

Utah State University

DigitalCommons@USU

---

Reports

Utah Water Research Laboratory

---

January 1993

## Development of Mountain Climate Generator and Snowpack Model for Erosion Predictions in the Western United States Using WEPP: Phase IV

David S. Bowles

U. Lall

D. G. Tarboton

E. Malek

B. Rajagopalan

T. Chowdhury

*See next page for additional authors*

Follow this and additional works at: [https://digitalcommons.usu.edu/water\\_rep](https://digitalcommons.usu.edu/water_rep)



Part of the [Civil and Environmental Engineering Commons](#), and the [Water Resource Management Commons](#)

---

### Recommended Citation

Bowles, David S.; Lall, U.; Tarboton, D. G.; Malek, E.; Rajagopalan, B.; Chowdhury, T.; and Kluzak, E., "Development of Mountain Climate Generator and Snowpack Model for Erosion Predictions in the Western United States Using WEPP: Phase IV" (1993). *Reports*. Paper 590.

[https://digitalcommons.usu.edu/water\\_rep/590](https://digitalcommons.usu.edu/water_rep/590)

This Report is brought to you for free and open access by the Utah Water Research Laboratory at DigitalCommons@USU. It has been accepted for inclusion in Reports by an authorized administrator of DigitalCommons@USU. For more information, please contact [digitalcommons@usu.edu](mailto:digitalcommons@usu.edu).



---

**Authors**

David S. Bowles, U. Lall, D. G. Tarboton, E. Malek, B. Rajagopalan, T. Chowdhury, and E. Kluzak

# Development of Mountain Climate Generator and Snowpack Model for Erosion Predictions in the Western United States using WEPP

**Research Completion Report for the funding period January 1, 1994 to January 31, 1995 of Phase IV**

Research Joint Venture Agreement INT-92660-RJVA, Amend. #3

## Submitted to:

U.S. Department of Agriculture, Forest Service  
Intermountain Research Station  
Forestry Sciences Laboratory  
1221 South Main Street  
Moscow, Idaho 83843

February 27, 1995 (DRAFT)

## Submitted by:

David S. Bowles, Gail E. Bingham, Upmanu Lall, David G. Tarboton, Ragagopalan Balaji, Esmail Malek, Erik Kluzak



UTAH WATER RESEARCH LABORATORY  
UTAH STATE UNIVERSITY  
LOGAN, UT 84322-8200

**Utah State**  
UNIVERSITY

99354.2

# CHAPTER 1

## INTRODUCTION

### 1.1 OBJECTIVE

The overall objective of this project is to develop a procedure for generating representative historical and synthetic climate sequences at ungaged locations throughout the mountainous Western United States. As a secondary objective, we are also developing a snowpack simulation model. The Utah Water Research Laboratory (UWRL) is conducting this project under a research joint venture agreement with the U.S. Forest Service (USFS) as part of the USFS Watershed and Erosion Prediction Project (WEPP).

This work is part of a USFS research and development effort and, as such, must provide a usable product within the project schedules established by the USFS. The MCLIGEN, which is being developed by the UWRL, will furnish climate inputs to the WEPP with the goal that acceptably accurate erosion predictions are provided for design and planning purposes. The representation of climate in mountainous areas is a major challenge because climatological data are scarce and meaningful interpolation of climate variables is difficult in complex terrain. The project is using existing techniques which provide adequate climate inputs, adapting existing procedures where appropriate, and developing new procedures within the constraints of available data and project resources.

Although MCLIGEN is being developed under the WEPP project for erosion prediction, it will have many other applications of interest to the USFS and other resource agencies. For example, the climate sequences would be useful for driving ecological models or resource assessment procedures which require climate inputs. MCLIGEN will also have the capability to be run under climate change scenarios.

### 1.2 USER REQUIREMENTS

The WEPP user will need these "climate sequences" accessible in three "event forms."

- *Selected representative historical events or sequences* (e.g., average, dry, and wet). This capability would enable users to make erosion estimates for climate sequences based on historical events either as observed or as outputs from the climate modeling system. In the first case, the user could select a recorded event or sequence of data from a station or stations which the user considers best represents the conditions at the site which is under evaluation. In the second case, the modeling system would adjust recorded historical events to be representative of ungaged locations.
- *Continuous simulation* of climate for up to 20-year periods using stochastic methods. This will be particularly useful in assessing the erosion potential from timber harvest areas, and it could be used to estimate a probability distribution of erosion potential, average potentials, or perhaps high or low extreme climate cases. High cases could be useful for design of sediment control measures, such as detention basins.
- *Design events* associated with various occurrence frequencies or return periods.



