

Runoff Modeling in Drainage Networks Using Probabilistic Graphs

E. Verstegen

Runoff Modeling in Drainage Networks Using Probabilistic Graphs

by

Emiel Verstegen

to obtain the degree of Master of Science
at the Delft University of Technology,
to be defended publicly on Friday August 5, 2016 at 13:30.

Student number: 1308408

Thesis committee:	Dr. ir. G. H. W. Schoups	TU Delft
	Prof. dr. ir. N. C. van de Giesen	TU Delft
	Prof. dr. ir. A. W. Heemink	TU Delft

An electronic version of this thesis is available at <http://repository.tudelft.nl/>.

Abstract

Conventional hydrological models use a deterministic approach. One could think of it like a black box, having an input, parameters, relations and an output. The parameters are calibrated by comparing the model output with observations of the system response (for example river runoff). When assessing the uncertainty in models the focus is often on the parameter uncertainty and all other variables are considered to be true. Not taking into account all sources of uncertainty, results in unreliable knowledge about the parameter uncertainty.

A method to include more sources of uncertainty is applying a probabilistic model. Using this model, all variables are described by probability distributions. In a spatially distributed model this allows for a spatial estimation of variables and their uncertainty in all model components. Due to the large amount of variables in a distributed model, the complexity of the exact solution of the probabilistic model increases. To be able to efficiently calculate an approximate solution, the probabilistic model is factorized and structured in a factor graph, a form of a probabilistic graph. A factor graph contains factors and variables, where the factors represent the relations between variables and physical knowledge while the variables represent the belief about the data. A factor graph is bipartite, which means that factors are only connected with variables and vice versa. Information propagates through the graph using message passing. This is a process where a factor updates each connected variable, based on a function of the other connected variables. If the graph has a tree structure, message passing starts at the highest level of the tree, progressing downward. This ensures that all information reaches the root of the tree. The process is then reversed in an upward sweep of message passing, propagating all gained knowledge also upstream. When approximations are used or when the tree contains cycles, multiple iterations (downward and upward sweeps) are needed for the variables to converge to their final value. The result is a posterior distribution of every variable in every cell.

In this research, a probabilistic graph is applied on a distributed runoff accumulation model. In each cell of the model a local runoff is calculated based on the precipitation, evaporation and an unknown bias term (initiated by bias parameters), which are all represented by a Gaussian distribution. Using flow paths, derived from a Digital Elevation Model, the local runoff is accumulated into accumulated runoff. A physical positivity constraint is added to the accumulated runoff and forcing data to prevent negative values. Multiple runoff observations (Gaussian distributed) are added resulting in spatial estimations of accumulated runoff, local runoff and bias, and an updated belief about the precipitation and evaporation. The bias parameters which initiate the bias in each cell contain uncertainty as well, allowing the parameters to be updated given the data received from the model. After the solution has converged, each dataset has been updated using the model structure (physical knowledge and constraints) and the prior knowledge from all other data sets.

By looking at the spatial distribution of the bias, conclusions can be drawn about the quality of the data and the water balance as a representation of the hydrological processes. Areas with a high posterior bias either have a mismatch between forcing data and runoff observations, the water balance does not represent the reality well, or the data conflicts with the positivity constraints. The model is applied on the Volta basins, where 3 areas are identified where the bias is higher than in other areas of the basin, mainly influenced by

the positivity constraints. These regions are next to a large river, next to lake Volta and the delta area at the mouth of the river. They are characterized by a negative prior local and accumulated runoff, indicating net evaporation while there is no water available. It is very likely that not the forcing data is faulty in these areas, but that an important hydrological process has not been incorporated. Given the large water bodies in the vicinity, ground water flow from the water body to the constrained cells is most probably the important hydrological process which is missing in the model. By applying probabilistic graphs on other (more complex) hydrological models allows for a better spatial estimation of both the variable values and their uncertainty, as well as a spatial evaluation of the performance of the model structure.

Preface

This thesis is the product of a very educative period of graduation. Combining my master track in water management and my background and interest in computer science, I decided to spend my graduation rethinking the way modeling was taught me. This process led me to better understand the concept of uncertainty in data, uncertainty propagation, Bayesian statistics and message passing. As knowledge does not always come easy, I thank my supervisor Gerrit Schoups for the many hours of discussion about the subject which taught me a lot. Furthermore I would like to thank the other committee members for their feedback during my presentations.

I would like to use this opportunity to thank my parents, who always showed support in the decisions I've made, bringing me to where I am today. I owe thanks to my girlfriend who, next to her loving support, challenged me to think with a different viewpoint and for the feedback on my report and presentations. Finally I thank my colleague graduation students for a nice atmosphere to work in, and the regular game of table tennis.

Emiel Verstegen,
Delft, July 22, 2016

Contents

Abstract	iii
Preface	v
List of Figures	ix
List of Algorithms	xi
List of Tables	xi
Nomenclature	xiii
1 Introduction	1
2 A probabilistic graphical model of spatially distributed runoff	5
2.1 Describing the probabilistic model	5
2.1.1 Deterministic model structure	5
2.1.2 Adding uncertainty to the model.	6
2.1.3 Describing the joint and posterior distribution	8
2.2 Translating into a graphical model	9
2.3 Solving the graphical model: calculating the marginal posterior	10
2.3.1 Message passing	11
2.3.2 Calculating a marginal	12
2.3.3 Calculating the incoming message.	12
2.3.4 Calculating the outgoing message	13
2.3.5 Scheduling factors: exploiting the tree structure	15
3 Model application using test data	17
3.1 Data description	18
3.2 Effect of different model setups	19
3.2.1 Adding observations	19
3.2.2 Adding constraints.	21
3.2.3 Adding a bias term	22
3.2.4 Effects of combining components	23
3.3 Influence of prior model parameters	23
3.3.1 Influence of data uncertainty	23
3.3.2 Influence of bias parameters.	24
3.4 Effect of model on posterior bias parameters	25
3.5 Convergence of different model structures	25
3.5.1 Factor scheduling	26
3.5.2 Positivity constraints.	26
3.5.3 Bias with parameter uncertainty.	27

4	Model application using real data	29
4.1	Data description	29
4.1.1	Data sources	29
4.1.2	Data processing	30
4.2	Uncertainty quantification	32
4.2.1	Precipitation uncertainty	32
4.2.2	Evaporation uncertainty	33
4.2.3	Runoff observations uncertainty	33
4.3	Implementation of the model	34
4.3.1	Prior	34
4.3.2	Model 2	35
4.3.3	Model 6	37
4.3.4	Model 7	37
4.4	Bias uncertainty influence	40
4.4.1	Model without positivity constraints	40
4.4.2	Model with positivity constraints	40
5	Conclusions	43
6	Recommendations	45
	References	47
A	Basin description	A1
B	Different model setups: Factor Graphs	B1
C	Different model setups: Results of testdata	C1
C.1	Model 1	C2
C.2	Model 2	C3
C.3	Model 3	C4
C.4	Model 4	C5
C.5	Model 5	C6
C.6	Model 6	C7
C.7	Model 7	C8
D	Results on Volta: figures	D1

List of Figures

1.1	The four components of a model	1
2.1	Translation from Digital Elevation Model (DEM) to flow path.	6
2.2	An example factor graph.	9
2.3	One grid cell in the factor graph of the distributed model	10
2.4	Message passing between a factor and a variable	11
2.5	Constraining a variable to be positive	14
3.1	Drainage direction of all cells within the test data grid	17
3.2	Influenced areas due to observations	19
3.3	Effects of adding an observation on the accumulated runoff	20
3.4	Effects of adding an observation on the local runoff	20
3.5	Effects of adding the constraints on the accumulated and local runoff	21
3.6	Mean value of the local bias and accumulated runoff, without parameter uncertainty	22
3.7	Mean value of the local bias and accumulated runoff, with parameter uncertainty (model 6).	22
3.8	Mean value of the local bias and accumulated runoff, with parameter uncertainty and constraints.	23
3.9	Mean value of the local bias under changing bias precision	24
3.10	Posterior value of the mean bias precision after each iteration for the closed (implied bias of zero) and non closed water balance.	25
3.11	Convergence of the unscheduled (Infer.NET model) and the scheduled (manual model 7), shown by plotting the mean accumulated runoff at the root cell.	26
3.12	Applying a positivity constraint multiple iterations on the same variable.	27
3.13	Difference in outlet runoff compared to the previous iteration (model 3).	27
3.14	Difference in outlet runoff compared to the previous iteration (model 6).	27
4.1	Re-sampling grid cells using the nearest neighbor algorithm	31
4.2	Calculated stream does not always match the river	32
4.3	River cross-over due to DEM uncertainties	32
4.4	Flow accumulation in a lake	32
4.5	Prior mean values for accumulated runoff, local runoff, precipitation and evaporation	35
4.6	Difference in mean value with respect to the prior for all the variables after adding observations (model 2)	36
4.7	Difference in mean value with respect to the prior for all the variables after adding observations and bias (model 6)	38
4.8	Difference in mean value with respect to the prior for all the variables after adding observations, bias and positivity constraints (model 7)	39
4.9	Results of adding bias uncertainty and overfitting.	41
A.1	Geographical map of the Volta basin	A1

B.1	Different factor graph structures represent all informed models (model 2-7)	B2
C.1	Model 1, uninformed, unconstrained, unbiased	C2
C.2	Model 2, informed, unconstrained, unbiased	C3
C.3	Model 3, informed, constrained, unbiased	C4
C.4	Model 4, informed, unconstrained, bias without parameter uncertainty	C5
C.5	Model 5, informed, constrained, bias without parameter uncertainty	C6
C.6	Model 6, informed, unconstrained, bias with parameter uncertainty	C7
C.7	Model 7, informed, constrained, bias with parameter uncertainty	C8
D.1	Mean accumulated runoff	D2
D.2	Standard deviation accumulated runoff	D3
D.3	Mean local runoff	D4
D.4	Standard deviation local runoff	D5
D.5	Mean precipitation runoff	D6
D.6	Standard deviation precipitation	D7
D.7	Mean evaporation	D8
D.8	Standard deviation evaporation runoff	D9
D.9	Bias Mean	D10
D.10	Bias SD	D10

List of Algorithms

1	Message passing (without making use of the model structure)	11
2	Function UpdateMarginals	11
3	Message Passing (with scheduling)	15

List of Tables

2.1	Infer.NET methods used to calculate outgoing messages	16
3.1	Prior values [m ³ /s] for precipitation, evaporation and local runoff	18
4.1	Model parameters	34
A.1	General information of the Volta basin	A2
A.2	Discharge at the mouth of the river by combining different data products	A2

Nomenclature

Symbol	Description	Unit ¹
General		
μ	Mean	[x]
σ	Standard Deviation (SD)	[x]
σ^2	Variance	[x] ²
$\tau = 1/\sigma^2$	Precision	[x] ⁻²
$CV = \sigma/\mu $	Coefficient of Variation	[-]
Model variables		
P_i	Precipitation in cell i (Gaussian distribution)	[m ³ /s]
E_i	Evaporation in cell i (Gaussian distribution)	[m ³ /s]
B_i	Local bias in cell i (Gaussian distribution)	[m ³ /s]
$Q_{loc,i}$	Local runoff in cell i (Gaussian distribution)	[m ³ /s]
$Q_{acc,i}$	Accumulated runoff in cell i (Gaussian distribution)	[m ³ /s]
Bias parameters		
μ_B	Bias mean (Gaussian distribution)	[mm/y]
τ_B	Bias precision (Gamma distribution)	[(mm/y) ⁻²]
Model parameters		
CV_P	Coefficient of Variation of precipitation data	[-]
CV_E	Coefficient of Variation of evaporation data	[-]
CV_{Obs}	Coefficient of Variation of runoff observation data	[-]
μ_μ	Mean of bias mean parameter μ_B	[mm/y]
σ_μ^2	Variance of bias mean parameter μ_B	[(mm/y) ²]
μ_τ	Mean of bias precision parameter τ_B	[(mm/y) ⁻²]
σ_τ^2	Variance of bias precision parameter τ_B	[(mm/y) ⁻⁴]
Other		
N	Set of all cells in the model	{...}
$M \subseteq N$	Set of all cells containing runoff observations	{...}
$u(i) \subset N$	Set of upstream neighbor cells of cell i	{...}
\mathcal{N}	Gaussian distribution given μ and σ^2	
\mathcal{G}	Gamma distribution given α and β	
$\alpha = \mu^2/\sigma^2$	Shape parameter	[-]
$\beta = \mu/\sigma^2$	Rate parameter	[x] ⁻¹
A_i	Surface area of cell i	[m ²]
C	Conversion factor 1: $3.17e^{-11}$	[$\frac{m}{mm} \frac{y}{s}$]

¹Unit [x] represents the unit of the variable, described by the mean, standard deviation, etc.

Introduction

When looking back on research of the last 50 years in Water Resources, it is clear that the topics have changed over time. Before the topics were mainly focusing on the hydrologic mechanisms and processes, nowadays the focus lies more on *uncertainty in modeling* and *data assimilation* (Rajaram et al. 2015). At the same time the amount of data sources is growing. Nowadays, data from distributed measurement networks, radar and satellites is published with a higher temporal and spatial resolution. The trend towards sharing data (open data) is also growing accordingly. It is of importance that the hydro-informatics follow these trends and make use of the available data. By combining many different data sources a higher temporal and spatial resolution can be achieved with an uncertainty that can be minimized (Chen & Han 2016).

Current practice in hydrological modeling

Current models can often be considered as a black box with four components. These components are 1) *forcing data*, 2) *a set of deterministic relations*, 3) *parameters*, and 4) *a model output*.

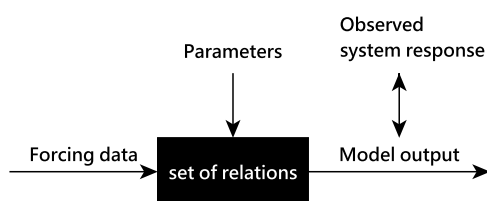


Figure 1.1: The four components of a model

Each of these components has an uncertainty connected with it. The *forcing data* is often an estimate based on indirect measurements or an extrapolation over space and/or time. In the case of rainfall, satellites provide estimates based on other physical phenomenon which have a relation with the rainfall, for example water vapor in the atmosphere or the cloud top temperature. These estimates are subject to uncertainties (AghaKouchak et al. 2009).

The *set of deterministic relations* that is used to calculate the model output is almost always a simplification of the reality, thus induce uncertainty in the model predictions (Refsgaard et al. 2006). An example of model uncertainty is neglecting water withdrawal for industry or groundwater flows which are not modeled.

The *model output* is then compared with the *observed system response* (often river runoff) in order to calibrate the model parameters. A common method to measure the (daily) river runoff is by using rating

curves, a relation between the water level and the discharge. Uncertainty in rating curves have multiple sources. The rating curve itself has uncertainties as the measurements of water level and discharge to create the rating curve are prone to errors. The morphology of the river can change, out-dating the rating curve. Often rating curves are created without taking extreme high or low runoffs into account. (Tomkins 2014, Di Baldassarre & Montanari 2009, Domeneghetti et al. 2012)

The *parameters* are entered into a model, based on measurements or expert knowledge. Often the parameters are calibrated based on the output of the model and an observed quantity. If uncertainties of other model components are not dealt with, all the discrepancy between the model and observation propagates into the parameter uncertainty (Moradkhani & Sorooshian 2009).

According to McMillan et al. (2011) and Kavetski et al. (2002), the uncertainty in forcing data, model structure, and observations is usually neglected during calibration and uncertainty assessment, and only parameter uncertainty is taken into account. McMillan further pleads for the use of error propagation of input data, claiming that no model can produce accurate prediction when forced with erroneous forcing data. One way to improve the model results is to reduce the error in the forcing data. Although the spatial and temporal resolution and coverage of the rainfall data have improved a lot over the last decades, uncertainty of this data will remain important in the close future due to the high variability of precipitation in both space and time (Kavetski et al. 2006).

Besides a better incorporation of uncertainty, Liu & Gupta (2007) plead for a better understanding, quantification and reduction of uncertainty.

Recent advances

One way of integrating uncertainty in a model is by defining the model with probabilistic/stochastic constraints, and the forcing data and observed system response with probability density functions. Recently, a conceptual study was conducting describing a spatially distributed model, accumulating runoff in a drainage network with a joint distribution (Schoups 2015). The joint distribution is a function dependent of all variables (3 million variables, due to the large amount of cells in the model). The posterior can be calculated by integrating the joint distribution over all the unknown variables, which is almost impossible considering the vast amount of variables. By factorizing the joint distribution, smaller problems are created which can be visualized by a factor graph. By solving these problems (factors) one by one an approximation of the posterior distribution of all variables can be achieved.

By implementing a hydrological model as a probabilistic graph it can 1) take uncertainty in both forcing data and system response into account, 2) combine different datasets to update knowledge of all of them, 3) show a spatial distribution of uncertainty and 4) calculate the posterior result with limited resources by exploiting the graph structure.

In this study the concept is applied to a real world example, extended with precipitation and evaporation data, runoff observations and physical constraints. It presents a the framework to build and solve a probabilistic hydrological model, incorporating multiple sources of uncertainty. The posterior distribution of the data is used to draw conclusions about the spatial estimate of the runoff and the spatial uncertainty of the model in terms of data and water balance. For this purpose the following research question is formed:

How to provide a *spatial estimate of the runoff* and the *uncertainty in the water balance* using a *probabilistic graph*?

Report outline

The probabilistic model structure which is used in this research and the translation to a graphical model is presented in Chapter 2. The model is then implemented using test data in Chapter 3 to analyze the influence of the different model components and parameters. In Chapter 4, the model is applied on a hydrological basin to research what the practical implications of the model are and how to draw conclusions from the model result. Conclusions regarding graphical modeling in general and the model result are presented in Chapter 5. Recommendations for future research are addressed in Chapter 6.

2

A probabilistic graphical model of spatially distributed runoff

In this chapter the basic model setup is introduced, and how the probabilistic component is added to the model. First the difference between deterministic models and probabilistic models is discussed in Section 2.1. The translation to factor graphs, which function as a graphical representation of the model, is outlined in Section 2.2. Finally, the mathematics and application of message passing and the algorithm for solving the model is explained in Section 2.3.

2.1. Describing the probabilistic model

Many hydrological models are based on *deterministic* relations and data. They assume the data to be true and the relations to be exact. *Probabilistic* models on the other hand, assume the model outcome to be a joint distribution over random variables. In this section the structure of the model and the probabilistic notation is described.

2.1.1. Deterministic model structure

The model setup is a distributed model. A distributed model divides the drainage basin in grid cells and calculates the runoff of each cell. Using a Digital Elevation Model (DEM), the flow direction is calculated to route the flow towards the mouth of the basin as is shown in Figure 2.1. Using this method the flows only converge, which results in a tree structured flow path (further explained in Chapter 4). In each cell of the grid, the local contribution to the flow is calculated using a water balance, which is described in Equation 2.1.

$$Q_{loc,i} = P_i - E_i + B_i \tag{2.1}$$

where P_i is the precipitation [m^3/s], E_i is the actual evaporation [m^3/s], B_i is a local bias term [m^3/s] accounting for the model structure uncertainty, and $Q_{loc,i}$ is the local runoff in cell $i \in N$ in [m^3/s], where $N = \{1, 2, \dots, n\}$ is the set of all cells in the model. By using long term averages for the forcing data and river runoff observations, change in storage can be neglected and a transport model does not have to be incorporated.

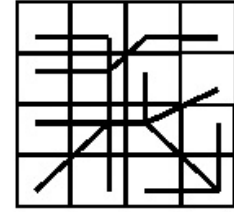
By means of the flow direction map the local runoff is accumulated downstream into the accumulated runoff, which is described by Equation 2.2.

$$Q_{acc,i} = Q_{loc,i} + \sum_{j \in u(i)} Q_{acc,j} \quad (2.2)$$

where $Q_{acc,i}$ is the accumulated runoff [m^3/s] flowing out of a cell, and $Q_{acc,j}$ is the inflow of accumulated runoff from upstream cells [m^3/s], where $j \in u(i)$ are the upstream neighbors of cell i .

45	38	44	45
40	34	50	60
58	31	30	53
50	45	32	22

DEM



Flow path

Figure 2.1: From Digital Elevation Model (DEM) to flow path. On the left the average elevation of each cell is shown. This is translated into a flowpath, shown on the right. Source: (Olivera et al. 1998)

The local bias term

The local bias term B_i in Equation 2.1 was added in order to represent (un)known processes that are not taken into account in the model structure. Examples of these processes that could cause a bias are:

- Lateral (ground)water flows, where the flow direction is other than indicated by the Digital Elevation Model
- Discharge or withdrawal of water for consumption, industry or irrigation
- Divergence of flows (human influence like canals or natural divergence in delta areas)
- Long term change of storage due to construction of dams

Positivity constraints

More physical knowledge is added to the model using constraints. In this model the precipitation P_i , evaporation E_i and the accumulated runoff $Q_{acc,i}$ are constrained to be positive. The first two are constrained because 'negative precipitation' can be considered evaporation and 'negative evaporation' can be considered precipitation. The accumulated runoff is constrained to be positive because a negative value would indicate a river flowing in the upstream direction. When the accumulated runoff is zero there is no water available to withdraw from the flow.

2.1.2. Adding uncertainty to the model

In the model presented in this research, variables and relations are not described by deterministic variables but by stochastic variables. This way the uncertainty of each component of the model can be incorporated. This is done by defining each variable as a probability distribution, which follows from the Bayesian approach. Most variables are incorporated using Gaussian distributions to describe the variables. The mean μ of the distribution represents the actual measured or estimated data point and the standard deviation σ is an expression for the uncertainty. The Gaussian distribution was chosen for the following reasons:

- The *central limit theorem* states the Gaussian distribution is considered a good representation of a random variable when this variable is an average of many smaller random variables, no matter their distribution (Castrup 2001, Lyon 2014). In this research the variables are long term averages of observations, which makes the Gaussian distribution suitable regardless of the distribution of the uncertainty on a single observation.

- The Gaussian distribution is closed under convolution (Vinga & Almeida 2004, supplementary material). This makes it possible to sum variables with the result being a Gaussian distribution. No approximations in the flow accumulation are needed.
- The product and ratio of two Gaussian Probability Density Functions (PDF's) is proportional to another Gaussian PDF (Bromiley 2003). These operations are needed to add or subtract information to the model, further explained in Section 2.3.

All uncertainty in the prior data is described relative to the magnitude of the mean, using the Coefficient of Variation (CV) [-], which is defined as:

$$CV = \frac{\sigma}{\mu} \quad (2.3)$$

Forcing data uncertainty

The *forcing data* consists of a precipitation variable P_i and evaporation variable E_i in each cell. These variables are based on the observed precipitation $\mu_{P,i}$ [m^3/s] and evaporation $\mu_{E,i}$ [m^3/s]. The *model parameters* CV_P and CV_E determine the standard deviation of the Gaussian distribution.

$$\begin{aligned} \text{Precipitation: } P_i &\sim \mathcal{N}(\mu_{P,i}, \sigma_{P,i}^2) \\ \text{where } \sigma_{P,i} &= CV_P * \mu_{P,i} \end{aligned} \quad (2.4)$$

$$\begin{aligned} \text{Evaporation: } E_i &\sim \mathcal{N}(\mu_{E,i}, \sigma_{E,i}^2) \\ \text{where } \sigma_{E,i} &= CV_E * \mu_{E,i} \end{aligned} \quad (2.5)$$

Model structure uncertainty

The *relations* that the model uses in order to calculate the output (Equation 2.1) can also contain uncertainty. Due to the simplification of the reality, not all hydrological processes in the water balance are captured. Therefore a *local bias* is added to the water balance. The prior local bias will be a Gaussian distribution given the mean μ_B [mm/y] and precision τ_B [$(\text{mm}/\text{y})^{-2}$], which are called the *bias parameters*. These variables are not local (for one cell only) but apply to all the cells in the model.

$$\begin{aligned} \text{Bias: } B_i &\sim \mathcal{N}(\mu_B, \sigma_B^2) \\ \text{where } \sigma_B^2 &= \frac{1}{\tau_B} \end{aligned} \quad (2.6)$$

Parameter uncertainty

Each local bias is described by the *bias parameters*: μ_B and τ_B . These parameters would normally be chosen using expert knowledge, and then optimized for better results. In this research the uncertainty of these parameters is added by describing the parameters with probability distributions. The bias mean is assumed to have a Gaussian distribution (Equation 2.7), described by two *model parameters*: μ_μ [mm/y] and σ_μ^2 [$(\text{mm}/\text{y})^2$]. The bias precision (reciprocal of the variance) is assumed to have a Gamma distribution (Equation 2.8). Using two *model parameters*, μ_τ [$(\text{mm}/\text{y})^{-2}$] and σ_τ^2 [$(\text{mm}/\text{y})^{-4}$], the shape parameter α_τ [-] and rate parameter β_τ [$(\text{mm}/\text{y})^2$] of the Gamma distribution are calculated.

$$\text{Bias mean: } \mu_B \sim \mathcal{N}(\mu_\mu, \sigma_\mu^2) \quad (2.7)$$

$$\text{Bias precision: } \tau_B \sim \mathcal{G}(\alpha_\tau, \beta_\tau) \quad (2.8)$$

$$\text{where } \alpha_\tau = \frac{\mu_\tau^2}{\sigma_\tau^2} \text{ and } \beta_\tau = \frac{\mu_\tau}{\sigma_\tau^2} \quad (\text{Dekking et al. 2005})$$

Because the bias parameters are uncertain, information from the model can influence the posterior distribution of the bias parameters. Information from one part of the drainage basin influences the parameters after which the parameters influence the other part of the basin. Because of this uncertainty in parameters a spatial correlation of local biases B_i is established. If both σ_μ and σ_τ approach zero, there is no uncertainty in the parameters and the parameters will not be updated by the information from the model. The prior distribution of the local bias B_i is in this case defined by Equation 2.9.

$$B_i \sim \mathcal{N}(\mu_\mu, 1/\mu_\tau) \quad \text{when } \sigma_\mu \rightarrow 0 \text{ and } \sigma_\tau \rightarrow 0 \quad (2.9)$$

System response uncertainty

The output of this model is compared with the *system response* (in this model the observed accumulated runoff). This quantity has an uncertainty which is also described by a Gaussian distribution. The observed value is the mean of the distribution and the uncertainty is defined as the standard deviation. The standard deviation is calculated relative to the mean value, using the *model parameter* CV_{obs} (Equation 2.10).

$$\begin{aligned} \text{Runoff observations: } Q_{acc,k} &\sim \mathcal{N}(\mu_{acc,k}, \sigma_{acc,k}^2) \\ \text{where } \sigma_{acc,k} &= CV_{obs} * \mu_{acc,k} \end{aligned} \quad (2.10)$$

where k is a cell in the subset of set M , which contains cells that have an observed accumulated runoff.

2.1.3. Describing the joint and posterior distribution

The variables can be divided in observed variables (the data, Equation 2.11) and the unknown variables (Equation 2.12):

$$\text{Observed variables: } \mathbf{D} = \{\boldsymbol{\mu}_P, \boldsymbol{\sigma}_P, \boldsymbol{\mu}_E, \boldsymbol{\sigma}_E, \mu_\mu, \sigma_\mu, \alpha_\tau, \beta_\tau, \boldsymbol{\mu}_{acc}, \boldsymbol{\sigma}_{acc}\} \quad (2.11)$$

$$\text{Unknown variables: } \mathbf{X} = \{\mathbf{P}, \mathbf{E}, \mathbf{B}, \mu_B, \tau_B, \mathbf{Q}_{loc}, \mathbf{Q}_{acc}\} \quad (2.12)$$

where the bold symbols are vectors of variables over all cells in the model (except $\boldsymbol{\mu}_{acc}$ and $\boldsymbol{\sigma}_{acc}$, which are vectors over runoff observations).

