# A Statistical Graphical Model of the California Reservoir System

**A. Taeb[1]** , **J. T. Reager[2]** , **M. Turmon[2]** , and **V. Chandrasekaran[3]**

[1]Department of Electrical Engineering, California Institute of Technology, Pasadena, CA, USA, [2]Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA, [3]Department of Computing and Mathematical Sciences and Department of Electrical Engineering, California Institute of Technology, Pasadena, CA, USA

**Abstract** The recent California drought has highlighted the potential vulnerability of the state's water management infrastructure to multiyear dry intervals. Due to the high complexity of the network, dynamic storage changes in California reservoirs on a state-wide scale have previously been difficult to model using either traditional statistical or physical approaches. Indeed, although there is a significant line of research on exploring models for single (or a small number of) reservoirs, these approaches are not amenable to a system-wide modeling of the California reservoir network due to the spatial and hydrological heterogeneities of the system. In this work, we develop a state-wide statistical graphical model to characterize the dependencies among a collection of 55 major California reservoirs across the state; this model is defined with respect to a graph in which the nodes index reservoirs and the edges specify the relationships or dependencies between reservoirs. We obtain and validate this model in a data-driven manner based on reservoir volumes over the period 2003–2016. A key feature of our framework is a quantification of the effects of external phenomena that influence the entire reservoir network. We further characterize the degree to which physical factors (e.g., state-wide Palmer Drought Severity Index (PDSI), average temperature, snow pack) and economic factors (e.g., consumer price index, number of agricultural workers) explain these external influences. As a consequence of this analysis, we obtain a system-wide health diagnosis of the reservoir network as a function of PDSI.

## 1. Introduction

### 1.1. Motivation

The state of California depends on a complex water management system to meet wide-ranging water demands across a large, hydrologically diverse domain. As part of this infrastructure, California has constructed 1,530 reservoirs having a collective storage capacity equivalent to a year of mean run off from California rivers (Graf, 1999). The purpose of this system is to create water storage capacity and extend seasonal water availability to meet agricultural, residential, industrial, power generation, and recreational needs.

Major state-wide California precipitation deficits during the years 2012–2015 rivaled the most intense 4 year droughts in the past 1,200 years (Griffin & Anchukaitis, 2014). The drought was punctuated by low snow pack in the Sierra Nevada, declining groundwater storage, and fallowed agricultural lands, in addition to significantly diminished reservoir levels (AghaKouchak et al., 2014; Famiglietti, 2014; Howitt et al., 2014). This sensitivity of the California reservoir network to external conditions (e.g., temperature, precipitation) has implications for state-wide water and agricultural security. In this paper, we seek a characterization of the relationships among the major California reservoirs and their sensitivity to state-wide physical and economic factors, with a view to investigating and quantifying the likelihood of systemic catastrophes such as the simultaneous exhaustion of multiple large reservoirs.

Such an analysis has been difficult to carry out on a system-wide scale due to the size and complexity of the reservoir network. In one direction, a body of work has focused on characterizing the behavior of a small collection of reservoirs using *physical laws* (e.g., Christensen & Lettenmaier, 2004; Christensen et al., 2006; Nazemi & Wheater, 2015). Such approaches quickly become intractable in settings with large numbers of reservoirs whose complex management is based on multiple economic and sectoral objectives (Howitt et al., 2014). The hard-to-quantify influence of human operators and the lack of system closure have made the modeling and prediction of reservoir network behavior using physical equations challenging in

hydrology and climate models (Solander et al., 2016). In a different direction, numerous works have developed *empirical techniques* for modeling the behavior of a small number of reservoirs (e.g., Ashaary et al., 2015; Barnett & Pierce, 2008; Bazartseren et al., 2003; Chen & Liu, 2015; Cheng et al., 2015; Hoerling & Eischeid, 2007; Kuria & Vogel, 2015; Linares-Rodriguez, et al., 2015; Liu & Chung, 2014; Marton, et al., 2015; Nash & Gleick, 1991, 1993; Phatafod, 1989; Revelle & Waggoner, 1983; Wisser et al., 2010; Yang et al., 2016, 2017; Zhang et al., 2017). However, these methods are not directly applicable to modeling a large reservoir network, as the water levels of major reservoirs in California exhibit complex interactions and are statistically correlated with one another (as is demonstrated by our analysis). This necessitates a proper quantification of the complex dependencies among reservoirs in determining the systemic characteristics of the reservoir network.

The focus of this work is to develop a state-wide model over the California reservoir network that addresses the following scientific questions:
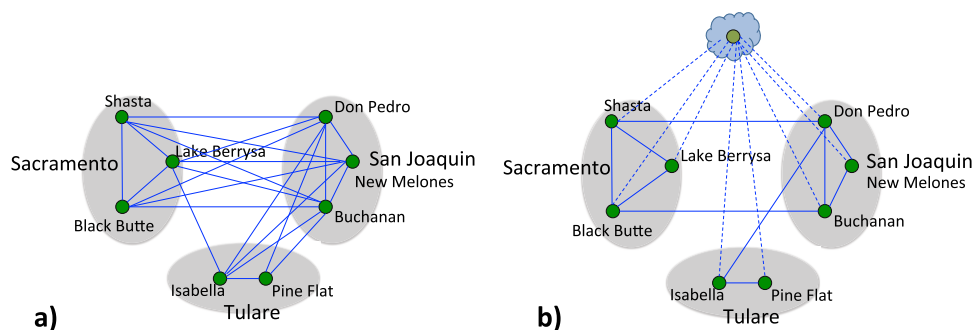
1. What are the interactions or dependencies among reservoir volumes? In particular, how correlated are major reservoirs in the system?
2. Are there common external factors influencing the network globally? Could these external factors cause a system-wide catastrophe?

To the best of our knowledge, this work is the first that attempts such a state-wide characterization of the California reservoir network. The state-wide external factors that we consider in our analysis include physical factors such as state-wide PDSI and average temperature, and economic factors such as the consumer price index and the number of agricultural workers. The focus on these state-wide external influences is driven by the global nature of our analysis; indeed, an exciting direction for further research is to complement our global model with local reservoir-specific factors to obtain an integrated picture of both systemic as well as local risks to the reservoir network.

Answering these questions for the California reservoir system raises a number of challenges, and it is important that any modeling framework that we consider addresses these challenges. First, reservoirs with similar hydrological attributes (e.g., altitude, drainage area, spatial location) tend to behave similarly. As an example, a pair of reservoirs that is approximately at the same altitude or in the same hydrological zone are more likely to have a stronger correlation than those in different altitudes/zones. Therefore, we seek a framework that ably models the complex heterogeneities in the reservoir system. A second challenge, which is in some sense in competition with the first one, is that compactly specified models are much more preferable to less succinct models, as concisely described models are often more interpretable and avoid problems associated with *over-fitting*. Finally, it is crucial that models with both of the preceding attributes have the additional feature that they can be identified in a computationally efficient manner.

### 1.2. Approach and Results

Gaussian *graphical models* offer an appealing and conceptually powerful framework with all the attributes just described. Graphical modeling is a prominent multivariate analysis technique that has been successfully employed in domains as varied as gene regulatory network analysis, social networks, speech recognition, and computer vision (see Jordan, 2004, for a survey on graphical modeling). These models are defined with respect to graphs, with nodes of a graph indexing variables and the edges specifying statistical dependencies among these variables. In a reservoir modeling context, the nodes of the graph correspond to reservoirs and an edge between two reservoirs would describe the strength of the interaction between the levels of those reservoirs. Formally, the strength of an edge specifies the degree of conditional dependence between the corresponding reservoirs; in other words, this is the dependence between two reservoirs conditioned on all the other reservoirs in the network. Informally, an edge in a graphical model denotes the extent to which two reservoirs remain correlated even after accounting for the influence of all the other reservoirs in the network. We illustrate these points using a toy example of a graphical model over a collection of eight reservoirs, shown in Figure 1a. (This figure is purely for explanatory purposes rather than a factual representation of the complex dependencies among reservoirs, which we obtain in section 3). One can imagine that the reservoir volume of Shasta (which is at a high elevation in northern California in the Sacramento hydrological zone) is independent of the reservoir Pine Flat and the reservoir Isabella (which are in southern California in the Tulare hydrological zone) after conditioning on volumes of reservoirs in the central portion of the state (e.g., Black Butte, Lake Berrysa, New Melones, Buchanan, and Don Pedro). These

**Figure 1.** Graphical structure between a collection of eight reservoirs (a) without latent variables and (b) with latent variables. Green nodes represent reservoirs (variables) and the clouded green node represents latent variables. Solid blue lines represent edges between reservoirs and dotted edges between reservoirs and latent variables. The reservoirs have been grouped according to hydrological zones.

relationships are encoded in a graphical model of Figure 1a. In particular, note that Shasta has an edge linking it to each of the reservoirs {Black Butte, Lake Berrysa, Don Pedro, New Melones, Buchanan}, but does not have an edge connecting it to the reservoirs {Pine Flat, Isabella}. Figure 1a is, of course, a cartoon demonstration of a graphical modeling framework. In practice, identifying conditional dependencies between pairs of reservoirs in large networks such as the one considered in our work is a challenging problem, and we describe tractable approaches to learning such a graphical structure underlying the complex California reservoir system in a completely data-driven manner in section 3. To the best of our knowledge, this is the first work that applies graphical modeling techniques to model reservoirs or other water resources.

The graphical modeling framework provides a common lens for viewing two frequently employed statistical techniques. On the one hand, a classical approach for obtaining a multivariate Gaussian distribution over reservoir volumes is via a maximum likelihood estimator. This estimator has been widely used in various domains in the geophysical sciences for multivariate analysis of a collection of random variables (Wackernagel, 2003). The model obtained by this maximum likelihood estimator is specified by a completely connected graphical structure, where all reservoirs are conditionally correlated given all other reservoirs. On the other hand, an independent reservoir model analyses the behavior of an individual reservoir independently of the other reservoirs in the network. This model results in a fully disconnected graphical model. In this paper, we learn a statistical graphical model over the reservoir network in a data-driven manner based on historical reservoir data. This model yields a sparse (yet connected) graphical structure describing the network interactions. We demonstrate that this model outperforms the model obtained via unregularized maximum likelihood estimator and an independent reservoir model. Thus, the reservoir behaviors are not independent of one another but can be specified with a moderate number of interactions. We demonstrate that a majority of these interactions are between reservoirs that are in the same basin or hydrological zone, and among reservoirs that have similar altitude and drainage area.

A natural question is whether some dependencies specified by the graphical model are due to a small number of external phenomena (drought, agricultural usage, Colorado river discharge, precipitation, etc.). For example, water held by a collection of nearby reservoirs might be influenced by a common snow pack variable. Without observing this common variable, all reservoirs in this set would appear to have mutual links, whereas if snow pack is included in the analysis, the common behavior is explained by a link to the snow pack variable. Accounting for latent structure removes these *confounding* dependencies and leads to *sparser and more localized* interactions between reservoirs. Figure 1b illustrates this point. Latent variable graphical modeling offers a principled approach to quantify the effects of *external phenomena* that influence the entire reservoir network. In particular, this modeling framework uses observational data to (1) identify the number of global factors (e.g., latent variables) that summarize the effect of external phenomena on the reservoir network, and (2) identify the residual reservoir dependencies after accounting for these global factors. Our experimental results demonstrate that the reservoir network at a monthly resolution has two distinct global factors, and residual dependencies persist after accounting for these global factors.

Latent variable graphical modeling obtains a mathematical representation of the external phenomena influencing the reservoir network. One is naturally interested in linking these mathematical objects to real-world signals (e.g., state-wide Palmer Drought Severity Index, snow pack, consumer price index). We present an approach for associating semantics to these latent variables. We find that the state-wide Palmer Drought Severity Index (PDSI) is highly correlated ($\rho \approx 0.88$) with one of the latent variables. PDSI is then included as a covariate in the *next* iteration of the graphical modeling procedure to learn a joint model over reservoirs and PDSI. Using this model, we characterize the system-wide behavior of the network to hypothetical drought conditions. In particular, we find that as PDSI approaches $-5$, there is a probability greater than 50% of simultaneous exhaustion of multiple large reservoirs. We further present an approach for identifying specific reservoirs in the network that are at high risk of exhaustion during extreme drought conditions. We find that the Buchanan and Hidden Dam reservoirs are at high risk and describe the stringent water management policies that were enforced to prevent exhaustion.

## 2. Data Set and Model Validation

Our primary data set consists of monthly averages of reservoir volumes, derived from daily time series of volumes downloaded from the California Data Exchange Center (CDEC). We also used secondary data for some covariates.

### 2.1. Reservoir Time Series

As described in section 1, there are 1,530 reservoirs in California. In this work, we perform statistical analysis on the largest 60 reservoirs in California. We apply our analysis on a subset of the reservoirs as they have a large amount of historical data available. Our technique can be extended to a larger collection of reservoirs given sufficient data. For these 60 reservoirs, daily volume data are available during the period of study (January 2003–November 2016). We excluded five reservoirs with more than half of their values undefined or zero, leaving 55 reservoirs. This list of daily values was inspected using a simple continuity criterion and approximately 50 specific values were removed or corrected. Corrections were possible in six cases because values had misplaced decimal points, but all other detected errors were set to missing values. The most common error modes were missing values that were recorded as zero volume, and a burst of errors in the Lyons reservoir during late October 2014 that seems due to a change in recording method at that time.

The final set of 55 reservoir volume time series spans 5,083 days over the 167 months in the study period. It contains two full cycles of California drought (roughly, 2007–2008 and 2012–2015) and three cycles of wet period (2004–2006, 2009–2011, 2016). Four California hydrological zones are represented, with 25, 20, 6, and 4 reservoirs in the Sacramento, San Joaquin, Tulare, and North Coast zones, respectively.

