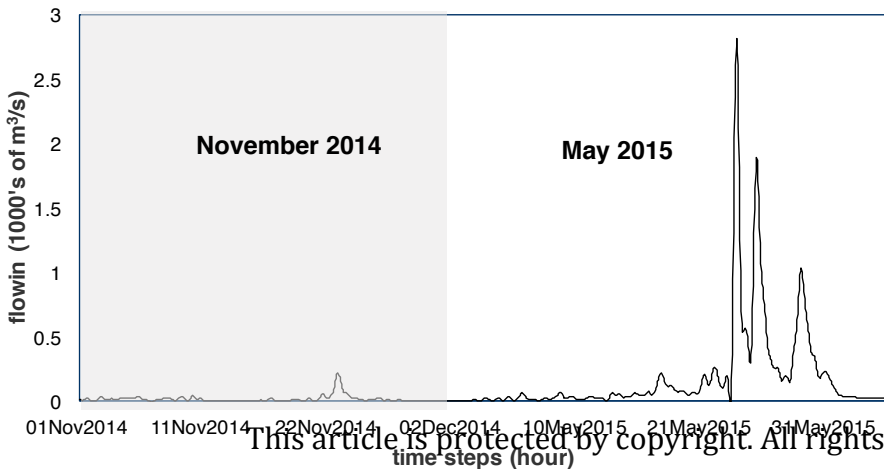
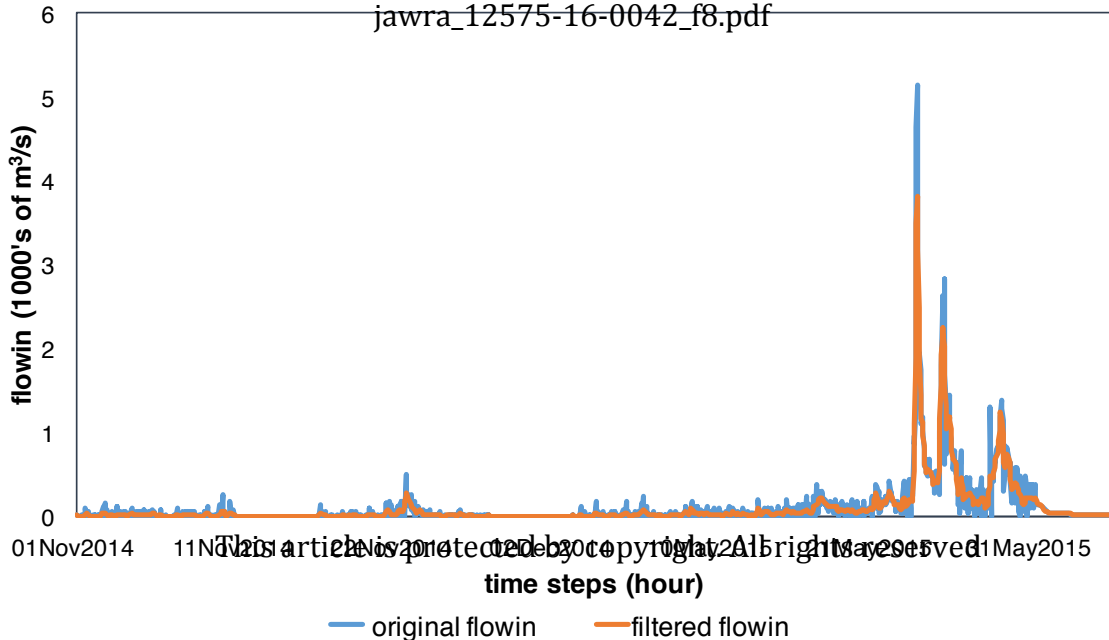


Figure 7.c. Reservoir Inflow at Level 3



This article is protected by copyright. All rights reserved.