The *joint distribution* $p(\mathbf{X}, \mathbf{D})$ can be written as a product of factorized relations (factors) between the variables (Koller & Friedman 2009, pp. 50), assuming the variables are independent. Apart from the four sources of uncertainty and the local and accumulating water balances, the positivity constraints are represented by a factor. The factorized *joint distribution* is defined as:

$$\begin{aligned} p(\mathbf{X}, \mathbf{D}) &= p(\mathbf{P}, \boldsymbol{\mu}_P, \boldsymbol{\sigma}_P, \mathbf{E}, \boldsymbol{\mu}_E, \boldsymbol{\sigma}_E, \mathbf{B}, \mu_B, \sigma_\mu, \tau_B, \alpha_\tau, \beta_\tau, \mathbf{Q}_{loc}, \mathbf{Q}_{acc}, \boldsymbol{\mu}_{acc}, \boldsymbol{\sigma}_{acc}) = \\ &\quad \mathcal{N}(\mu_B | \mu_\mu, \sigma_\mu^2) \mathcal{G}(\tau_B | \alpha_\tau, \beta_\tau) \times \quad \text{parameter uncertainty} \\ &\quad \prod_{i \in N} \left[\mathcal{N}(P_i | \mu_{P,i}, \sigma_{P,i}^2) \mathcal{N}(E_i | \mu_{E,i}, \sigma_{E,i}^2) \times \quad \text{forcing data uncertainty} \right. \\ &\quad \mathcal{N}(B_i | \mu_B, \sigma_B^2) \times \quad \text{model structure uncertainty} \\ &\quad \delta(Q_{loc,i} - P_i + E_i - B_i) \delta(Q_{acc,i} - Q_{loc,i} - \sum_{j \in u(i)} Q_{acc,j}) \times \quad \text{water balances} \\ &\quad \left. H(P_i) H(E_i) H(Q_{acc,i}) \right] \times \quad \text{positivity constraints} \\ &\quad \prod_{k \in M} \mathcal{N}(Q_{acc,k} | \mu_{acc,k}, \sigma_{acc,k}^2) \quad \text{system response uncertainty} \end{aligned} \quad (2.13)$$

The *dirac delta function* δ and the *unit step function* H describe respectively the water balances and positivity constraints.

The *posterior distribution* $p(\mathbf{X}|\mathbf{D})$ is a function of the *joint distribution* $p(\mathbf{X}, \mathbf{D})$ divided by the *marginal likelihood* of the data $p(\mathbf{D})$:

$$p(\mathbf{X}|\mathbf{D}) = \frac{p(\mathbf{X}, \mathbf{D})}{p(\mathbf{D})} \propto p(\mathbf{X}, \mathbf{D}) \quad (2.14)$$

or in terms of the variables of the model:

$$p(\mathbf{P}, \mathbf{E}, \mathbf{B}, \mu_B, \tau_B, \mathbf{Q}_{loc}, \mathbf{Q}_{acc} | \boldsymbol{\mu}_P, \boldsymbol{\sigma}_P, \boldsymbol{\mu}_E, \boldsymbol{\sigma}_E, \mu_\mu, \sigma_\mu, \alpha_\tau, \beta_\tau, \boldsymbol{\mu}_{acc}, \boldsymbol{\sigma}_{acc}) = \frac{p(\mathbf{P}, \boldsymbol{\mu}_P, \boldsymbol{\sigma}_P, \mathbf{E}, \boldsymbol{\mu}_E, \boldsymbol{\sigma}_E, \mathbf{B}, \mu_B, \mu_\mu, \sigma_\mu, \tau_B, \alpha_\tau, \beta_\tau, \mathbf{Q}_{loc}, \mathbf{Q}_{acc}, \boldsymbol{\mu}_{acc}, \boldsymbol{\sigma}_{acc})}{p(\boldsymbol{\mu}_P, \boldsymbol{\sigma}_P, \boldsymbol{\mu}_E, \boldsymbol{\sigma}_E, \mu_\mu, \sigma_\mu, \alpha_\tau, \beta_\tau, \boldsymbol{\mu}_{acc}, \boldsymbol{\sigma}_{acc})} \quad (2.15)$$

The model results is the *marginal posterior* of each variable. The marginal posterior is a representation of a variable's value and uncertainty, in terms of mean and variance, after incorporating all the information from the model and the other variables. Let Y be the variable of interest, the marginal posterior of this variable can then be written as the integral over all the variables except Y :

$$p(Y|\mathbf{D}) = \int \dots \int_{x \in \mathbf{X} \setminus Y} p(\mathbf{X}, \mathbf{D}) dx \quad (2.16)$$

The model as defined by the joint distribution in Equation 2.13 has 5 unknown variables per cell. With an increasing number of cells, the complexity of this integral increases drastically. In order to solve the model efficiently with a large number of cells, another method is used to approximate the marginal posterior.

2.2. Translating into a graphical model

When the *posterior* distribution becomes very complex as the number of variables increase, *graphical models* can be used in order to exploit the structure within a complex distribution (Koller & Friedman 2009, pp. 3). There are two types of graphical models, Bayesian and Markov networks, respectively directed and undirected networks. A *factor graph*, used in this research, is capable of representing both those networks. In this study a mixture between directed and undirected dependencies is used.

The factor graph notation is used to transfer Equation 2.13 to a graph which is solved step by step. Frey et al. (1998) defined a factor graph as “a bipartite graph that expresses how a global function of several variables factors into a product of local functions”.

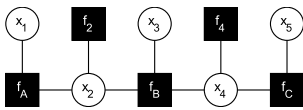


Figure 2.2: An example factor graph. Source: Frey et al. (1998)

An example of a factor graph is given in Figure 2.2. This example contains 5 variables $x_1 \dots x_5$ and 5 factors f_A, f_B, f_C, f_2, f_4 , connected with *edges*. The graph is *bipartite* which means that every variable is only connected with factors, and every factor is only connected with variables. The variables are given a prior probability distribution and the factors represent a local relation or function.

The global function illustrated in Figure 2.2 is:

$$g(x_1, x_2, x_3, x_4, x_5) = f_A(x_1, x_2) f_B(x_2, x_3, x_4) f_C(x_4, x_5) f_2(x_2) f_4(x_4) \quad (2.17)$$

where the value of the set of variables is the product of functions, depending on the set of variables.

All the components from Equation 2.13 are structured in a factor graph as can be seen in Figure 2.3, where the graph of one cell of the model is displayed. The circle nodes in the figure are the variables from the joint distribution. The grey circles are the observed variables (data \mathbf{D}), while the white circles are the unknown variables (\mathbf{X}) of which we want to know the posterior distribution. The factorized relations of Equation 2.13 are represented in the graph as squares. In the following sections all factors of the model are covered.

Gaussian from mean and variance: f_P , f_E and f_{obs}

These factors initialize variables to have a Gaussian distribution, derived from the deterministic variables mean μ and variance σ^2 . Every cell has an observed precipitation and evaporation, but only a couple of cells (in subset M) contain an observation of accumulated runoff.

Water balances: f_{loc} and f_{acc}

The water balance factor f_{loc} implements Equation 2.1, calculating the local runoff $Q_{loc,i}$ from the forcing variables P_i and E_i and the bias B_i . The factor also works in the other direction, calculating a new value for P_i , E_i or B_i when new information about the local runoff is received.

The factor f_{acc} implements Equation 2.2, calculating the outgoing accumulated runoff $Q_{acc,i}$ as a function of the incoming accumulated runoff from upstream cells $Q_{acc,j}$ and the local runoff $Q_{loc,i}$ from within the cell. Just as with factor f_{loc} , the factor f_{acc} can calculate a new value for all connected (unknown) variables.

Constraints: C_P , C_E , C_{acc}

The constraining factors assign a new value to the connected variable, as a function of the variable itself. If the connected variable is well above zero (positive mean and low CV), the new value will be equal the old value and the constraint has no influence. The exact constraining procedure is discussed in Section 2.3.4.

Bias: f_μ , f_τ and f_B

The factors f_μ and f_τ initialize the variables μ_B and τ_B to have respectively a Gaussian and a Gamma distribution. The factor f_B updates the local bias B_i based on the variables μ_B and τ_B . At the same time factor updates the parameters, based on the local bias.

The bias parameters μ_B and τ_B are defined in mm/y and $(mm/y)^{-2}$ respectively, not to be influenced by the surface area of a cell. The data used later on in Chapter 4 has a resolution in degree latitude and longitude, resulting a different surface area of a cell depending on the position of the cell. This means that the values of μ_B and τ_B are converted to m^3/s before applying factor f_B .

2.3. Solving the graphical model: calculating the marginal posterior

As described in the previous section, the posterior result of the graph can be calculated by solving it factor by factor. This section describes how each factor can be solved using *Message Passing*.

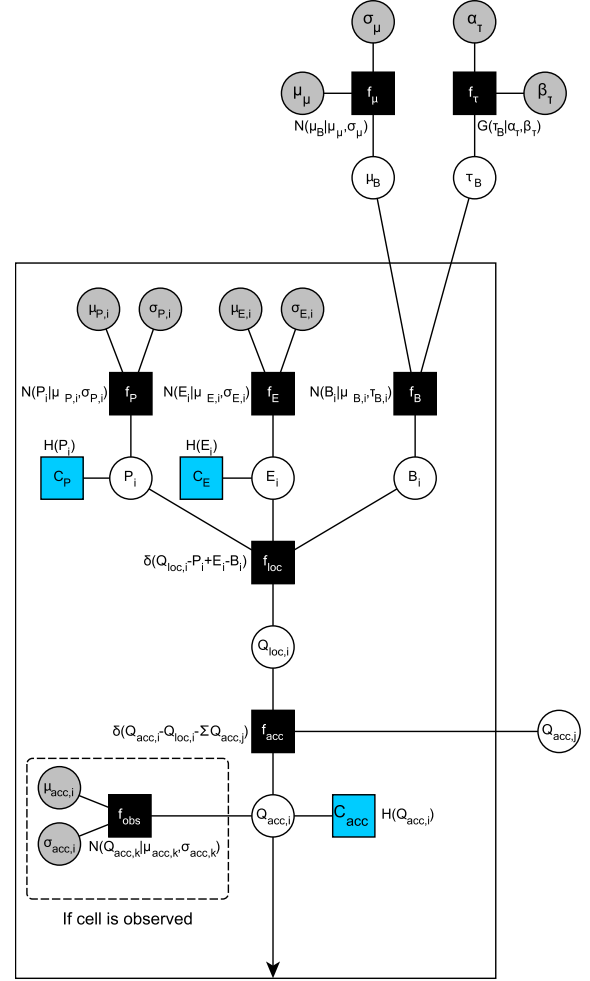


Figure 2.3: One grid cell in the factor graph of the distributed model. Observed variables (data) are shown as grey circles, the unknown variables as white circles. The squares represent factors.

2.3.1. Message passing

The message passing algorithm (Algorithm 1) uses *Marginal Updating* (Algorithm 2), where the marginals of the variables connected with a factor are updated after every iteration. For the next sections, consider one factor f and its connected variable x , as depicted in Figure 2.4

The *outgoing* message $m_{f \rightarrow x}$ (from the factor to the variable) is a probability distribution which contains all the information from the right side of the graph. The outgoing message to variable x is a function of the incoming messages from all connected variables, except x itself ($m_{y_1 \rightarrow f} \dots m_{y_n \rightarrow f}$). The function calculating the outgoing message depends on the type of factor.

The *incoming* message $m_{x \rightarrow f}$ contains all the information from the left side of the graph and is equal to the variable marginal distribution divided by the outgoing message $m_{f \rightarrow x}$. Division removes information because the marginal distribution of the variable is the product of the incoming and outgoing message (information from left and right side of the graph). By dividing with the outgoing message the information of the right side of the graph is removed from the marginal, resulting in a distribution representing the knowledge from the left side of the graph.

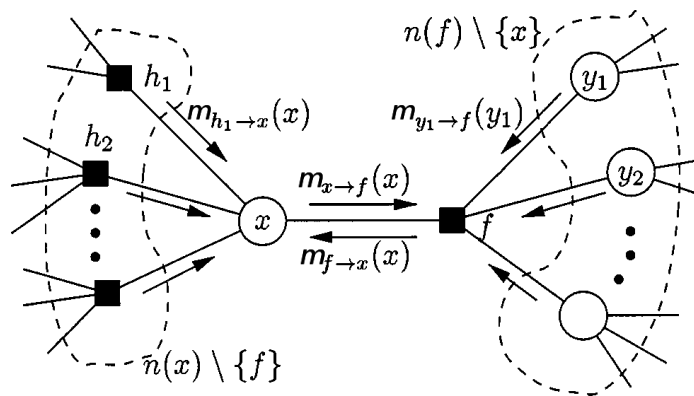


Figure 2.4: Message passing between a factor and a variable (variable x , factor f , messages m and neighbors of both nodes). Source: (Kschischang et al. 2001, p. 502), edited to comply with symbol conventions.

Algorithm 1: Message passing (without making use of the model structure)

```

for iteration  $i$  in NumberOfIterations do
  foreach Factor  $f$  in the factor graph do
     $\lfloor$  function UpdateMarginals( $f$ );

```

Algorithm 2: Function UpdateMarginals

```

Function UpdateMarginals (Factor  $f$ )
  1) foreach Connected variable  $v$  to factor  $f$  do
     $\lfloor$  incoming message = variable marginal / outgoing message;
  2) foreach Connected variable  $v$  to factor  $f$  do
     $\lfloor$  outgoing message = factor operation on (array of) incoming message(s);
  3) foreach Connected variable  $v$  to factor  $f$  do
     $\lfloor$  variable marginal = incoming message * outgoing message;

```

Initialization of the model

Before using the algorithm, initialization of the model needs to take place. This is done by setting values of:

- Unobserved *variables* (defined by their connected observed variables):
 - Define precipitation P_i and evaporation E_i marginals as a Gaussian distribution.
 - Set bias B_i as a Uniform distribution, μ_B as a Gaussian distribution and τ_B as a Gamma distribution.
 - Set local runoff $Q_{loc,i}$ marginals as a Uniform distribution.
 - Set observed accumulated runoff $Q_{acc,i}$ marginals as Gaussian distribution, in the cells without observation it is set to a Uniform distribution.
- Set all *messages* between factors and variables to a Uniform distribution.

2.3.2. Calculating a marginal

The marginal belief of variable x , written as $b(x)$, is calculated by multiplying the incoming message (information from the other factors connected to the variable) and the outgoing message (information from the factor in question).

$$b(x) = m_{x \rightarrow f} * m_{f \rightarrow x} \quad (2.18)$$

Given that the incoming and outgoing messages are Gaussians, the updated marginal will also be a Gaussian with the mean and variance (Bromiley 2003):

$$\begin{aligned} \frac{1}{\sigma_x^2} &= \frac{1}{\sigma_{x \rightarrow f}^2} + \frac{1}{\sigma_{f \rightarrow x}^2} \\ \mu_x &= \left(\frac{\mu_{x \rightarrow f}}{\sigma_{x \rightarrow f}^2} + \frac{\mu_{f \rightarrow x}}{\sigma_{f \rightarrow x}^2} \right) \sigma_x^2 \end{aligned} \quad (2.19)$$

The updated marginal belief has a mean which is an average of the mean of both messages, weighted on the precision (reciprocal of the variance). The precision of the updated marginal is a sum of the precisions of the messages. The next two sections describe how the incoming and outgoing messages are calculated.

2.3.3. Calculating the incoming message

A variable marginal is a product of all the information it received. The incoming message $m_{x \rightarrow f}$ is calculated by removing the outgoing message $m_{f \rightarrow x}$ from the marginal distribution $b(x)$.

$$m_{x \rightarrow f} = \frac{b(x)}{m_{f \rightarrow x}} \quad (2.20)$$

Given a Gaussian marginal distribution and outgoing message, the incoming message $m_{x \rightarrow f}$ is a Gaussian distribution as well. The mean and variance of the incoming message are calculated in Equation 2.21

$$\begin{aligned} \frac{1}{\sigma_{x \rightarrow f}^2} &= \frac{1}{\sigma_x^2} - \frac{1}{\sigma_{f \rightarrow x}^2} \\ \mu_{x \rightarrow f} &= \left(\frac{\mu_x}{\sigma_x^2} - \frac{\mu_{f \rightarrow x}}{\sigma_{f \rightarrow x}^2} \right) \sigma_{x \rightarrow f}^2 \end{aligned} \quad (2.21)$$

2.3.4. Calculating the outgoing message

The outgoing message $m_{f \rightarrow x}$ is calculated as a function of all incoming messages from the connected variables of f , except variable x itself (Figure 2.4, set $y \in n(f) \setminus \{x\}$). This method is called Expectation Propagation (Minka 2001a,b) and is a parametric approximation. This means that beliefs (marginals) are not exactly propagated, but their expectations (mean and variance) are. The general formula for calculating an outgoing message using Expectation Propagation is (Winn & Minka 2007):

$$m_{f \rightarrow x} = \frac{\text{proj} \left[m_{x \rightarrow f} \int_{y \in n(f) \setminus \{x\}} (f(y, x) m_{y \rightarrow f}) dy \right]}{m_{x \rightarrow f}} \quad (2.22)$$

where the outgoing message is a function of the factor, the incoming messages and a projection. If the formula between the projection brackets does not yield a proper Gaussian distribution, the projection creates a Gaussian distribution which matches the result in terms of expected mean and variance. This process is called M-projection (Koller & Friedman 2009, pp. 274).

In case the message does not need projection ($\text{proj}[x] = x$), exact beliefs are propagated instead of approximated beliefs (expectations). This specific case of Expectation Propagation is called Belief Propagation and the outgoing message to a variable $m_{f \rightarrow x}$ is not dependent anymore on the incoming message from that same variable $m_{x \rightarrow f}$, which is shown in Equation 2.23.

$$\begin{aligned} m_{f \rightarrow x} &= \frac{\text{proj} \left[m_{x \rightarrow f} \int_{y \in n(f) \setminus \{x\}} (f(y, x) m_{y \rightarrow f}) dy \right]}{m_{x \rightarrow f}} \\ &= \frac{m_{x \rightarrow f} \int_{y \in n(f) \setminus \{x\}} (f(y, x) m_{y \rightarrow f}) dy}{m_{x \rightarrow f}} \\ &= \int_{y \in n(f) \setminus \{x\}} f(y, x) m_{y \rightarrow f} dy \end{aligned} \quad (2.23)$$

In the model, outgoing messages are calculated using methods from the Infer.NET library, developed by Microsoft Research (Minka et al. 2014). Table 2.1, at the end of this chapter, gives an overview of the Infer.NET methods that are used for the calculation of outgoing messages. Next, two examples of these methods are shown, using Belief and Expectation propagation.

Example: local water balance f_{loc}

The calculation of the outgoing message sent from the local water balance factor f_{loc} to the variable Q_{loc} is one of the simplest, and uses Belief Propagation. This means that the message is only dependent on the incoming messages from the variables P , E and B , as is explained in Equation 2.23. Substituting the variables used in this research into the equation, it becomes a triple integral over the three incoming messages (\mathbf{m}), shown in Equation 2.24.

$$\begin{aligned} m_{f_{loc} \rightarrow Q_{loc}} &= \iiint \delta(Q_{loc} - P_i + E_i - B_i) \times m_{P_i \rightarrow f_{loc}} \times m_{E_i \rightarrow f_{loc}} \times m_{B_i \rightarrow f_{loc}} d\mathbf{m} \\ &\sim \mathcal{N}(\mu_{f_{loc} \rightarrow Q_{loc}}, \sigma_{f_{loc} \rightarrow Q_{loc}}^2) \end{aligned} \quad (2.24)$$

When solving this integral it shows that the result is a Gaussian distribution of which the mean follows the water balance, and all the variances are summed:

$$\begin{aligned} \mu_{f_{loc} \rightarrow Q_{loc}} &= \mu_{P_i \rightarrow f_{loc}} - \mu_{E_i \rightarrow f_{loc}} + \mu_{B_i \rightarrow f_{loc}} \\ \sigma_{f_{loc} \rightarrow Q_{loc}}^2 &= \sigma_{P_i \rightarrow f_{loc}}^2 + \sigma_{E_i \rightarrow f_{loc}}^2 + \sigma_{B_i \rightarrow f_{loc}}^2 \end{aligned}$$

This factor is an example of Belief Propagation as the parametric approximation equals the exact belief propagation.

Example: constraining the accumulated runoff C_{acc}

The constraining factor is only connected with one variable and is the most evident example of Expectation Propagation. It is responsible for updating the variable for it to be positive. This factor works different from the previous factor in the sense that it is only connected with one variable. This means that Equation 2.22 simplifies to:

$$m_{C_{acc} \rightarrow Q_{acc}} = \frac{\text{proj}[m_{Q_{acc} \rightarrow C_{acc}} \int H(Q_{acc}) dm_{Q_{acc} \rightarrow C_{acc}}]}{m_{Q_{acc} \rightarrow C_{acc}}} \quad (2.25)$$

$$b(Q_{acc}) = m_{C_{acc} \rightarrow Q_{acc}} * m_{Q_{acc} \rightarrow C_{acc}} = \text{proj}[m_{Q_{acc} \rightarrow C_{acc}} \int H(Q_{acc}) dm_{Q_{acc} \rightarrow C_{acc}}] \quad (2.26)$$

In practice this means that the updated belief is a projection (parametric approximation) of the truncated belief. The expected mean and variance for a truncated Gaussian are shown in Equation 2.27 (Greene 2003, pp. 759).

$$\begin{aligned} E[x|x > a] &= \mu + \sigma \lambda(\alpha) \\ \text{Var}[x|x > a] &= \sigma^2 (1 - \delta(\alpha)) \end{aligned} \quad (2.27)$$

where:

$$\alpha = \frac{a - \mu}{\sigma}, \quad a = 0$$

$$\lambda(\alpha) = \frac{\phi(\alpha)}{1 - \Phi(\alpha)}$$

$\phi(\alpha)$: standard normal density function

$\Phi(\alpha)$: standard normal cumulative distribution function

The updated marginal belief is a Gaussian distribution with always has a higher mean and lower variance than the original marginal belief, which can be observed in Figure 2.5. Another observation is that there is still a probability on a negative value of the variable. When using multiple iterations, the constraints will be applied multiple times. The consequence of this is further discussed in Section 3.5.

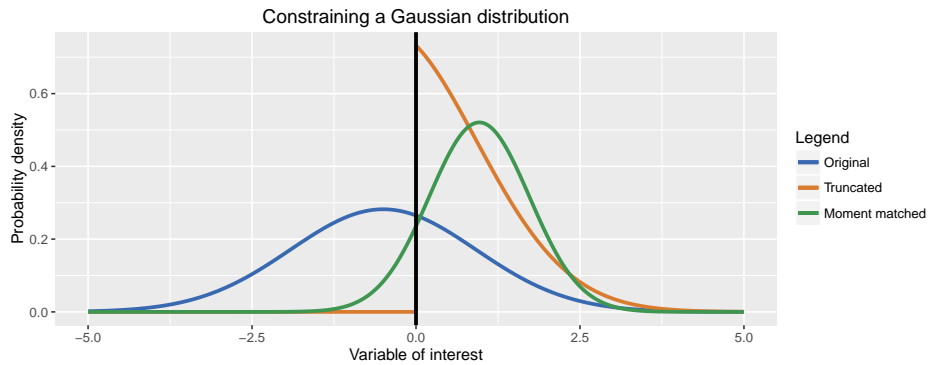


Figure 2.5: Constraining a variable to be positive

2.3.5. Scheduling factors: exploiting the tree structure

As mentioned before, a drainage network built from a DEM has a tree structure. Each cell in the network can have multiple upstream neighbors, but only one downstream neighbor. This means that eventually all cells flow towards one basin outlet (the root of the tree). The factor graph within one cell also has a tree structure (see Figure 2.3). All together, this results in the full factor graph having a tree structure as well.

Marginal Updating is a local operation which only updates the marginals of the connected variables of one factor at a time. That is why it is important to schedule the order in which factors update their variable marginals in order to minimize the amount of iterations needed to solve the model. Per iteration, two sweeps are needed to propagate the information throughout the model (Algorithm 3).

The *downwards sweep* starts updating the factors at the *leaf cells* (cells at the highest level of the flow tree, having no upstream neighbors). Once all these factors updated the connected variables, factors of one level lower in the tree are updated. This continues until the root of the tree (mouth of the river) has been reached, and all the marginals contain upstream information. Because down the line information is added due to observations and positivity constraints, an *upward sweep* is also needed. This sweep will start at the root of the tree with updating marginals, moving its way up the tree level by level. To determine the order in which the factors are solved, the Breadth-First algorithm is applied to assign every cell in the grid a distance from the root cell.

Multiple iterations of these sweeps have to be performed to converge to a solution. This is due to the fact that cycles in the tree have been introduced by means of bias parameters (variables which are connected to all cells) and approximations are made by re-projecting the constrained variables.

Algorithm 3: Message Passing (with scheduling)

```

for iteration  $i$  in NumberOfIterations do
  General pass:
  foreach Cell  $c$  in All Cells do
    UpdateMarginals(Factor Constraint Precipitation  $C_P$ );
    UpdateMarginals(Factor Constraint Evaporation  $C_E$ );
    UpdateMarginals(Factor Bias  $f_B$ );
    UpdateMarginals(Factor Local runoff  $f_{loc}$ );
  Downward pass:
  foreach Cell  $c$  in CellOrder do
    UpdateMarginals(Factor Accumulated runoff  $f_{acc}$ );
    UpdateMarginals(Factor Constraint Accumulated runoff  $C_{acc}$ );
  Upward pass:
  foreach Cell  $c$  in inverted CellOrder do
    UpdateMarginals(Factor Accumulated runoff  $f_{acc}$ );
    UpdateMarginals(Factor Constraint Accumulated runoff  $C_{acc}$ );

```

Table 2.1: Infer.NET methods used to calculate outgoing messages. Documentation: [Microsoft Research Cambridge](#)

Factor	Direction	Method
f_{loc}	$Q_{loc,i}$	FastSumOp.SumAverageConditional()
f_{loc}	P_i, E_i, B_i	FastSumOp.ArrayAverageConditional()
f_{acc}	$Q_{acc,i}$	FastSumOp.SumAverageConditional()
f_{acc}	$Q_{loc,i}, Q_{acc,j}$	FastSumOp.ArrayAverageConditional()
C_P, C_E, C_{acc}	C_P, C_E, C_{acc}	IsPositiveOp.XAverageConditional()
f_B	B_i	GaussianOp.SampleAverageLogarithm()
f_B	μ_B	GaussianOp.MeanAverageLogarithm()
f_B	τ_B	GaussianOp.PrecisionAverageLogarithm()
f_P, f_E, f_μ	P_i, E_i, μ_B	Gaussian.FromMeanAndVariance()
f_τ	τ_B	Gamma.FromShapeAndRate()

3

Model application using test data

In this chapter the influence of adding the observation, constraint and bias components to the model is discussed. Seven different model setups using a set of test data were analyzed. The factor graph structure of the model setups enumerated below are illustrated in Appendix B. Some results are highlighted using figures in this chapter. A complete overview of the results can be found in Appendix C.