Users will choose the form of climate input which they use. UWRL work is focused on the first and second forms listed above. The generator will have the capability of providing climate inputs based on locational information (such as latitude, longitude, elevation, slope, and aspect).

### **1.3 TECHNICAL APPROACH**

MCLIGEN will depend on historical, physically based interpolations of weather sequences from a mesoscale-climate modeling system which is comprised of four nested layers:

1. An existing synoptic scale forecast model (300 x 300 km);
2. A regional scale climate model (50 x 50 km);
3. A local scale climate model (10 x 10 km); and
4. A specific point climate predictor, referred to as "ZOOM."

Two additional MCLIGEN components are:

5. A local scale stochastic climate generator; and
6. A point energy balance snowmelt model.

We will provide the USFS with a database for the Western U.S. consisting of 13.4 years of climate data on a 50-km grid from the second layer model. It is anticipated that this database will be maintained at the regional level of the USFS, and that 13.4-year sequences of 10 x 10 km data could be provided for each forest using the local scale climate model. The 10 x 10 km datasets will be used for input to ZOOM, which will provide climate data at any point on a watershed subarea for which erosion analyses are being made. These climate data will be inputs to the hydrologic, snowpack, and erosion prediction components of WEPP. ZOOM, the local scale climate simulator, and the 60-km resolution climate database are the project deliverables from this part of our work. The RegCM2 model and the ECMWF synoptic data set can also be provided if desired.

Stochastically generated sequences of climate variables are also required to run WEPP. We are developing and testing a nonparametric wet/dry spell stochastic model for daily precipitation. Initially, historical-climate data are being modeled, but as generated climate data from the climate modeling system become available they will be modeled via the input variables to ZOOM.

The snowpack simulation model will need to meet some requirements which other snowpack models have not been demonstrated to meet. Specifically, the snowpack simulation model should provide snowmelt predictions which represent the characteristics of snowmelt events which are important for accurate erosion predictions. Also, the model should be a transportable model which will work well for the variety of conditions which exist across the Western U.S., and the model should be physically based so that it will require a minimum of calibration for local conditions. This last requirement is particularly important because WEPP users may not always have the expertise to calibrate the snowpack model. The snowpack simulation model is being developed by combining the best components of existing energy balance models with improvements which we are developing at the UWRL.

### **1.4 PROJECT STATUS**

Three developmental phases were defined in the work plan submitted to the USFS on September 8, 1989:

Phase I: Climate data evaluation and generator design.

Phase II: MCLIGEN coding and evaluation at representative sites.

Phase III: Generalization to the entire Western United States.

Work undertaken during this funding period, beginning January 1, 1994 and January 31, 1995, is part of Phase II.

Tasks remaining to complete work are presented in the last sections of Chapters 2, 3, and 4. The interrelationship between these tasks is illustrated in Figure 1--1. Tasks are divided into three increments according to funding availability. Funding for increment I was secured, and work was conducted in the period ending January 31, 1995. The schedule for increment II tasks is based on a January 31, 1996 completion.

Papers presented on work conducted under this project during the reporting period are listed below:

## **1.5 OUTLINE OF REPORT**

The report is divided into four chapters and an executive summary. Chapters 2, 3, and 4 address the three major parts of work: climate modeling system, snowpack modeling, and stochastic modeling. Each chapter includes a literature review, discussion of the proposed methodology, and description of work plan.