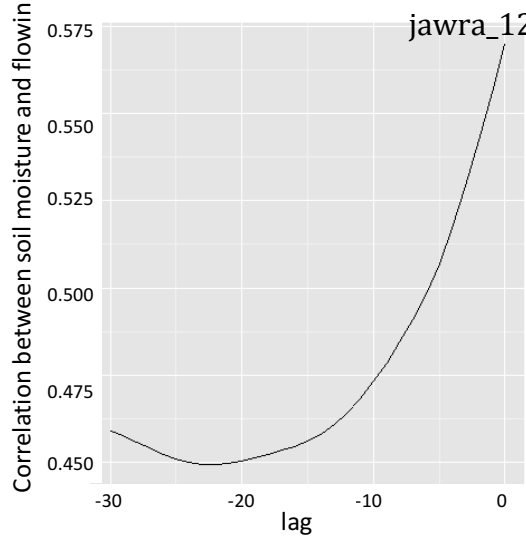


FIGURE 9.a. Correlation between Soil Moisture and *Flowin*

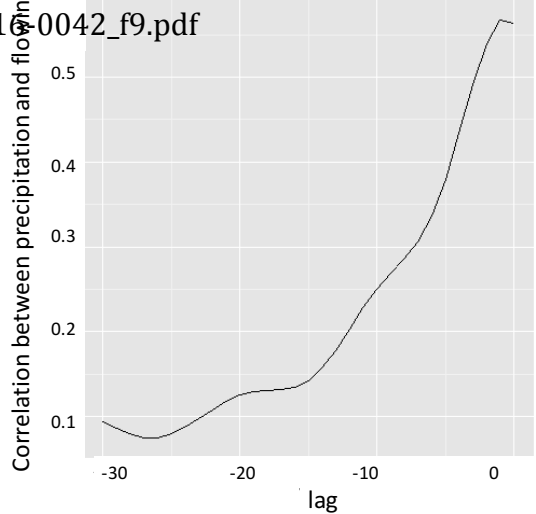


FIGURE 9.b. Correlation between Precipitation and *Flowin*

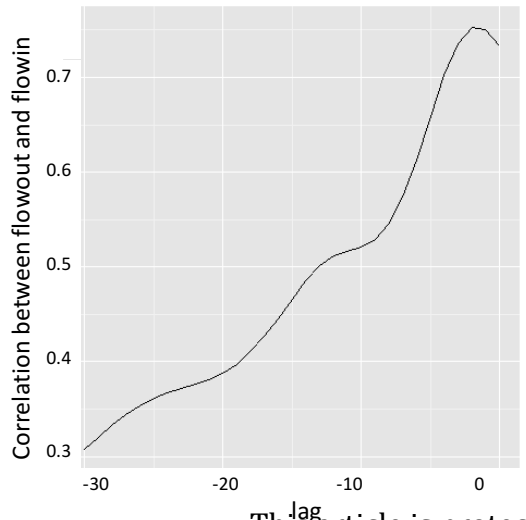
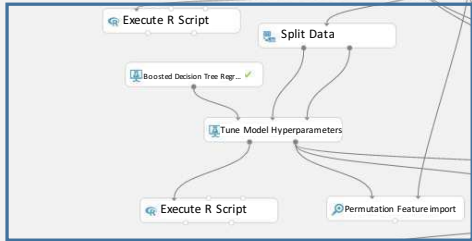


FIGURE 9.c. Correlation between *Flowout* of Upstream Reservoir and *Flowin*

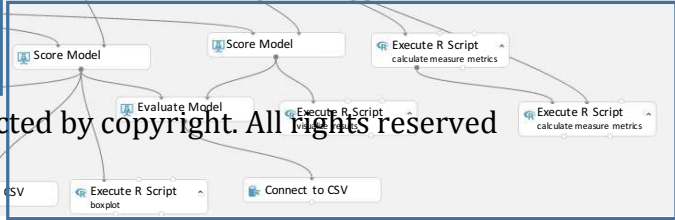
This article is protected by copyright. All rights reserved