Model 1 includes forcing data, excludes runoff observations, bias term and positivity constraints

Model 2 includes forcing data, runoff observations, excludes bias term and positivity constraints

Model 3 includes forcing data, runoff observations and positivity constraints, excludes bias term

Model 4 includes forcing data, runoff observations and bias term (without parameter uncertainty), excludes positivity constraints

Model 5 includes forcing data, runoff observations, bias term (without parameter uncertainty) and positivity constraints

Model 6 includes forcing data, runoff observations and bias term (with parameter uncertainty), excludes positivity constraints

Model 7 includes forcing data and runoff observations, bias term (with parameter uncertainty) and positivity constraints

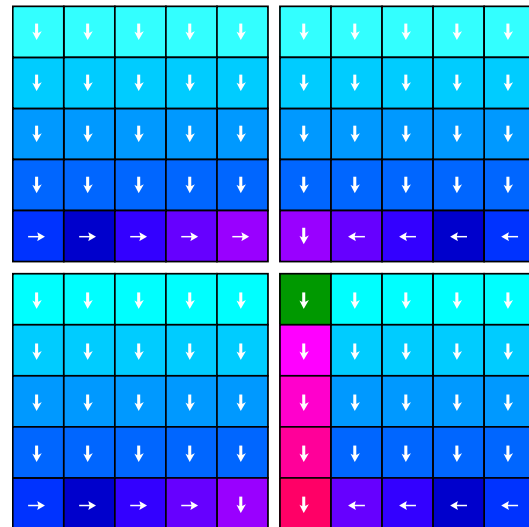


Figure 3.1: Drainage direction of all cells within the test data grid. The color indicates the flow accumulation (increasing accumulation from lightblue to red). The green cell contains a runoff observation.

In the next sections the test data is described and the influence of adding each component to the model is discussed.

3.1. Data description

The test data consists of a 10 x 10 grid, subdivided in 4 basins of 5 x 5 cells (Figure 3.1). The drainage direction of every cell is indicated by an arrow and the flow accumulation (amount of upstream cells) by the cell color. The green cell in Figure 3.1 is the location of the runoff observation (used in models 2 to 7). The basins have different prior mean values for precipitation and evaporation, resulting in different prior local runoff in each basin. The uncertainty of the data is described relative to the magnitude by the coefficient of variation (CV , Equation 2.3).

The CV in this hypothetical example is 20% for precipitation, 30% for evaporation and 5% for the runoff observation. The prior distributions of these variables per basin, and the resulting local runoff, can be found in Table 3.1.

Table 3.1: Prior values [m^3/s] for precipitation, evaporation and local runoff

	Basin	Precipitation		Evaporation		Local runoff		
		μ	σ	μ	σ	μ	σ	CV
1	Upper left	10	2	8	2.4	2	3.12	1.56
2	Upper right	20	4	10	3	10	5.00	0.5
3	Bottom left	20	4	10	3	10	5.00	0.5
4	Bottom right	50	10	10	3	40	10.44	0.26

The basins differ from each other on the value of the forcing data and on the position with respect to the runoff observation (upstream or downstream of the observation).

- Basin 1 has a low precipitation and evaporation, but in the same order of magnitude. This leads to a local runoff with a low mean value μ , a low absolute uncertainty σ but a high relative uncertainty CV . This basin is positioned upstream of the observed cell.
- Basin 2 has a precipitation which is much higher than the evaporation. This results in a higher runoff than in basin 1, a higher absolute uncertainty but a lower relative uncertainty. Also basin 2 is upstream of the runoff observation.
- Basin 3 has the same input data as in basin 2, but basin 3 is not at all connected with the rest of the basins. This in order to indicate the transfer of information purely by means of bias (models 4 to 7).
- Basin 4 has a very high precipitation compared to the other basins. This leads to an even higher local runoff and absolute uncertainty but a very low relative uncertainty. This basin is connected to basin 1 and 2 downstream of the runoff observation.

As mentioned above, the relative uncertainty of local runoff is high when the precipitation and evaporation are in the same order of magnitude. The mean local runoff is the difference between the mean values of precipitation and evaporation, while the variance σ^2 of the local runoff is sum of the variance of the precipitation and evaporation. This leads to a high value for the coefficient of variation

For all *informed* models (model 2-7), an accumulated runoff observation ($\mu_{obs} = 100 m^3/s$ and $\sigma_{obs} = 5 m^3/s$) is added in the observed cell. This observation is much lower than the calculated accumulated runoff in model 1 ($\mu_{acc} = 340 m^3/s$ and $\sigma_{acc} = 31.3 m^3/s$), to be able to clearly asses its effects.

3.2. Effect of different model setups

In this section more information is added to the model. First the accumulated runoff observation is added. The positivity constraints preventing negative variables are implemented next. Finally the bias is added to the model, first without bias parameter uncertainty and later with bias parameter uncertainty.

The spatial influence of adding information is expressed in absolute change in the mean and standard deviation compared with the prior (model 1), as well as in the relative information gain which is measured by the Kullback-Leibler divergence (KL divergence). The KL divergence is a measure of relative difference between the prior and the posterior distribution. The general formulation of the KL divergence is (Bishop 2006, p. 55):

$$KL(p||q) = - \int p(x) \log q(x) dx - \left(- \int p(x) \log p(x) dx \right) \quad (3.1)$$

where q is the prior distribution and p is the posterior distribution. Let us denote the distributions as Gaussian distributions $p(x) = \mathcal{N}(\mu_1, \sigma_1)$ and $q(x) = \mathcal{N}(\mu_2, \sigma_2)$. The KL divergence for two Gaussian distributions is (Belov & Armstrong 2009):

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \quad (3.2)$$

3.2.1. Adding observations

The influence of adding the observation can be divided into three areas (Figure 3.2). There are the cells *upstream of the observations* (blue), the *observed cell itself* (green) and the *stream downstream of the observed cell* (red). Cells on (lateral) side-branches of the tree are not influenced (black).

After one iteration (a downward pass and an upward pass) the result has converged. Only one iteration is needed because the graph does not contain any cycles due to absence of the bias parameters and approximations (positivity constraints).

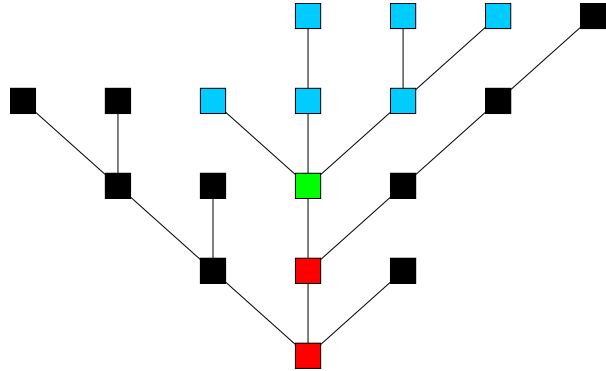


Figure 3.2: Observation influence spreads to the upstream cells (blue), the observed cell itself (green), and the downstream cells (red). The lateral cells (black) are not influenced by the observation.

Influence on accumulated and local runoff

In the *observed cell* the mean accumulated runoff has dropped from $340 \text{ m}^3/\text{s}$ to $106 \text{ m}^3/\text{s}$ and the standard deviation from $31.2 \text{ m}^3/\text{s}$ to $4.9 \text{ m}^3/\text{s}$ (Figure 3.3). This is a direct result of multiplying the prior accumulated runoff by the observation, resulting in a posterior accumulated runoff (posterior variables are denoted with an asterisk *). The resulting mean is an average of the prior and the observations, weighted on both uncertainties (as described in Equation 2.19). The posterior uncertainty is always smaller than the smallest prior uncertainty, as more knowledge is added.

The information propagates to the cell downstream of the observed cell. To keep the model consistent, the flow accumulation (Equation 2.2) is applied again in the downstream cell with the new knowledge of the observed cell. As the other inflows (other upstream neighbors and the local runoff) have not changed, the accumulated runoff in the downstream cell will change with the same magnitude as the observed cell (in terms of mean and variance). The local runoff downstream of the observation will not be influenced, which can be seen in Figure 3.4.

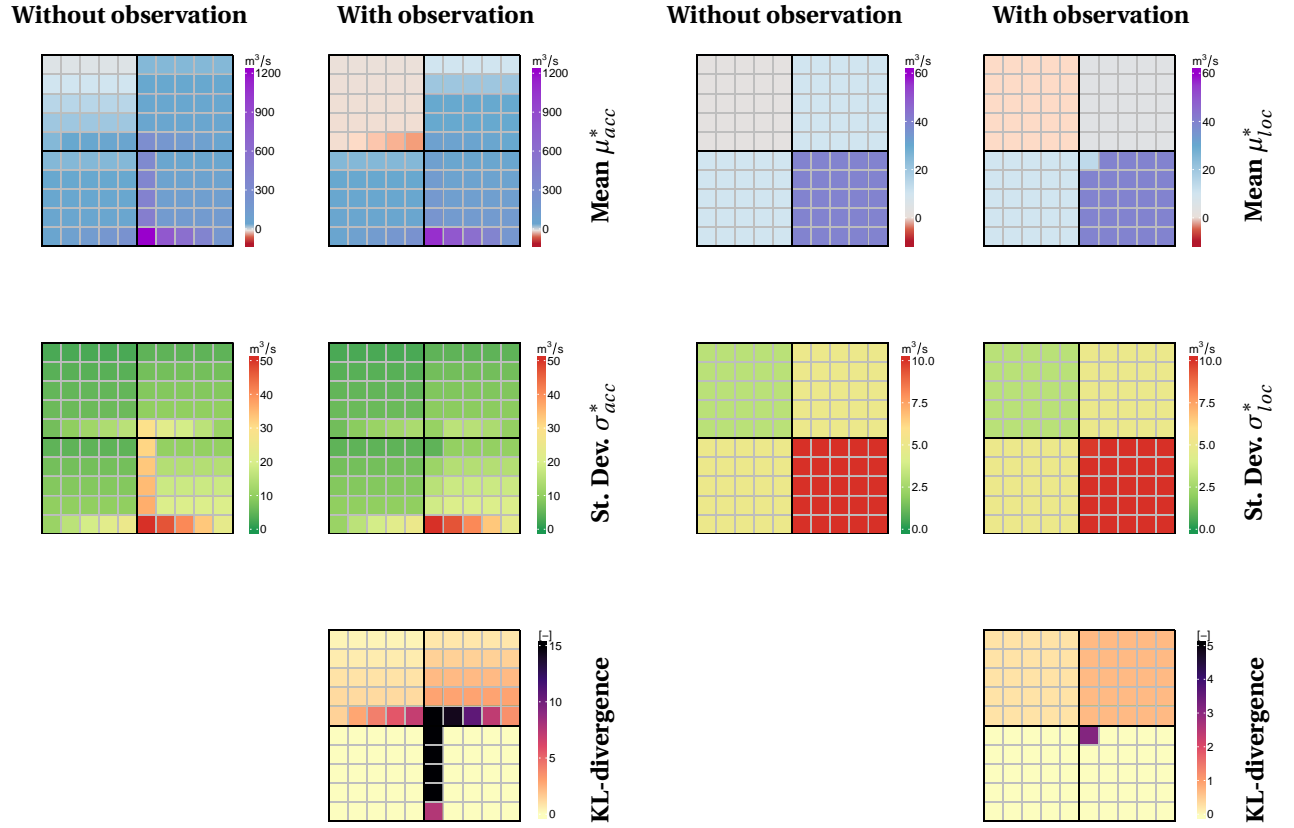


Figure 3.3: Effects of adding an observation on the accumulated runoff. Model 1 (left) and model 2 (right)

Figure 3.4: Effects of adding an observation on the local runoff. Model 1 (left) and model 2 (right)

The gained information also propagates upstream. This means that the summed accumulated runoff from the upstream cells of the observed cell plus the local runoff must equal the posterior accumulated runoff in the observed cell as stated in Equation 2.2. The propagation of this information depends on the uncertainty of the inflows. Because the local runoff has a small absolute uncertainty compared to the inflowing accumulated runoff, the mean of the local runoff will decrease only slightly (in absolute sense) compared with the mean of the accumulated inflow.

This process is repeated with the upstream neighbor cells. These cells now have an updated accumulated runoff that influences their local runoff, and in turn their upstream neighbors' accumulated runoff. With every cell upstream a bit of the information dissipates into the local runoff, while the rest propagates further upstream. The variance of the local runoff does not vary within the basins, which causes the decrease in mean local runoff to be equal throughout the cells in each basin. This dissipation can be observed in the KL divergence (defined by Equation 3.2), indicating information gain.

In cells where the mean value of the accumulated runoff was already low, the posterior mean can have a negative value. This can be prevented with a positivity constraint, which is discussed in Section 3.2.2 and will be added in the next section.

Not only the mean of the variables is influenced by the addition of the observation. As stated before, the uncertainty in the observed cell has decreased by the addition of information. The accumulated runoff close to the runoff observation also decreases, but with distance from the observed cell the effect dissipates as it gets absorbed by the local runoff. The decrease of the standard deviation in the local runoff is very limited. In Figure 3.4 the effect on the standard deviation is not visible as the difference is only in the range of decimals.

Precipitation & Evaporation

The decreased local runoff in, and upstream of, the observed cell propagates into the forcing data. This causes the mean precipitation to decrease and the mean evaporation to increase (Figure C.2a and C.2b in Appendix C). The standard deviation in basin 1 is higher for the evaporation (Figure C.2f) than for precipitation (Figure C.2e), therefore most of the change in mean local runoff propagates towards the evaporation dataset. In basin 2 the precipitation has a higher standard deviation, causing propagation of information mainly towards the precipitation.

3.2.2. Adding constraints

In the previous section, one of the mentioned side effects of adding the observation is that the mean accumulated runoff μ_{acc}^* becomes negative in some areas of the model. In order to prevent the accumulated runoff and forcing data to become negative, positivity constraints are added besides the observation.

After applying the constraints, it is evident that there are no more negative mean values in the accumulated runoff (right top in Figure 3.5). The effect in a constrained cell is twofold, as explained in Section 2.3.4. Firstly, the mean value in those cells will increase, as can be seen in the top figures in Figure 3.5. The other effect is a slight decrease in the standard deviation of the accumulated runoff in constrained cells due to truncation.

These two effects occurs throughout the model where the accumulated runoff is negative or low compared to the uncertainty (high CV). In fact, all constrained distributions with a probability on negative values are truncated and a new moment matched distribution is assigned to the variable. This new distribution might still contain probabilities in the negative domain. During a next iteration the variable constraints are applied again. Especially in leaf cells, where the accumulated runoff is equal to the local runoff, the constraints influence the result. This is because the variance is high compared to the mean value. The mean value of accumulated runoff in these cells is the difference in precipitation and evaporation ($\mu_{acc} = \mu_{loc} = \mu_P - \mu_E$) while the variance is the sum of variances ($\sigma_{acc}^2 = \sigma_{loc}^2 = \sigma_P^2 + \sigma_E^2$).

The observation of the accumulated runoff still has to be matched. The accumulated runoff has increased in the leaf cells of basin 1, so it needs to decrease somewhere else. This can be observed downstream in basin 1, where the local runoff becomes negative (bottom figures in Figure 3.5). This decreases the accumulated runoff in order to match the observation. The location where the local runoff will be negative will depend on the water availability (accumulated runoff) and uncertainty of the local runoff. Regardless of the negative local runoff in the downstream part of basin 1, the outflow of this basin has increased with the addition of the positivity constraints. This increase will be compensated by a decrease of runoff from basin 2, which is less effected by the positivity constraints (there is enough water available to decrease).

Unlike the model without constraints, this model cannot be solved using only 1 iteration. During the downward pass the constraints do not have any effect as the mean value is well above zero. As the information of the observation propagates upstream during the upward pass, the mean decreases and the constraints start

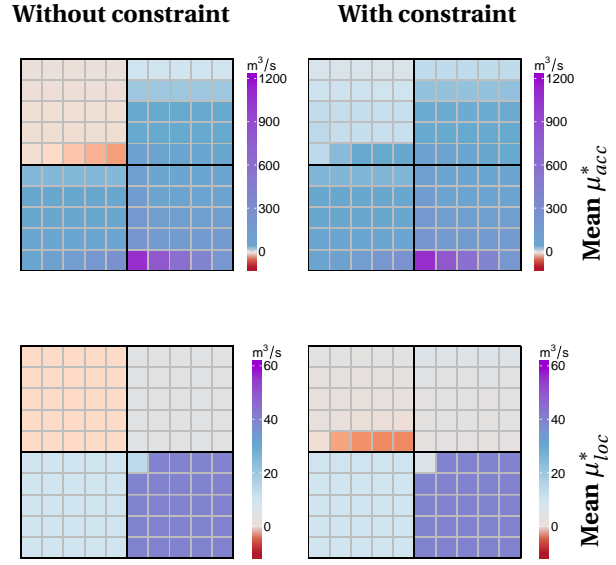


Figure 3.5: Effects of adding the constraints on the accumulated runoff (top) and local runoff (bottom). Model 2 (left) and model 3 (right)

effecting the variables. During the next iterations, information from the constraints propagates throughout the model until consistency is reached.

The KL values are not necessarily higher compared to the observed model without constraints, even though more information has been added. This is because the KL value depends on both change in mean and variance. The constraints raise the mean values throughout the model, while the observation contradicts this information and lowers the mean values. The uncertainty values are lowered by both the constraints and observation.

3.2.3. Adding a bias term

The bias term accounts for unknown or unmodeled processes and represents the model structure uncertainty. In this section a bias term is added to the observed model, without constraints and bias parameter uncertainty (model 4). The influence of adding uncertainty to the bias parameter (model 5) is considered next.

Bias in observed model, without parameter uncertainty, $\sigma_\mu, \sigma_\tau \rightarrow 0$

In the first application the bias term is added to the observed model without parameter uncertainty. This means that posterior values of the bias parameters are equal to the prior values. In other words, the model will not influence the parameters. Effectively this means that the factor f_B is removed from the factor graph (model 4, Appendix B).

The local bias B_i will have a prior value of $\mathcal{N}(\mu_\mu, \frac{1}{\mu_\tau}) = \mathcal{N}(0, 100)$ and thus can be influenced. In Figure 3.6a it can be observed that the area upstream of the observation has a posterior local bias

with a negative mean. This is because the observation is much lower than the prior accumulated runoff. The local bias in the rest of the basin is not influenced and the mean remains zero.

The effect on the accumulated runoff (Figure 3.6b) is a higher decrease in mean value. Because more uncertainty is added to the local runoff with the addition of the bias, the uncertainty in the accumulated runoff increases. This gives more weight to the runoff observation compared with the prior accumulated runoff.

Bias in observed model, with parameter uncertainty

Uncertainty is added to the bias parameters by setting σ_μ and σ_τ to positive, non zero values (model 6, Appendix B). The posterior parameters can now differ from the prior parameter, as added information from the observation propagates through the local bias to the bias parameters. The updated bias parameters now influence the local bias in all the cells of the model. That the information about the bias gathered in the upstream cells spreads out over the whole basin, which can be observed in Figure 3.7a.

This effects the accumulated runoff is now not only upstream and downstream of the observation, but also in the lateral connected cells in basin 4 and the completely disconnected basin 3 (Figure 3.7b).

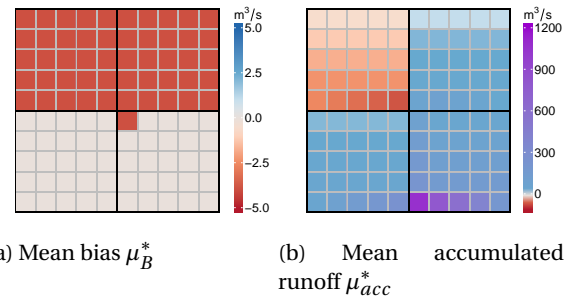


Figure 3.6: Mean value of the local bias and accumulated runoff, without parameter uncertainty (model 4). Due to the observation a bias exists in the upstream area.

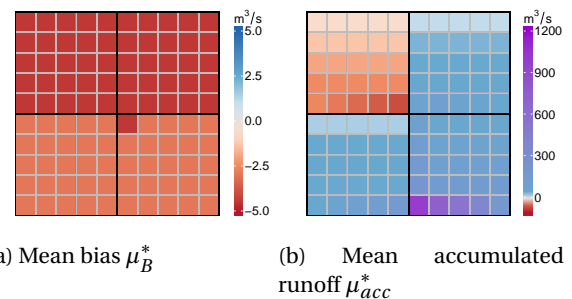


Figure 3.7: Mean value of the local bias and accumulated runoff, with parameter uncertainty (model 6).

3.2.4. Effects of combining components

In this section the interaction between observations, bias with parameter uncertainty and positivity constraints is analyzed. The observation decrease the runoff in the upstream area, causing a negative local bias in parts of the upstream area (Figure 3.8a). This local bias influences the bias parameters, causing a negative mean bias parameter μ_μ . The bias decreases the mean accumulated runoff and increases the uncertainty in the runoff. Both these effects cause constraints to have an effect in a bigger area of the model (as opposed to Section

3.2.2, where the effect was mainly limited to basin 1). In basin 1, 2 and 3 the effect of constraints is noticeable at the leaf cells where the constraints on the accumulated bias propagate to the bias, which is less negative or even positive in these cells. The accumulated runoff shows a smoothed result of the accumulated runoff in the area upstream of the observation (Figure 3.8b). This is due to the fact that the accumulated runoff is constrained by the positivity constraints on one side, and the observation on the other side. Because the bias has enough spatial variability, it adopts a value so that the accumulated runoff complies to all its constraints.

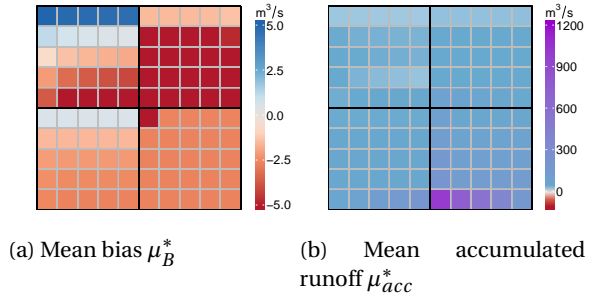


Figure 3.8: Mean value of the local bias and accumulated runoff, with parameter uncertainty and constraints.

3.3. Influence of prior model parameters

Before running the model, the prior input has to be determined. The prior data uncertainty is set to be relative to the mean value of the data. In the next sections the influence of the relative uncertainty and bias uncertainty are discussed.

3.3.1. Influence of data uncertainty

The mean values for precipitation, evaporation and runoff observations will be from data sources in the next chapter. The uncertainty of the data is entered in the model using 3 model parameters, the Coefficient of Variance for the three data sets:

- Precipitation: CV_P
- Evaporation: CV_E
- Runoff observation: CV_{acc}

Several conclusions can be made with respect to the value of these parameters. The most obvious conclusion is that a higher value of the CV (high prior uncertainty) results in a higher posterior uncertainty of that same dataset. Secondly, a high prior uncertainty on one dataset causes this dataset to adapt more to added information. For example, if the precipitation has a much higher absolute uncertainty than evaporation, the precipitation will adapt more to the observation than the evaporation in terms of mean value. The uncertainty of the observed accumulated runoff determines to which extent the model will adapt to these observations, opposed to the calculated runoff from the forcing data. All these uncertainties are relative to each other, if all the uncertainties increase, not much may happen to the result.

Finally when using constraints, a high uncertainty in the forcing data can cause a high accumulated runoff because the constraint in each leaf cell re-projects the truncated Gaussian causing a relatively high mean value for the accumulated runoff.

3.3.2. Influence of bias parameters

In this section the influence of the model parameters is covered. The model parameters - enumerated below - define the prior bias parameters according to Equation 2.7 and 2.8.

- Mean of the bias mean: μ_μ
- Variance of the bias mean: σ_μ^2
- Mean of the bias precision: μ_τ
- Variance of the bias precision: σ_τ^2

Mean of the bias mean μ_μ

The mean of the bias mean parameter is used to indicate prior knowledge about the bias. For example when there is an overall increase in storage in the basin this parameter should be negative to compensate for the storage term which is not modeled. If there is no knowledge about the general bias, this should be set to zero.

The data implies a bias (discrepancy between $P-E$ and Q_{acc} over the whole basin), defining the *implied bias* (if the local bias B_i equals the implied bias, the prior accumulated runoff matches the runoff observations). The posterior μ_μ^* has a value somewhere between the prior μ_μ and the implied bias. A prior μ_μ closer to the implied bias strengthens the belief and reduces the posterior parameter uncertainty σ_μ^* .

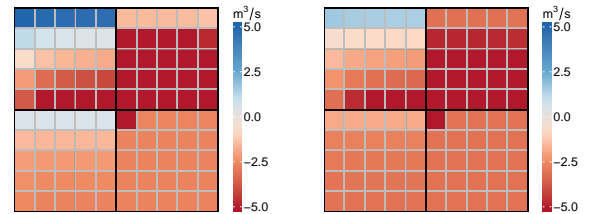
Uncertainty of the bias mean σ_μ

The uncertainty of the bias mean indicates how sure one is about the μ_μ . If this parameter is high, the posterior μ_μ^* can deviate more from the prior μ_μ , because it is more sensitive to the data (implied bias). When multiple observations are added, the implied bias in different parts of a basin is different. A higher uncertainty of the mean will cause a higher spatial variation of the posterior local bias B_i^* .

Mean of the bias precision μ_τ

The mean of the bias precision has a large influence on the uncertainty of the local bias B_i . A lower mean bias precision μ_τ results in a higher posterior variance of the local bias $\sigma_{B,i}^{2*}$. A higher variance of the local bias has as a result that the mean value of the local bias is more sensitive to the gap in the local water balance. In Figure 3.9 the mean bias is shown with low and high mean bias precision μ_τ .

A low mean bias precision results in a higher spatial variability of the mean local bias μ_B^* (and a better fit to the other data). The average mean bias in a cell deviates more from the prior μ_μ because it is more influenced by the model. A high precision results in a mean bias μ_B^* which is spatially more smooth and closer to μ_μ .



(a) Mean local bias μ_B^* with low μ_τ (high uncertainty) (b) Mean local bias μ_B^* with high μ_τ (low uncertainty)

Figure 3.9: Mean value of the local bias under changing mean bias precision in model 7 (observed, constrained, bias with parameter uncertainty)

Uncertainty of the bias precision σ_τ

This parameter indicates the uncertainty of the μ_τ parameter and the extent in which the parameter can change in the posterior with respect to the prior. As variances only decrease with the addition of information, the posterior μ_τ^* will only increase compared with its prior. The σ_τ indicates to which extent the uncertainty will decrease based on the added data.

3.4. Effect of model on posterior bias parameters

The prior bias parameters have a large influence on the posterior parameters, but also the model structure and data has an effect on the posterior value of the parameters. In this section the influence of the runoff observation and constraints on the bias parameters is explained.

The posterior mean of the bias mean μ_μ^* - in the case of a basin with one observation - is rather straightforward, assuming its prior value μ_μ is zero. If the observation is lower than the prior result, the μ_μ^* is negative. With a higher observation than the prior result, the μ_μ^* is positive. The extend of deviation from zero depends on the implied bias and the prior variance of the mean σ_μ^2 . A higher implied bias and variance allows for a higher deviation from zero. With increasing uncertainty of the forcing data or observation data (higher CV values), there will be a higher degree of adaptation of these datasets. A smaller bias is needed to close the water balance, resulting in a small μ_μ^* . When multiple observations are added to the model, different areas can have different implied biases. These areas will influence the μ_μ^* , considering the size of the area and the uncertainty of the implied bias.