We are interested in long-term reservoir behavior and thus model reservoir volumes at a monthly time scale. In particular, we average the data from daily down to 167 monthly observations. The reservoir data exhibit strong seasonal components. Hence, a seasonal adjustment step is performed to remove these predictable patterns, so that we can model deviations from the underlying trend in the reservoir behavior. With the exception of the Farmington reservoir (which has volume less than $10^8$ m$^3$), the joint volume anomalies of the remaining 54 reservoirs are well-approximated by a multivariate Gaussian distribution. This is demonstrated by a Q-Q plot in Figure S1 of the supporting information. Since a large amount of historical data is available for the Farmington reservoir, we have included it in our analysis. These observed properties suggest that the reservoir data is amenable to the multivariate Gaussian models we employ in this paper. Before being used in the fitting algorithms, each time series is also rescaled by its standard deviation so that each series has unit variance. We note that our statistical approach identifies correlations between reservoir volumes. Since the correlation between two random variables is normalized by their respective variances, this transformation is appropriate.

### 2.2. Covariate Time Series

Latent variable graphical modeling identifies a mathematical representation of the global factors influencing the reservoir network. We link these global factors to real-world signals using ancillary data, i.e., *covariates*, which are observable variables, exogenous to the model, that may affect a large fraction of reservoirs. The particular covariates that we use are temperature (averaged values over California downloaded from NOAA), Palmer Drought Severity Index (averaged values over California downloaded from NOAA), hydroelectric power

generation of California (downloaded from U.S. Energy Information and Administration), Colorado river discharge (averaged values downloaded from United States Geological Survey), and Sierra Nevada snow pack covariate (manually averaged in the Sierra Nevada region where the elevation is over 100 m, gridded observations downloaded from NOAA). Note that since we are interested in state-wide covariates that exert influence over the entire network, these hydrological indicators were averaged over the state of California (or in the case of snow pack and Colorado river discharge, averaged over a large region in the Sierra Nevada and Colorado river, respectively). In addition to these hydrological indicators, we use the following economic factors: state-wide number of agricultural workers (downloaded from State of California Employment Development Department) and state-wide consumer price index (downloaded from Department of Industrial Relations).

For each of the seven covariates, we obtain averaged monthly observations from 2003 to 2016. We apply a time lag of 2 months to the covariates temperature, snow pack, Colorado river discharge, and Palmer Drought Severity Index (the reason for a 2 months lag is explained in section 4.4). As with the reservoir time series, we remove seasonal patterns with a per-month average.

### 2.3. Model Validation

To ensure that the model of the reservoirs is representative of reservoir behavior, we perform model validation using a technique known as *holdout validation* (Hastie et al., 2009). The objective of this technique is to produce models that are not overly tuned to the idiosyncrasies of observational reservoir data, so that these models are representative of future reservoir behavior. In a holdout validation framework, the available data are partitioned into a training set, and a disjoint validation set. The training set is used as input to a fitting algorithm to identify a model. The accuracy of this model is then validated by computing the average log-likelihood of the validation set with respect to the distribution specified by the model. Here, larger values of log-likelihood are indicative of better fit to data. For our experiments, we set aside monthly observations of reservoir volumes and covariates from January 2004–December 2013 as a training set ($n_{train} = 120$) and monthly observations from January 2003–December 2003 and January 2014–November 2016 as a (disjoint) validation set ($n_{test} = 47$). Both the training and validation observations contain a significant amount of annual and interannual variability.

## 3. Dependencies Underlying the Reservoir Network

### 3.1. Method: Graphical Modeling

A common approach for fitting a graphical model to reservoirs is to choose the simplest model, that is, the sparsest network that adequately explains the observational data. Easing this task, for Gaussian graphical models, the graphical structure is encoded in the sparsity pattern of the precision matrix (inverse covariance matrix) over the variables. Specifically, zeros in the precision matrix of a multivariate Gaussian distribution indicate absent edges in the corresponding graphical model. Thus, the number of edges in the graphical model equals the number of nonzeros of the precision matrix $\Theta$. As an example, consider the toy graphical model in Figure 1a. Suppose that the precision matrix $\Theta$ of size $8 \times 8$ is indexed according to the ordering {Shasta, Black Butte, Lake Beryssa, Isabella, Pine Flat, Don Pedro, New Melones, and Buchanan}. Then, $\Theta$ has the following structure:

$$\Theta = \begin{pmatrix} \star & \star & \star & 0 & 0 & \star & \star & \star \\ \star & \star & \star & 0 & 0 & \star & \star & \star \\ \star & \star & \star & \star & 0 & \star & \star & \star \\ 0 & 0 & \star & \star & \star & \star & \star & \star \\ 0 & 0 & 0 & \star & \star & \star & 0 & \star \\ \star & \star & \star & \star & \star & \star & \star & \star \\ \star & \star & \star & \star & 0 & \star & \star & \star \\ \star & \star & \star & \star & \star & \star & \star & \star \end{pmatrix},$$

where $\star$ denotes a nonzero value. The intimate connection between a graphical structure and the precision matrix implies that fitting a sparse Gaussian graphical model to reservoir observational data is equivalent to estimating a sparse precision matrix $\Theta$. Hence, the reservoirs are modeled according to the distribution $y \sim \mathcal{N}(0, \Theta^{-1})$, where $\Theta$ is sparse. Note that the preprocessing to remove climatology causes the mean to be zero. A natural technique to fit such a model to observational data is to minimize the negative log-

likelihood (e.g., maximum likelihood estimation) of data while controlling the sparsity level of $\Theta$. The log-likelihood function of the training observations $\mathcal{D}_{\text{train}} = \{y^{(i)}\}_{i=1}^{n_{\text{train}}} \subset \mathbb{R}^{55}$ (after removing some additive constants and scaling) is given by the concave function

$$\ell(\Theta; \mathcal{D}_{\text{train}}) = \log \det (\Theta) - \text{tr}[\Theta \cdot \Sigma_n] \quad, \tag{1}$$

where $\Sigma_n = \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} y^{(i)} y^{(i)\prime}$ is the sample covariance matrix. Thus, fitting a graphical model to $\mathcal{D}_{\text{train}}$ translates to searching over the space of precision matrices to identify a matrix $\Theta$ that is sparse and also yields a small value of $-\ell(\Theta; \mathcal{D}_{\text{train}})$. This formulation, however, is a computationally intractable combinatorial problem. Recent work (Friedman et al., 2008; Yuan & Lin, 2007) has identified a way around this road block by using a convex relaxation:

$$\hat{\Theta} = \underset{\Theta \in \mathbb{S}^{55}}{\arg \min} -\ell(\Theta; \mathcal{D}_{\text{train}}) + \lambda \, ||\Theta||_1,$$

$$\text{s.t.} \quad \Theta \succeq 0 \quad. \tag{2}$$

The notation $\mathbb{S}^{55}$ denotes the set of symmetric $55 \times 55$ matrices. The constraint $\succ 0$ imposes positive definiteness so that the joint distribution of reservoirs is nondegenerate. The regularization term $|| \cdot ||_1$ denotes the $L_1$ norm (element-wise sum of absolute values) that promotes sparsity in the matrix $\Theta$. The $L_1$ penalty, and more broadly, regularization techniques, are widely employed in inverse problems in data analysis to overcome ill-posedness and avoid problems such as *over-fitting* to moderate sample size (see the textbooks/monographs Bühlmann & van de Geer, 2011; Wainwright, 2014; and the references therein). These regularization approaches have proved to be valuable in many applications, including cameras (Duarte et al., 2008), magnetic resonance imaging (Lustig et al., 2008), gene regularity networks (Zhang & Kim, 2014), and radar (Herman & Strohmer, 2009).

The regularization parameter $\lambda$ in (2) provides overall control of the trade-off between the fidelity of the model to the data and the complexity of the model. In particular, the program (2) with $\lambda = 0$ yields the familiar maximum likelihood covariance estimator. This estimator has a well-known closed form solution $\hat{\Theta} = \Sigma_n^{-1}$. Generally, $\Sigma_n^{-1}$ will not contain any zeros. This implies that the estimated graphical structure is fully connected with close fit to the training data $\mathcal{D}_{\text{train}}$. However, as explored in section 3.2, this model may be overtuned to the idiosyncrasies of the training observations $\mathcal{D}_{\text{train}}$ and will not generalize to future behavior of reservoirs (a phenomenon known as *over-fitting*). Larger values of $\lambda$ yield a sparser graphical model with very large $\lambda$ resulting in a completely disconnected graphical model where the reservoirs are independent of one another. Importantly, for any choice of $\lambda > 0$, equation (2) is a convex program with a unique optimum, and can be solved efficiently using general purpose off-the-shelf solvers (Toh et al., 2006). Further theoretical support of this estimator is presented in (Ravikumar et al., 2011).

We select the regularization parameter $\lambda$ by *holdout validation*. In particular, for any choice of $\lambda$, we supply the training observations $\mathcal{D}_{\text{train}}$ to (2) to learn a graphical model and compute the average log-likelihood of this model on the validation set $\mathcal{D}_{\text{test}} = \{y^{(i)}\}_{i=1}^{n_{\text{test}}} \subset \mathbb{R}^{55}$. We sweep over all values of $\lambda$ to choose the model with the best validation performance. Let the selected model (after holdout validation) be specified by the precision matrix $\hat{\Theta}$. As discussed earlier, the matrix $\hat{\Theta}$ specifies the structural properties of the graphical model of the network. An edge between reservoirs $r$ and $r'$ is present in the graph if and only if $\hat{\Theta}_{r,r'} \neq 0$, with larger magnitudes indicating stronger interactions. We denote the strength of an edge as the normalized magnitude of the precision matrix entry, that is,

$$s(r, r') = |\hat{\Theta}_{r,r'}| / (\hat{\Theta}_{r,r} \hat{\Theta}_{r',r'})^{1/2} \geq 0. \tag{3}$$

The quantity $s(r, r')$ can be viewed as the partial correlation between reservoirs $r$ and $r'$, given all other reservoirs. In particular, a large $s(r, r')$ indicates that reservoirs $r$ and $r'$ are highly correlated even after accounting for the influence of all the other reservoirs in the network. A small value of $s(r, r')$ indicates that the reservoirs $r$ and $r'$ are weakly correlated conditioned on all the reservoirs. Finally, $s(r, r') = 0$ indicates that reservoirs $r$ and $r'$ are independent conditioned on all the remaining reservoirs.

### 3.2. Results: Graphical Model of Reservoir Network
In this section, we explore the properties of a graphical model over the reservoir network. As described in section 3.1, we learn a graphical model by specifying a regularization parameter $\lambda$ and supplying

observations $\mathcal{D}_{\text{train}}$ to the convex program (2). We vary $\lambda$ from 0 to 1 to identify a collection of graphical models. For $\lambda \geq 1$, the graphical model is completely disconnected and not of interest. For each graphical model, we measure the training performance as the log-likelihood of training observation $\mathcal{D}_{\text{train}}$ and the validation performance as the log-likelihood of validation observations $\mathcal{D}_{\text{test}}$. Figure S2 in the supporting information shows the training and validation performance of graphical modeling across $\lambda$. Recall that $\lambda = 0$ yields the unregularized maximum likelihood (ML) estimate. This model has a training performance of $-23.91$ and validation performance of $-1140.4$. Large $\lambda$ (here, $\lambda = 1$) yields an independent reservoir model, where the graphical structure is disconnected. This model has a training performance of $-82.23$ and validation performance of $-101.95$. To obtain a graphical model over the reservoir network, we choose $\lambda = 0.23$ where the validation performance is maximized (i.e., the choice of $\lambda$ using holdout validation). This model has a training performance of $-62.38$ and validation performance of $-\mathbf{85.43}$. Supporting information Table S1 summarizes the training and validation performances of these three models. Results of supporting information Table S1 and Figure S2 show that the training performance is a decreasing function of $\lambda$: smaller values of $\lambda$ lead to a closer fit to training observations. However, small values of $\lambda$ yield a high complexity model that fits the idiosyncrasies of the training data and thus suffers from overfitting. This is evident from the poor validation performance of the unregularized ML estimate (when $\lambda = 0$). The specified graphical model is the superior model since it has a better validation performance than the unregularized ML estimate and an independent reservoir model. Thus, the reservoir behaviors are not independent, but can be characterized by a moderate number of dependencies. In the supporting information, we characterize the sensitivity of the graphical model to the choice of the regularization parameter $\lambda$.
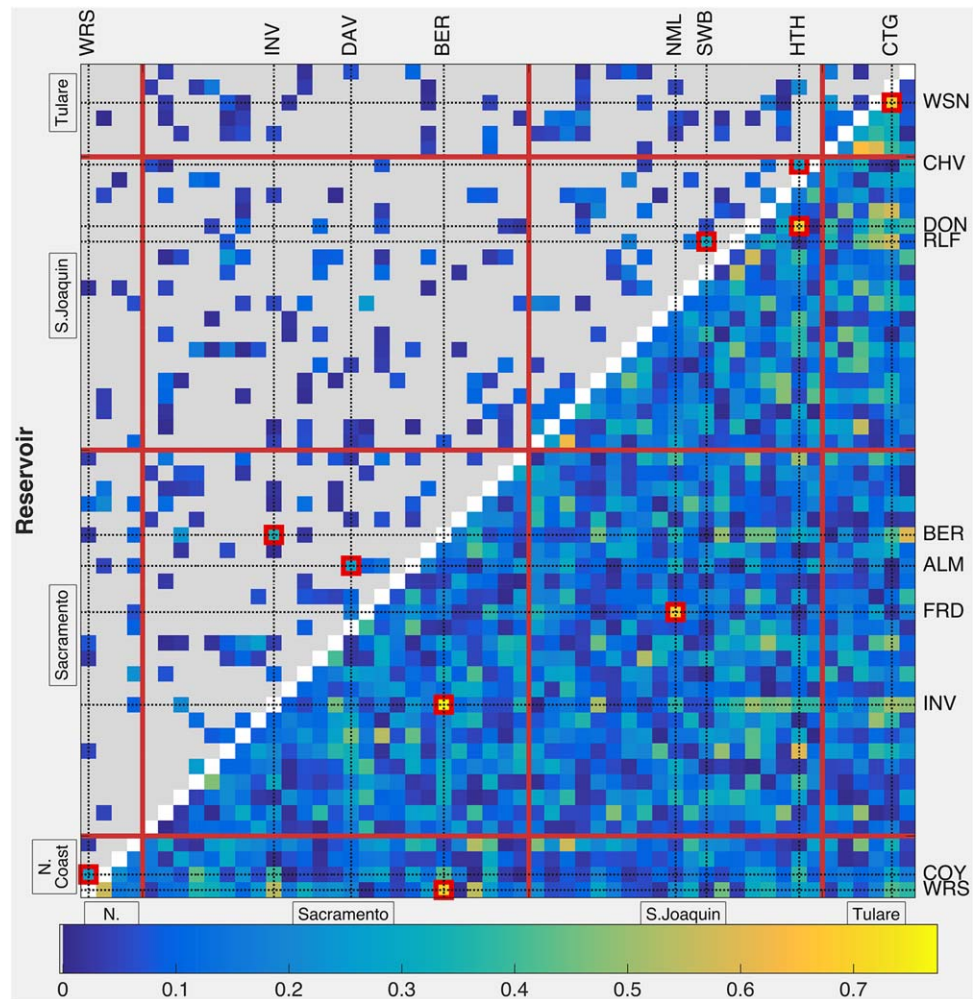
We further explore the properties of the specified graphical model, consisting of 285 edges. Using relation (3), we compute the strength of the connections in the graphical structure. The upper triangle of Figure 2 shows the dependence relationships between reservoirs in this graphical model. The five strongest edges in this graphical structure are between reservoirs Relief—Main Strawberry, Cherry—Hetch Hetchy, Invisible Lake—Lake Berryessa, Almanor—Davis, and Coyote Valley—Warm Spring. We show the geographical location of these pairs of reservoirs in Figure 3. The presence of these strong edges is sensible: each such edge is between reservoirs in the same hydrological zone, and four of these five edges are between pairs of reservoirs fed by the same river. The five most connected reservoirs in order Folsom Lake, Antelope river, Black Butte River, New Exchequer, and French Meadows, all of which are large reservoirs (volume $\geq 10^8$ m$^3$). We show the five strongest connections to Folsom lake in Figure 3, all of which are either connected or are in close proximity to the Sacramento River. As a point of comparison, the lower triangle of Figure 2 shows the graphical structure of the unregularized maximum likelihood estimate. This model yields a fully connected network.