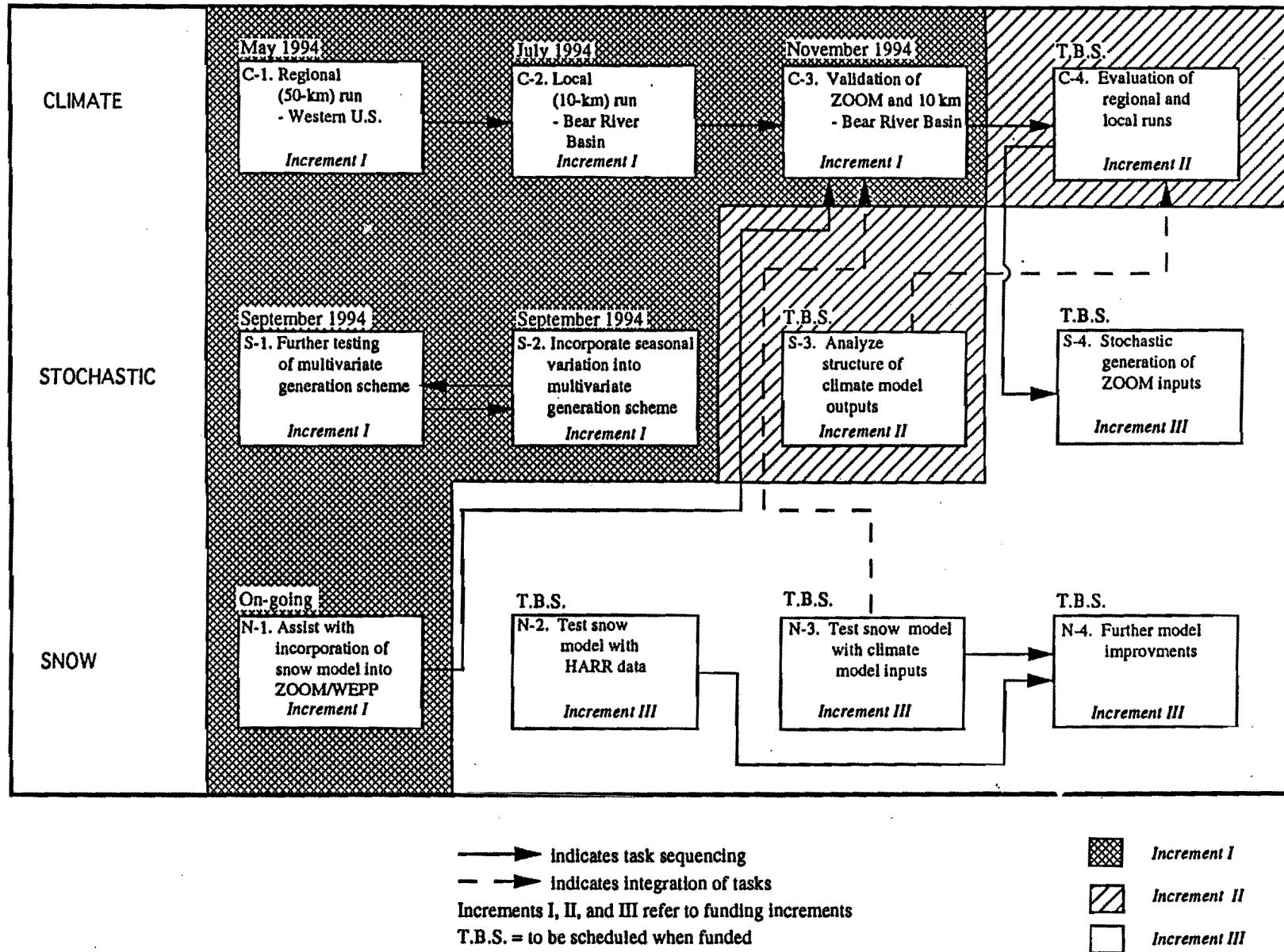


Figure 1-1. Interrelationship of third-year tasks.

# CHAPTER 2

## CLIMATE MODELING SYSTEM

### 2.1 OBJECTIVE FOR THIS REPORTING PERIOD

The objective established for the current reporting period was to complete the development of the RegCM2 all-season model for use in developing the MCLIGEN model, to test the RegCM2 model to the western U.S. Forest Service study area, to refine the model to generate valid climate parameters at 10 km resolution and to test it at 10-km resolution over a 360 x 370 km application area.