The posterior mean bias precision μ_τ^* is mainly influenced by the spatial variation of the bias. When the bias is more or less uniform over the whole basin, the precision will be high. In case a model has multiple runoff observations, implying different biases per sub basin will decrease the mean bias precision μ_τ^* . Positivity constraints can cause very local but high biases. This results in a high spatial variation in the bias, causing the mean bias precision μ_τ^* to be low. In Figure 3.10 the value of μ_τ^* is shown after each iteration for the test model with the original observation and an observation which closes the water balance (observation matches the prior accumulated runoff, causing a posterior mean local bias μ_B^* of zero throughout the model). Given an observation closing the water balance, results in a higher precision because our data enforces our prior belief of the bias.

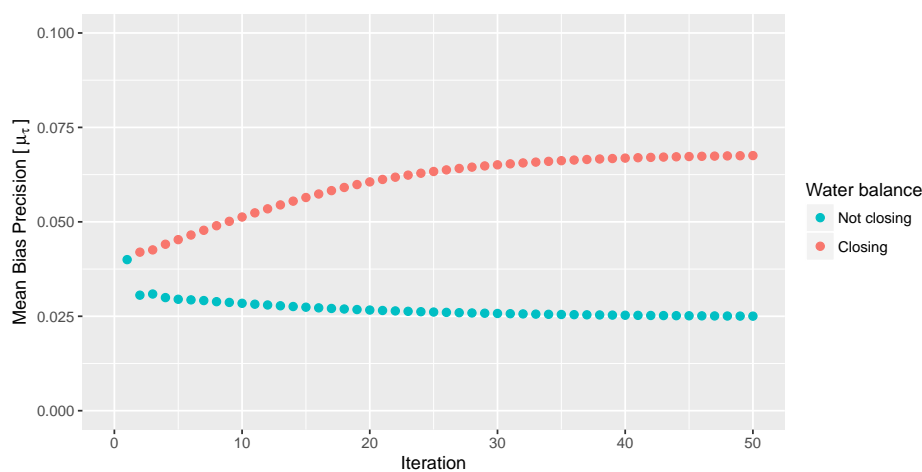


Figure 3.10: Posterior value of the mean bias precision after each iteration for the closed (implied bias of zero) and non closed water balance.

3.5. Convergence of different model structures

In this section the convergence of the different models is analyzed. The importance of factor scheduling is addressed in the first section. The subsequent sections deal with the convergence using scheduled factors. Using factor scheduling, models 1, 2 and 4 (forcing data, forcing data + observations, forcing data + observations + bias without parameter uncertainty) converge on the first iteration because those models do not use approximations and there are no cycles in their factor graph.

3.5.1. Factor scheduling

Scheduling of factors in a downward sweep and an upward sweep (discussed in the previous chapter) has a big influence on the speed of convergence of the model.

The model (model 7, including positivity constraints and bias with parameter uncertainty) is first executed in an unscheduled manner. This is done by using the Microsoft Infer.NET package, in which the relations between variables and factors are defined. The software package then executes all the message passing automatically, but does not use the tree structure of the factor graph. The model starts solving one type of factor in the first cell (left top of the model), going down the rows to the last cell (right bottom of the model), before progressing to the next type of factor. This limits the propagation of information in the upward and left direction, because those cells were already solved during the iteration and will only propagate the information further in the next iteration.

Next, the model is executed without using the Infer.NET package but by manually implementing the message passing and scheduling (downward and upward sweep). Figure 3.11 shows the accumulated runoff at the outlet of the basin. The model with scheduling has almost converged after 10 iterations, while the unscheduled model needs several iterations only to propagate any information to the accumulated runoff at the outlet of the model. After the first information reached the outlet, it takes another 80 iterations to approach the point of convergence.

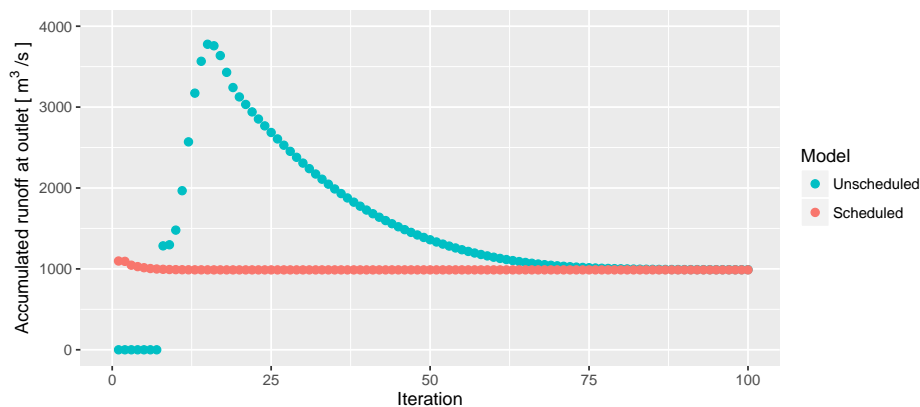


Figure 3.11: Convergence of the unscheduled (Infer.NET model) and the scheduled (manual model 7), shown by plotting the mean accumulated runoff at the root cell.

With upscaling of the model the difference in needed iterations only grows. The model with factor scheduling still propagates all information through the model in one iteration. The unscheduled model needs more iterations just to propagate knowledge through the model before the solution can start converging.

3.5.2. Positivity constraints

As explained in the previous chapter, applying a positivity constraint does not always result in a variable which has no probability on a negative value. In a next iteration of the model, the variable will be again subject to the positivity constraint. The effect of this is observed in Figure 3.12, where the mean value of the variable increases with each iteration, while the variance decreases.

This process starts with a rapid convergence but it slows down as the variable nears its convergence point. The convergence process can take very long, but the difference with each iteration will become so small that for the use in this model it can be considered converged. The influence of the constraints on the accumulated runoff is depicted in Figure 3.13, which shows the absolute difference of the outlet runoff with respect to the previous iteration on a logarithmic scale.

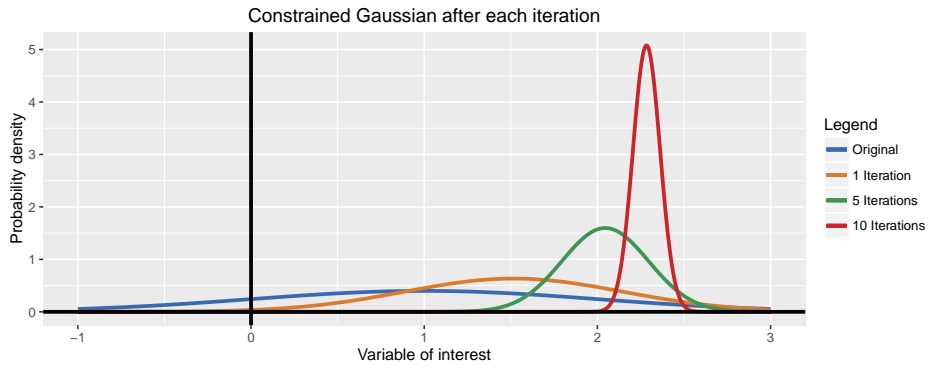


Figure 3.12: Applying a positivity constraint multiple iterations on the same variable.

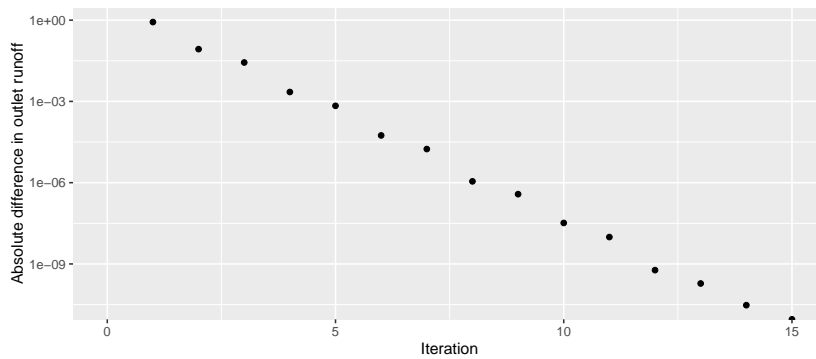


Figure 3.13: Difference in outlet runoff compared to the previous iteration (model 3).

3.5.3. Bias with parameter uncertainty

Each cell in the model updates the bias parameters, which in turn updates the local bias in all other cells. The update of a local bias due to the updated bias parameters will propagate towards the accumulated runoff. This results in the updating of the local bias elsewhere in the model, which will update the bias parameters again. This cyclic behavior results in a need for multiple iterations to converge to a solution. Just as with the positivity constraints, the convergence speed is high at the beginning and slows down with every iteration. Figure 3.14 shows the difference in outlet runoff compared to the previous iteration. After a rapid decrease of the influence, a logarithmic decrease of the difference is observed. Note that the addition of bias parameter uncertainty results in a slower convergence compared to the addition of the positivity constraints.

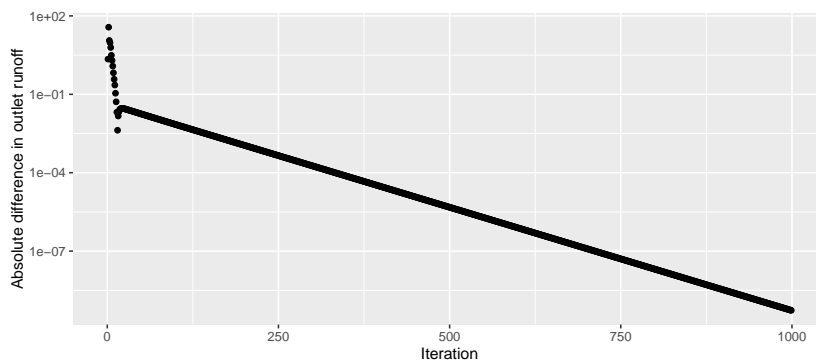


Figure 3.14: Difference in outlet runoff compared to the previous iteration (model 6).

4

Model application using real data

In the previous chapter, different model structures and the influence of model parameter influences was analyzed. In this chapter, previously gained knowledge is applied on the Volta basin (Appendix A). The runoff in the basin is monitored at 10 runoff locations.

The data used for the model and the data processing methods are outlined in Section 4.1. The quantification of uncertainty of the different data sets is described in Section 4.2. The results of implementing different model components is described in Section 4.3, followed by the influence of bias uncertainty in Section 4.4.

4.1. Data description

The datasets which are used in this research are described in Section 4.1.1. In order to use the data several processing steps are performed, which are described in Section 4.1.2.

4.1.1. Data sources

For the setup of the model as described in Chapter 2, 5 data sources are needed. The first two data sources are the spatially distributed forcing data precipitation and evaporation. The third dataset contains point measurements of the river runoff. The fourth data source is the Digital Elevation Model, in order to accumulate the local runoff downstream and connect the previously described datasets. Finally a land cover map is used to identify open water.

Precipitation: CHIRPS

The **C**limate **H**azards Group **I**nfra**R**ed **P**recipitation with **S**tation data (CHIRPS) is a precipitation data product combining long-time averages, infra-red observations, ground precipitation observations, remote sensing precipitation data (Tropical Rainfall Measuring Mission, TRMM), and atmospheric model data (Funk et al. 2014). Its coverage is quasi-global (50°S–50°N), with a spatial resolution of 0.05° (≈ 5km at the equator). The temporal resolution ranges from daily (for Africa 6 hourly) to annual, containing data from the year 1981 to present.

Actual evaporation: CMRSET

The CMRSET dataset is a scaled dataset derived from the MOD16 ET algorithm dataset by Mu et al. (2011) for actual evaporation (Guerschman et al. 2009). It contains monthly estimates for Actual Evapotranspiration

(AET).

River runoff: GRDC

The river runoff observations are retrieved from the Global River Discharge Centre (GRDC) archive. The GRDC operates under the World Meteorological Organisation (WMO) and is hosted by the German Federal Institute of Hydrology ([Bundesanstalt für Gewässerkunde 2014](#)). Its international data archive contains over 200 years of river discharge data, a total of more than 9000 stations collecting daily or monthly data.

Digital Elevation Model: GMTED2010

In order to derive the flowpaths, the Digital Elevation Model (DEM) *Global Multi-resolution Terrain Elevation Data 2010* (GMTED2010) is used, a product from the United States Geological Survey (USGS) and National Geospatial-Intelligence Agency (NGA) . It is available in 30, 15 and 7.5 arc-second resolution and is derived from 11 datasets of which the SRTM Digital Terrain Elevation Data is the most important ([Danielson & Gesch 2011](#)). For this study the 30 arc-second resolution ($\approx 1\text{ km}$ at the equator) is used as this is already a higher resolution than the precipitation and evaporation data. An even higher resolution would increase the amount of cells in the model and the amount of stations that needs correction (Section 4.1.2).

Land cover

For the identification of lakes (which will be used further on in this chapter) the *0.5 km MODIS-based Global Land Cover Climatology* ([Broxton et al. 2014](#)) is used. This global map with a resolution of 15 arc seconds classifies each grid cell as one of 17 classes, of which open water is one.

4.1.2. Data processing

In order to prepare the data for use in the model several operations are performed. The following data processing steps are executed using 3 open source software packages: GRASS ([GRASS Development Team 2012](#)), QGIS ([QGIS Development Team 2015](#)), and R ([R Core Team 2015](#)).

Averaging of Precipitation, Evaporation and Discharge measurements

The precipitation and evaporation datasets contain monthly estimates. Per year the monthly data is summed resulting in a yearly precipitation and evaporation. This is then averaged over the period from the year 2000 up to and including the year 2012. This time frame was chosen because all datasets contain sufficient data in this period. The local water balance in this model uses the precipitation and evaporation variables in m^3/s . Therefore the data products which are in mm/y are converted accordingly:

$$X_{\text{m}^3/\text{s}} = X_{\text{mm}/\text{y}} * A_i * C \quad (4.1)$$

where $X_{\text{m}^3/\text{s}}$ is the variable of interest in volume per second [m^3/s], $X_{\text{mm}/\text{y}}$ is the variable of interest in depth per year [mm/y], A_i is the surface area of a cell [m^2] (cells have different surface areas depending on their geographical coordinates as the resolution is in radian degrees), and C is the conversion constant [$\frac{\text{m}}{\text{mm}} \frac{\text{y}}{\text{s}}$] that converts the millimeters to meters and the year to seconds.

The river runoff measurements contain monthly values per station. More data gaps exist in this data. An average of the available months is used in order to calculate a long term average. Depending on the measurement station, 31 to 86 months (out of 156 months) were used to calculate the average yearly value.

From digital elevation model to flow direction map

In order to build a drainage direction map from the digital elevation model, the flow direction of each cell is determined. This is done by using the D8 algorithm. For every cell, the neighbor (one of the 8 adjacent

cells) with the lowest elevation is considered the receiving cell and all the water from the current cell will flow towards that neighbor. If a cell is surrounded only by neighbor cells which all have a higher elevation, this cell is called a pit. The pit is a local depression and all the water flowing towards this cell would disappear from the catchment. In order to prevent this, the cell is raised to the level of its lowest neighbor and it will drain in the direction of this neighbor. The result will be a raster map containing the flow direction of every cell, starting at 1 for the east direction, going counterclockwise until 8 for southeast.

Re-sampling grids

The precipitation and evaporation dataset have a resolution of 0.05° , the DEM has a 30 arc-second resolution ($= 0.0083333^\circ \approx 1\text{km}$ at equator) and the land cover has a 15 arc-second resolution. The data is re-sampled in order for the grids to have the same resolution and origin as the DEM data. This is done by using the Nearest Neighbor Resampling algorithm. This results for every cell of precipitation and evaporation in 36 cells containing the same value, slightly shifted with perspective to the larger original cell as is depicted in Figure 4.1.

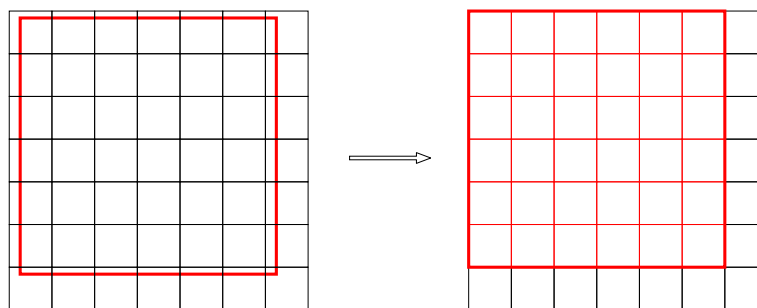


Figure 4.1: Re-sampling grid cells using the nearest neighbor algorithm. On the left: large red cell belongs to the precipitation/evaporation grid, small black cells belong to drainage direction grid. Origins and resolution do not match. On the right: red cell was shifted and split into smaller cells to match the raster of the drainage direction

Correction of observation stations

Not all measurement station the GRDC archive have well specified locations. Often the latitude and longitude coordinates are documented with only 2 decimals which can induce an error of about 1 kilometer. Moreover, errors in the stream network (from the DEM) result in a difference between the calculated position and the real position of the river. This can be observed in Figure 4.2, where the calculated flow approximates the real position of the river, but the measurement station (red circle) is not positioned in the calculated flow. Fortunately, for each station the name of the river which is measured is documented, allowing for manual correction using OpenStreetMap as a background map.

Using the D8 algorithm, flows can only converge. In reality flows do sometimes diverge naturally or by man made structures. Canals can withdraw water from the river and discharge it again further downstream or into a different river. If discharge measurement stations clearly observe only part of the flow (a side canal), they should be excluded in order to get the best results.

When two rivers are in each others proximity there is a chance that the D8 algorithm makes a false cross-over, as is shown in Figure 4.3. As a consequence, the accumulated runoff is assigned to the wrong river. The error induced by this effect is only resolved at the point where the rivers converge. Runoff observations between the cross-over and the point where these two rivers converge would introduce errors and should be excluded from the model.

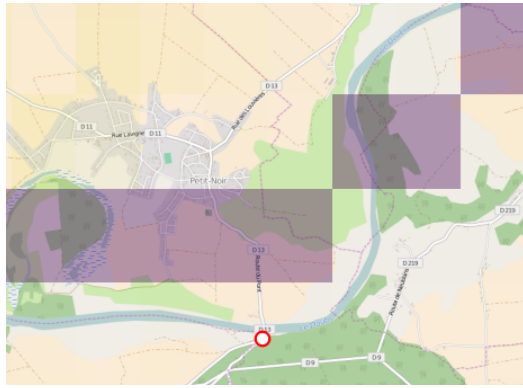


Figure 4.2: The calculated stream does not exactly match the river, causing an error if the location of the observation station (red circle) is not moved accordingly. Background: ©OpenStreetMap contributors



Figure 4.3: The calculated stream does not follow the stream from the right-bottom corner, but at the bend in the river it makes a cross-over to another river. Background: ©OpenStreetMap contributors

Not constraining accumulated runoff in lakes

The stream calculated by the D8 algorithm is only one cell wide. When a river enters a lake (or any other flat plain), a one cell wide path to the outlet of the lake is created. The other cells covering the lake are lateral inflows to this stream, as is shown in Figure 4.4.

The evaporation in lakes is often higher than precipitation because of the open water and water availability, which often causes a negative local runoff in the cells covering a lake. Considering that the lateral streams in the lake contain cells with a negative local runoff, the accumulated runoff of these branches will also be negative. By constraining these cells the lakes' withdrawal (due to the negative local runoff) of the stream becomes a discharge.

To prevent this effect from happening an exception to the accumulated runoff constraint is added to cells which are classified as open water in the land use dataset. The lateral inflows will be able to have a negative accumulated runoff it becomes possible to withdraw water from the stream.

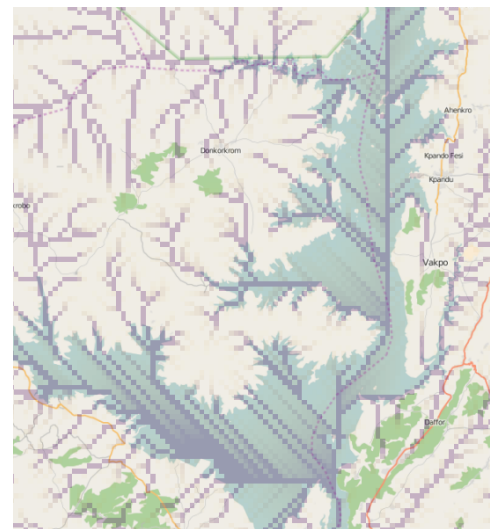


Figure 4.4: Flow accumulation in a lake. Background: ©OpenStreetMap contributors

4.2. Uncertainty quantification

Documented quantification of uncertainty of the datasets is often ambiguous, as sometimes an absolute range is given (*mm/month*, as is the case with the TRMM data), sometimes relative to the measured value. In many cases these values indicate the uncertainty in the 95% confidence interval, but this is not always communicated.

4.2.1. Precipitation uncertainty

The TRMM dataset includes an absolute error (in mm/month). As the CHIRPS dataset (which does not include an error estimation) is based on the TRMM, the TRMM error estimation will be used as an indicator for the magnitude of uncertainty in the CHIRPS dataset.

For the Volta basin the average relative error ($CV = \frac{\sigma}{\mu}$) is 3.2%. In recent studies by Karimi (Karimi &

Bastiaanssen 2015, Karimi et al. 2015) the Mean Average Percentage Error (MAPE) on precipitation products was found to be 18.5%. Although MAPE is not directly related to the standard deviation, it indicates a higher uncertainty than the TRMM uncertainty suggests. The validation papers used in the literature study of Karimi report a deviation of TRMM data between 0 and 64%. In other words, there is not much certainty about the uncertainty.

4.2.2. Evaporation uncertainty

The CMRSET dataset for evaporation does not include a quantification of uncertainty either. Research has been conducted for validation of the MOD16 data product, of which CMRSET is a scaled version. It is assumed that the uncertainty of the two data sets does not differ much. Velpuri et al. (2013) reports an uncertainty of 25% in the MOD16 actual evaporation dataset. The literature is not clear about the definition of uncertainty, therefore it is assumed to be the deviation from the mean at a 95% confidence interval. Assuming a normal distribution the Coefficient of Variance can be calculated:

$$\begin{aligned}\mu \pm 1.96\sigma &= \mu \pm 0.25\mu \\ 1.96\sigma &= 0.25\mu \\ CV = \frac{\sigma}{\mu} &= \frac{0.25}{1.96} \approx 0.13\end{aligned}\tag{4.2}$$

According to the literature study of Karimi & Bastiaanssen (2015), on average the accuracy of the evaporation products is higher than those of precipitation products. The MODIS product has a MAPE ranging from 0.6% to 18% with a mean of 6%. Again, as is the case with the precipitation, there is no real consensus about the uncertainty of the evaporation dataset.

4.2.3. Runoff observations uncertainty

The same principle is applied to the runoff observations. The exact method for the collection of the data in the GRDC river runoff dataset is not known. Quantification of uncertainty of individual discharge measurements is in the range of 5-20% (Hersch 2009). Rating curves are often used to determine runoff based on water level observations. Di Baldassarre & Montanari (2009) investigated the uncertainty induced by these rating curves using a 1D model of the Po river in Italy. The global uncertainty (combining measurement and rating curve uncertainty) was found to be 25.6% at the 95% confidence interval on average. Applying equation 4.2 again, it leads to a CV of 0.13.

4.3. Implementation of the model

In this section the model is applied on the Volta basin and the results are shown. The influence of different components of the model is discussed. The estimations of data uncertainty discussed in the previous section are used, together with an estimate of the bias parameters, all of which can be found in Table 4.1. The model parameters concerning the bias are chosen for the following reasons:

- μ_μ : no information about a prior bias is available
- σ_μ^2 : allows posterior bias mean to deviate from prior, but it should not account for too much of the adaptation due to observations
- μ_τ : adds uncertainty to the local water balance but it should not be the main source of uncertainty. Still small compared to forcing data uncertainty
- σ_τ^2 : causes the shape parameter of the Gamma distribution of the precision to be 1, allowing the bias precision to approach zero

Table 4.1: Model parameters

Coefficient of Variation			Bias parameters			
Precipitation	Evaporation	Observation	$\mu_\mu[\frac{mm}{y}]$	$\sigma_\mu^2[\frac{mm^2}{y}]$	$\mu_\tau[\frac{mm^{-2}}{y}]$	$\sigma_\tau^2[\frac{mm^{-4}}{y}]$
3.2%	13%	13%	0	100	0.01	0.0001

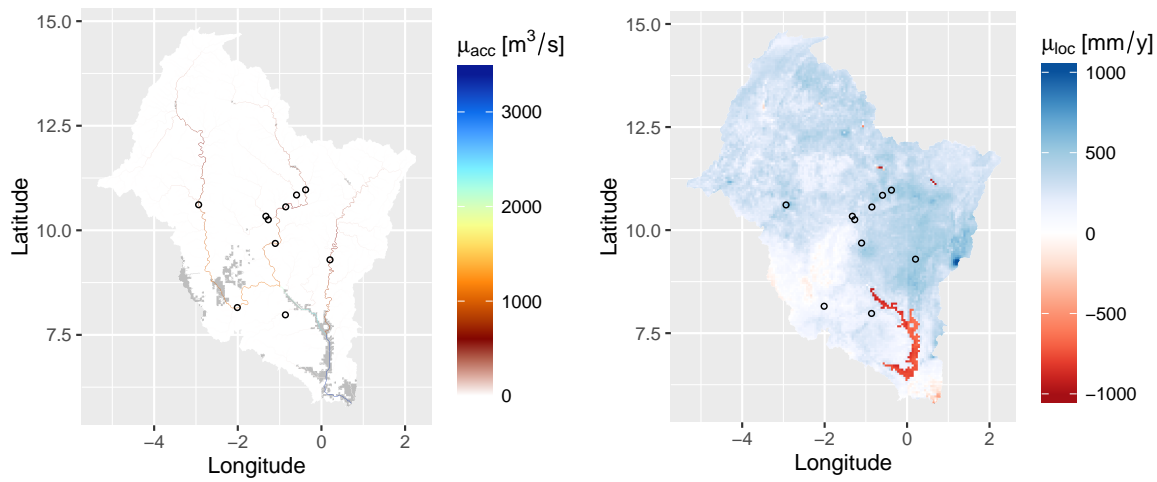
In the first section the prior results are calculated, adding forcing data but no observations, bias and positivity constraints (model 1). In the following sections, the effect on the results caused by the different model components is shown. Note that in those sections, the graphs show the difference in mean value of the variable in respect to the prior results. An overview and additional graphs regarding the standard deviation can be found in Appendix D. The results for the local runoff, precipitation, evaporation and bias are converted to a more comprehensible mm/y .