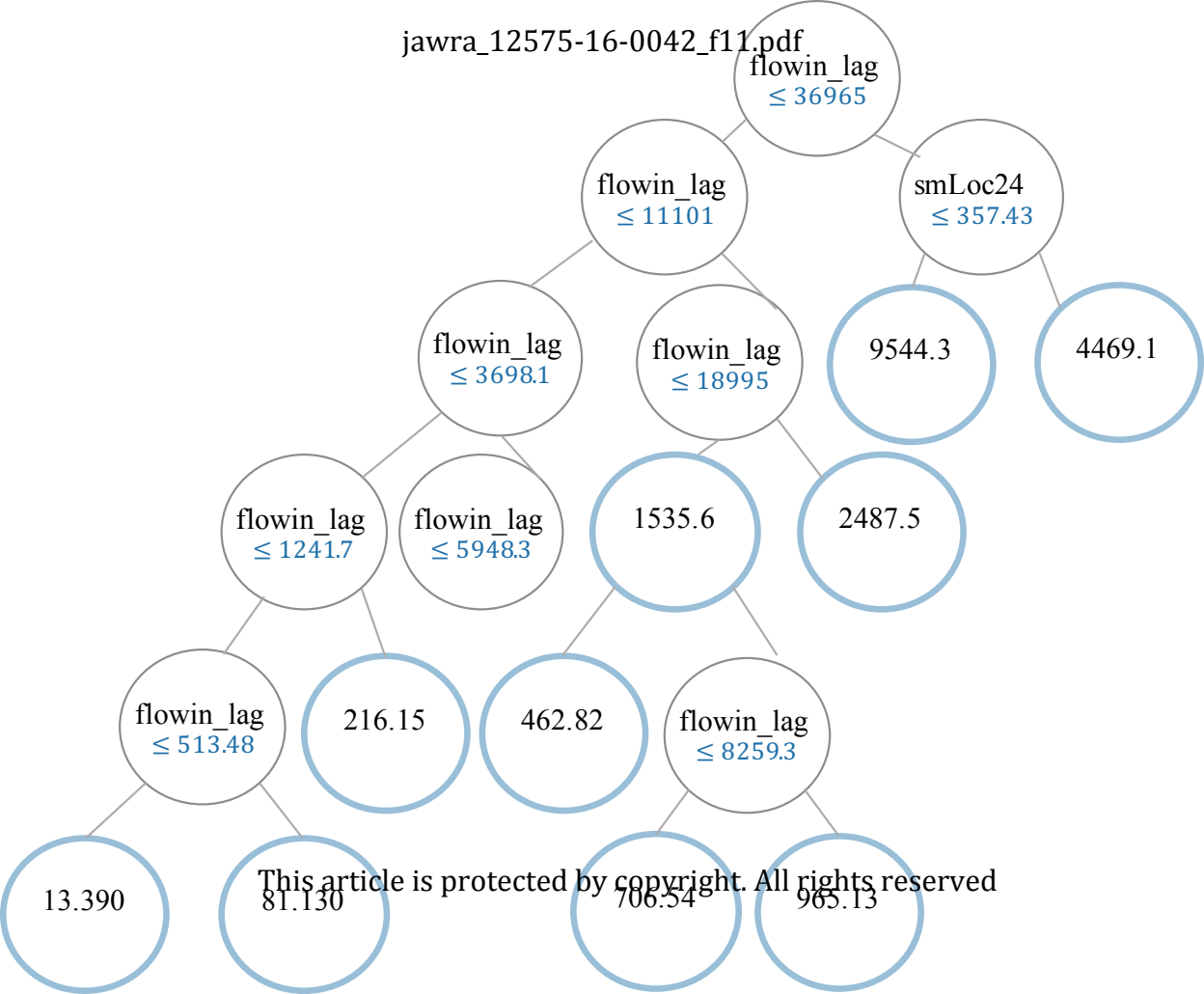
data flow

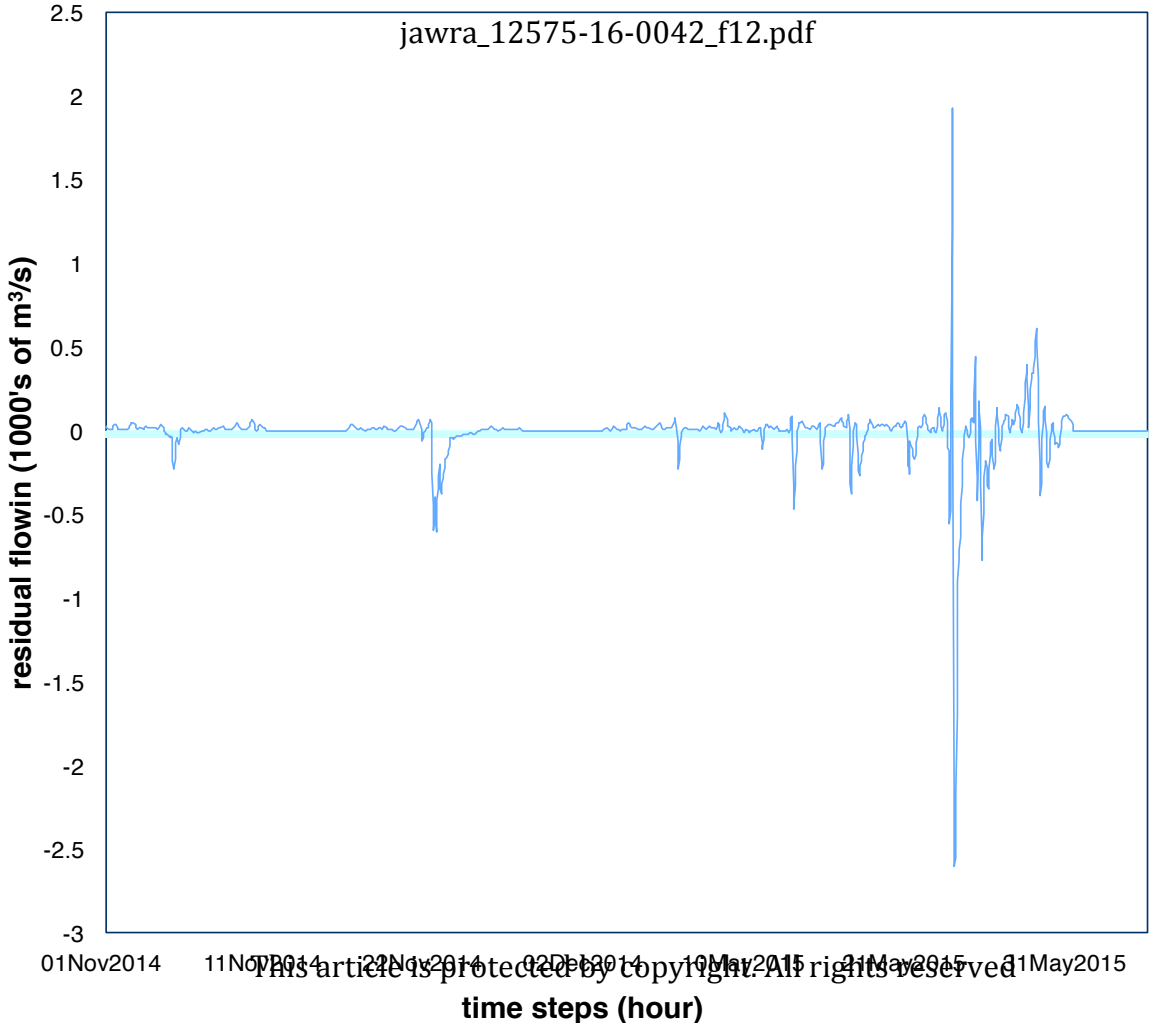