Furthermore, we observe that a majority of interactions in this graphical model are among reservoirs that have similar drainage area (e.g., land where water falls off into reservoirs) and elevation. Figure 4a shows a plot of the ratios of drainage areas between pairs of reservoirs connected via an edge and the strength of the connections. Figure 4b shows a plot of the ratios of altitudes between pairs of connected reservoirs and the strength of the connections. As a point of comparison, Figures 4c and 4d show similar metrics for the unregularized ML estimate. Examining Figure 4, we observe that graphical modeling removes (or weakens) dependencies between reservoirs of vastly different drainage area or elevation. This is expected since reservoirs with substantially different drainage area or elevation are less likely to have similar variability.

We observe that a large portion of the strong interactions occur between reservoirs in the same hydrological zone, here denoted $h(r)$. To quantify this observation, we consider

$$\kappa = \frac{\sum_{r,r' \text{ and } h(r)=h(r')} s(r, r')}{\sum_{r,r'} s(r, r')} \quad , \tag{4}$$

the ratio of within-zone edge strength to total edge strength. The model we fit has $\kappa = 0.85$, so 85% of the total edge strength is between reservoirs in the same hydrological zones. In comparison, $\kappa = 0.46$ for an unregularized ML estimate. Nevertheless, we notice some surprising connections between reservoirs that are geographically far apart. In the next section, we propose a framework to quantify the influence of external phenomena on the reservoir network. We further explore the effect of these external phenomena to remove the confounding relationships between geographically distant reservoirs.

**Figure 2.** Linkages between reservoir pairs in the graphical model (top triangle) compared with those of the unregularized maximum likelihood estimate (bottom triangle). Connection strength $s(r, r')$ is shown in the image map, with unlinked reservoir pairs drawn in gray. The four hydrological zones are separated by red lines. Red boxes surround the five strongest connections in each model.
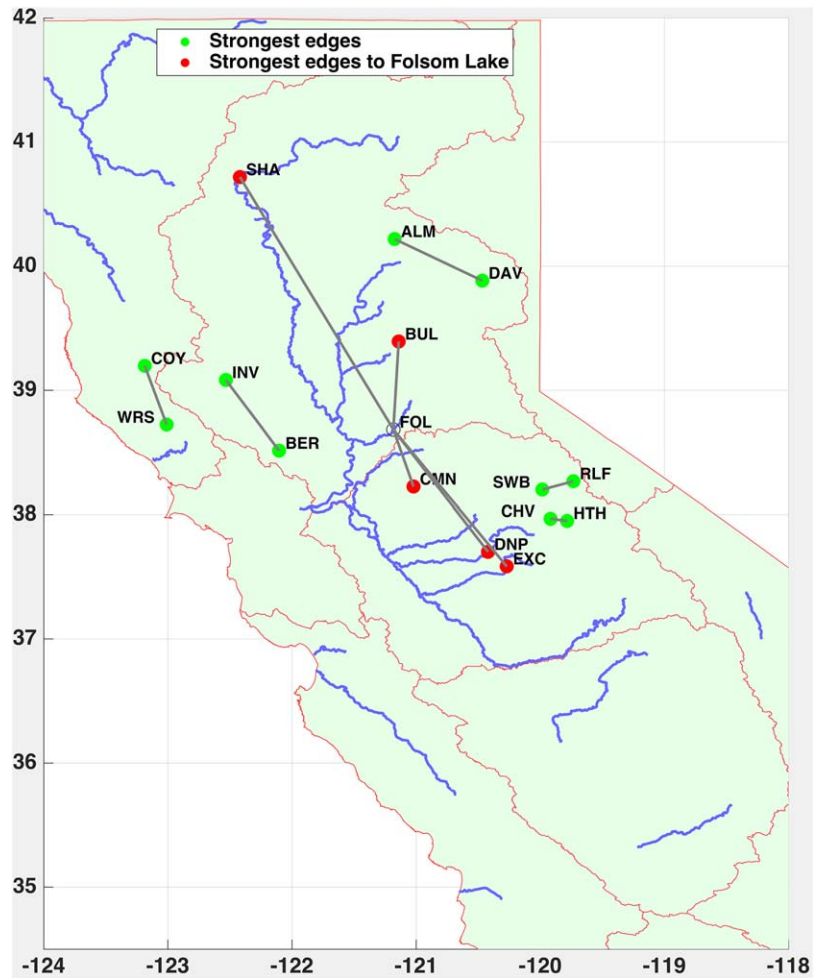
## 4. Global Factors of the Reservoir Network

We identified a graphical model over California reservoirs. Could some of these dependencies specified by the graphical model be due to external phenomena (e.g., global factors)? In this section, we describe an approach, known as *latent variable graphical modeling*, that identifies the number and effect of global factors influencing the reservoir network. Since these global factors are not directly observed (although we later discuss an approach to link global factors to real-world signals), we also denote them as *latent variables*.

### 4.1. Method: Latent Variable Graphical Modeling

As shown by Chandrasekaran et al. (2012), fitting a latent variable graphical model corresponds to representing the precision matrix of the reservoir volumes $\Theta$ as the difference $\Theta = S - L$, where $S$ is sparse and $L$ is a low rank matrix. The matrix $L$ accounts for the effect of external phenomena, and its rank is equal to the number of global factors; these global factors summarize the effect of external phenomena on the reservoir network. The matrix $S$ specifies the residual conditional dependencies among the reservoirs after extracting the influence of global factors. Moreover, the sparsity pattern of $S$ encodes the residual graphical structure among reservoirs. As an example, consider the toy model shown in Figure 1b. Suppose that the matrix $S$ is
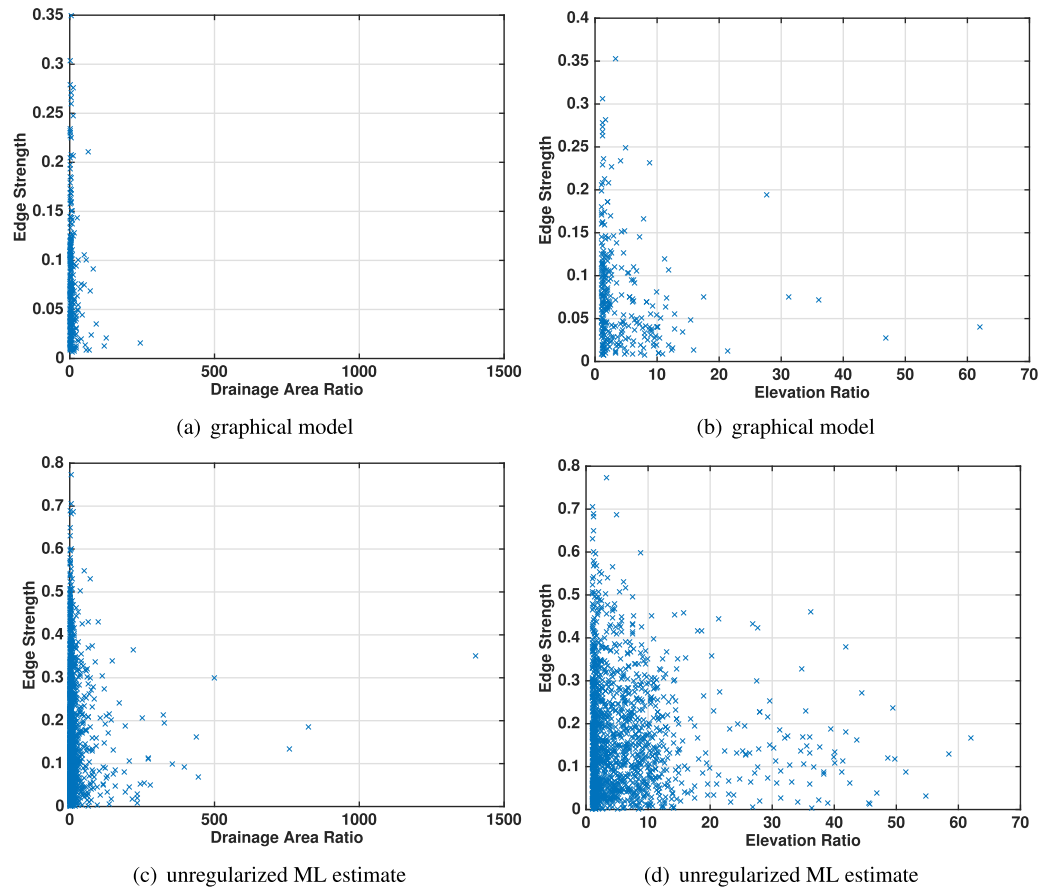
**Figure 3.** A schematic of California and its river network with some reservoir connections. Green nodes represent the five pairs of reservoirs with strongest edge strength in the graphical model. The red nodes represent the five strongest edges to Folsom Lake, which is the most connected reservoir in the network. The acronyms for the reservoirs are: WRS, Wishon; COY, Coyote Valley; INV, Indian Valley; BER, Lake Berryessa; SHA, Shasta; BUL, Bullards Bar; FOL, Folsom Lake; CMN, Camanche; DNP, Don Pedro; EXC, New Exchequer; ALM, Almanor Lake; DAV, Lake Davis; SWB, Main Strawberry; RLF, Relief; CHV, Cherry Valley; HTH, Hetch-Hetchy.

indexed according to the ordering {Shasta, Black Butte, Lake Berrysa, Isabella, Pine Flat, Don Pedro, New Melones, and Buchanan}. Then $S$ has the structure:

$$
S = \begin{pmatrix}
\star & \star & \star & 0 & 0 & 0 & 0 & \star \\
\star & \star & \star & 0 & 0 & \star & 0 & 0 \\
0 & 0 & 0 & \star & \star & 0 & 0 & 0 \\
0 & 0 & 0 & \star & \star & 0 & 0 & \star \\
0 & 0 & 0 & \star & \star & 0 & 0 & 0 \\
0 & \star & 0 & 0 & 0 & \star & \star & \star \\
0 & 0 & 0 & 0 & 0 & \star & \star & \star \\
\star & 0 & 0 & \star & 0 & \star & \star & \star
\end{pmatrix},
$$

where $\star$ denotes a nonzero entry. Fitting a latent variable graphical model to reservoir volumes is to identify the simplest model, e.g., smallest number of global factors and sparsest residual network, that adequately explains the data. In other words, we search over the space of precision matrices $\Theta$ that can be

(a) graphical model      (b) graphical model

(c) unregularized ML estimate      (d) unregularized ML estimate

**Figure 4.** (a) Ratios of drainage areas between pairs of reservoirs connected with an edge and their corresponding edge strengths in a graphical model. (b) Ratios of elevations of pairs of reservoirs connected with an edge and their corresponding edge strengths in a graphical model. (c) Ratios of drainage areas between pairs of reservoirs connected with an edge and their corresponding edge strengths in an unregularized maximum likelihood (ML) estimate. (d) Ratios of elevations of pairs of reservoirs connected with an edge and their corresponding edge strengths in an unregularized maximum likelihood (ML) estimate.

decomposed as $\Theta = S - L$ to identify a matrix $S$ that is sparse, a matrix $L$ that has a small rank, and also yields a small negative log-likelihood $-\ell(\mathcal{D}_{\text{train}}, S - L)$. As with the case of graphical modeling, this formulation is a computationally intractable combinatorial problem. Based on a recent work by Chandrasekaran et al. (2012), a computationally tractable estimator is given by:

$$(\hat{S}, \hat{L}) = \arg\min_{S, L \in \mathbb{S}^{55}} -\ell(S - L; \mathcal{D}_{\text{train}}) + \lambda(||S||_1 + \gamma \text{tr}(L)),$$

$$\text{s.t.} \quad S - L \succ 0, \ L \geqslant 0 \quad . \tag{5}$$

The constraint $\succ 0$ imposes positive definiteness on the precision matrix estimate $S - L$, so that the joint distribution of reservoirs is nondegenerate. The constraint $\geqslant 0$ imposes positive semidefiniteness on the matrix $L$ (see Chandrasekaran et al., 2012, for an explanation of this constraint). Here, $\hat{L}$ provides an estimate for the low-rank component of the precision matrix (corresponding to the effect of latent variables on the reservoir volumes), and $\hat{S}$ provides an estimate for the sparse component of the precision matrix (corresponding to the residual dependencies between reservoirs after accounting for the latent variables).