- Model 2: Forcing data + observations
- Model 6: Forcing data + observations + bias (with parameter uncertainty)
- Model 7: Forcing data + observations + bias (with parameter uncertainty) + positivity constraints

4.3.1. Prior

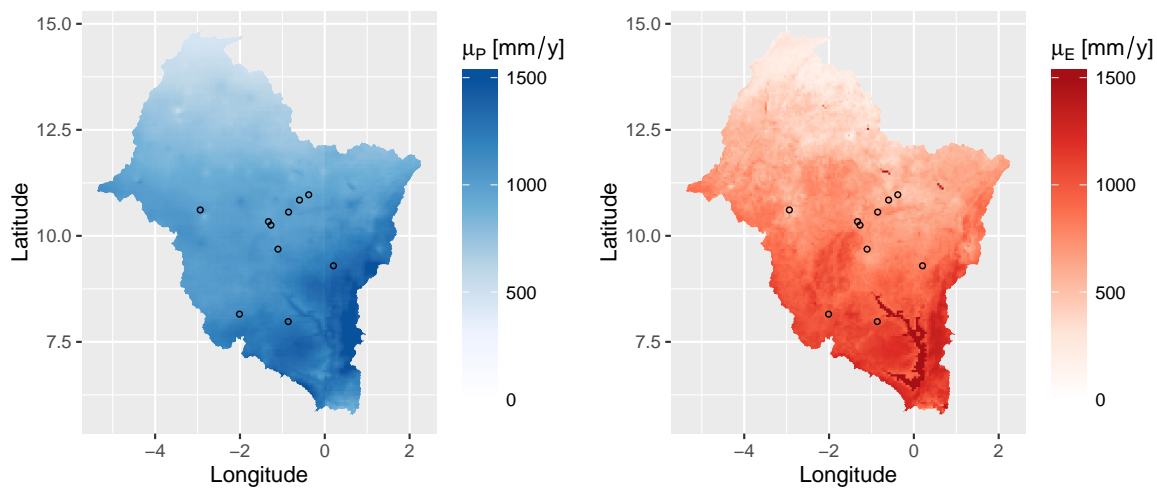
The prior values of each variable are shown in Figure 4.5. The prior is the result of the model where only forcing data is added to the model and propagated to the local runoff and accumulated runoff. The locations of the runoff observations, that will be added in the next models, are indicated with black circles.

In the mean local runoff (Figure 4.5b), Lake Volta in the south can be clearly identified, due to its high negative local runoff due to high values of open water evaporation (see Figure A.1 for a geographical map of the Volta basin). Just north-west and south-east of Lake Volta there are two areas with a negative local runoff due to higher levels of evaporation than precipitation. This leads to a negative accumulated runoff (Figure 4.5a). When applying positivity constraints in model 7, these areas will be prone to high changes in variable values.



(a) Prior mean accumulated runoff, grey areas indicate a negative value

(b) Prior mean local runoff



(c) Prior mean precipitation

(d) Prior mean evaporation

Figure 4.5: Prior mean values for a) accumulated runoff, b) local runoff, c) precipitation, and d) evaporation. The circles indicate the position of the runoff observation stations.

4.3.2. Model 2

In this model the observations are added, which are lower than the prior data suggests. In the accumulated runoff the influence of these observations results in a decrease in the mean value in the whole upstream area (although the effect mainly occurs in the bigger streams as the uncertainty there is larger) and the river downstream (Figure 4.6a). The decrease is rather small, as the uncertainty on the accumulated runoff at the observation stations is lower than the uncertainty of the runoff observation itself. Due to the addition of information by means of runoff observations, the uncertainty in the accumulated runoff decreases slightly. In the local runoff and forcing data the influence only occurs upstream of the observation (Figure 4.6b).

The effect on precipitation and evaporation (Figures 4.6c and 4.6d) is of different extend in each region upstream of an observation depending on the runoff observation. Also within each upstream sub catchment,

which has to do with the fact that in each cell the variance is relative to the value. A higher variance results in a higher adaptation to the observation. This is also the reason why the evaporation adapts much more to the lower observations than the precipitation, as the variance on the evaporation is higher than on the precipitation.

The standard deviation σ in these variables hardly decreases. This is because the small decrease in variance σ_{acc}^2 in accumulated runoff is divided over all upstream cells, resulting in an even smaller decrease in standard deviation.

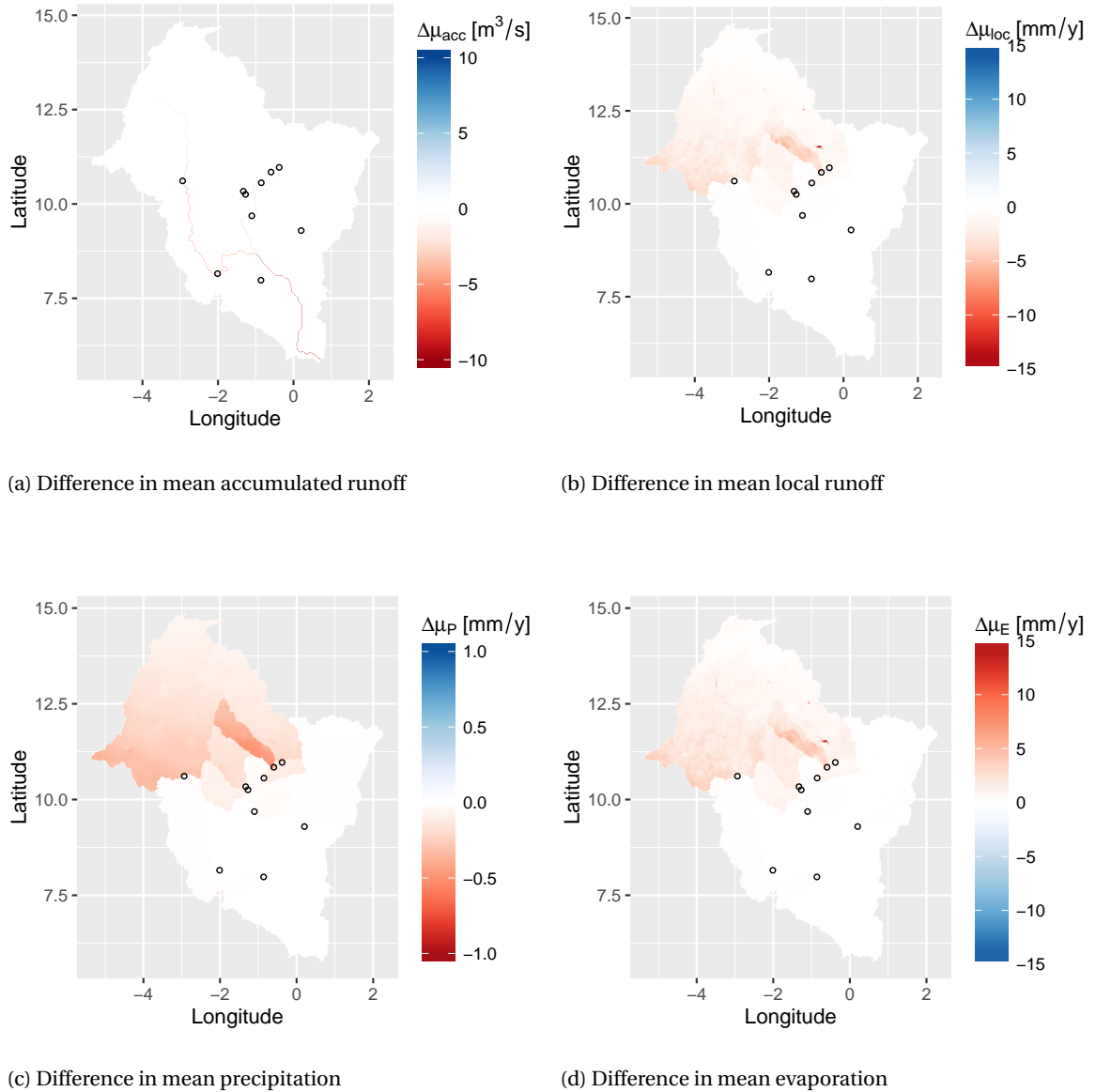


Figure 4.6: Difference in mean value with respect to the prior for all the variables after adding observations (model 2)

4.3.3. Model 6

The bias is added to the local water balance. In essence, the bias adds uncertainty to the local water balance and consequently will increase the uncertainty in the local runoff and accumulated runoff. This results in a higher decrease of those variables, compared with model 2 (Figures 4.7a and 4.7b). At the local water balance this adaptation to the observations is partly absorbed by the bias, resulting in a smaller difference in the precipitation and evaporation (Figures 4.7c and 4.7d).

In the cells upstream of the observations, the local bias becomes negative as the local runoff decreases. This updates the knowledge about the bias parameters. Due to the updated bias parameters, the local bias in the area downstream of the observations will also become negative (Figure 4.7e) resulting in a decrease of the local runoff in that area.

4.3.4. Model 7

Finally the positivity constraint is added on top of the observations and bias. In contrast with the previous results, an increase in accumulated and local runoff is observed compared with the prior results. This effect is caused at the cells in which the positivity constraint increases the accumulated runoff. The increase in accumulated runoff propagates into the increase of the local runoff. Further propagation into the local bias takes place, resulting in updated bias parameters. During the model run the bias parameter μ_μ becomes positive, which will influence the local bias in other cells. This results in an overall increase of local runoff (although the influence of the observations can still be seen, Figure 4.8b) and therefore also an increase in accumulated runoff (Figure 4.8a). Although the bias absorbs some of the increase in local runoff, the precipitation and evaporation adapt to a high degree to the change in local runoff 4.8c and 4.8d.

When taking a closer look on the spatial distribution of the bias (Figure 4.8e), a higher bias is observed in the regions where the prior local runoff was negative, and around Lake Volta. Apparently, the data suggests a negative local runoff where there is no water available. This could be the result of faulty precipitation and evaporation data, or suggest that the used water balance does not represent the processes in a good way. The areas with a negative local runoff according to the precipitation and evaporation data could be fed through ground water flows or irrigation, where water is pumped from the ground water or nearby rivers or lakes causing a higher evaporation than naturally would be possible.

The high bias due to constraints around lake Volta has two explanations. The first reason is that the original evaporation data has a coarser resolution than the digital elevation model. This means that the original pixel of the evaporation product overlaps model cells which represent open water but also cells which represent land. In the cells representing land, the evaporation is still in the same order of magnitude as the evaporation above the water. Even outside of the influence of these coarser cells constraining is observed through a higher bias, which can be explained by ground water flows from the lake to the cells that cover the flat plains around the lake.

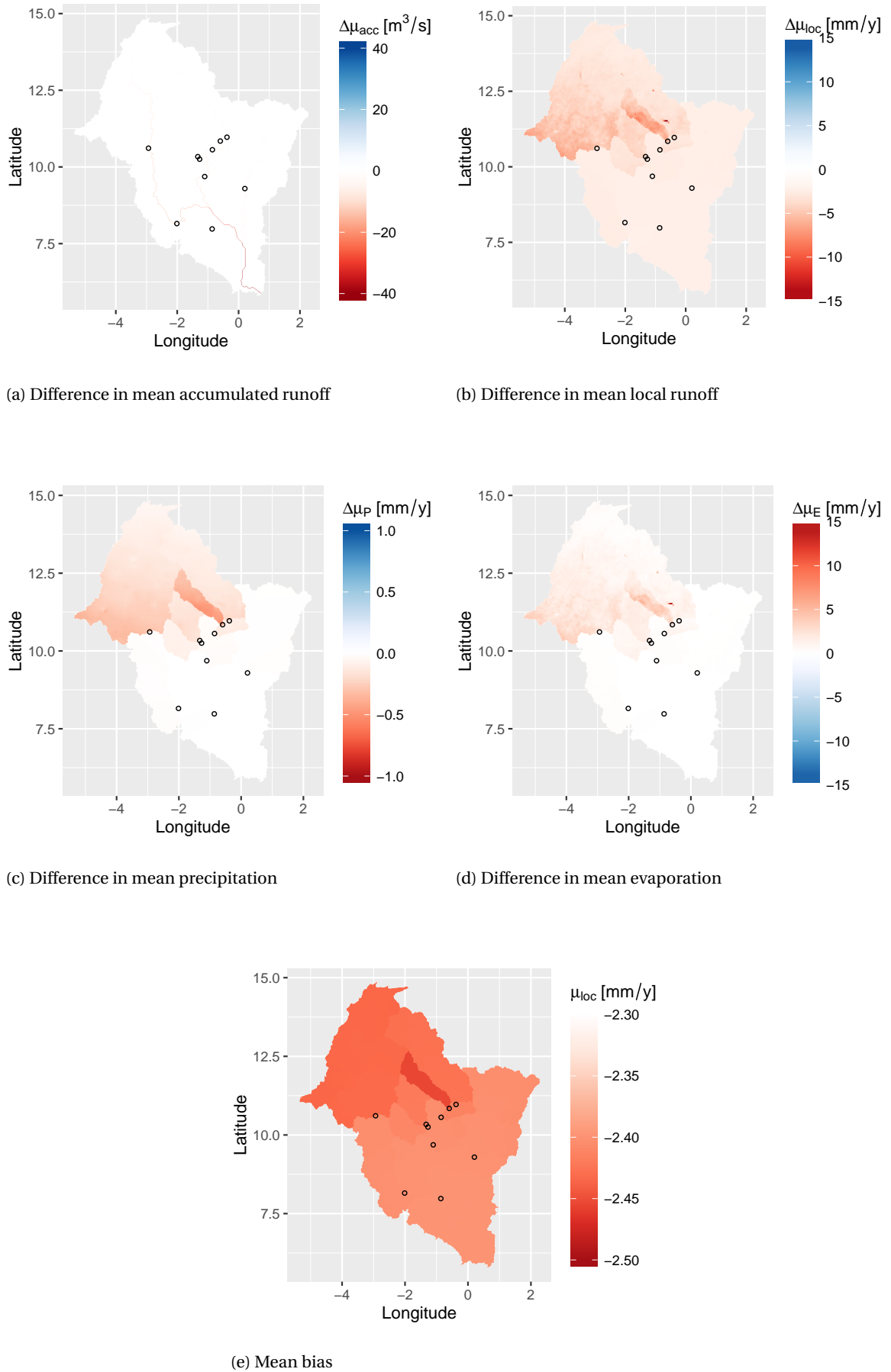
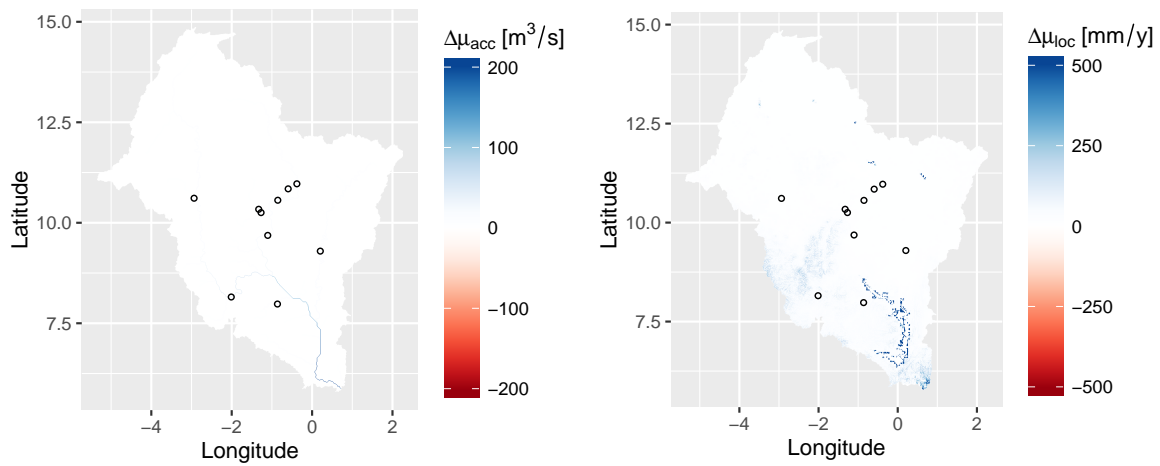
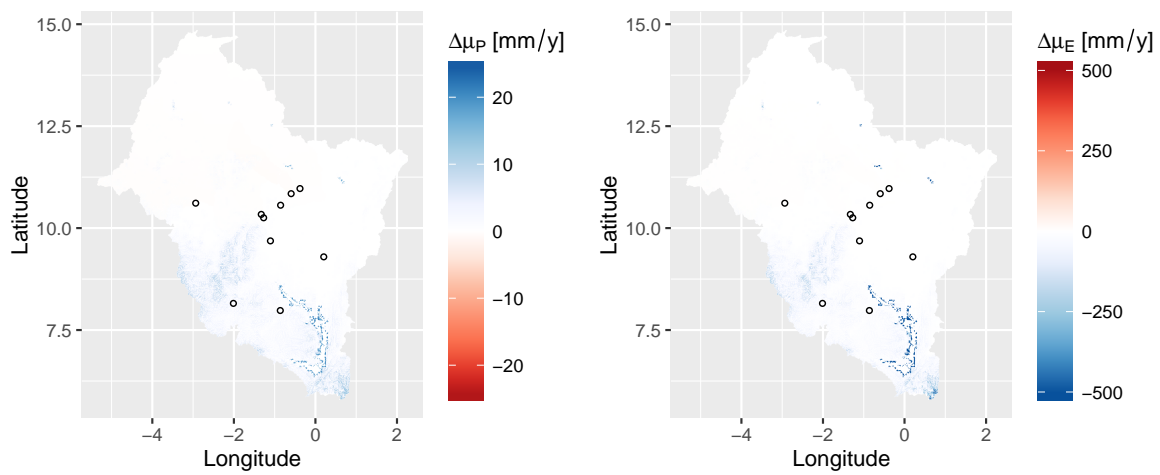


Figure 4.7: Difference in mean value with respect to the prior for all the variables after adding observations and bias (model 6)



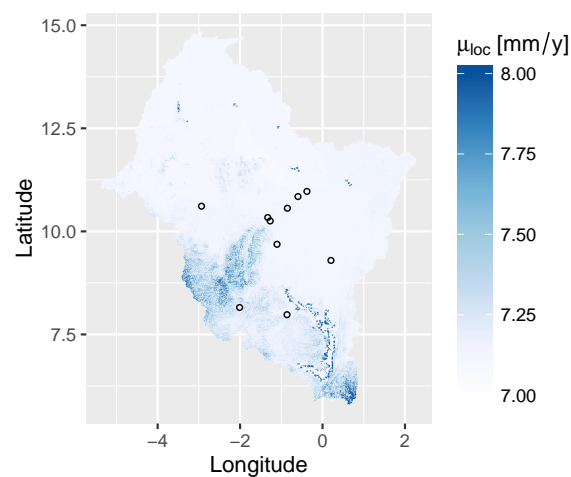
(a) Difference in mean accumulated runoff

(b) Difference in mean local runoff



(c) Difference in mean precipitation

(d) Difference in mean evaporation



(e) Bias

Figure 4.8: Difference in mean value with respect to the prior for all the variables after adding observations, bias and positivity constraints (model 7)

4.4. Bias uncertainty influence

As described in the previous chapter, the mean bias precision parameter μ_τ is an indicator for spatial variation of the local bias. A lower precision results in a higher variance of the local bias. A higher variance allows the local bias to be more influenced by the data received from the model. First the influence of this parameter on the model without positivity constraints is considered, followed by the model containing the positivity constraints.

4.4.1. Model without positivity constraints

The prior mean bias parameter μ_μ is set to zero and the observations are lower than the prior data. A small bias uncertainty of the bias results in:

- A posterior μ_μ^* close to its prior.
- A small spatial variation of the mean local bias.
- Overall a better match between the model results and the runoff observations compared to the model without bias. Due to limited spatial variation, the match between the model result and observed runoff may worsen. Example: when all but one runoff station observe a lower runoff than the prior, the μ_μ parameter will be negative. Because of a low mean bias precision, the local bias will also become negative even though the runoff observations suggests otherwise. As a result the accumulated runoff decreases compared to the prior result, causing a worse fit to the runoff observation.

Increasing the bias uncertainty results in:

- A posterior μ_μ^* further from its prior.
- A higher spatial variation of the mean local bias.
- The result at some observation stations match better with the observation due to a higher (negative) bias. Some stations which were already matching well at a lower bias uncertainty now match less well. This is because the increase in mean bias can not be compensated by the increased spatial variation.

Increasing the bias uncertainty even further causes overfitting:

- A posterior μ_μ^* close to its prior.
- A very big spatial variation of the bias.
- Results at all observation station match very well because the local bias will adapt to a high degree to the runoff observations due to the high uncertainty.

4.4.2. Model with positivity constraints

When the constraints are added, the bias variance influences the results through another mechanism. This has to do with the moment matching of the constrained variables. A higher variance in accumulated runoff means that constraints increase the mean value more. This is illustrated using the following examples.

Using a low bias uncertainty, effects on the constrained model results are:

- Even though the observations are lower than the prior result, the posterior μ_μ^* is positive. This is a result of the constraints which increase the local bias, which in turn updates the bias parameters.

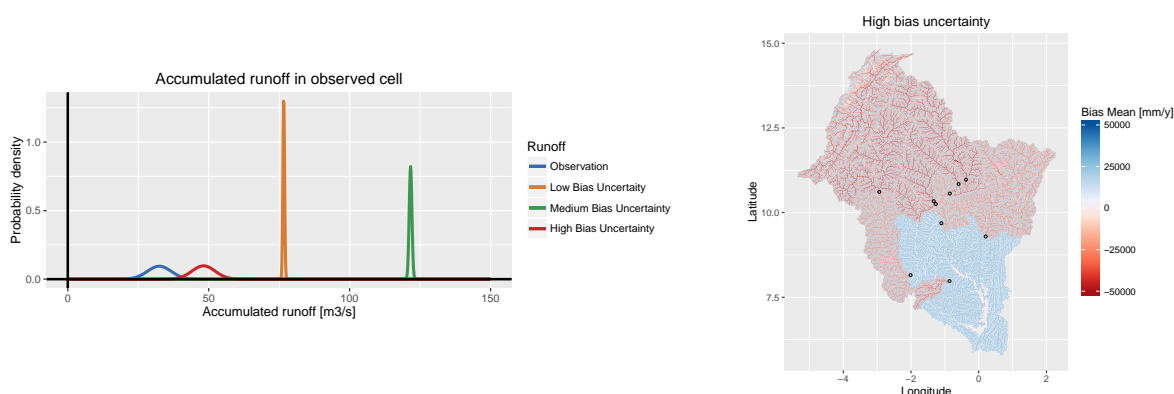
- Throughout the model the local bias is positive, which increases the local and accumulated runoff.
- Observations match less well compared with the unconstrained model which is a result of the increase in accumulated runoff.

Adding more bias uncertainty:

- A lower bias precision μ_τ causes a higher variance in the local bias and therefore in the local runoff and the accumulated runoff. Due to the higher variance in the accumulated runoff, the constraints have a higher increasing effect on the mean of the constrained variable. This results in a higher local bias and in turn in a higher posterior μ_μ^* .
- Accumulated runoff matches even worse with the observed runoff at most observation locations.

Overfitting due to too high bias uncertainty:

- Very high uncertainty in the local bias leads to an extreme effect of the constraints in all cells that have a low accumulated runoff, resulting in an extremely high mean local bias.
- The high bias uncertainty allows for a high spatial differentiation of the bias. Because the accumulated runoff has to match with the observations, the bias in the larger flows has an extreme negative value decreasing the accumulated runoff.
- The posterior μ_μ^* has decreased again, compared to the previous result.
- At the runoff observation locations, the accumulated runoff matches well with the observed runoff.



(a) Observation and accumulated runoff at a runoff observation station for different values of the bias precision parameter.

(b) Too high bias uncertainty leads to overfitting.

Figure 4.9: Results of adding bias uncertainty and overfitting.

Although the accumulated runoff matches better with the observed runoff when the bias uncertainty is high (Figure 4.9a), the bias uncertainty parameter is now so large that the posterior results of the mean bias are unrealistic as is shown in Figure 4.9b. This results in a very chaotic accumulated runoff which is very high, far away from runoff observations, but decreases towards these stations to match the runoff observation. In order to prevent overfitting a validation step can be conducted. This can be done by removing a runoff observation and executing the model. If there is no overfitting, the accumulated runoff should still more or less match with the removed runoff observations. When overfitting occurs, the accumulated runoff will not match with the removed runoff observation.

5

Conclusions

Over the course of the previous chapters insight is gained on how probabilistic modeling works and how some hydrological processes can be described. In this chapter the conclusions regarding the research question '*how to provide a spatial estimate of the runoff and the uncertainty in the water balance using a probabilistic graph?*' are discussed.

Spatial estimation of the accumulated runoff

The structure of a factor graph represents the model structure which consists of physical knowledge. By supplying this graph with estimations of forcing data and runoff observations (observed variables which are Gaussian distributed) and solving it by means of message passing, the unobserved variables can be determined. The result is a posterior estimation for the true accumulated runoff in each cell of the distributed model. The posterior estimation includes all prior data (forcing data and runoff observations) as well as the physical knowledge (water balances, positivity constraints).

In the streams in the vicinity of runoff observations, the uncertainty of the estimation decreases due to the addition of knowledge from the observations. In the example applied in this research this influence is small and dissipates fast as the distance to the runoff observation increases.

Including sources of uncertainty

In this research, uncertainty was successfully added to the four components of the model (forcing uncertainty, model structure uncertainty, parameter uncertainty, validation data uncertainty). The introduction of uncertainty in the forcing data enables the model to make use of different data sources, regardless of their quality, as long as the uncertainty is well assessed.

Moreover, the model can be easily 'calibrated' using different measurements. The uncertainty of runoff observations depends on the method used to estimate the runoff. More precise methods of runoff estimation will have a higher influence on the model result. Given the increase of available hydrological data, probabilistic modeling is the way forward in order to incorporate more of the available data.

Propagation of information to all data

The focus in this research was to achieve a spatial estimation of the accumulated runoff, which was done by combining forcing data, runoff observations and physical knowledge (water balance and constraints). Not only was the accumulated runoff estimated, the belief about the precipitation and evaporation data was updated at the same time. This is a form of data assimilation and can improve the knowledge of this data, by using the physical knowledge combined with other data sources.

The bias as a spatial estimation of the water balance error

The posterior bias gives an indication of how well the prior data fits with the observations and the positivity constraints. The closer the posterior mean bias in a cell is to the prior mean bias (zero in this case), the better the applied water balance represent the processes occurring in that cell. Upon investigation of the bias in the Volta basin, a higher bias is found in three regions. These regions are either relative dry and in the vicinity of a large river, in the flat planes around lake Volta and in the delta area close to the mouth of the river. These local differences in bias are induced by the positivity constraints. This means that the model suggests that not enough water is available to evaporate, causing a higher value for the bias. This can be interpreted as faulty data or a bad representation of the processes.

In this case, both are the case. The resolution of the evaporation data is rather course compared with the elevation dataset. This causes that a pixel of the evaporation overlaps both water and land, causing high evaporation values over land. But also outside the influence of the courser resolution we see a higher bias due to positivity constraints because of available water limitations. The most probable explanation for the high bias is the absence of ground water flows in the model. All regions with a high bias are close to a large water body from which ground water can flow to the cells that are currently subjected to positivity constraints. Including ground water flows in the model could reduce the bias, improving the model.

Spatial distribution of uncertainty

Another strength of this model is the spatial distribution of uncertainty. By combining multiple observations and physical constraints, an insight in the quality of data in different areas of the basin can be gained. If the prior accumulated runoff matches well with the runoff observation, uncertainty in the upstream forcing data will decrease more, compared to the case where the accumulated runoff does not match well with the observation. The spatial distribution of uncertainty can provide us with knowledge about the performance of data products in different regions, for example if a satellite precipitation product performs well in a mountainous area compared to flatlands. Unfortunately only a small influence of the runoff observations on the uncertainty of other data was found, especially with respect to the influence of the positivity constraints, which reduce uncertainty greatly when applied.