build training model

build test model



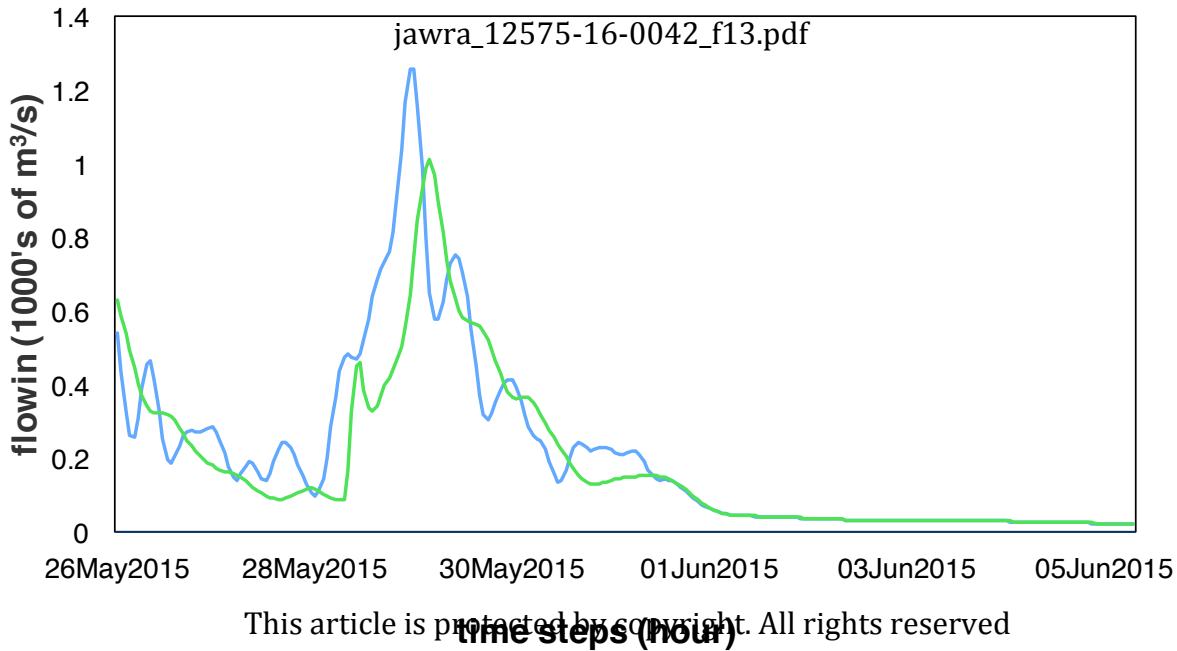




01Nov2014 11Nov2014 21Nov2014 