The regularization parameter $\gamma$ provides a trade-off between the graphical model component and the latent component. In particular, for very large values of $\gamma$, the convex program (5) produces the same estimates as the graphical model estimator (2) (that is, $\hat{L} = 0$ so that no latent variables are used). As $\gamma$ decreases, the number of latent variables increases and correspondingly the number of edges in the residual graphical structure decreases; this is because latent variables account for a global signal common to all reservoirs.

The regularization parameter $\lambda$ provides overall control of the trade-off between the fidelity of the model to the data and the complexity of the model.

As before, the function $|| \cdot ||_1$ denotes the $L_1$ norm that promotes sparsity in the matrix $S$. The role of the trace penalty on $L$ is to promote low-rank structure (Fazel, 2002). As before, for $\lambda, \gamma \geq 0$, equation (5) is a convex program with a unique optimum that can be solved efficiently. Theoretical support for this estimator is presented in Chandrasekaran et al. (2012).

Similar to the graphical model setting, we use the *holdout validation* technique to determine the number of global latent variables and edges in the graphical structure between reservoirs. Concretely, for a particular choice of $\lambda, \gamma$, we supply $\mathcal{D}_{\text{train}}$ as input to the program (5) to learn a latent variable graphical model and compute the average log-likelihood of this model on the validation set $\mathcal{D}_{\text{test}}$. We sweep over all possible choices of $\gamma, \lambda$ and choose a set of parameters that yield the best validation performance.

Let the selected model (after holdout validation) be specified by the parameters $(\hat{S}, \hat{L})$. The matrix $\hat{L}$ denotes the effect of $k = \text{rank}(\hat{L})$ latent variables on the reservoir network. The matrix $\hat{S}$ encodes the residual graphical structure between reservoirs after incorporating $k$ latent variables. We can quantify the strength of the edges of this graphical structure using the relation (3) with $\hat{\Theta}$ replaced with $\hat{S}$. Finally, we quantify the portion of the variability of the network explained by the latent variables as follows: the model estimates the covariance matrix of reservoirs as $(\hat{S} - \hat{L})^{-1}$ so that $y \sim \mathcal{N}(0, (\hat{S} - \hat{L})^{-1})$. Given that the variance of a reservoir $r$ is $[(\hat{S} - \hat{L})^{-1}]_{r,r}$, we denote the overall variance of the network as $\sum_{r=1}^{55}[(\hat{S} - \hat{L})^{-1}]_{r,r}$. The variance of reservoir $r$, conditioned on $k$ latent variables, is given by $(\hat{S}^{-1})_r$. We thus denote the variance of the network conditioned on $k$ latent variables by $\sum_{r=1}^{55}[\hat{S}^{-1}]_{r,r}$. Furthermore, we define the ratio

$$\delta(k) = \frac{\sum_{r=1}^{55}[(\hat{S} - \hat{L})^{-1} - \hat{S}^{-1}]_{r,r}}{\sum_{r=1}^{55}[(\hat{S} - \hat{L})^{-1}]_{r,r}}, \tag{6}$$

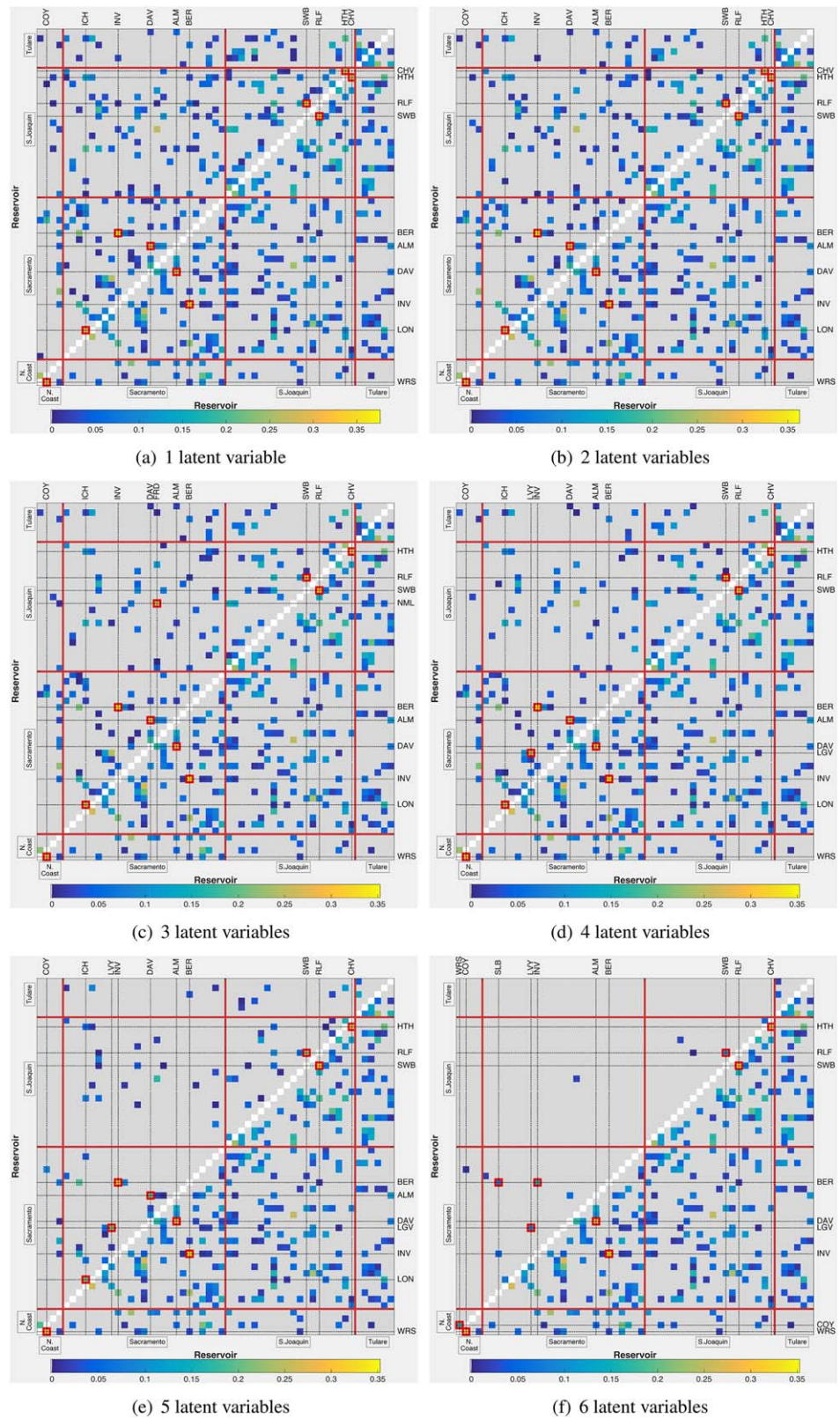as the portion of the variability of the network explained by $k$ latent variables.

### 4.2. Results: Accounting for Global Factors of the Reservoir Network

We first explore the effect of global factors on the connectivity of the reservoir network. Using observations $\mathcal{D}_{\text{train}}$ as input to the convex program (5), we vary the regularization parameters $(\lambda, \gamma)$ to learn a collection of latent variables graphical models. Figure 5 shows the residual conditional graphical structure corresponding to each model. We observe that an increase in the number of latent variables leads to sparser structures and stronger inner-zone connections. Indeed, the ratios of inner zone edge strengths to total edge strength are $\kappa = 0.90$, $\kappa = 0.91$, $\kappa = 0.93$, $\kappa = 0.94$, $\kappa = 0.97$, and $\kappa = 0.99$ for models with 1, 2, 3, 4, 5, and 6 latent variables, respectively. These results support the idea that latent variables extract global features that are common to all reservoirs, and incorporating them results in more localized interactions. The residual dependencies that persist (even after including several latent variables) can be attributed to unmodeled local variables.

Further, appealing to relation (6), the portion of the variability of the network explained by 1, 2, 3, 4, 5, and 6 latent variables is given by $\delta(1) = 0.23$, $\delta(2) = 0.25$, $\delta(3) = 0.28$, $\delta(4) = 0.31$, $\delta(5) = 0.32$, $\delta(6) = 0.40$, respectively. Thus, the effect of latent variables on the network increases as we incorporate more of them in the model. Nonetheless, even six latent variables explain less than 50% of the reservoir variability, with the other portion attributed to residual conditional dependencies between reservoirs. Furthermore, this experiment suggests that both the influence of global latent variables and residual dependencies among reservoirs are important factors of the reservoir network variability.

We now focus on one of these latent variables. In particular, we choose the parameters $(\gamma, \lambda)$ via holdout validation with the validation set $\mathcal{D}_{\text{test}}$ to learn a latent variable graphical model consisting of two latent variables together with a residual graphical model (conditioned on the latent variables) having 171 edges. This is the model corresponding to Figure 5b. Thus, the reservoir network consists of two global factors, and some residual dependencies persist after accounting for their influence. The training and validation performance of this model (in terms of log-likelihood) are given by $-62.11$ and $-85.87$, respectively.

The conditional dependency relationships between reservoir pairs in this residual graphical structure are shown in the upper triangle of Figure 5b. Comparing this graphical structure with the graphical structure without any latent variables (lower triangle of Figure 5b), accounting for the global factors weakens or removes many

**Figure 5.** Linkages between reservoir pairs in the latent-variable sparse graphical model (top triangle) with varying number of latent variables compared with those of the ordinary sparse graphical model (bottom triangle). Connection strength $s(r, r')$ is shown in the image map, with unlinked reservoir pairs drawn in gray. The four hydrological zones are separated by red lines. Red boxes surround the five strongest connections in each model.

connections between reservoirs: 134 are removed and 252 are weakened. Of the 134 edges removed, 94 are between reservoirs in different hydrological zones. Further, the latent variable graphical model has comparable model complexity and training/testing performance to the graphical model without latent variables. We conclude that many of the connections in the graphical model (without latent variables) are due to unmodeled global factors and accounting for these variables leads to fewer remaining conditional dependencies.

Finally, of the 55 reservoirs in our system, 35 are used for sourcing hydroelectric power. In the graphical structure without latent variables, there are 154 pairwise edges between reservoirs that are used for generating hydroelectric power. Once the latent variables are incorporated, all but 15 of these edges are weakened or removed. This suggests that hydroelectric power is strongly correlated to one of the global factors. We verify this hypothesis in the next section.

### 4.3. Method: Interpreting Latent Variables Via Correlation Analysis

Latent-variable graphical modeling identifies a mathematical representation of the global factors of the reservoir network. Naturally, one is interested in linking these mathematical variables to real-world signals to aid understanding of factors that globally affect the reservoir network. We propose an approach to give physical interpretations to the estimated global factors. The high-level intuition of this approach is to identify a space of all possible latent variable data termed *the latent space*. Then we compute the correlation of external covariates (the covariates we consider are in section 2.2) with this space. Candidate covariates with high correlation are variables that globally influence the reservoir network.

Suppose, we identified a latent variable graphical model with estimates $(\hat{S}, \hat{L})$ and $k = \text{rank}(\hat{L})$. Let $z \in \mathbb{R}^k$ denote the latent variables (i.e., $k$ global variables influencing the reservoir network) and $y \in \mathbb{R}^{55}$ denote reservoir volumes; further, partition the joint precision matrix of $(y, z)$ as $\tilde{\Theta} = \begin{pmatrix} \tilde{\Theta}_y & \tilde{\Theta}'_{zy} \\ \tilde{\Theta}_{zy} & \tilde{\Theta}_z \end{pmatrix}$. A natural approximation for the observations of $z$ given observations $\mathcal{D}_{\text{train}}$ is the conditional mean:

$$\tilde{z}^{(i)} = \mathbb{E}[z^{(i)} | y^{(i)}] = -\tilde{\Theta}_z^{-1} \tilde{\Theta}_{zy} y^{(i)}. \tag{7}$$

If $\tilde{\Theta}_z$ and $\tilde{\Theta}_{zy}$ were explicitly known, the length $n_{\text{train}}$ observations $\{\tilde{z}^{(i)}\}_{i=1}^{n_{\text{train}}} \subset \mathbb{R}^k$ would provide an estimate of the latent variables given observations $\mathcal{D}_{\text{train}}$. As discussed in Chandrasekaran et al. (2012), the low-rank component in the decomposition of the marginal precision matrix of $y$ is $\hat{L} = \tilde{\Theta}'_{zy} \tilde{\Theta}_z^{-1} \tilde{\Theta}_{zy}$. However, even though we have $\hat{L}$, this does not uniquely identify $\tilde{\Theta}_z^{-1} \tilde{\Theta}_{zy}$. Indeed, for any nonsingular $A \in \mathbb{R}^{k \times k}$, one can transform $\tilde{\Theta}_z \to A \tilde{\Theta}_z A'$ and $\tilde{\Theta}_{zy} \to A \tilde{\Theta}_{zy}$ without altering $\hat{L}$. In terms of $z$, these observations imply that for any nonsingular $A$, $\{A^{-1} \tilde{z}^{(i)}\}_{i=1}^{n_{\text{train}}}$ is an equivalent realization of the latent variable data: $z$ is recoverable only up to a nonsingular transformation.