6

Recommendations

During the research several problems were encountered which can have an influence on the model quality. This chapter addresses these issues and the recommendations for future research.

Positivity constraints and representation of constrained variables

The positivity constraints used in this research did not always show the intended effects. When constraining, the result is a symmetrical Gaussian with an increased mean value and a decreased variance in order to prevent the probability on a negative value. This causes that the probability on a constrained variable close to zero is very low. While in the used basin this poses no problems for precipitation and evaporation (they do not approach zero), it does so for the accumulated runoff in regions with a low local runoff. An area with a prior mean precipitation and evaporation which are equal implies that there is no runoff from this area, but by constraining this variable the probability of zero runoff is almost non-existent.

A better choice for the representation of those variables would be a Gamma distribution. This distribution has zero probability in the negative domain and does not need to be constrained. A Gamma distribution can have its highest probability at zero, decreasing in probability as the variable value goes up. When the variable mean is close to zero it will show an asymmetric distribution, but as the mean value increases compared to the variance, it will approach the Gaussian distribution.

The problem with the Gamma distribution is that accumulation of independent Gamma distributions do not provide a new Gamma distribution but can be approximated (Murakami 2015). This research was conducted to explore the possibilities of the framework of graphical models in runoff modeling, therefore the less complicated and exact Gaussian distribution was used.

Quantification of uncertainty

During the literature study it became clear how hard it is to gather information about the uncertainty of data sets. Some data sets do not report any uncertainty, others report it in many different quantities which can not easily be compared with one another. More research should be conducted for a better quantification of the data uncertainty.

In this research the uncertainty was considered to be relative to the mean value of the variable. The strength of this model is that it can include a spatially distributed uncertainty, which can be used to improve

the model results. Satellite precipitation estimates in mountains can be more uncertain (due to high spatial variation of precipitation) than estimates in flat plains where point measurements are used to calibrate and validate the data. This spatial variation of uncertainty can be incorporated using the model used in this research.

Convergence with increasing graph complexity

The introduction of approximations and cycles into the model increases the amount of iterations needed to converge. Since this model uses mostly linear relations (local water balance and flow accumulation), the amount of iterations is still manageable. When using non-linear relations, the model will apply more approximations resulting in a slower convergence.

Many hydrological models use time steps to calculate time series, where variables in time step $t + 1$ are dependent on variables in time step t . This could be visualized by one factor graph per time step, plus factors connecting the dependent variables in both graphs. If there are multiple connections between the two time steps, this causes additional cycles in the factor graph. These extra cycles will require the model to iterate more often in order to converge.

Uncertainty in the Digital Elevation Model

One source of uncertainty which is not included in this research is the uncertainty of the Digital Elevation Model. The DEM is the base for the structure of the factor graph and its uncertainty can not be included in the model. If uncertainties in the DEM lead to errors in the flow direction map this means that the knowledge is not propagated in the right way.

Examples are the cross-over of the flow path to a different river and divergence of flows as discussed in Chapter 4. The former can be corrected manually, the latter could also be implemented manually, however the proportionality of the divergent flows has to be estimated, which introduces a new source of uncertainty.

References

- AghaKouchak, A., Nasrollahi, N. & Habib, E. (2009), 'Accounting for uncertainties of the TRMM satellite estimates', *Remote Sensing* **1**(3), 606–619.
- Andah, W. E. I., van de Giesen, N. & Biney, C. a. (2003), 'Water, climate, food, and environment in the Volta Basin', *Adaptation strategies to changing environments. Contribution to the ADAPT project* <http://www.weap21.org/downloads/ADAPTVolta.pdf>.
- Belov, D. I. & Armstrong, R. D. (2009), 'Distributions of Kullback – Leibler Divergence and Its Application for the LSAT'.
URL: <http://www.lzac.org/lzacresources/research/rr/pdf/rr-09-02.pdf>
- Bishop, C. M. (2006), *Pattern Recognition and Machine Learning*, Vol. 4.
URL: <http://www.library.wisc.edu/selectedtocs/bg0137.pdf>
- Bromiley, P. A. (2003), Products and Convolutions of Gaussian Distributions, Technical Report 3.
URL: <http://tina.wiau.man.ac.uk/docs/memos/2003-003.pdf>
- Broxton, P. D., Zeng, X., Sulla-Menashe, D. & Troch, P. A. (2014), 'A global land cover climatology using MODIS data', *Journal of Applied Meteorology and Climatology* **53**(6), 1593–1605.
- Bundesanstalt für Gewässerkunde (2014), 'Global Runoff Data Centre'.
URL: http://www.bafg.de/GRDC/EN/Home/homepage_node.html
- Castrup, H. (2001), 'Distributions for uncertainty analysis', *Proceedings of the International Dimensional Workshop* **12**(May 2004).
URL: http://www.isgmax.com/articles_papers/distributions_for_uncertainty_analysis_-_revised.pdf
- Chen, Y. & Han, D. (2016), 'Big data and hydroinformatics', *Journal of Hydroinformatics* pp. 1–16.
URL: <http://jh.iwaponline.com/cgi/doi/10.2166/hydro.2016.180>
- Danielson, J. J. & Gesch, D. B. (2011), Global Multi-resolution Terrain Elevation Data 2010, Technical report.
- Dekking, F., Kraaikamp, C., Lopuhaä, H. & Meester, L. (2005), *A Modern Introduction to Probability and Statistics*, Springer.
URL: <http://amstat.tandfonline.com/doi/pdf/10.1198/tech.2007.s502>
- Di Baldassarre, G. & Montanari, a. (2009), 'Uncertainty in river discharge observations: a quantitative analysis', *Hydrology and Earth System Sciences Discussions* **6**(1), 39–61.
- Domeneghetti, a., Castellarin, a. & Brath, a. (2012), 'Assessing rating-curve uncertainty and its effects on hydraulic model calibration', *Hydrology and Earth System Sciences* **16**(4), 1191–1202.
- Frey, B. J., Kschischang, F. R., Loeliger, H.-A. & Wiberg, N. (1998), 'Factor graphs and algorithms', *Proceedings 35th Allerton Conference on Communications, Control, and Computing* pp. 666—680.
URL: http://www.psi.toronto.edu/psi/pubs2/1999_and_before/134.pdf

- Funk, C. C., Peterson, P. J., Landsfeld, M. E., Pedreros, D. H., Verdin, J. P., Rowland, J. D., Romero, B. E., Husak, G. J., Michaelsen, J. C. & Verdin, a. P. (2014), 'A quasi-global precipitation time series for drought monitoring', *U.S. Geological Survey Data Series* **832**, 4.
- GRASS Development Team (2012), 'Geographic Resources Analysis Support System (GRASS GIS) Software'.
URL: <http://grass.osgeo.org>
- Greene, W. H. (2003), *Econometric Analysis*, Pearson Education, Inc., New Jersey.
- Guerschman, J. P., Van Dijk, A. I. J. M., Mattersdorf, G., Beringer, J., Hutley, L. B., Leuning, R., Pipunic, R. C. & Sherman, B. S. (2009), 'Scaling of potential evapotranspiration with MODIS data reproduces flux observations and catchment water balance observations across Australia', *Journal of Hydrology* **369**(1-2), 107–119.
URL: <http://dx.doi.org/10.1016/j.jhydrol.2009.02.013>
- Herschey, R. W. (2009), *Streamflow Measurement*, third edit edn, Taylor & Francis.
- Karimi, P. & Bastiaanssen, W. G. M. (2015), 'Spatial evapotranspiration, rainfall and land use data in water accounting – Part 1: Review of the accuracy of the remote sensing data', *Hydrology and Earth System Sciences Discussions* **19**(1), 507–532.
URL: <http://www.hydrol-earth-syst-sci-discuss.net/11/1073/2014/>
- Karimi, P., Bastiaanssen, W. G. M., Sood, A., Hoogeveen, J., Peiser, L., Bastidas-Obando, E. & Dost, R. J. (2015), 'Spatial evapotranspiration, rainfall and land use data in water accounting – Part 2: Reliability of water accounting results for policy decisions in the Awash basin', *Hydrology and Earth System Sciences Discussions* **19**(1), 533–550.
URL: <http://www.hydrol-earth-syst-sci-discuss.net/11/1125/2014/hessd-11-1125-2014.html>
- Kavetski, D., Franks, S. W. & Kuczera, G. (2002), 'Confronting input uncertainty in environmental modelling', *Science* **6**, 49–68.
- Kavetski, D., Kuczera, G. & Franks, S. W. (2006), 'Bayesian analysis of input uncertainty in hydrological modeling: 1. Theory', *Water Resources Research* **42**(3), 1–9.
- Koller, D. & Friedman, N. (2009), *Probabilistic Graphical Models: Principles and Techniques*.
- Kschischang, F. R., Frey, B. J. & Loeliger, H. a. (2001), 'Factor graphs and the sum-product algorithm', *IEEE Transactions on Information Theory* **47**(2), 498–519.
- Liu, Y. & Gupta, H. V. (2007), 'Uncertainty in hydrologic modeling: Toward an integrated data assimilation framework', *Water Resources Research* **43**(7), 1–18.
- Lyon, A. (2014), 'Why are normal distributions normal?', *British Journal for the Philosophy of Science* **65**(3), 621–649.
- McMillan, H., Jackson, B., Clark, M., Kavetski, D. & Woods, R. (2011), 'Rainfall uncertainty in hydrological modelling: An evaluation of multiplicative error models', *Journal of Hydrology* **400**(1-2), 83–94.
URL: <http://dx.doi.org/10.1016/j.jhydrol.2011.01.026>
- Microsoft Research Cambridge (n.d.), 'Infer.NET code documentation'.
URL: http://research.microsoft.com/en-us/um/cambridge/projects/infernet/codedoc/html/R_Project_Infer.htm

- Minka, T. P. (2001a), A family of algorithms for approximate Bayesian inference, PhD thesis.
URL: papers2://publication/uuid/37D3C7DD-C308-4279-86AC-52057DE5CB29
- Minka, T. P. (2001b), Expectation Propagation for Approximate Bayesian Inference, in 'Proceedings of the Seventeenth Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-01)', Vol. 17, pp. 362–369.
URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.86.1319&rep=rep1&type=pdf>
- Minka, T. P., Winn, J., Guiver, J., Webster, S., Zaykov, Y., Yangel, B., Spengler, A. & Pronskill, J. (2014), 'Infer.NET 2.6'.
URL: <http://research.microsoft.com/infernet>
- Moradkhani, H. & Sorooshian, S. (2009), 'General Review of Rainfall-Runoff Modeling : Model Calibration , Data Assimilation , and Uncertainty Analysis', *Review Literature And Arts Of The Americas* **63**, 1–24.
URL: <http://www.springerlink.com/index/q04565241t811818.pdf>
- Mu, Q., Zhao, M. & Running, S. W. (2011), 'Improvements to a MODIS global terrestrial evapotranspiration algorithm', *Remote Sensing of Environment* **115**(8), 1781–1800.
URL: <http://dx.doi.org/10.1016/j.rse.2011.02.019>
- Murakami, H. (2015), 'Approximations to the distribution of sum of independent non-identically gamma random variables', *Mathematical Sciences* **9**(4), 205–213.
URL: <http://link.springer.com/10.1007/s40096-015-0169-2>
- Olivera, F., Reed, S. & Maidment, D. (1998), HEC-PrePro v. 2.0: An ArcView Pre-Processor for HEC's Hydrologic Modeling System, in 'ESRI User's Conference', University of Texas at Austin - Center for Research in Water Resources, Austin, Texas.
URL: <http://www.crrw.utexas.edu/gis/gishyd98/runoff/webfiles/esri98/p400.htm>
- QGIS Development Team (2015), 'QGIS 2.8 Geographic Information System User Guide'.
URL: <http://qgis.osgeo.org>
- R Core Team (2015), 'R: A Language and Environment for Statistical Computing'.
URL: <http://www.r-project.org/>
- Rajaram, H., Bahr, J., Blöschl, G., Cai, X., Scott Mackay, D., Michalak, A. M., Montanari, A., Sanchez-Villa, X. & Sander, G. (2015), 'A reflection on the first 50 years of Water Resources Research', *Water Resources Research* (March 1965), n/a–n/a.
URL: <http://doi.wiley.com/10.1002/2015WR018089>
- Refsgaard, J. C., van der Sluijs, J. P., Brown, J. & van der Keur, P. (2006), 'A framework for dealing with uncertainty due to model structure error', *Advances in Water Resources* **29**(11), 1586–1597.
- Schoups, G. H. W. (2015), Belief propagation for inference in linear-Gaussian models with applications in hydrology.
- Tomkins, K. M. (2014), 'Uncertainty in streamflow rating curves: Methods, controls and consequences', *Hydrological Processes* **28**(3), 464–481.
- Velpuri, N., Senay, G., Singh, R., Bohms, S. & Verdin, J. (2013), 'A comprehensive evaluation of two MODIS evapotranspiration products over the conterminous United States: Using point and gridded FLUXNET and

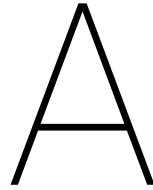
water balance ET', *Remote Sensing of Environment* **139**, 35–49.

URL: <http://linkinghub.elsevier.com/retrieve/pii/S0034425713002253>

Vinga, S. & Almeida, J. S. (2004), 'Renyi continuous entropy of DNA sequences', *Journal of Theoretical Biology* **231**(3), 377–388.

Winn, J. & Minka, T. P. (2007), 'Expectation Propagation & Variational Message Passing'.

URL: http://videolectures.net/abi07_winn_ipi/



Basin description

In this research the Volta basin is used to investigate how the model performs using a real basin and which conclusions can be drawn from it. The basin is located in West Africa and spans almost 400,000 km^2 over six countries (Mali, Burkina Faso, Benin, Togo, Ivory Coast and Ghana). There are 10 runoff observation station located throughout the basin that collect monthly data of the discharge.

An important water body in the basin is Lake Volta, located in the south. The lake can have a change of water storage over the long term data, but as all runoff observations are positioned upstream of the lake it should have no effect on the model.

The Volta has bi-modal rainfall pattern climate, semi-arid sub-humid savanna (Andah et al. 2003). This causes some streams to dry up during parts of the year. Because long term averages are used this does not pose a problem for the model.

Several data products for precipitation and actual evaporation were examined on the average value in the Volta basin (Table A.1). The examined precipitation products do not differ that much, but the evaporation products do. When applying the water balance with different data products, and calculating the total runoff from the basin, the results show a large variation in total runoff (Table A.2)

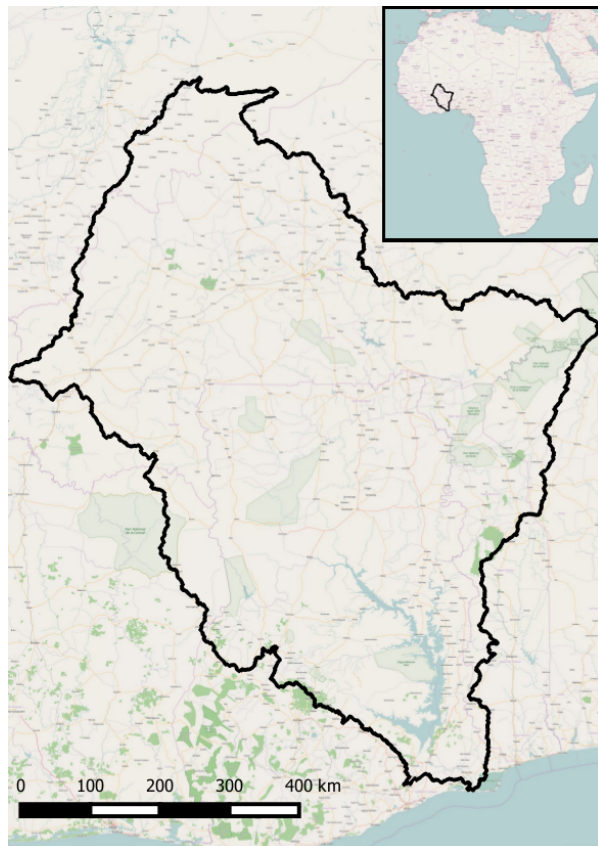


Figure A.1: Geographical map of the Volta basin ©OpenStreetMap contributors

Table A.1: General information of the Volta basin

Area [km ²]	394,196
Amount of discharge stations [-]	10
Station density [1/km ²]	39,420
Precipitation [mm]	
CHIRPS	995
TRMM	1,010
Evaporation [mm]	
CMRSET	744
FAO	856
MOD16	961

Table A.2: Discharge at the mouth of the river by combining different data products

Discharge at mouth [m³/s]			
		Precipitation	
		CHIRPS	TRMM
Evaporation	CMRSET	3,137	3,325
	FAO	1,737	1,925
	MOD16	425	612

B

Different model setups: Factor Graphs

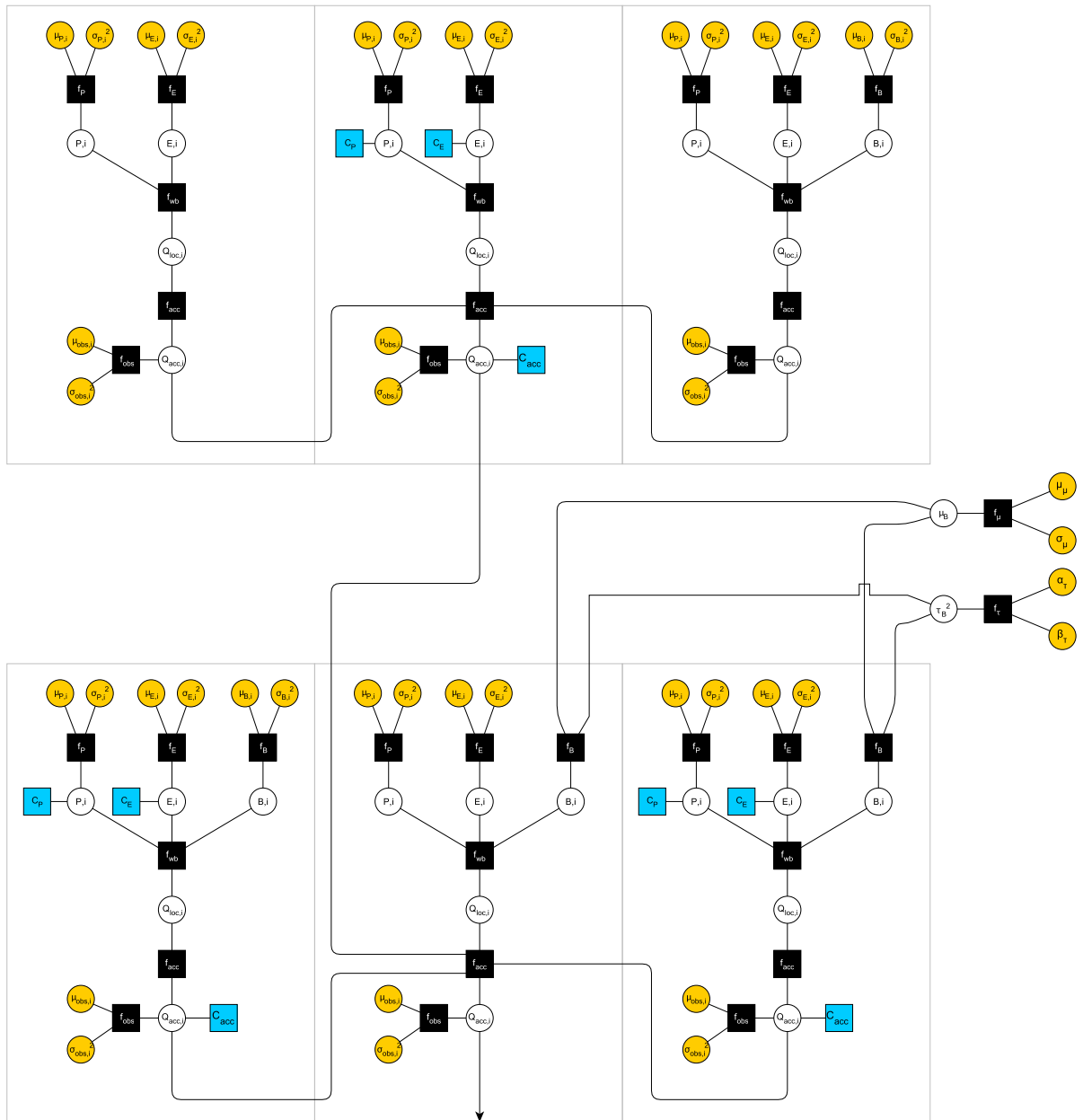


Figure B.1: 6 connected cells with each a different factor graph, representing a different model structure. In the top row from left to right model 2, 3 and 4. On the bottom row model 5, 6 and 7.