01Dec2014 10Mar2015 21May2015 31May2015

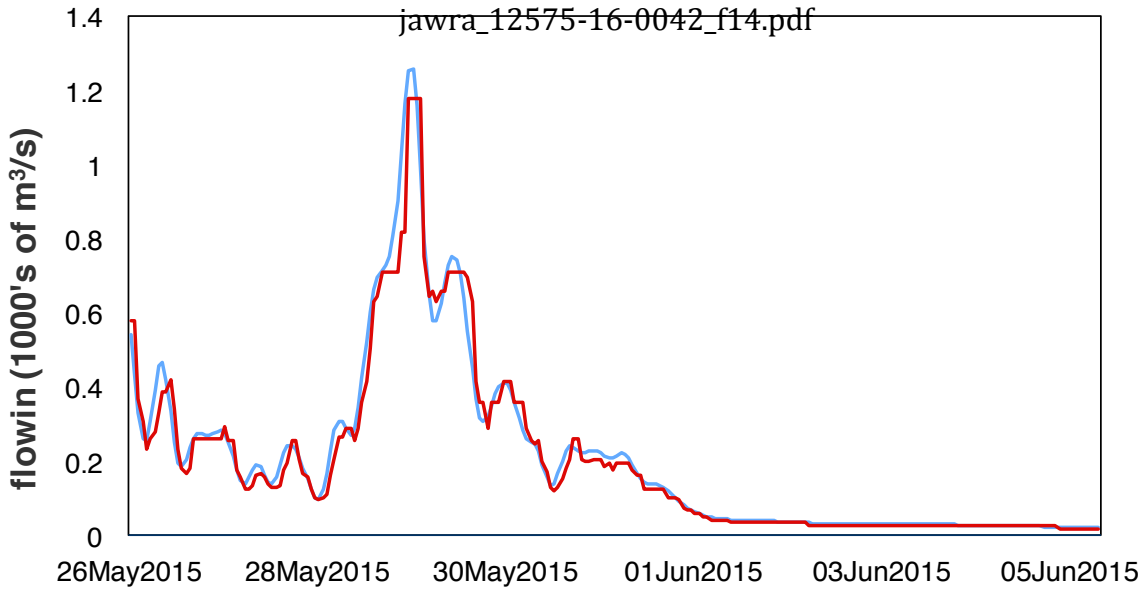
time steps (hour)