Nevertheless, the structure of the low-rank matrix $\hat{L}$ places a constraint on the effect of the latent variables $z$ on $y$. Let $\tilde{Z} \in \mathbb{R}^{n \times k}$ denote a (nonunique) realization of latent variable observations. As we have seen, $\tilde{Z} A^{-1}$ is an equivalent realization. The key *invariant* is the column-space of $\tilde{Z}$, a $k$-dimensional linear subspace of $\mathbb{R}^{n_{\text{train}}}$. We thus, define the *latent space* to be the column-space of $\tilde{Z}$. We recover the latent space as follows: Let $Y \in \mathbb{R}^{n_{\text{train}} \times 55}$ denote observations of reservoir volumes, (7) becomes $\tilde{Z} = Y \tilde{\Theta}'_{zy} \tilde{\Theta}_z^{-1}$. Since the column-space of $Y \tilde{\Theta}'_{zy} \tilde{\Theta}_z^{-1}$ is equal to the column-space of $Y \hat{L}$, the basis elements of the latent space are given by the $k$ left singular vectors of the matrix $YL$, which can be readily computed. We interpret the underlying latent variables by correlating each covariate with this latent space. The manner in which we compute these correlations is presented in the supporting information. A covariate with a large correlation has a strong influence over the entire network.

Suppose we have identified a particular covariate with a large correlation. As described in the supporting information, we can appropriately modify our technique to identify other covariates that are correlated with the latent space after taking away the effect of the specified covariate. Taking this effect away from further analysis is important since the covariates may be dependent on one another (e.g., PDSI and temperature). A covariate that has a high correlation is another global factor. We can repeat this procedure to identify all the $k$ global factors of the reservoir network.

We make two remarks. First, the observations $\{y^{(i)}\}$ used in (7) to characterize the latent space need not be the same as the data employed to identify a latent variable graphical model using the estimator (5). In particular, to quantify the correlation of a covariate with the global factors, we use observations $\{y^{(i)}\}$ in (7)

that are of the same time scale and period as the data that is available for the covariate. As an example, if data for a particular covariate is only available from January 2005–January 2016 at a monthly scale, we use monthly observations of y during the same time period in (7) to characterize the latent space, and subsequently link the observations of the covariate to this space. Second, we note that a subset of the authors of the present paper have proposed an alternate approach for giving physical interpretation to the global factors. This procedure is different that the one proposed in this paper and is based on solving a convex optimization program (Taeb & Chandrasekaran, 2017).

### 4.4. Results: Semantics for Global Factors of the Reservoir Network

The latent variable graphical model identified two global factors influencing the reservoir network. As described in section 4.3, this yields a two-dimensional latent space corresponding to all possible observations of the global factors. To obtain real-world representation of these two global factors, we link the two-dimensional *latent space* to the seven covariates described in section 2.2. Recall from section 2.2 that the covariates PDSI, Colorado river discharge, temperature, and snow pack had a time lag of 2 months. The time lag for each of these covariates was selected to maximize their correlation with the latent space.

We find that the covariates PDSI and hydroelectric power have the largest correlations with $\rho = 0.88$ and $\rho = 0.80$, respectively. Secondary covariate influences are due to consumer price index, Colorado river discharge, Sierra Nevada snow pack (their correlations values are all less than $\rho = 0.5$) with little influence from the number of agricultural workers and temperature. We deduce that PDSI, being computed from variables like precipitation and temperature that control mass balance, is a forcing function on system-wide reservoir levels, while correlation of water levels with aggregate hydropower generation is a system-wide response to high reservoir levels across the network. We then take the effect of PDSI away from the latent space to find the correlation of the modified latent space with the remaining six covariates. We notice that the correlation of CPI (consumer price index) and Colorado river discharge with the latent space do not change very much, since they are unlikely to be structurally connected to PDSI. On the other hand, the correlation of number of agricultural workers, Sierra Nevada snow pack, hydroelectric power, and temperature are significantly reduced as they are largely dependent on PDSI. Nevertheless, all the six covariates have less than 0.5 correlation with the modified latent space. Further tests with additional covariates could yield candidates with strong influence over the reservoir network. The complete list of each covariate and its correlation with the latent space before and after removing PDSI is shown in supporting information Table S2.

In the subsequent section, we describe an approach for incorporating PDSI as a covariate in the next iteration of graphical modeling to learn a joint distribution over reservoir volumes and PDSI. Since we identified one of the two global factors influencing the network, we account for the presence of residual latent variables in the modeling framework.

## 5. Systemic Dependency of the Network to Global Factors

The previous experiment confirmed that the state-wide PDSI signal is a strong forcing function on the entire reservoir network. For purposes of full generality, suppose that using the approach described in section 4.3, we discovered a collection of covariates that are the global factors of the reservoir network. We can extend our modeling framework to incorporate these covariates and characterize the behavior of the network subject to extreme values of these covariates.

### 5.1. Method: Conditional Latent Variable Graphical Modeling

Let $x \in \mathbb{R}^q$ be a collection of covariates that are global factors of the reservoir network (in our setting, $q = 1$ and $x$ is the PDSI variable). Since $x$ can account for the effect of some of the global factors, the distribution of $y$ given $x$ may still depend on a few *residual latent variables*. Therefore, we fit a latent variable graphical model to the conditional distribution of $y|x$. We term this modeling framework as *conditional latent variable graphical modeling*.

Let $\Sigma$ be the join covariance matrix of $(y, x) \in \mathbb{R}^{55+q}$ and $\Theta = \Sigma^{-1}$ be the corresponding joint precision matrix partitioned as $\Theta = \begin{pmatrix} \Theta_y & \Theta_{yx} \\ \Theta'_{yx} & \Theta_x \end{pmatrix}$. The conditional precision matrix of $y$ given $x$ is equal to the submatrix $\Theta_y$. Following the description of section 4.1, fitting a latent variable graphical model to the distribution

of $y$ given $x$ corresponds to decomposing the submatrix $\Theta_y$ as the difference $S_y - L_y$. The matrix $L_y$ is the effect of residual latent variables on the reservoirs after regressing on the covariates $x$, and its rank is equal to the number of residual latent variables. The matrix $S_y$ specifies the residual dependencies among reservoirs after accounting for $x$ and residual latent variables. The sparsity pattern of $S_y$ encodes the residual graphical structure among reservoirs.

Let $\mathcal{D}_{\text{train}}^+ = \{(y^{(i)}; x^{(i)})\}_{i=1}^{n_{\text{train}}} \subset \mathbb{R}^{55+q}$ be the training set of reservoir volumes augmented with covariate data and let $\mathcal{D}_{\text{test}}^+ = \{(y^{(i)}; x^{(i)})\}_{i=1}^{n_{\text{test}}} \subset \mathbb{R}^{55+q}$ be the corresponding validation set. A natural approach for fitting a conditional latent variable graphical model is to choose the simplest model, e.g., the smallest number of residual latent variables and sparsest residual graphical model, that adequately explains the data. Following a similar line of reasoning as the case of latent variable graphical modeling, we arrive at the following estimator for fitting a conditional latent variable graphical model to the observations $\mathcal{D}_{\text{train}}^+$ (A. Taeb & V. Chandrasekaran, Sufficient dimension reduction and modeling responses conditioned on covariates: An integrated approach via convex optimization, 2015, arXiv:1508.03852, hereinafter referred to as A. Taeb & V. Chandrasekaran, online report, 2015):

$$(\hat{\Theta}, \hat{S}_y, \hat{L}_y) = \underset{\substack{\Theta \in \mathbb{S}^{55+q} \\ S_y, L_y \in \mathbb{S}^{55}}}{\arg\min} -\ell(\Theta; \mathcal{D}_{\text{train}}^+) + \lambda(\|S_y\|_1 + \gamma \text{tr}(L_y))$$

$$\text{s.t.} \quad \Theta \succ 0, \quad \Theta_y = S_y - L_y, \ L_y \succeq 0. \tag{8}$$

The term $\ell(\Theta; \mathcal{D}_{\text{train}}^+)$ is the Gaussian log-likelihood function over the variables $(y, x)$, which after removing constants terms and scaling is given by

$$\ell(\Theta; \mathcal{D}_{\text{train}}^+) = \log \det (\Theta) - \text{tr}\left[\Theta \cdot \Sigma_n^+\right] \quad,$$

where $\Sigma_n^+ = \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \begin{pmatrix} y^{(i)} \\ x^{(i)} \end{pmatrix} \begin{pmatrix} y^{(i)} \\ x^{(i)} \end{pmatrix}'$ is the sample covariance matrix of reservoirs and covariates. The

program (8) with $\lambda = 0$ is the unregularized multivariate maximum likelihood estimator of reservoirs and covariates. For $\lambda, \gamma \geq 0$, the regularized maximum likelihood estimator (8) is a convex program with a unique optimum and can be solved efficiently, similar to estimators (1) and (5). Theoretical support for this estimator is presented in A. Taeb and V. Chandrasekaran (online report, 2015). We note that a conditional graphical model could also be obtained using other techniques, such as the convex program proposed by Frot et al. (B. Frot, L. Jostins, & G. McVean, Latent variable model selection for Gaussian conditional random fields, 2017, arXiv:1512.06412).

We select the regularization parameters $\lambda, \gamma$ in (8) via holdout validation with the testing set $\mathcal{D}_{\text{test}}^+$. Concretely, for a particular choice of $\lambda, \gamma$, we supply $\mathcal{D}_{\text{train}}^+$ as input to the program (8) to obtain a conditional latent variable graphical and validate the performance on the validation set $\mathcal{D}_{\text{test}}^+$. We perform this procedure as we vary $\lambda, \gamma$, and choose the model with the best validation performance.

Suppose we obtain a conditional latent variable graphical model over $(y, x) \in \mathbb{R}^{55 \times q}$ with estimates $(\hat{\Theta}, \hat{S}_y, \hat{L}_y)$. We use this model to characterize the behavior of the network in response to the covariates $x$ in the month of November (the analysis can be done for any month). Our metric for the behavior of the network is the *probability of simultaneous exhaustion*: the probability that the volumes of a collection of reservoirs drop below zero. Letting $\hat{\Sigma} = \hat{\Theta}^{-1}$, the composite variable $(y, x) \in \mathbb{R}^{55+q}$ is distributed as $(y, x) \sim \mathcal{N}(0, \hat{\Sigma})$. (Preprocessing to remove climatology causes the mean to be zero.) To determine the behavior of a collection of $K$ reservoirs $\mathbf{r} = \{r_1, r_2, \ldots, r_K\}$ as the covariates $x$ vary, we extract the $(K+q) \times (K+q)$ block of $\hat{\Sigma}$ corresponding to $y_\mathbf{r} \in \mathbb{R}^K$ and $x$, and recall that

$$y_\mathbf{r}|x \sim \mathcal{N}(\hat{\Sigma}_{y_\mathbf{r},x}\hat{\Sigma}_x^{-1}x, \ \hat{\Sigma}_{y_\mathbf{r}} - \hat{\Sigma}_{y_\mathbf{r},x}\hat{\Sigma}_x^{-1}\hat{\Sigma}_{x,y_\mathbf{r}}) \quad, \tag{9}$$

an instance of the standard expressions for the conditional mean and variance of these jointly Gaussian variables. Let the November climatology, subtracted during preprocessing, for reservoir volume $y_r$ ($r \in \mathbf{r}$) be $\mu_{y_r}$, and the November climatology of $x$ be $\mu_x \in \mathbb{R}^q$. Let the scaling used to make the time series of $y_r$ have unit variance be $a_{y_r}$ and the scaling matrix used to make the time series of each covariate to have unit variance be $a_x \in \mathbb{R}^{q \times q}$. Then, for $x = u$, the probability that at least $k$ of $K$ reservoirs have their volume drop below zero in November is:

$$P(A_K(k)|x=a_x(u-\mu_x)), \tag{10}$$

where $A_K(k)$ is the event that $y_r \leq -\mu_{y_r} a_{y_r}$ for at least $k$ of the $K$ reservoirs. The probability in (10), or that of any system-wide event, can be computed using Monte Carlo draws from the joint conditional distribution.

We can further use the model to identify "weak nodes" of the network: reservoirs that are at high risk of exhaustion. In particular, we compute the probability of each reservoir conditioned on PDSI, namely,

$$P(y_r < -\mu_{y_r} a_{y_r}|x=a_x(u-\mu_x)) \quad, \tag{11}$$

by applying equation (10) with $K = 1$.

### 5.2. Results: Network Behavior Under Drought

To obtain a system-wide response to drought, we follow the approach described in section 5.1 to compute the probability of exhaustion of a collection of reservoirs conditioned on particular PDSI. We obtain this probability by learning a conditional latent variable graphical model over reservoir volumes and PDSI. This probability is computed for the month of November, when reservoirs are typically at their lowest, but the same calculation applies to any month. Since we applied a time lag of 2 months to the PDSI time series, these probabilities are computed based on September PDSI.

To learn a joint distribution, let $x \in \mathbb{R}$ denote PDSI and consider a conditional latent-variable graphical model over $(y,x) \in \mathbb{R}^{55+1}$. Using observations $\mathcal{D}_{train}^+$ (consisting of 55 reservoir volumes and PDSI values) and appropriate choice of regularization parameters $\lambda, \gamma$ (using holdout validation), we fit a latent-variable graphical model to the conditional distribution $y|x$ via the estimator (8). The estimated model consists of one residual latent variable (e.g., rank $(\hat{L}_y)=1$). Recall that the reservoir network consists of two global factors. Evidently, by regressing away the effect of PDSI, we are left with one residual latent variable, which supports the observation that PDSI is a global factor of the reservoir network. It is plausible that a portion of the residual latent variable is due to management behavior.

The conditional latent-variable graphical modeling procedure also provides an estimate of a graphical model of the conditional distribution of $y$ conditioned on PDSI (e.g., the matrix $\hat{S}_y$)—this graphical model consists of 206 edges. The training and validation performance of this model is $-61.79$ and $-88.52$, respectively. We now compute the systemic response to drought based on the conditional latent variable graphical model over reservoirs and PDSI. Appealing to relation (10), we can compute the probability that at least $k$ of $K$ reservoirs have their volume drop below zero in November. Here, we consider those reservoirs having capacity of at least $10^8 m^3$ ($K = 31$). Of the 55 reservoirs in our data set, 22 have capacity below $10^8 m^3$. Two of the 33 remaining (Terminus and Success) are flood-control reservoirs: they are unique in that their volume routinely falls below 10% of capacity, independent of PDSI. Thus, we focus on the remaining 31 large reservoirs in what follows. We vary PDSI and compute (10) for selected values of $k$. Figure 6 indicates that with sustained precipitation deficits and a PDSI approaching $-5$, the probability that three or more of California's major reservoirs run dry is greater than 50%. This probability increase above 80% as PDSI drops to $-6$.