C

Different model setups: Results of testdata

C.1. Model 1

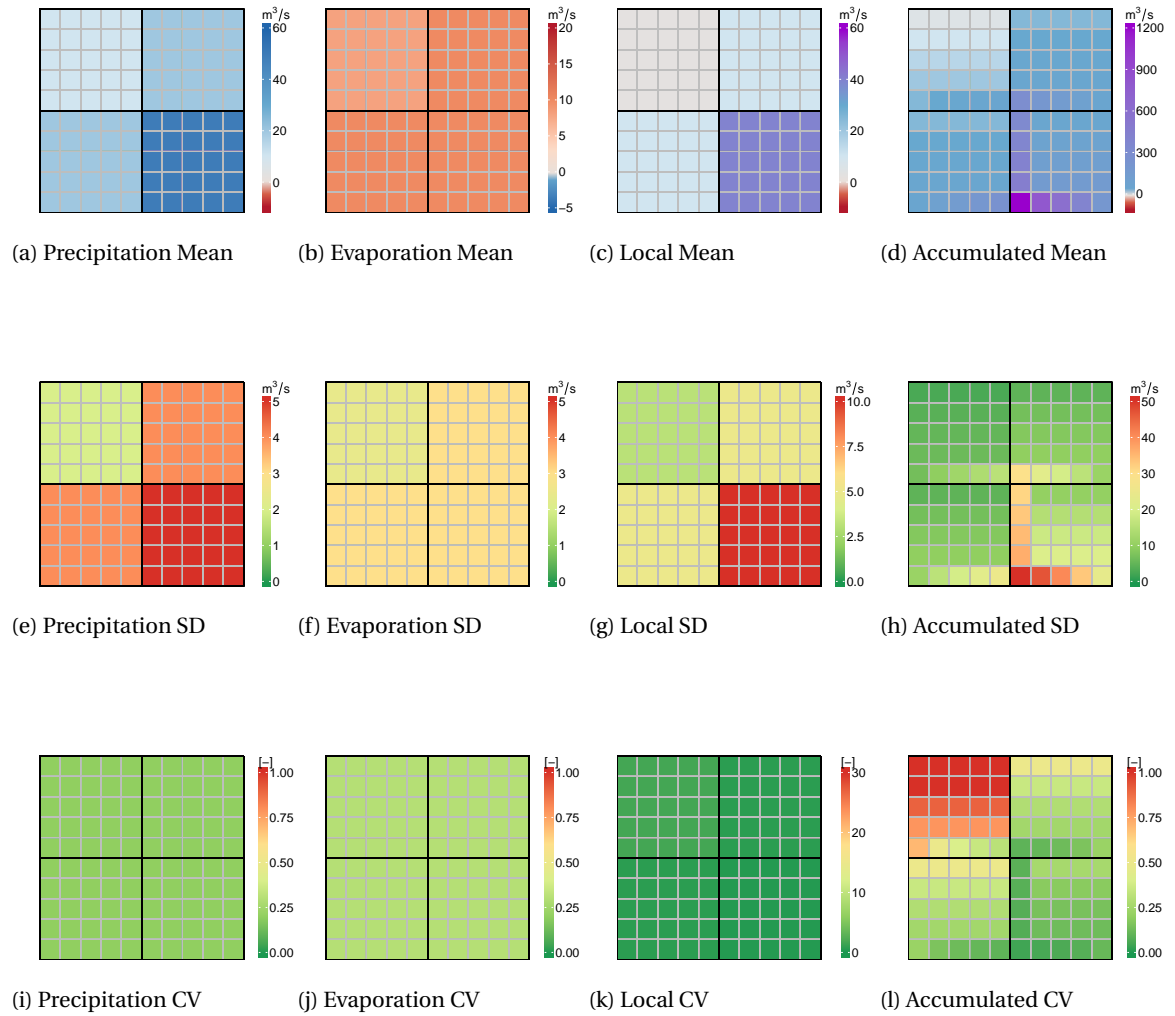


Figure C.1: Model 1, uninformed, unconstrained, unbiased

C.2. Model 2

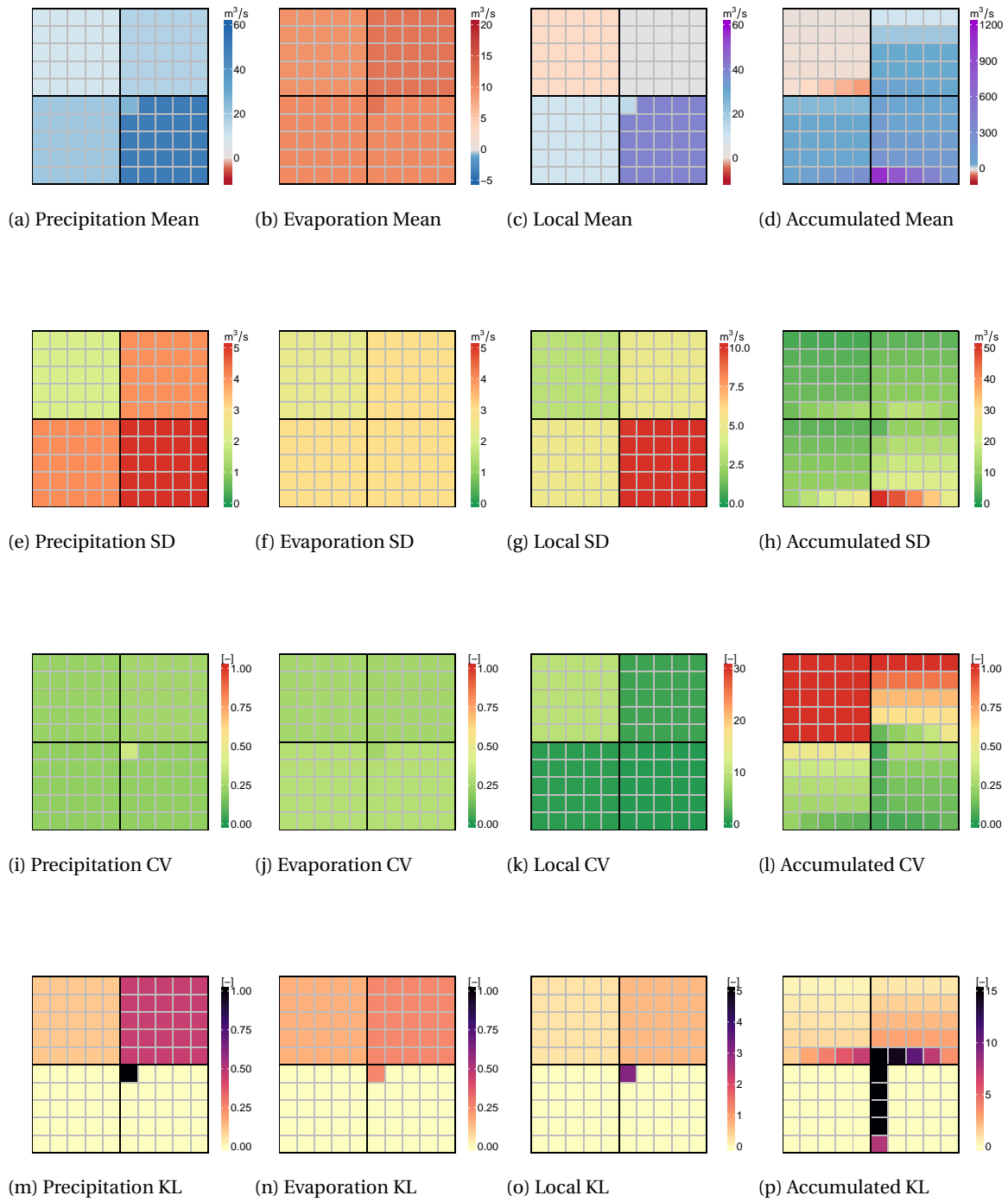


Figure C.2: Model 2, informed, unconstrained, unbiased

C.3. Model 3

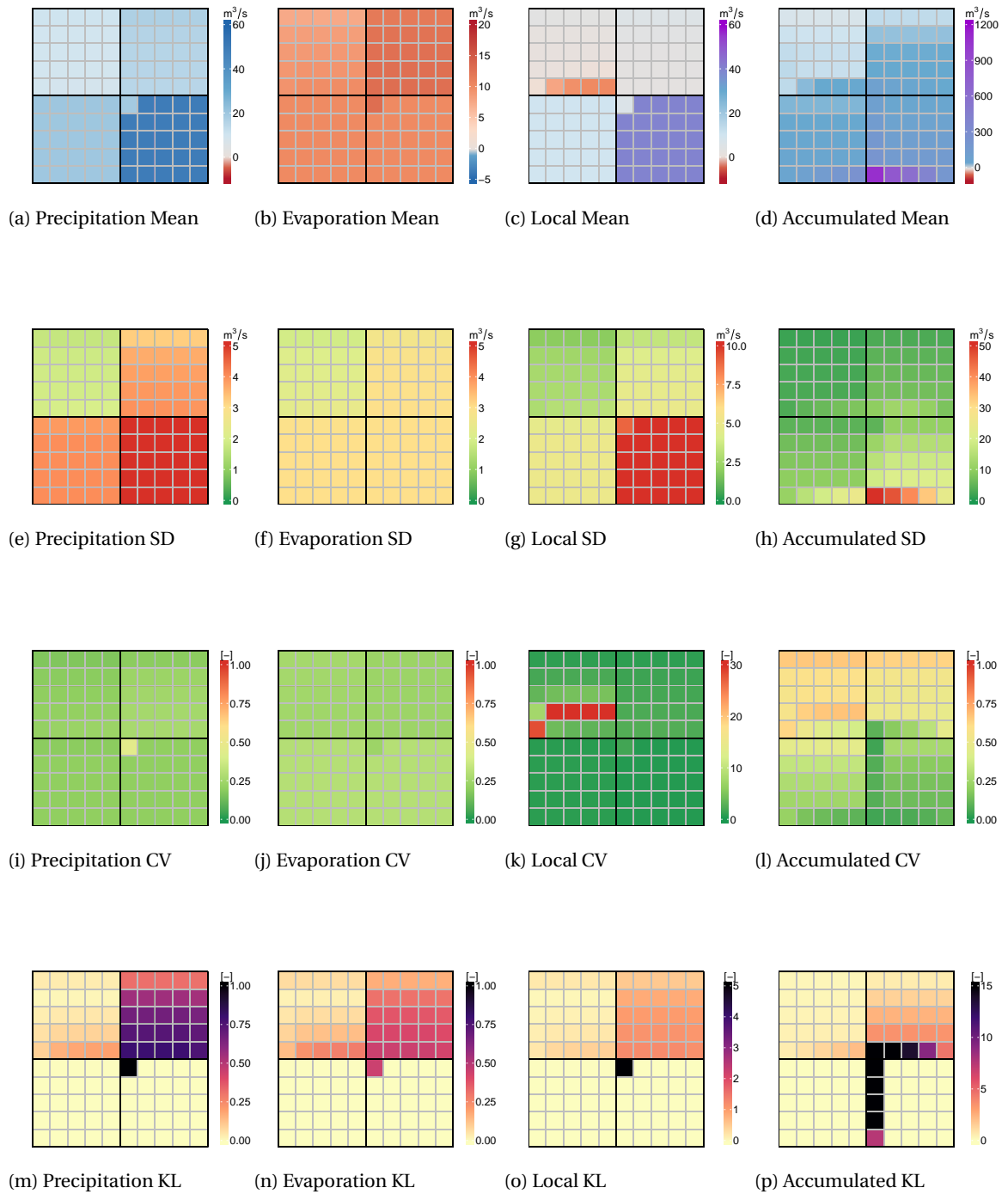


Figure C.3: Model 3, informed, constrained, unbiased

C.4. Model 4

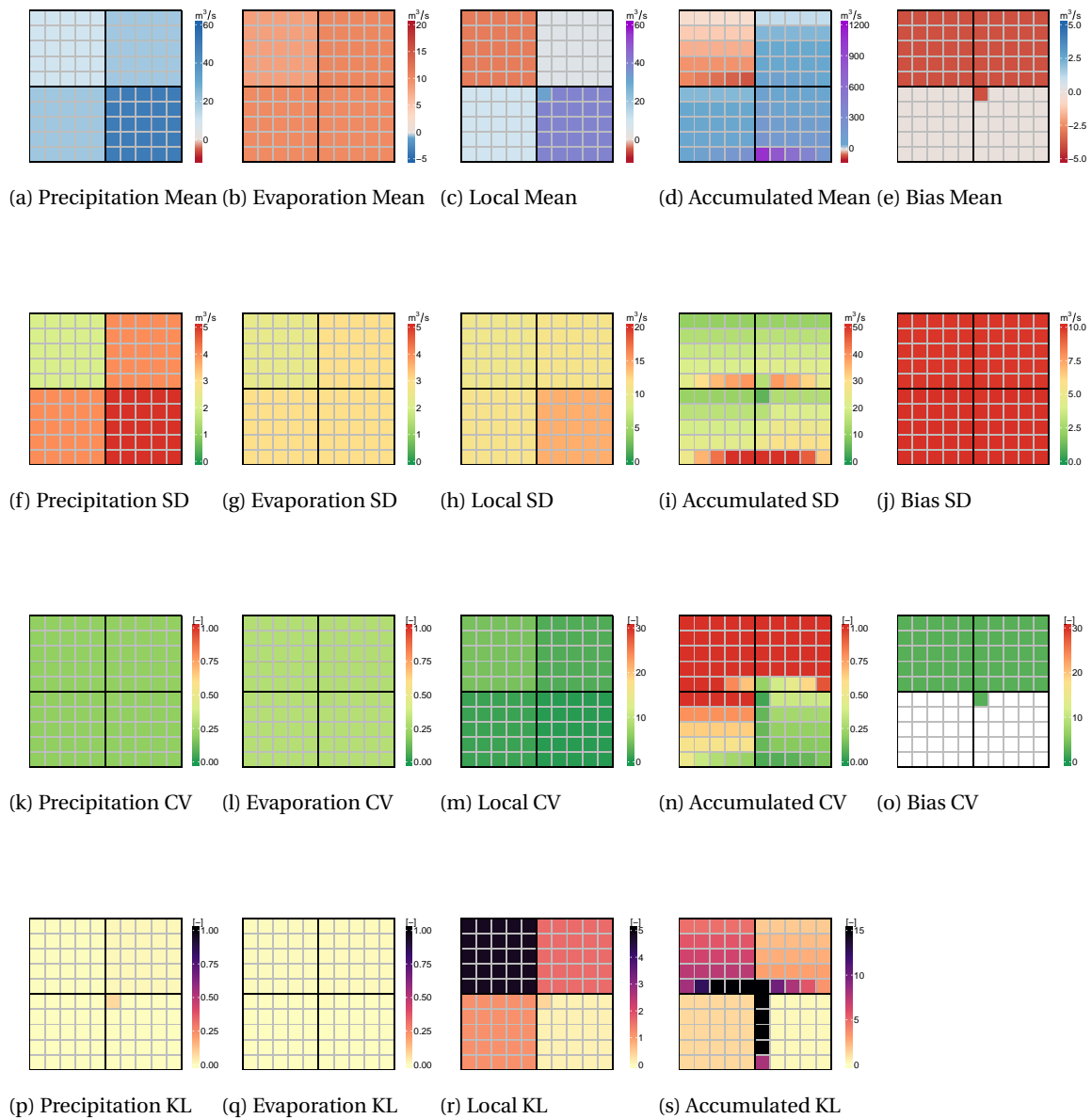


Figure C.4: Model 4, informed, unconstrained, bias without parameter uncertainty

C.5. Model 5

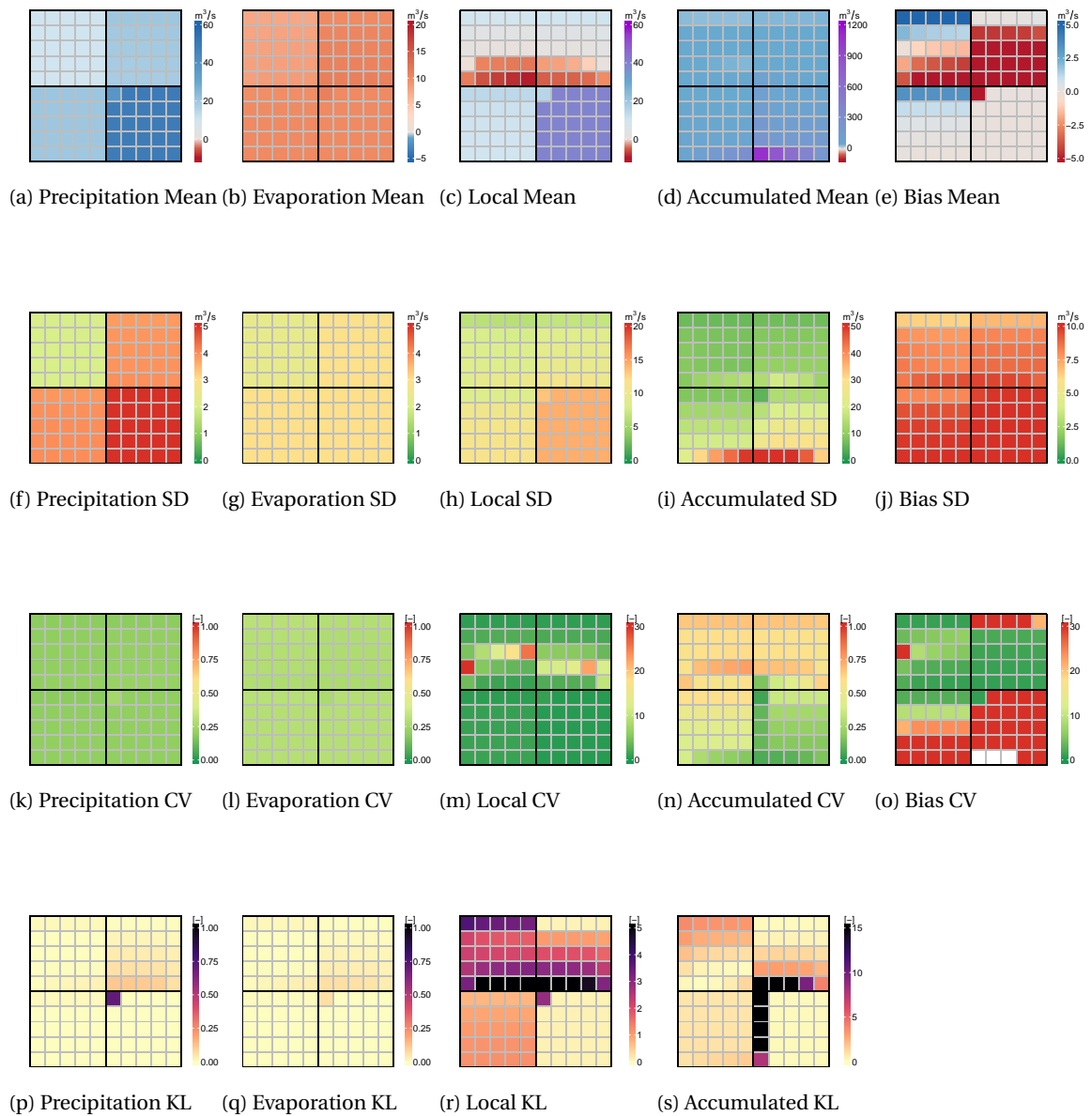


Figure C.5: Model 5, informed, constrained, bias without parameter uncertainty

C.6. Model 6

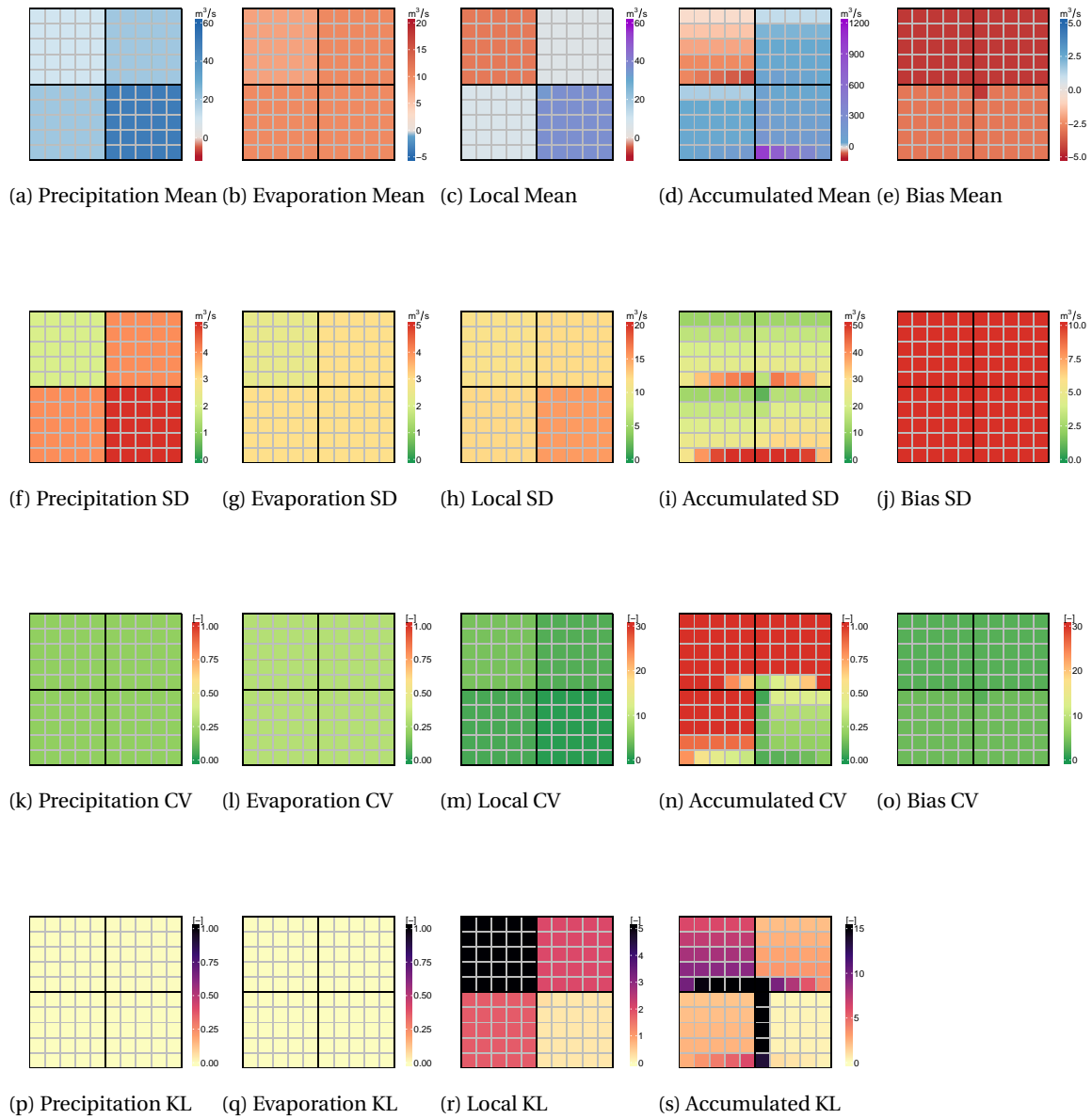


Figure C.6: Model 6, informed, unconstrained, bias with parameter uncertainty

C.7. Model 7

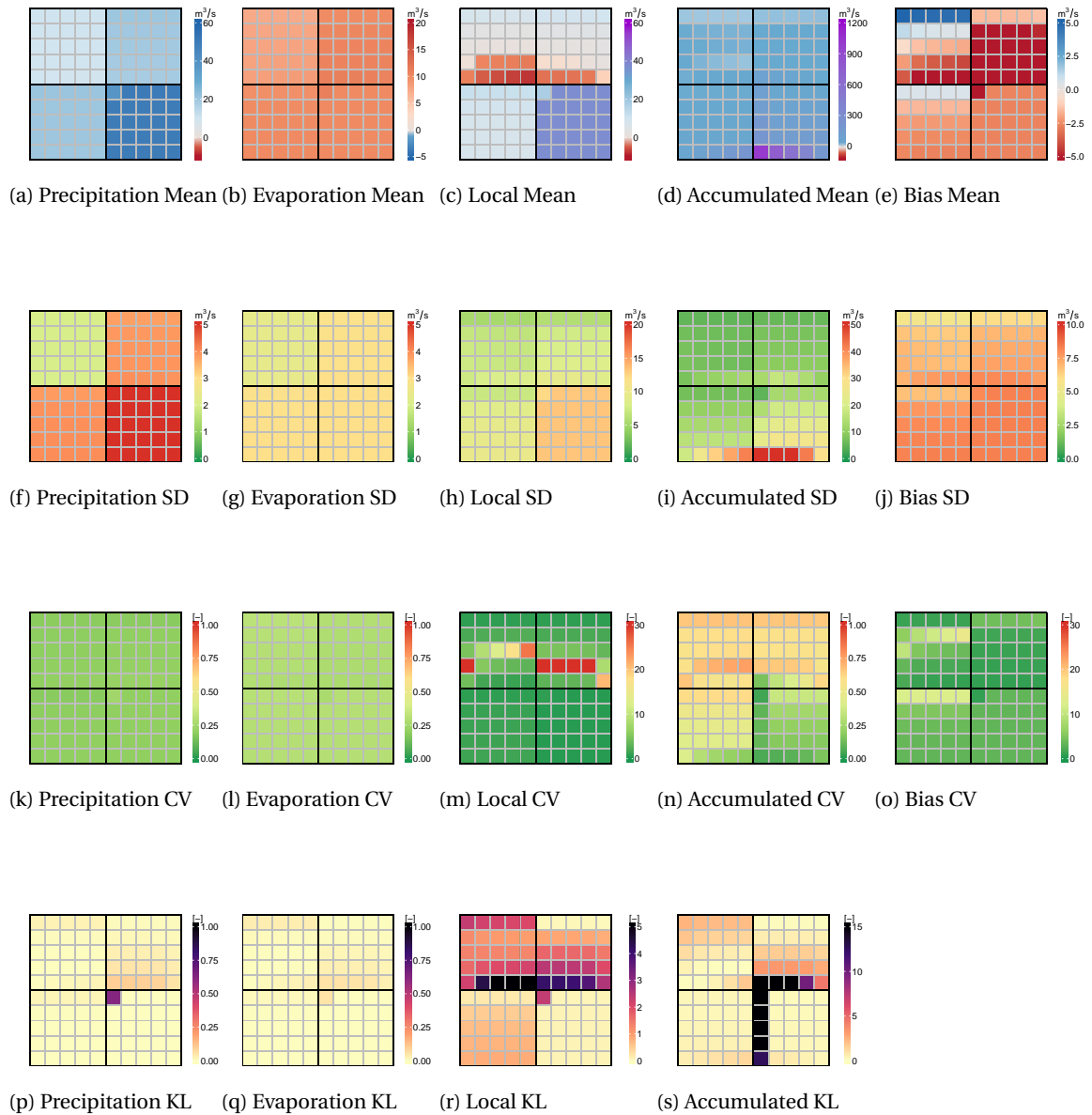
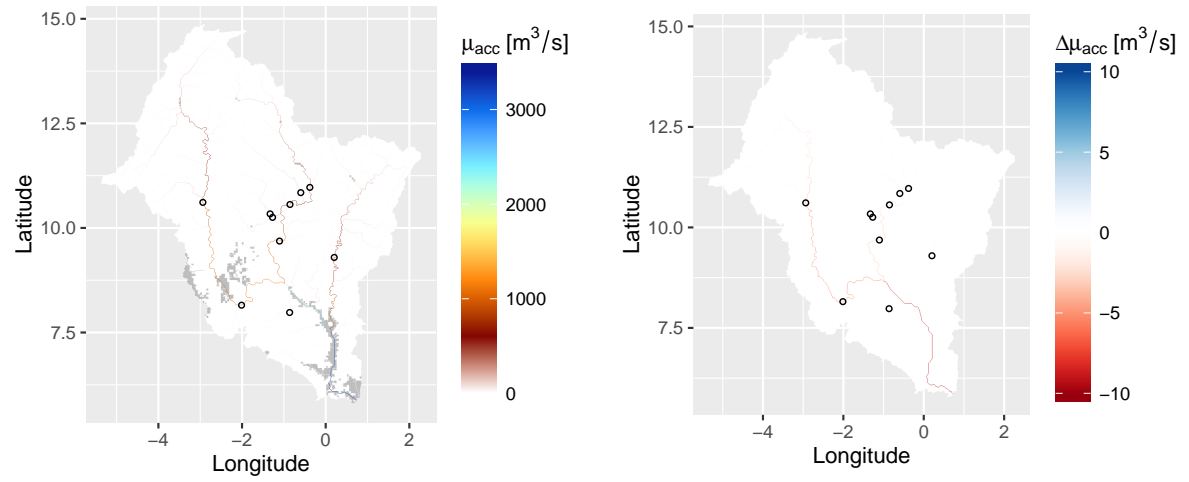


Figure C.7: Model 7, informed, constrained, bias with parameter uncertainty

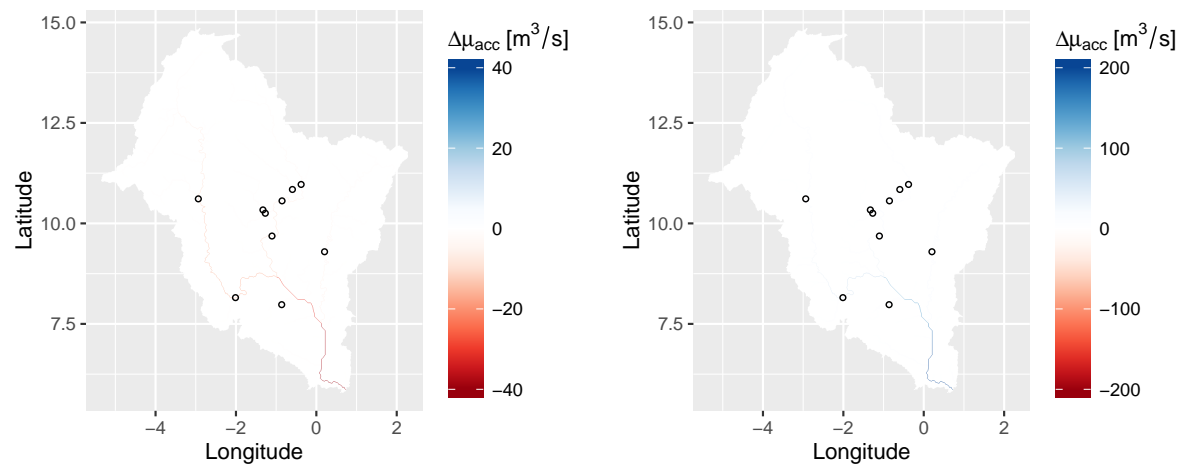
D

Results on Volta: figures



(a) Mean accumulated runoff model 1, grey areas indicate a negative value

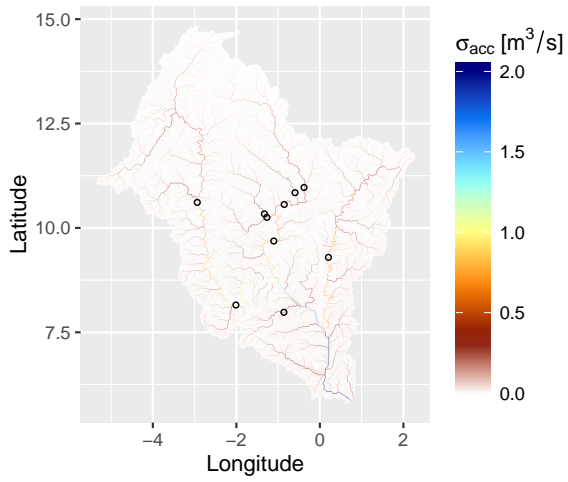
(b) Difference of mean accumulated runoff in model 2, compared with model 1



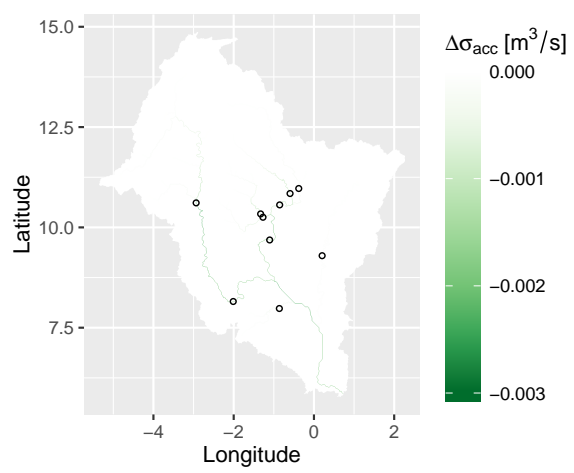
(c) Difference of mean accumulated runoff in model 6, compared with model 1