### 2.2 TASKS

The following tasks were established for the current reporting period:

1. Complete the regional RegCM2 (50 km) model run for the period, December 1978 to July 1992, for the Western United States. NCAR was to do the computing for the 10 year period between 1982 to 1992, and we were committed to provide the end periods.
2. Complete a local RegCM2 (10 km) model run for the period, December 1978 to July 1992, for the Bear River Basin in Utah, Idaho, and Wyoming.
3. Validate the performance of ZOOM and the 10 km Bear River Basin run through:
  - a.) Validating the 10 km model for the whole Bear River Basin.
  - b.) Incorporation of the WEPP soils model and the USU snow model into ZOOM.
  - c.) Intercomparison of the performance of the BATS soils and snow submodels with the WEPP soil model and the USU snow model to determine if the models significantly affect the calculation of local climate variables.
  - d.) Comparison of the BATS soil and snow submodels, the WEPP soils model, and the USUS snow model, each driven by ZOOM climate data, with the snow course, SNOTEL, RAWS, and satellite data for all seasons.
4. Evaluate the regional and local model output through comparisons with data from NOAA Coop Stations, SNOTEL, RAWS, NWS rawinsonde data, and satellite data. (Funding for this task was not received as part of the first increment of FY 94 funding.)

### 2.3 ACCOMPLISHMENTS AND PROBLEMS

Last year's report includes a good overview of the climate modeling system. Appendix D, a preprint from our presentation to the American Meteorological Society (AMS) meeting in Dallas TX, also presents an overview of the MCLIGEN climate modeling system, driving a distributed hydrological model.

The Regional run (50 km) is only about 60% completed. However, we have had to shoulder this entire task. At the beginning of the year, we anticipated doing the computing for only about 2.5 years of this 13 year period, while NCAR completed the rest of the run. After completing a significant portion of their run, NCAR discovered a mistake that rendered much of their effort invalid. Consequently, we have had to use a significant portion of our computing resources to complete this task. To date, we have completed over half of the regional (50 km) Western U.S. RegCM2 run, and roughly a third of the local scale (10 km) RegCM2 run (see sections 2.3.2 and 2.3.3 respectively). To allow us to complete the 50 km run in the next couple of months, we have broken it up into chunks, and we are running each segment on a separate machine.

Durning the year, we have improved and debugged much of the process (including RegCM2 itself and it's input and output support programs) (see section 2.3.1). We have developed a parallel version of RegCM2 that has significantly cut linear CPU time (see Appendix B). We have been confronted with some significant challenges in trying to maintain reasonable accuracy in the 10 km modeling (see sections 2.3.3 and 2.5.2 for details). Through out the year, we have visually quality controlled every byte of each days model run. In this way we have avoided serious mistakes, like the one that cost NCAR so dearly. Early on, we invested in the development of a visual data display system that lets us view each days output in a video format. Several significant model limitations have been identified and corrected using this technique. We have begun the task of actually comparing surface climate data to the modeling process (see section 2.3.5).

A summary of the significant developments is listed below:

Loss of the 10 years of 50 km model output being contributed by NCAR.

Because, of the long CPU run times of both the 50 and 10 km runs, we realized that we could not afford to wait until we had all the problems solved with RegCM2 to proceed. So we have had to make changes to the model while running. We have documented when these changes occur and their nature.

To effectively use our many available workstations, we have split the long runs up into shorter segments, and overlapped these segment to minimize startup transients in the runs.

Because, of the need to concentrate on getting the highest quality regional and local scale climate runs in a timely fashion, we have concentrated our resources on the RegCM2 runs. Paul Swetik, Charlie Luce, and Jeff Blatt have begun efforts on constructing WEPP/ZOOM and we detail this progress and ours in section 2.3.4.

### **2.3.1 MCLIGEN Climate Model System Development**

Over the past year we have worked on three general area's of development in the MCLIGEN modeling system. The first area is development of the RegCM2 model, the second is the development of a parallel version of RegCM2, and the third is creation/development of support programs for RegCM2. Developing a parallel version of RegCM2 has been critical in total wall clock time to complete the model runs in a reasonable amount of time. The changes to the model have been important in improving the performance when diffishencies were oberved in the output. The support programs have helped our understanding of problems in the model, validating model output, processing boundary condition data and in getting the model started.