### 5.3. Implications

The results of Figure 6 indicate that under severe drought conditions (e.g., small values of PDSI), there is a high risk of simultaneous exhaustion of multiple large reservoirs. To further investigate the implications of drought on reservoir conditions, we use (11) to compute the probability of exhaustion of each reservoir as a function of PDSI. As shown in the supporting information, our results indicate that the reservoirs Hidden Dam and Buchanan have the highest risk of exhaustion (among the 31 reservoirs with capacity $\geq 10^8 m^3$) under severe drought conditions. Stringent management practices; however, have prevented these reservoirs from running dry. Specifically, the Madera



**Figure 6.** System-wide response to drought in a conditional latent variable graphical model: probability that at least $k$ reservoirs out of 31 large reservoirs (with capacity $\geq 10^8 m^3$) will have volume fall to zero, for a range of PDSI; Dashed black line: average September PDSI (September 2004–2015). Dashed blue line: September 2014 PDSI. Dashed red line: September 2015 PDSI. Dashed green line: September 2016 PDSI.

Irrigation District, which owns the water rights of the Hidden Dam reservoir, allowed for the release of very small amount of water during the drought period of 2014–2015. This is because the reservoir volume had reached the minimum pool of 5,000 acre feet ($6.1 \times 10^6$ m$^3$, $\approx 5\%$ of the total capacity) required for recreational purposes. The Buchanan reservoir received a similar degree of stringent management. During the 2014–2015 period, the reservoir volume reached the minimum pool of 10,000 acre feet ($12.2 \times 10^6$ m$^3$, $\approx 6\%$ of the total capacity) required for recreational purposes. As a result, the Chowchilla Water District, which owns the water rights of the Buchanan reservoir, determined that no water will be released during the 2014–2015 period.

Thus, at low reservoir volumes, the stringent management that these reservoirs receive results in their behavior deviating from the predictions of our model. To further highlight this distinction, we examine the historical reservoir volumes of Buchanan and Hidden Dam as a function of PDSI. Suppose we restrict our attention to PDSI greater than $-3$. In this regime, the correlation of the Buchanan and Hidden Dam reservoirs with PDSI as obtained from our model is similar to the empirical historical average. On the other hand, for PDSI values less than $-3$, the empirical correlations are significantly reduced. Concretely, the empirical correlation of the Buchanan reservoir is a factor $\approx 6/100$ of the value estimated by our model. The empirical correlation of the Hidden Dam is a factor $\approx 2/5$ of the correlation estimated by our model (refer to the supporting information for further discussion). The significant reductions in these correlations for low PDSI values highlight the impact of the severe management practices. Our model is representative of the reservoir behavior in a "Business as Usual" (BAU) regime where heavy management practices have not been employed and therefore correlations of PDSI and reservoirs volumes are independent of PDSI value. Consequently, an alternative interpretation of our results is that Figure 6 provides an advanced guideline as to when strict reservoir management *needs* to be employed to leave the BAU regime—in effect breaking the correlation of PDSI and reservoir volumes—to prevent reservoir exhaustion. More specifically, we propose the following rule of thumb in situations where one may have advanced prediction of the PDSI value: if the exhaustion probabilities are low at the predicted value of PDSI, no heavy management effort is likely to be needed and the reservoir could be operated in a BAU setting. If these probabilities start to rise above 50%, this indicates trouble and that water managers *need* to prepare to leave the BAU regime.

To summarize, the proposed model characterizes the risk of exhaustion of large California reservoirs during extreme drought. The proposed methodology can be used to inform water managers of potential risks under typical management behavior. Additionally, the method used here can forecast other key events that precede reservoir exhaustion, such as when power generation is made impossible as water levels drop below turbine inlets, or when water levels reach the minimum pool for recreational purposes.

## 6. Discussion and Future Directions

The California reservoir system is summarized by a complex, dynamic network of correlated time series that respond to a diverse set of global and local factors, including both natural climate processes and human decision-making. Our objective was to develop the first statewide model of this complex network to address these scientific questions:

1. What are the interactions or dependencies among reservoir volumes?
2. Are there common external factors influencing the network globally? Could these external factors cause a system-wide catastrophe?

We appealed to a powerful modeling framework, known as graphical modeling, to address these questions. These models characterize the complex relationships among reservoirs, and can be learned efficiently based on solving a regularized maximum likelihood estimator. We identified a graphical model consisting of 285 edges over the reservoir network and demonstrated that $\approx 85\%$ of the dependencies are between reservoirs in the same hydrological zone. We observed that reservoirs with similar hydrological attributes (e.g., elevation and drainage area) tend to exhibit stronger dependencies. We further characterized Folsom Lake to be the most connected reservoir in the network, and demonstrated its strong dependencies with reservoirs connected to the Sacramento river. To address question 2, we quantified the influence of external phenomena on the network using an extension of the graphical modeling framework, known as *latent variable graphical modeling*. These models can be learned efficiently based on solving a generalization of the maximum likelihood estimator in the graphical modeling setting. Using historical reservoir data, we determined

two global factors influence the reservoir network at a monthly resolution, and proposed a novel methodology to obtain physical interpretation of these global factors. We found that PDSI was highly correlated ($\rho \approx 0.88$) with one of the global factors. We then used PDSI as a covariate in the *next* iteration of the graphical modeling procedure to characterize risks of system-wide catastrophe in response to hypothetical drought conditions. We also identified that Buchanan and Hidden Valley reservoirs are high susceptible to exhaustion.

The approach applied here to study reservoirs has the potential to be applicable across many complex data problems in the geosciences. The graphical modeling technique can be first used to model the complex network of variables. The model can be enhanced to account for global factors (latent variables) that influence the entire network. Then a latent space summarizing all possible configurations of latent variable data can be estimated by model optimization. Candidate external forcing data can be linked to this latent space to find matches. Once a best match is found, the effect of this covariate can be taken away and other covariates could be tested to identify all the factors of the global system variability. Then the latent variables could be included as covariates in a new iteration of the graphical modeling procedure to learn a joint model over the network variables and covariates. Using this model, the behavior of the network under extreme values of the global factors can be characterized. This procedure has the additional value of directing and prioritizing observational efforts.

There are several interesting directions for future research. The analysis of this paper was over a network of 55 major reservoirs in California. It would be interesting to obtain volumetric measurements of many more reservoirs (currently the amount of data available is insufficient for analysis on a larger set of reservoirs) and apply our procedure to obtain a model over this larger network; indeed, there is no other obstruction to carrying out a more extensive analysis with the methodology presented in this paper. Further, the statistical framework developed in this paper is focused on a global model of the reservoir network and the influence of state-wide variables. An exciting direction for future investigation is to complement our modeling framework to account for local variables (e.g., local temperature, local precipitation, etc.). Specifically, associated with each reservoir, we can include a collection of local variables and apply our framework to the reservoir volumes after regressing on the local variables. As described, this procedure would model the reservoir network at both local and global scales.

## References

AghaKouchak, A., Cheng, L., Mazdiyasni, O., & Farahmand, A. (2014). Global warming and changes in risk and concurrent climate extremes: Insights from the 2014 California drought. *Geophysical Research Letters*, 41, 8847–8852. https://doi.org/10.1002/2014GL062308

Ashaary, N., Ishak, W., & Ku-Mahamud, K. (2015). Forecasting model for the change of reservoir water level stage based on temporal pattern of reservoir water level. In *Proceedings of the 5th International Conference on Computing and Informatics* (Paper no. 203, pp. 692–697). Istanbul, Turkey: University Utara Malaysia.

Barnett, T., & Pierce, D. (2008). When will lake Mead go dry? *Water Resources Research*, 44, W03201. https://doi.org/10.1029/2007WR006704

Bazartseren, B., Hildebrandt, G., & Holz, K. (2003). Short-term water level prediction using neural networks and neuro-fuzzy approach. *Neurocomputing*, 55(3), 439–450. https://doi.org/10.1016/S0925-2312(03)00388-6

Bühlmann, P., & van de Geer, S. (2011). *Statistics for high dimensional data*, Springer Series in Statistics. Berlin, Germany: Springer. https://doi.org/10.1007/978-3-642-20192-9

Chandrasekaran, V., Parillo, P. A., & Willsky, A. S. (2012). Latent variable graphical model selection via convex optimization. *Annals of Statistics*, 40(4), 1935–1967. https://doi.org/10.1214/11-AOS949

Chen, W., & Liu, W. (2015). Water quality modeling in reservoirs using multivariate linear regression and two neural network models. *Advances in Artificial Neural Systems*, 6, 521721. https://doi.org/10.1155/2015/521721

Cheng, C., Feng, Z., Niu, W., & Liao, S. (2015). Heuristic methods for reservoir monthly inflow forecasting: A case study of Xinfengjiang reservoir in Pearl river, China. *Water*, 7(8), 4477–4495. https://doi.org/10.3390/w7084477

Christensen, N., & Lettenmaier, D. (2004). A multimodel ensemble approach to climate change impacts on the hydrology and water resources of the Colorado River Basin. *Hydrology and Earth System Sciences*, 3, 1–44. https://doi.org/10.5194/hess-11-1417-2007

Christensen, N., Wood, A., Voisin, N., & Lettenmaier, D. (2006). Effects of climate change on the hydrology and water resources of the Colorado basin. *Climate Change*, 62, 337–363. https://doi.org/10.1023/B:CLIM.0000013684.13621.1f

Duarte, M., Davenport, M. A., Takhar, J. N., Laska, J. N., Ting, S., Kelly, K. F., & Baraniuk, R. G. (2008). Single-pixel imaging via compressive sampling. *IEEE Signal Processing Magazine*, 25(2), 83–91. https://doi.org/10.1109/MSP.2007.914730

Famiglietti, J. (2014). The global groundwater crisis. *Nature Climate Change*, 4, 945–948. https://doi.org/10.1038/nclimate2425

Fazel, M. (2002). *Matrix rank minimization with applications* (PhD thesis).. Stanford, CA: Department of Electrical Engineering, Stanford University.

Friedman, J., Hastie, T., & Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3), 432–441. https://doi.org/10.1093/biostatistics/kxm045

Graf, W. (1999). Dam nation: A geographic census of American dams and their large-scale hydrologic impacts. *Water Resources Research*, 35(4), 1305–1311. https://doi.org/10.1029/1999WR900016

Griffin, D., & Anchukaitis, K. (2014). How unusual is the 2012–2014 California drought? *Geophysical Research Letters*, *41*, 9017–9023. https://doi.org/10.1002/2014GL062433

Hastie, T., Tibshirani, T., & Friedman, R. (2009). *The elements of statistical learning*. New York, NY: Springer.

Herman, M., & Strohmer, T. (2009). High-resolution radar via compressed sensing. *IEEE Transactions on Signal Processing*, *57*(6), 2275–2284. https://doi.org/10.1109/TSP.2009.2014277

Hoerling, M., & Eischeid, J. (2007). Past peak water in the West. *Southwest Hydrology*, *6*(1), 18–19.

Howitt, R., Medellín-Azuara, J., MacEwan, D., Lund, J. R., & Sumner, D. (2014). *Economic analysis of the 2014 drought for California agriculture*. Davis, CA: Center for Watershed Sciences, University of California.

Jordan, M. I. (2004). Graphical models. *Statistical Science*, *19*, 140–155. https://doi.org/10.1214/088342304000000026

Kuria, F., & Vogel, R. (2015). Uncertainty analysis for water supply reservoir yields. *Journal of Hydrology*, *529*(1), 257–264. https://doi.org/10.1016/j.jhydrol.2015.07.025

Linares-Rodriguez, A., Lara-Fanego, V., Pozo-Vazquez, D., & Tovar-Pescador, J. (2015). One-day-ahead stream flow forecasting using artificial neural networks and a meteorological mesoscale model. *Journal of Hydrologic Engineering*, *20*(9), 05015001. https://doi.org/10.1061/(ASCE)HE.1943-5584.0001163

Liu, W., & Chung, C. (2014). Enhancing and predicting accuracy of the water stage using a physical-based model and an artificial neural network-genetic algorithm in a river system. *Water*, *6*(6), 1642–1661. https://doi.org/10.3390/w6061642

Lustig, M., Donoho, D., Santos, J., & Pauly, J. (2008). Compressed sensing MRI. *IEEE Signal Processing Magazine*, *25*(2), 72–82. https://doi.org/10.1109/MSP.2007.914728

Marton, D., Starý, M., & Menšík, P. (2015). Analysis of the influence of input data uncertainties on determining the reliability of reservoir storage capacity. *Journal of Hydrology and Hydromechanics*, *63*(4), 287–294. https://doi.org/10.1515/johh-2015-0036

Nash, L., & Gleick, P. (1991). The sensitivity of stream flow in the Colorado Basin to climatic changes. *Journal of Hydrology*, *125*(3), 221–241. https://doi.org/10.1016/0022-1694(91)90030-L

Nash, L., & Gleick, P. (1993). The Colorado Basin and climate change. (Rep. EPA 230-R-93-009). Washington, DC: United States Environmental Protection Agency.

Nazemi, A., & Wheater, H. (2015). On inclusion of water resource management in earth system models–part 1: Problem definition and representation of water demand. *Hydrology and Earth System Sciences*, *19*, 33–61. https://doi.org/10.5194/hess-19-33-2015

Phatafod, R. (1989). Riverflow and reservoir storage models. *Mathematical Computational Modeling*, *12*(9), 1057–1077. https://doi.org/10.1016/0895-7177(89)90227-6

Ravikumar, P., Wainwright, M., Raskutti, G., & Yu, B. (2011). High-dimensional covariance estimation by minimizing $\ell_1$-penalized log-determinant divergence. *Electronic Journal of Statistics*, *5*, 935–980. https://doi.org/10.1214/11-EJS631

Revelle, P., & Waggoner, P. (1983). Effects of carbon dioxide-induced climatic change on water supplies in the western United States. In *Changing climate: Report of the carbon dioxide assessment committee* (pp. 419–432). Washington, DC: National Academy of Sciences, National Academy Press.