(d) Difference of mean accumulated runoff in model 7, compared with model 1

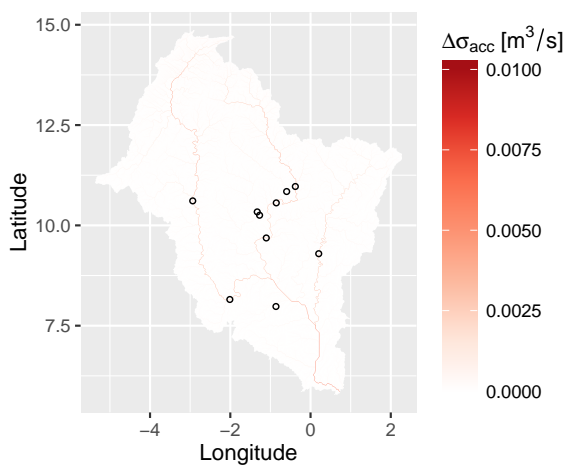
Figure D.1: Mean accumulated runoff



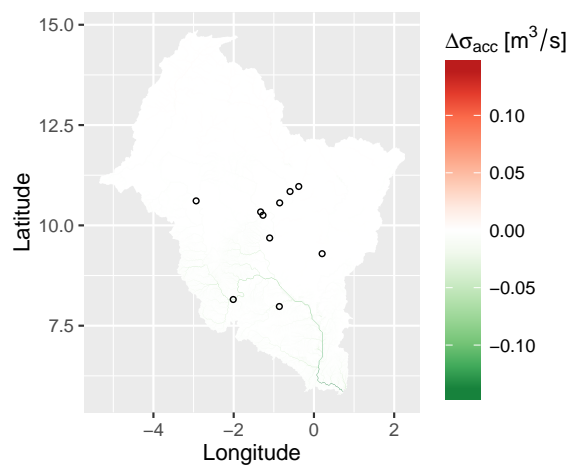
(a) SD accumulated runoff in model 1



(b) Difference of SD accumulated runoff in model 2, compared with model 1

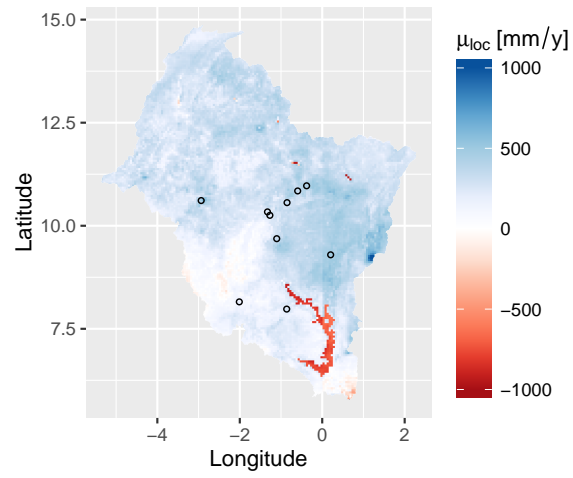


(c) Difference of SD accumulated runoff in model 6, compared with model 1

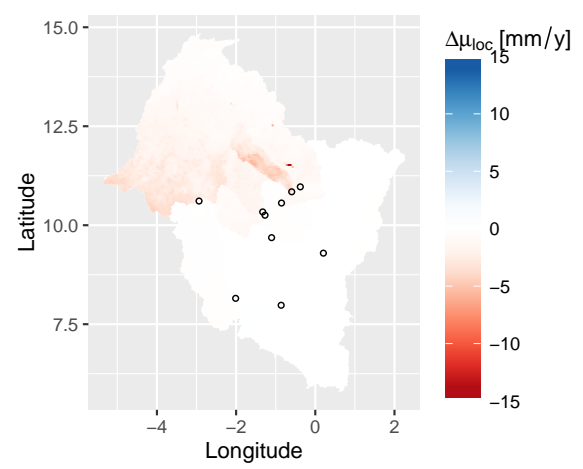


(d) Difference of SD accumulated runoff in model 7, compared with model 1

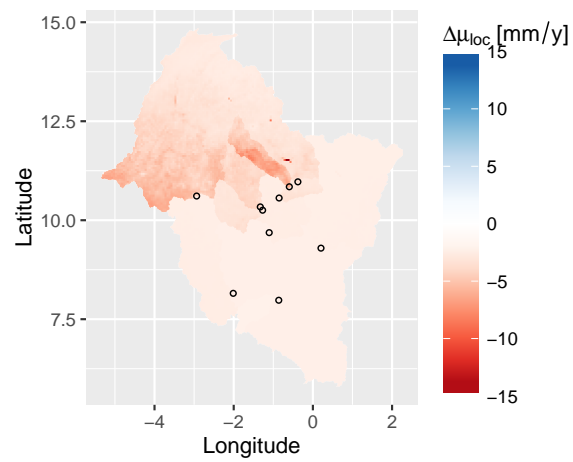
Figure D.2: Standard deviation accumulated runoff



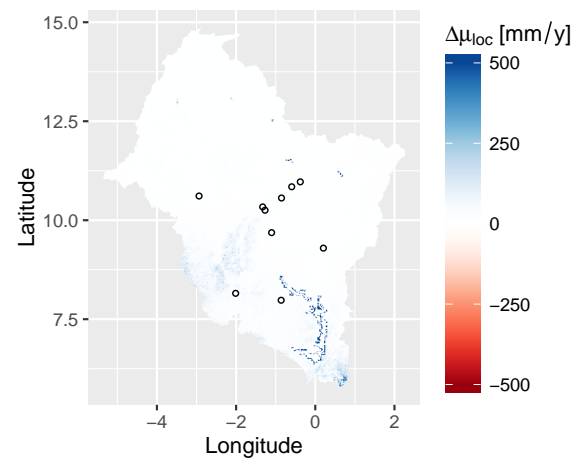
(a) Mean local runoff model 1



(b) Difference of mean local runoff in model 2, compared with model 1

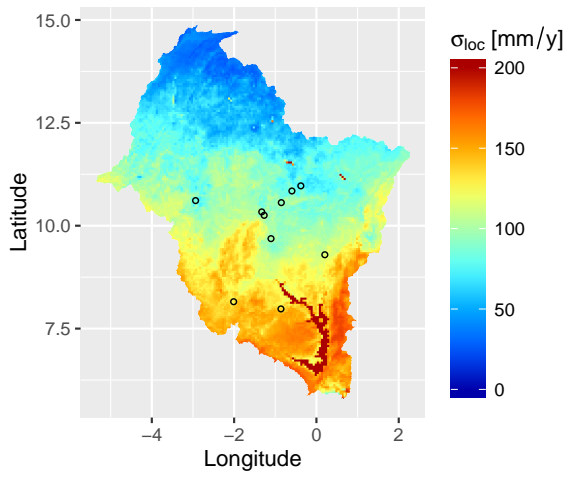


(c) Difference of mean local runoff in model 6, compared with model 1

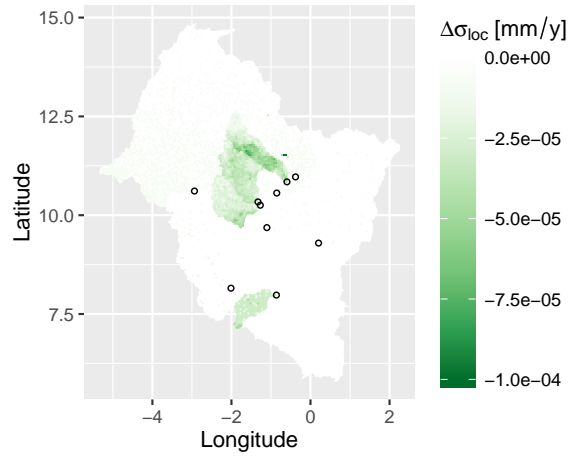


(d) Difference of mean local runoff in model 7, compared with model 1

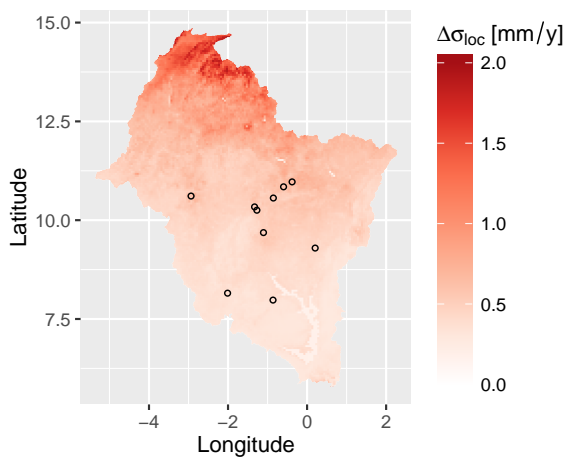
Figure D.3: Mean local runoff



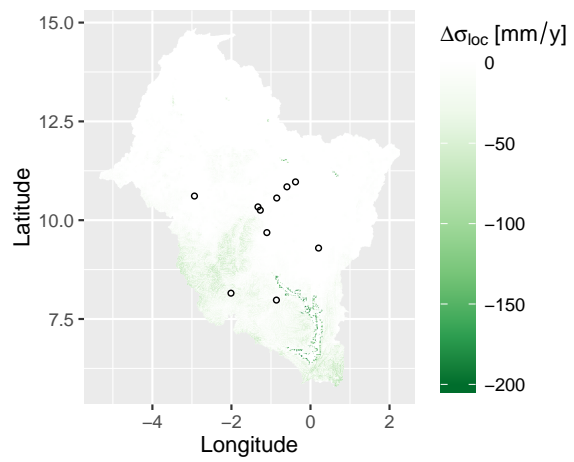
(a) SD local runoff in model 1



(b) Difference of SD local runoff in model 2, compared with model 1

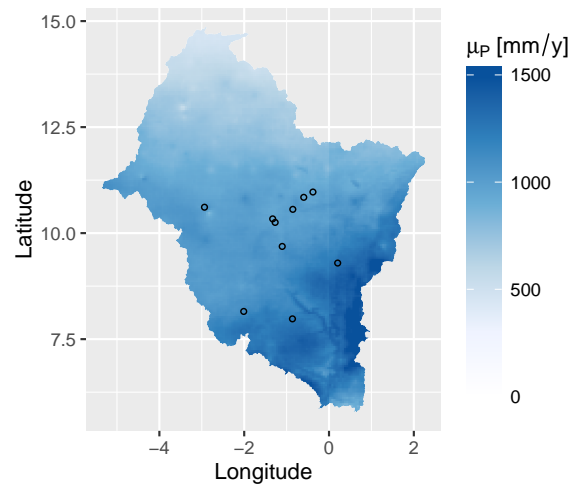


(c) Difference of SD local runoff in model 6, compared with model 1

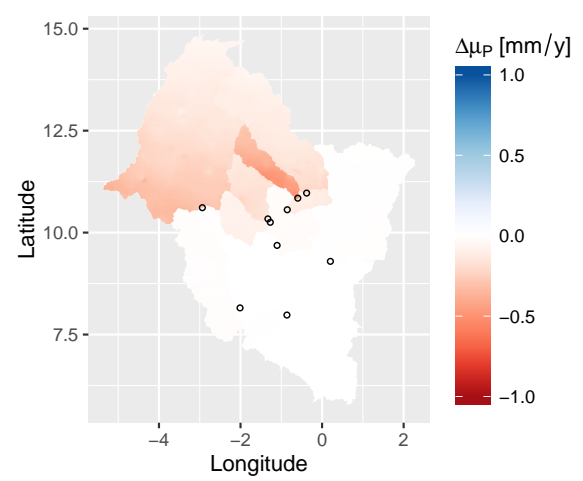


(d) Difference of SD local runoff in model 7, compared with model 1

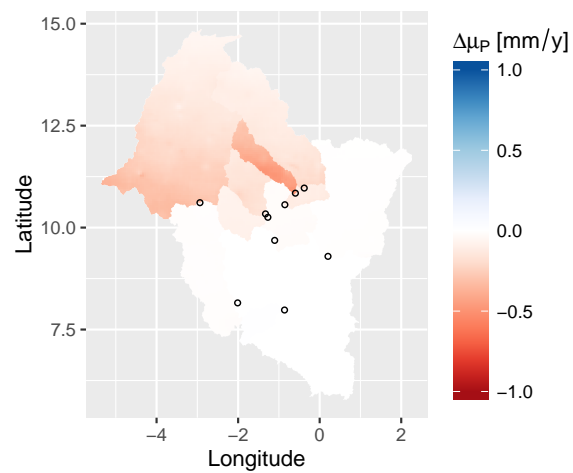
Figure D.4: Standard deviation local runoff



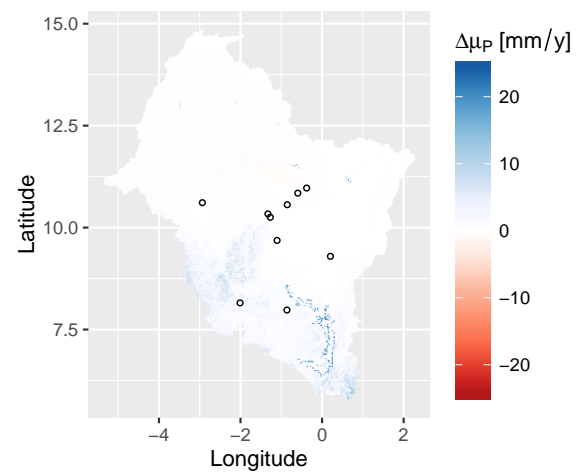
(a) Mean precipitation runoff model 1



(b) Difference of mean precipitation runoff in model 2, compared with model 1

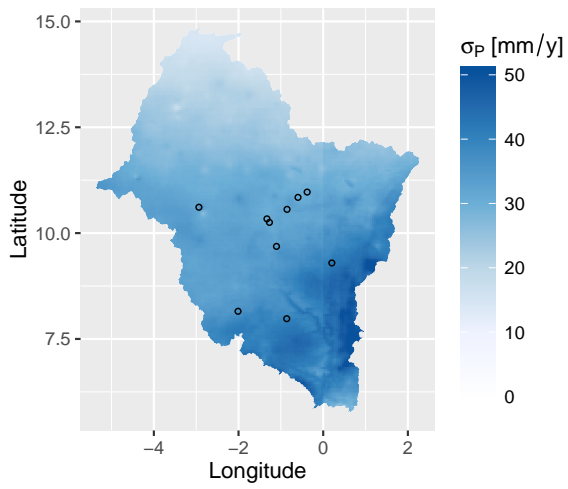


(c) Difference of mean precipitation runoff in model 6, compared with model 1

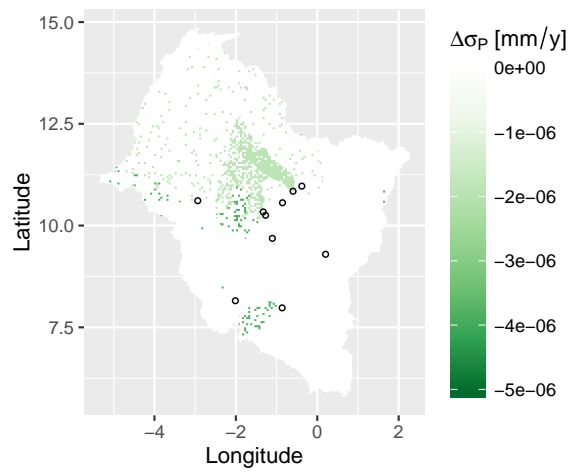


(d) Difference of mean precipitation runoff in model 7, compared with model 1

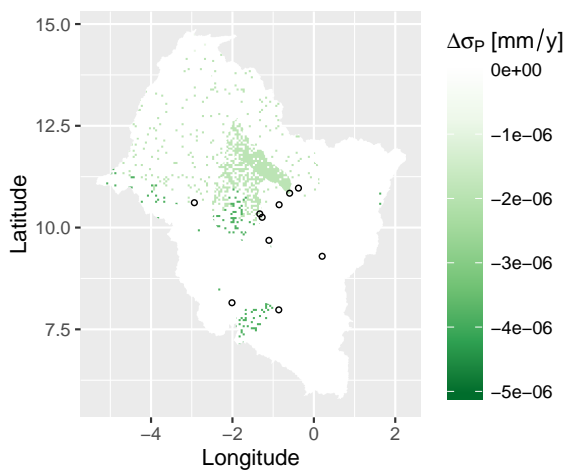
Figure D.5: Mean precipitation runoff



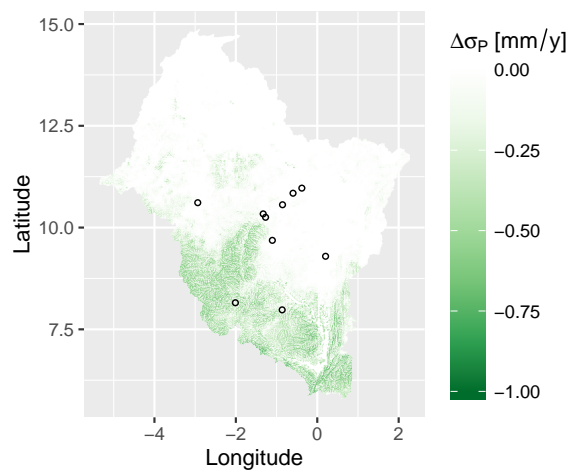
(a) SD precipitation in model 1



(b) Difference of SD precipitation in model 2, compared with model 1

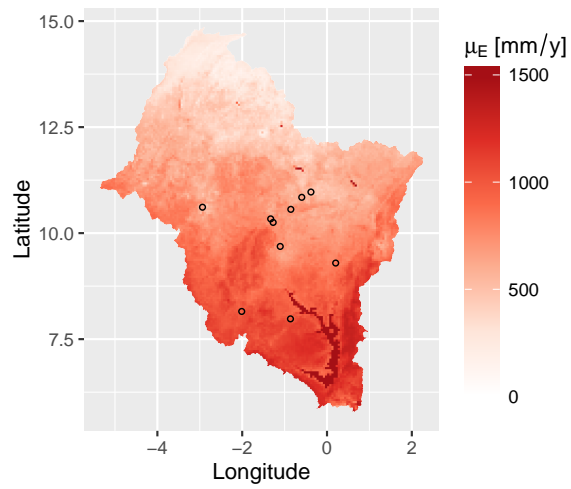


(c) Difference of SD precipitation in model 6, compared with model 1

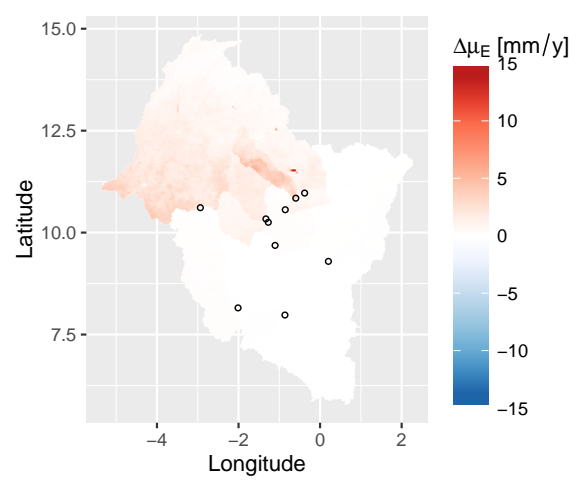


(d) Difference of SD precipitation in model 7, compared with model 1

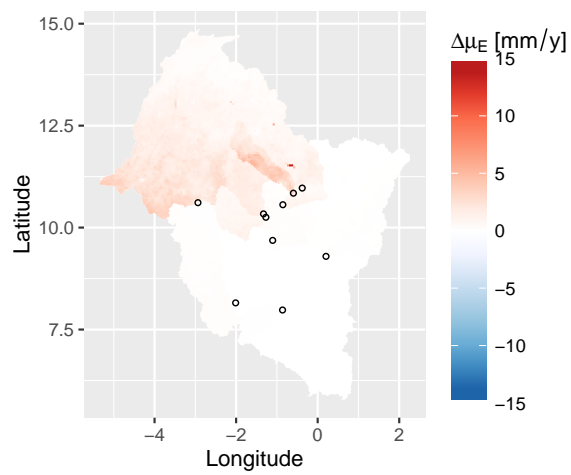
Figure D.6: Standard deviation precipitation



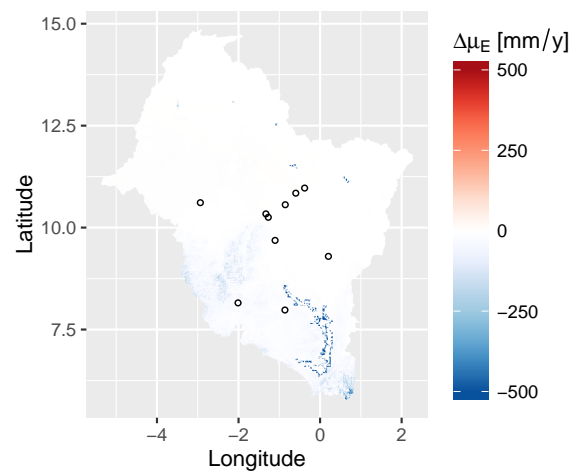
(a) Mean evaporation model 1



(b) Difference of mean evaporation in model 2, compared with model 1

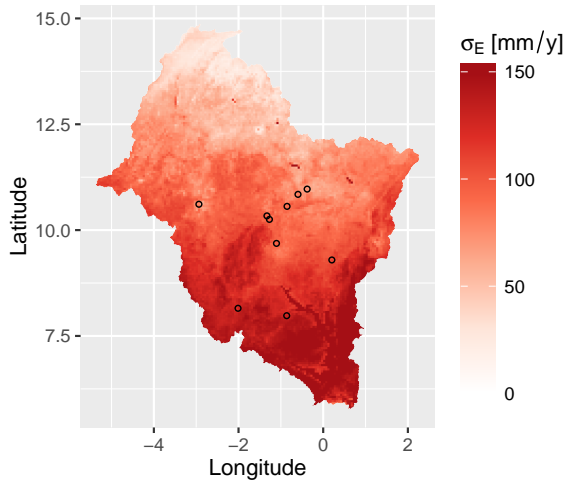


(c) Difference of mean evaporation in model 6, compared with model 1

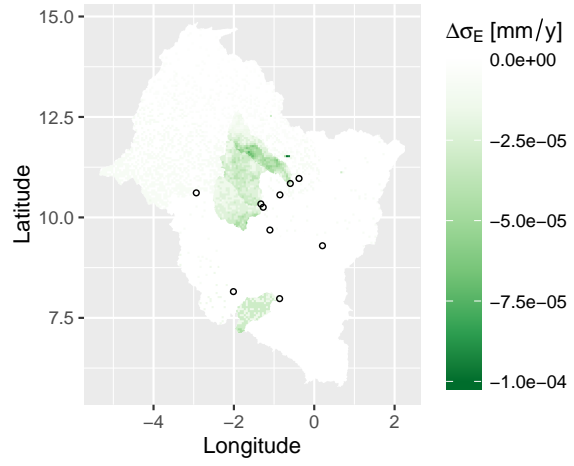


(d) Difference of mean evaporation in model 7, compared with model 1

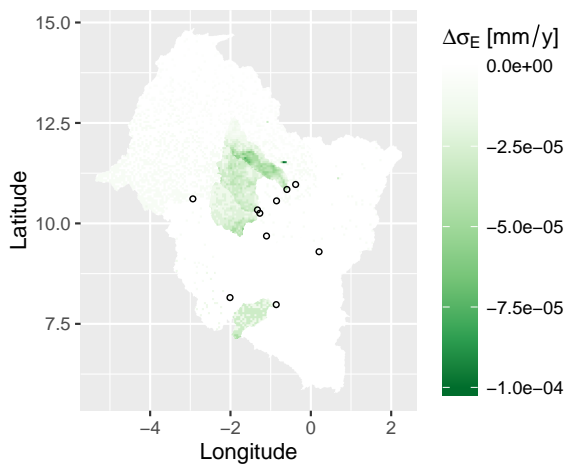
Figure D.7: Mean evaporation



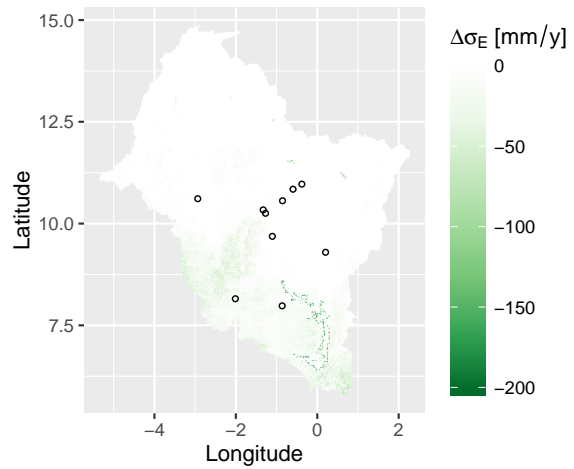
(a) SD evaporation runoff in model 1



(b) Difference of SD evaporation runoff in model 2, compared with model 1

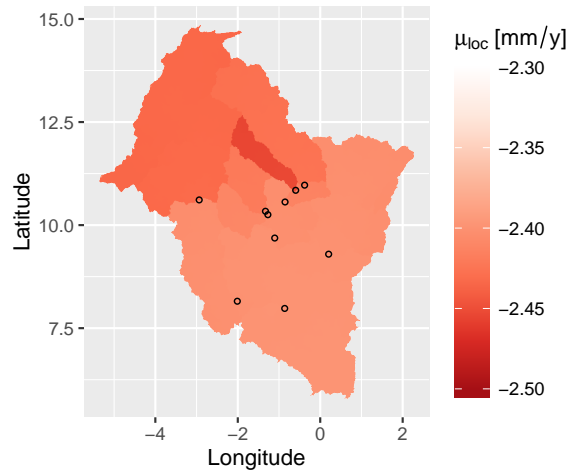


(c) Difference of SD evaporation runoff in model 6, compared with model 1

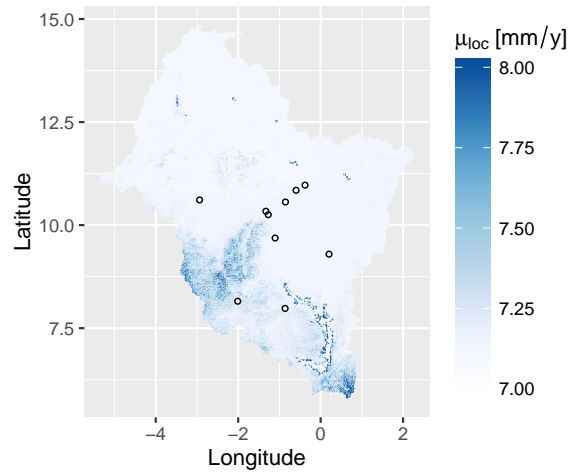


(d) Difference of SD evaporation runoff in model 7, compared with model 1

Figure D.8: Standard deviation evaporation runoff

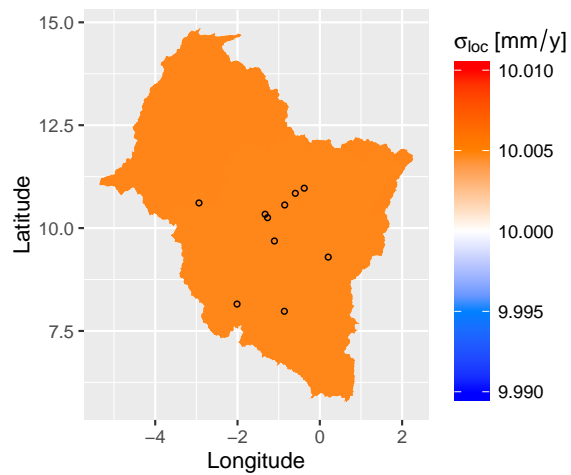


(a) Bias Mean model 6

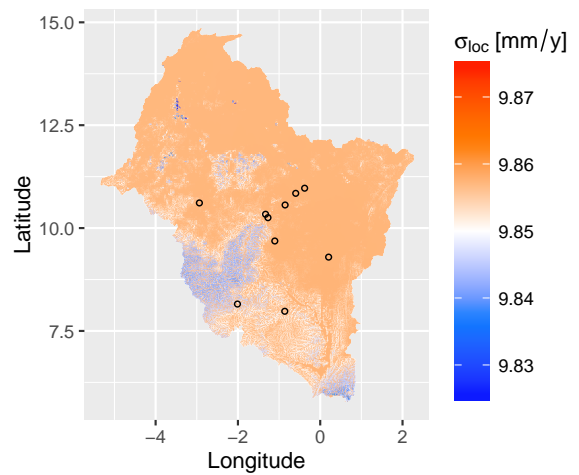


(b) Bias Mean model 7

Figure D.9: Bias Mean



(a) Bias SD model 6



(b) Bias SD model 7

Figure D.10: Bias SD