jawra_12575-16-0042_f13.pdf



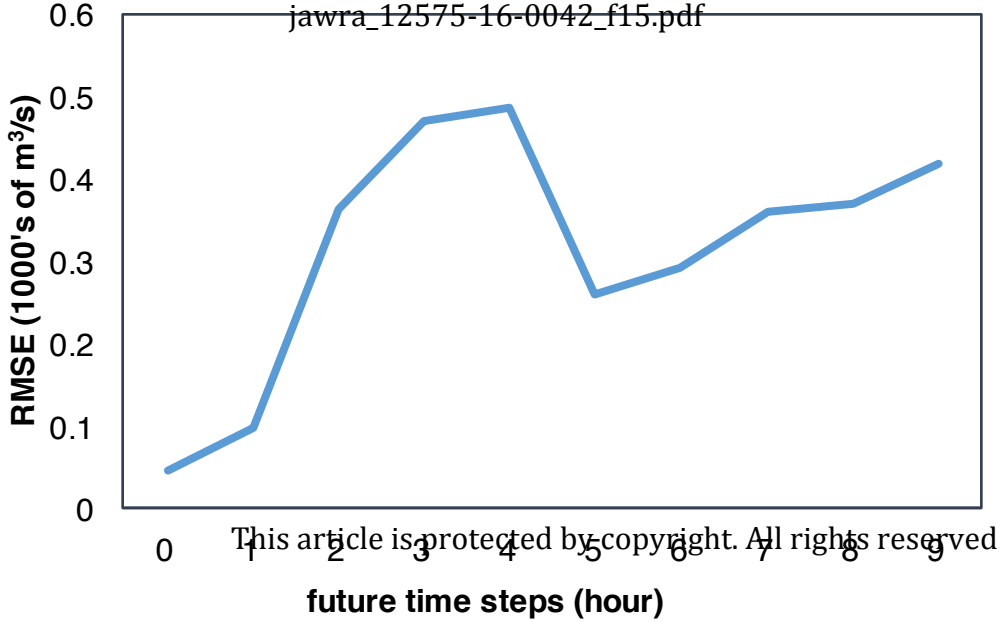
This article is protected by copyright. All rights reserved

— filtered observed flowin — predicted flowin



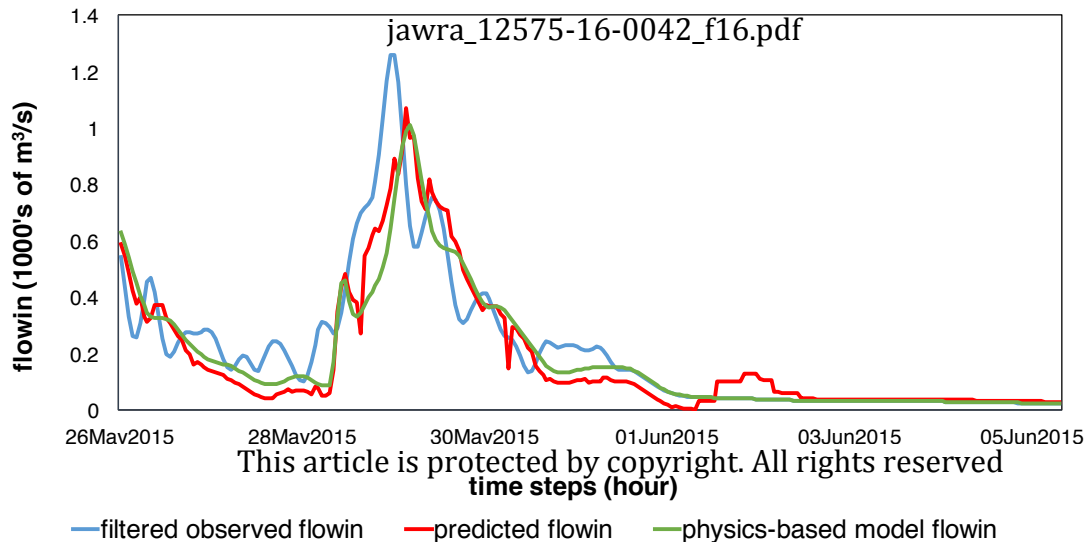
This article is protected by copyright. All rights reserved