Solander, K., Reager, J., Thomas, B., David, C., & Famiglietti, J. (2016). Simulating human water regulation: The development of an optimal complexity, climate-adaptive reservoir management model for an LSM. *Journal of Hydrometeorology*, *17*(3), 725–744. https://doi.org/10.1175/JHM-D-15-0056.1

Taeb, A., & Chandrasekaran, V. (2017). Interpreting latent variables in factor models via convex optimization. In *Mathematical programming*. Berlin, Germany: Springer. https://doi.org/10.1007/s10107-017-1187-7

Toh, K., Todd, M., & Tutuncu, R. (2006). *SDPT3–A MATLAB software package for semidefinite-quadratic-linear programming*. Retrieved from: http://www.math.nus.edu.sg/mattohkc/sdpt3.html accessed April 2016

Wackernagel, H. (2003). *Multivariate geostatistics: An introduction with applications*. Berlin, Germany: Springer.

Wainwright, M. (2014). Structured regularizers for high-dimensional problems: Statistical and computational issues. *Annual Review of Statistics and its Applications*, *1*, 233–253. https://doi.org/10.1146/annurev-statistics-022513-115643

Wisser, D., Fekete, B., Vörösmarty, C., & Schumann, A. H. (2010). Reconstructing 20th century global hydrography: A contribution to the Global Terrestrial Network-Hydrology (GTN-H). *Hydrology and Earth System Sciences*, *14*, 1–24. https://doi.org/10.5194/hess-14-1-2010

Yang, T., Asanjan, A., Welles, E., Gao, X., Sorooshian, S., & Liu, X. (2017). Developing reservoir monthly inflow forecasts using artificial intelligence and climate phenomenon information. *Water Resources Research*, *53*, 2786–2812. https://doi.org/10.1002/2017WR020482

Yang, T., Gao, X., Sorooshian, S., & Li, X. (2016). Simulating California reservoir operation using the classification and regression-tree algorithm combined with a shuffled cross-validation scheme. *Water Resources Research*, *52*, 1626–1651. https://doi.org/10.1002/2015WR017394

Yuan, M., & Lin, Y. (2007). Model selection and estimation in the Gaussian graphical model. *Biometrika*, *94*, 19–35. https://doi.org/10.1093/biomet/asm018

Zhang, L., & Kim, S. (2014). Learning gene networks under SNP perturbations using eQTL datasets. *PLOS Computational Biology*, *10*(4), e1003608. https://doi.org/10.1371/journal.pcbi.1003420

Zhang, W., Liu, P., Want, H., Chen, J., Lei, X., & Feng, M. (2017). Reservoir adaptive operating rules based on both of historical streamflow and future projections. *Journal of Hydrology*, *553*, 691–707. https://doi.org/10.1016/j.jhydrol.2017.08.031

# Supporting Information for A Statistical Graphical Model of the California Reservoir System

A. Taeb,[1] J.T. Reager,[2] M. Turmon,[2] V. Chandrasekaran,[1]

**Contents of this file**

**Preprocessing the Reservoir/Covariates Data:** Let $\{\bar{y}^{(i)}\}_{i=1}^{n_{\text{train}}} \subset \mathbb{R}^{55}$ and $\{\bar{y}^{(i)}\}_{i=1}^{n_{\text{test}}} \subset \mathbb{R}^{55}$ be the averaged monthly reservoir volumes in the training and validation set respectively. Focusing on a reservoir $r$ and the month of January, let $\mu_{\bar{y}_r}$ be the average reservoir level during January (obtained only from training observations). For each observation $i$ in January, we apply the transformation:

$$\tilde{y}_r^{(i)} = \bar{y}_r^{(i)} - \mu_{\bar{y}_r}.$$

————

[1]California Institute of Technology.

[2]Jet Propulsion Laboratory

We repeat the same steps for all months. Furthermore, letting $\sigma_r$ be the sample standard deviation of the training observations $\{\tilde{y}_r^{(i)}\}_{i=1}^{n_{\text{train}}}$, we produce unit variance observations with the transformation,

$$y_r^{(i)} = \frac{1}{\sigma_r^{1/2}} \tilde{y}_r^{(i)}.$$

We repeat the same steps for all reservoirs to obtain the preprocessed reservoir observations $\{y^{(i)}\}_{i=1}^{n_{\text{train}}}$ and $\{y^{(i)}\}_{i=1}^{n_{\text{test}}}$. Finally, the same steps are repeated to preprocess the covariates data.

**Checking Gaussianity:** We verify that the joint reservoir anomalies (after preprocessing steps) can be well-approximated by a multivariate Gaussian distribution. To check for the Gaussianity assumption, we use a commonly employed method known as Q-Q plot. This is a graphical procedure for comparing two probability distribution by plotting their quantiles against each other. In particular, we compare the quantiles of the reservoir observations with a multivariate normal distribution. Figure 1(a) shows the Q-Q plot for the 55 reservoirs. We notice that by removing the Farmington reservoir, the Q-Q plot shown in Figure 1(b) exhibits a strong linear relationship, suggesting that these 54 reservoirs are well-approximately jointly by a multivariate Gaussian distribution.

**Sensitivity of Graphical Model to $\lambda$:** As described in the main text, the regularization parameter $\lambda$ is varied from 0 to 1 to identify a collection of graphical models. For each graphical model, we measure the training and validation log-likelihood performances. Figure 2 illustrates the training and validation performances for different values of $\lambda$. Recall that $\lambda = 0$ corresponds to an unregularized maximum likelihood estimate and $\lambda = 1$ corresponds to independent reservoir model. We chose $\lambda = 0.23$ to obtain a graphical

model with the best validation performance. The training and validation performances of

these models are summarized in Table 1.

To demonstrate that the graphical model estimate does not vary significantly under small

perturbations to $\lambda$, we also obtain graphical model estimates with $\lambda = 0.26$ and $\lambda = 0.20$

(Recall that the edge strengths in a graphical model contain the relevant information of

the model). Figure 3(a) compares the edge strengths of the model with $\lambda = 0.23$ and

the model with $\lambda = 0.20$. Furthermore, Figure 3(b) compares the edge strengths of the

model with $\lambda = 0.23$ and the model with $\lambda = 0.26$. Evidently, strong edges persist across

all models, with a few weak edges removed or added as $\lambda$ is varied. The total number

of edges in the graphical model when $\lambda = 0.20$, $\lambda = 0.23$, and $\lambda = 0.26$ is 295, 285, and

279 respectively. Furthermore, the quantity $\kappa$ (defined in equation (4) of main paper) is

$0.852, 0.859$, and $0.862$ for $\lambda = 0.20$, $\lambda = 0.23$, and $\lambda = 0.26$. These results suggest that

our conclusions are not particularly sensitive to the choice of the regularization parameter,

although we chose $\lambda = 0.23$ as it leads to the best validation performance.

**Correlating Covariates to the Latent Space:** Latent variable graphical modeling

identifies a summarization of external phenomena influencing the reservoir network; these

influences are summarized by global latent variables. In the main paper, we introduced the

*latent space*, a space of all possible configurations of the latent variable time series. Here,

we describe the manner in which compute the correlation of a candidate covariate with

the latent space. Let $\mathcal{T} \subset \mathbb{R}^n$ with $\dim(\mathcal{T}) = k$ denote the latent space. Let $X_1 \in \mathbb{R}^{n_{\text{train}}}$

be the $n_{\text{train}}$ observations of the covariate $x_1$ (normalized to have unit variance). The

correlation of this covariate with the latent space is given by:

$$\text{corr}(x_1) = \left\| \mathcal{P}_{\mathcal{T}}(X_1) \right\|_{\ell_2},$$

where $\mathcal{P}_{\mathcal{T}}$ denotes the projection matrix onto the subspace $\mathcal{T}$. By definition, the quantity

$\text{corr}(x_1)$ is between 0 and 1 with large values indicating that the covariate $x_1$ has a strong

influence over the entire reservoir network.

Suppose we have identified a covariate $x_1$ that is highly correlated with the latent space.

We can modify our technique to identify other covariates that are correlated with the

latent space after taking away the effect of the covariate $x_1$.

Let $U_1 D_1 V_1'$ be the reduced SVD of $X_1$ where $U_1 \in \mathbb{R}^{n_{\text{train}}}$, $D_1 \in \mathbb{R}$ and $V_1 \in \mathbb{R}$. Let

$X_2 \in \mathbb{R}^{n_{\text{train}}}$ be the $n_{\text{train}}$ observations of the covariate $x_2$. The correlation of a covariate

$x_2$ with the latent space after taking away the effect of $x_1$ is given by:

$$\text{corr}_{x_1}(x_2) = \left\| (I - U_1 U_1') \mathcal{P}_{\mathcal{T}} (I - U_1 U_1')(X_2) \right\|_{\ell_2}.$$

If the quantity $\text{corr}_{x_1}(x_2)$ is large, then the covariate $x_2$ is strongly correlated to the second

global statewide variable. We can once again take away the effect of the covariates $x_1$

and $x_2$ from the latent space, and find its correlation with another covariate $x_3$. Let

$U_2 D_2 V_2'$ be the reduced SVD of $[X_1, X_2] \in \mathbb{R}^{n_{\text{train}} \times 2}$ where $U_2 \in \mathbb{R}^{n_{\text{train}} \times 2}$, $D_2 \in \mathbb{R}^{2 \times 2}$ and

$V_2 \in \mathbb{R}^{2 \times 2}$. Let $X_3 \in \mathbb{R}^{n_{\text{train}}}$ be the $n_{\text{train}}$ observations of the covariate $x_3$. The correlation

of a covariate $x_3$ with the latent space after taking away the effect of $x_1$ and $x_2$ is given

by:

$$\text{corr}_{x_1, x_2}(x_3) = \left\| (I - U_2 U_2') \mathcal{P}_{\mathcal{T}} (I - U_2 U_2')(X_3) \right\|_{\ell_2}.$$

Similarly, if the quantity $\text{corr}_{x_1, x_2}(x_3)$ is large, then the covariate $x_3$ is strongly correlated

to the third global driver. We can repeat this procedure to identify all the $k$ global drivers

influencing the reservoir network.

The latent variable graphical model identified two global drivers influencing the reservoir network. As described in the preceding paragraphs, this yields a two dimensional latent space corresponding to all possible observations of the global drivers. To obtain real-world representation of these two global drivers, we link the two dimensional *latent space* to the 7 covariates described in Section 2.2 (main paper). The correlation values of each covariate with the latent space are shown in the second column of Table 2. We then take the effect of PDSI away from the latent space to find the correlation of the modified latent space with the remaining 6 covariates. These correlation values are shown in the third column of Table 2.

**Identifying Reservoirs Most at Risk of Exhaustion:** As described in the main text, our modeling framework serves a powerful tool to identify reservoirs that are high risk of exhaustion so that appropriate preventive management practices could be employed. For each reservoir, we sweep over a range of PDSI and use (11) (main text) to compute probabilities of exhaustion. Figure 4 shows those reservoirs (among 31 large reservoirs with capacity greater than $10^8 m^3$) that were highly sensitive to PDSI. Evidently, these reservoirs are at high risk of exhaustion, and additionally, some have a greater sensitivity to small PDSI changes than others. We focus on two reservoirs with highest risk of exhaustion: Buchanan and Hidden Dam reservoir. We consider Figure 5 which demonstrates the historical volumes of these reservoirs in response to PDSI. Notice that as expected, there is a positive correlation between PDSI and reservoir volumes: smaller values of PDSI generally result in a lower volume. An interesting phenomenon seems to occur for very small values of PDSI (e.g. less than 3 corresponding to drought period 2014-2015). In

this range, changes to PDSI do not appear to substantially impact the reservoir volumes. In other words, the correlation between PDSI and reservoir volumes is significantly reduced as compared to the correlation during normal and wet periods. To provide concrete numbers on the reduction in this correlation, we focus on November volumes of Buchanan and Hidden Dam reservoirs and the corresponding September PDSI values. We further restrict to observations where PDSI is less than 3. We compute the Pearson Correlation Coefficient between PDSI and each reservoir during this period. This correlation for the Buchanan reservoir is a factor of $\approx 6/100$ of the value estimated by our model. Similarly, the correlation for the Hidden Dam is a factor $\approx 2/5$ of the correlation estimated by our model. As described in the main paper, the large drops in correlations are due to strict management. Figure 6 demonstrates the amount of water from precipitation into the Hidden Dam and Buchanan reservoirs, the total inflow, and the outflow as a consequence of the stringent management efforts. Examining Figure 6, notice that there was little to no outflow of water, which keeps the reservoir volumes mostly constant and prevents them from running dry.