— filtered observed flowin — predicted flowin

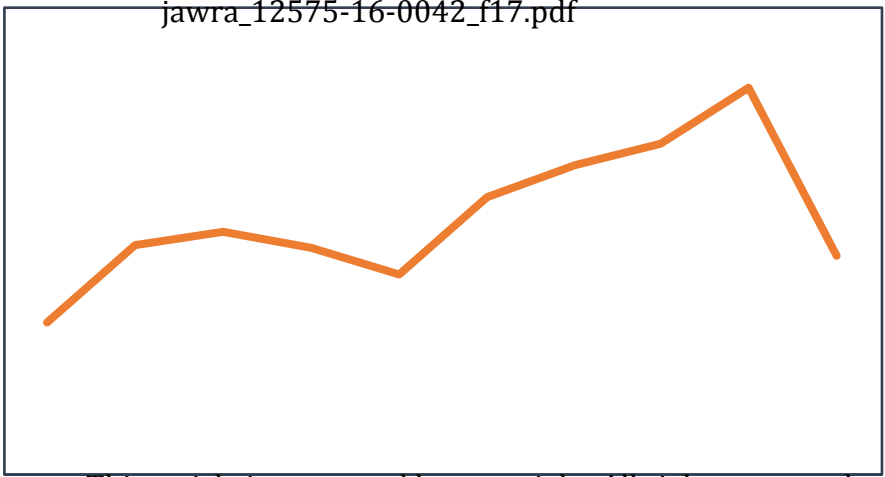


This article is protected by copyright. All rights reserved

future time steps (hour)



RMSE (1000's of m³/s)



This article is protected by copyright. All rights reserved

0 1 2 3 4 5 6 7 8 9

future time steps (hour)

flowout m³/s
flowin_lag m³/s
StartTime (YYYY-MM-DDThh)
EndTime (YYYY-MM-DDThh)

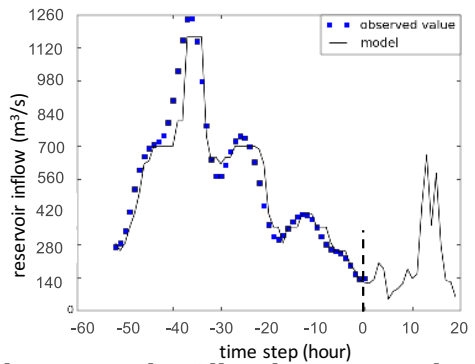
PredictStartTime 2015-05-31T04 (YYYY-MM-DDThh)
PredictStopTime 2015-05-31T04 hours

jawra_12575-16-0042_f18.pdf

Compute

Compute

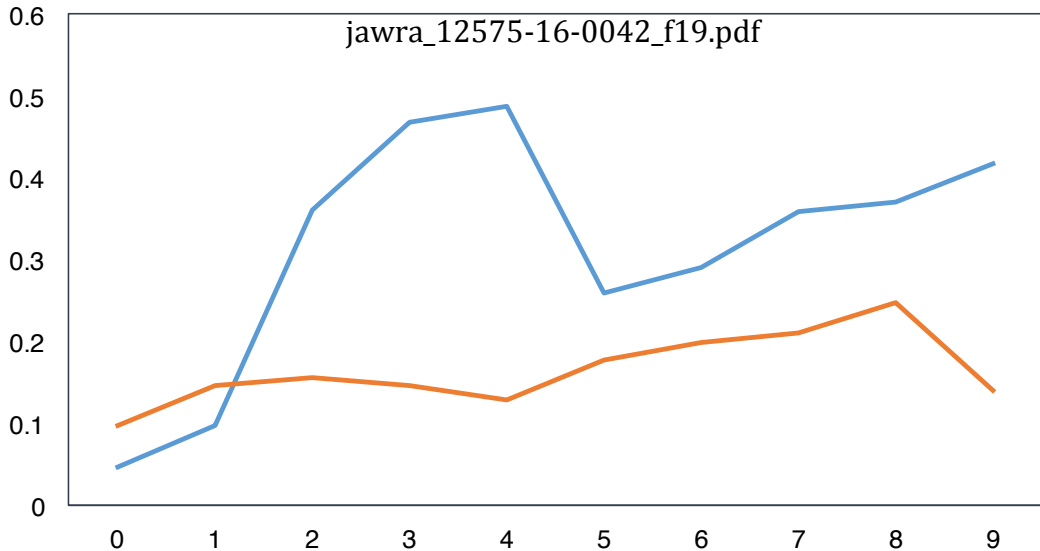
Reservoir inflow from time 2015-05-01T00 to 2015-05-01T01 is 787.0 m³/s.



This article is protected by copyright. All rights reserved

FIGURE 18.a. Flow Prediction Model to Calculate Reservoir Inflow

FIGURE 18.b. Reservoir Inflow Prediction

RMSE (1000's of m³/s)

This article is protected by copyright. All rights reserved

— performance by statistical model

— performance by hybrid model