(a)                                                                            (b)
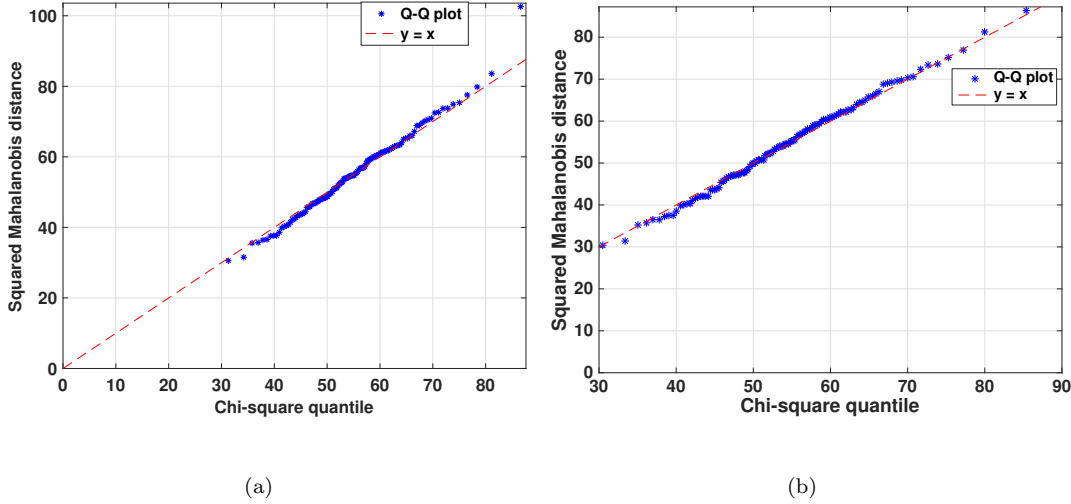
Figure 1: (a): Q-Q plot of the entire set of 55 reservoirs, (b): Q-Q plot of 54 reservoirs (excluding the Farmington reservoir). The Q-Q plots are against a multivariate Gaussian distribution. Notice that $y = x$ is a close approximation to the Q-Q plot in (b) implying that 54 reservoirs (excluding Farmington reservoir) is well approximated by a multivariate Gaussian distribution.

| Model | Training performance | Validation performance |
|---|---|---|
| unregularized ML estimate ($\lambda = 0$) | $-23.91$ | -1140.4 |
| independent reservoir model ($\lambda = 1$) | $-83.23$ | $-101.95$ |
| graphical model ($\lambda = 0.23$) | $-63.52$ | $\mathbf{-85.54}$ |

Table 1: Training and validation performances of unregularized maximum likelihood (ML) estimate, independent reservoir model, and graphical model. As larger values of log-likelihood are indicative of better performance, the graphical model is the superior model.
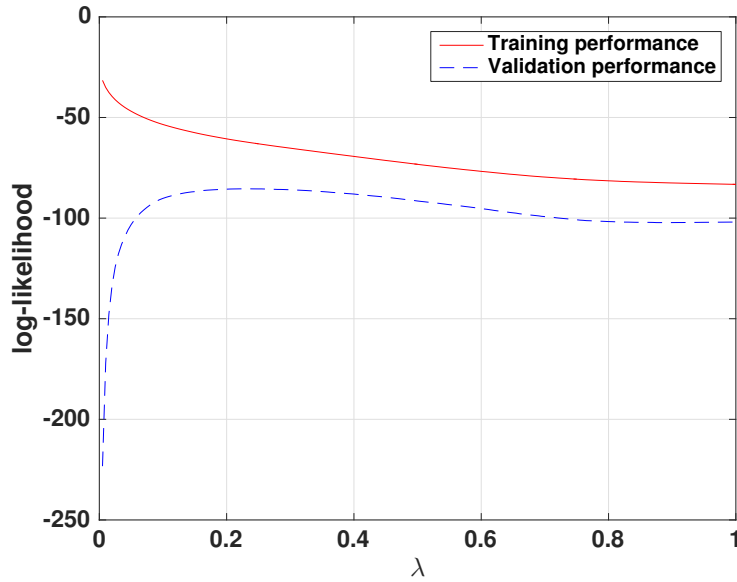
Figure 2: Training and validation performance of graphical modeling for different values of the regularization parameter $\lambda$. The training performance is computed as the average log-likelihood of training samples and the validation performance is computed as the average log-likelihood of validation samples.
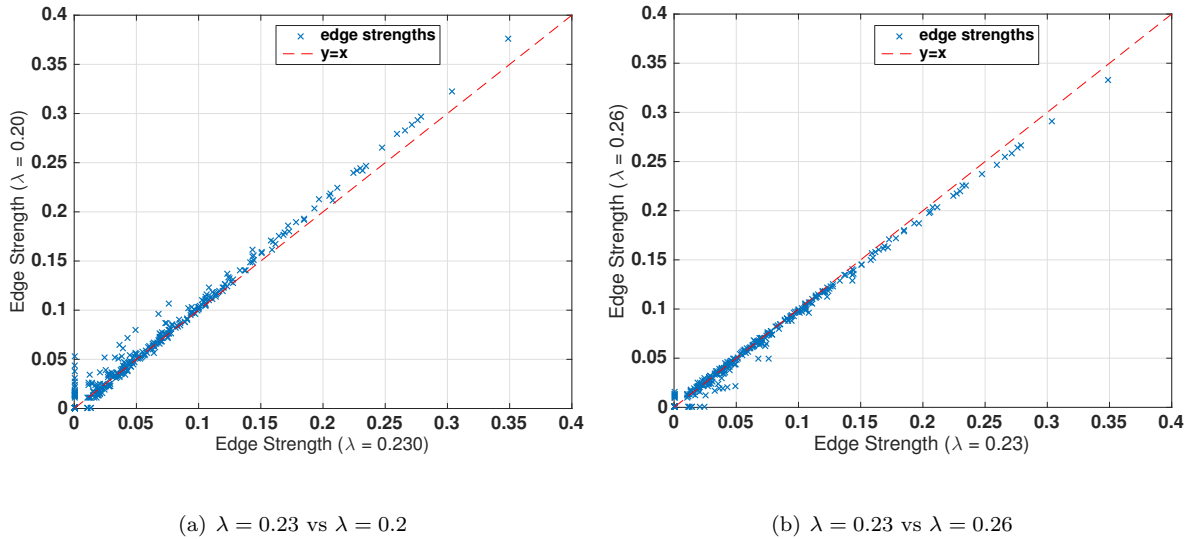


(a) $\lambda = 0.23$ vs $\lambda = 0.2$



(b) $\lambda = 0.23$ vs $\lambda = 0.26$

Figure 3: Sensitivity of the graphical model estimate to perturbations of $\lambda$ around the optimal value $\lambda = 0.23$ (this choice of $\lambda$ leads to optimal validation performance): we observe that strong edges in the original model are strong edges in the perturbed model (i.e. with perturbed $\lambda$) with approximately the same strength.

| Covariate | Correlation | Correlation after removing PDSI |
|---|---|---|
| Palmer Drought Severity Index (PDSI) | 0.88 | N/A |
| Hydroelectric power | 0.80 | 0.09 |
| Sierra Nevada snow pack | 0.50 | 0.32 |
| Consumer Price Index (CPI) | 0.33 | 0.25 |
| Colorado river discharge | 0.29 | 0.23 |
| Number of agricultural workers | 0.17 | 0.03 |
| Temperature | 0.10 | 0.04 |

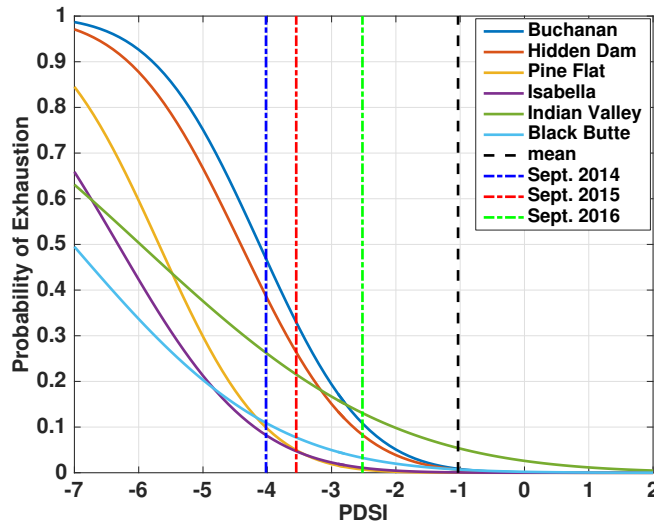Table 2: Covariates and correlations with the latent space before and after removing PDSI



Figure 4: Individual reservoir responses to drought in a conditional latent variable graphical model: probability that six most-at-risk reservoirs out of 31 large reservoirs (with capacity $\geq 10^8 m^3$) will have volume drop below zero; Dashed black line: average September PDSI (September 2004-September 2015). Dashed blue line: September 2014 PDSI. Dashed red line: September 2015 PDSI. Dashed green line: September 2016 PDSI.

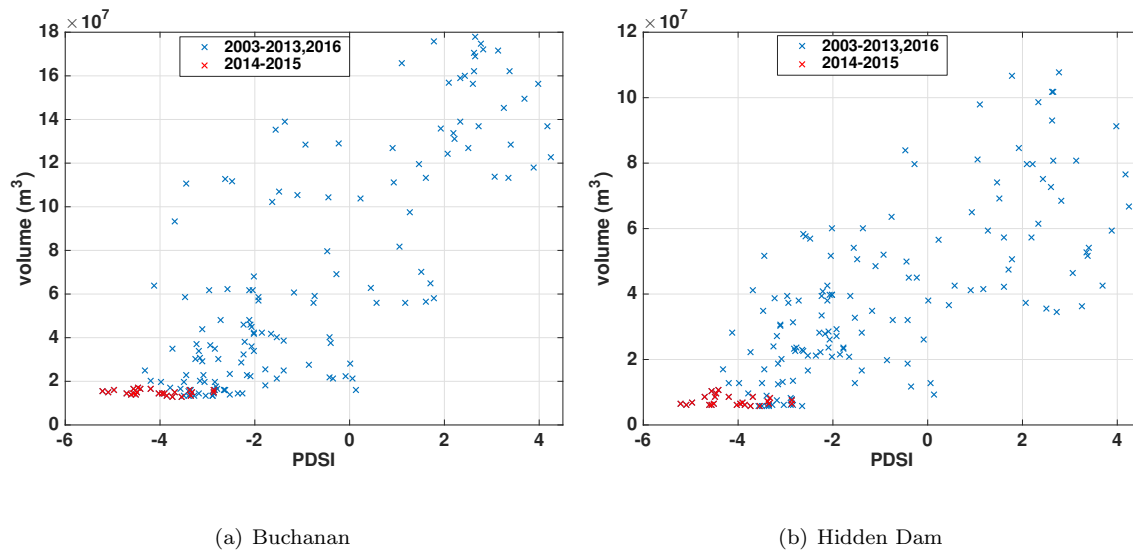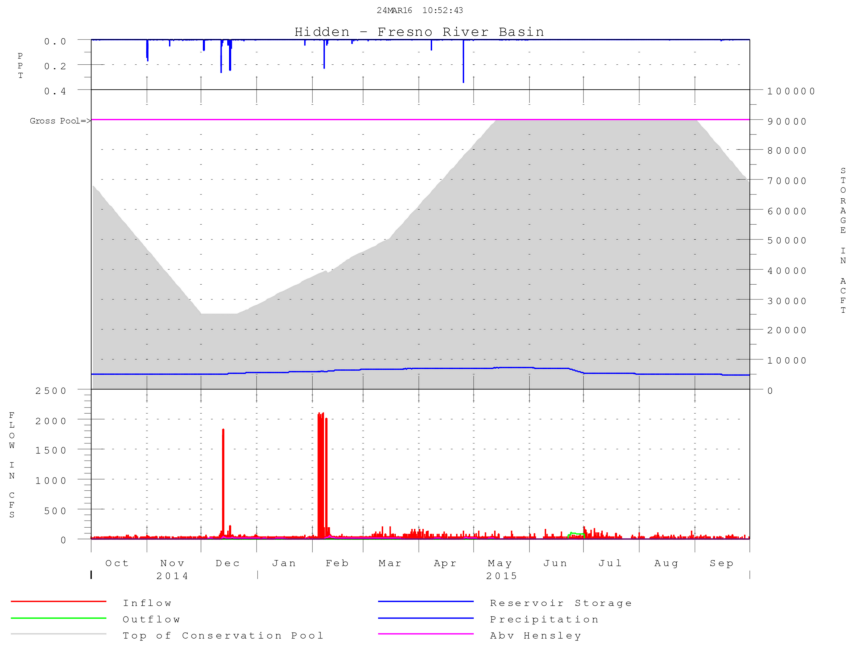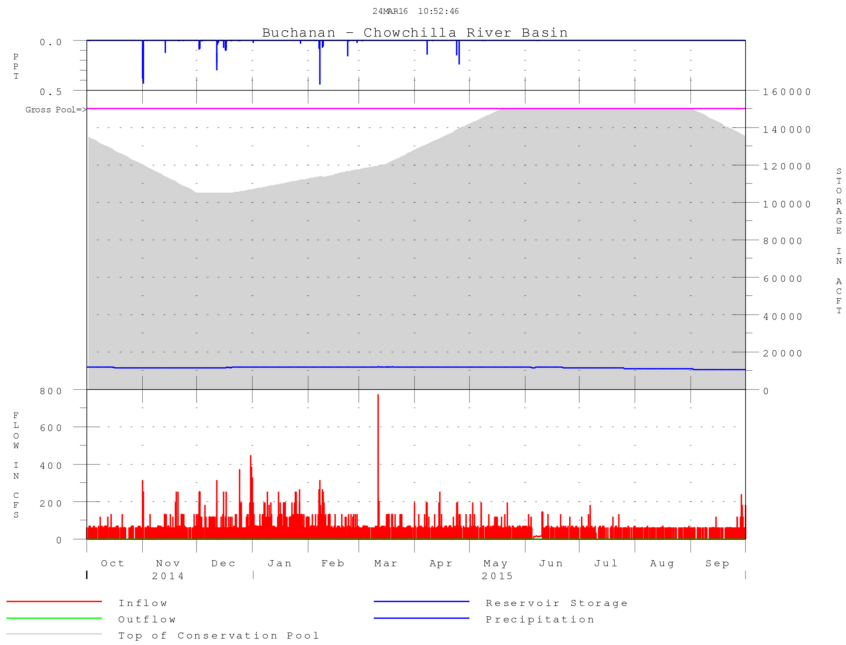(a) Buchanan                                    (b) Hidden Dam

Figure 5: PDSI vs reservoir levels for the Buchanan and Hidden Dam reservoirs during the period of study (i.e. January 2003 to November 2016). Notice a positive correlation between PDSI and the reservoir volumes: smaller values of PDSI generally lead to lower reservoir volumes. During the 2014-2015 drought period (shown in red), the correlation is substantially reduced as a result of stringent management efforts.

(a) Hidden Dam, 2014-2015



(b) Buchanan, 2014-2015

Figure 6: Inflows, outflows, precipitation, and water levels for the Buchanan and Hidden Dam reservoirs during the extreme drought period of 2014-2015. Notice that there was little precipitation, leading to marginal inflow of water into each reservoirs. Due to heavy management, there was little to no outflow of water from these reservoirs, preventing them from running dry. These figures are obtained from the Sacramento District Water Control Data System at http://www.spk-wc.usace.army.mil/plots/california.html.