1. Changes to RegCM2. There have been several changes necessary to the model as we have been using it. As there is not a need to go into the details here, we will just mention a few areas where changes were made. These include:
  - Double Precision RegCM2.
  - Changes in output files.
  - Changes in output BATS variables.
  - Changes as suggested by NCAR that impliment the Holtslag Planetary Boundary Layer (PBL) routine.
  - Change to fractional vegetation cover.
  - Change in leaf temperature calculation.
2. Parallel RegCM2. We have developed a preliminary parallel version of RegCM2 that has speeded up the processing by up to 10X the original rate. We presented this work at the American Geophysical Meeting (AGU) meeting in San Fransico and we reprint the details of our

AGU poster in Appendix B. We just point out that the current version of parallel RegCM2 is the first phase of the changes that are needed to optimise this process. We are now working on a fully parallel version of RegCM2. To make the model run efficiently on multiple processors we need to modularize the TEND routine, which currently has 80% of the code. This would allow us to split the job into smaller chunks than is now possible. The other note we make is that our work with parallel RegCM2 has caught the interest of Cray Research Corp. and Tony Meys of Cray Research has given us some free CPU time on a Cray C-90 to develop a parallel version of RegCM2 for the Cray. They have also promised us a few evening sessions on the C-90 where we would have the machine to ourselves. Preliminary estimates suggest in one session we may be able to run up to a half a years worth of climate runs.

## 2. Creation of Support programs for RegCM2.

- 2.1 Creation of video program/Color plotting. We spent considerable effort creating a program to animate the model data in a 3D color visualization format. This has been instrumental in understanding the model output, and fixing errors. It has also been cost-effective way for initial verification of the model output. Ensuring that the model is progressing and giving reasonable output qualitatively. This effort was also presented at the fall-94 AGU meeting (See Appendix C).
- 2.2 Creation of program to set Initial snow depth. We also have spent effort developing a program to set the initial snow depth for both the 50 km and 10 km runs. This program uses snow depth vs. elevation curves for various regions as well as using SNOTEL and NOAA/COOP station data to set values where data is available. This was important since the ECMWF data starts in December when snow depth is significant. Without setting the initial snow depth the first water year becomes useless, without some reasonable values for a starting snowpack.

### 2.3.2 Regional Scale Climate Model (50 km RegCM2 Model)

Goal: "Complete the regional RegCM2 (50 km) model run for the period, December 1978 to July 1992, for the Western United States."

Progress Summary. Currently we have processed over 2000 days of the Western U.S. 50 km run. This covers from Dec/1979 to Middle of 1984. NCAR has run from Aug/1985 to Aug/1986. We have also have over 200 days into a run that starts Aug/1/1988.

Support programs.

Changes in ECMWF to IN file program:

Creation of a NMC IN file program. As NCAR did not archive ECMWF observational analysis for Dec/1979 and ECMWF does not keep archives of it's old analysis, we had to create our IN files from NMC data.

Changes to RegCM2 effecting the regional runs:

Change in Sub-grid Scale precipitation routine

Changes in horizontal diffusion.

Semi Explicit Moisture Scheme

NCAR 50 km Multi-year model run. Our NCAR collaborators originally proposed running a 10 year 50 km run that would covering the entire continental U.S. starting Apr/1981. Our original plan was to use their 50 km run for the bulk of our needed 50 km runs. However, after running from Apr/1981 to 1984 they noticed that the Sea Surface Temperatures had been entered incorrectly. At this point they stopped

running so they could evaluate. At first they thought the effect was unimportant and only effected the coast. However, on closer examination they found several other problems with the runs. Precipitation in the summer was too high. Desert regions were not being modeled well.

Breakup of runs into sectional chunks. Because of the intensive CPU and run times that RegCM2 takes to run, we have decided that it is only prudent to break up the entire run into different time segments and run each segment on different machines. We have already started a segment starting in 1988. NCAR has run a section in 1985. Other segments that we plan to run are:

Jun/1/1986-Oct/1/1988,  
Aug/1/1990-Oct/1/1991,  
Aug/1/1991-Oct/1/1992, and  
Aug/1/1992-latest data.

The 1990 runs will take about 73 days of running on an indigo-2, the 1988 run will take 37 days to complete, the 1986 run 28 days, and the beginning run will take 25 days to complete. Thus we expect the 50 km runs to be finished in 2 months. The ECMWF observational analysis archives beyond 1989 are in a different format than the earlier datasets. Because, of the change in format we have been working on changing our codes to accomodate the new format. The newer datasets also have an archive with data stored 4X per day as well as the usual 2 X per day. We will use the the data stored the most often.

Creation of IN to INB file program. To help speed-up the spin-up time needed for each time segment we created a program that will use the soil moisture and temperature profile from a previous years run. This way the BATS fields will be able to reach equilibrium sooner.

### **2.3.3 Local Scale Climate Model (10 km RegCM2 Model)**

Goal: "Complete a local RegCM2 (10 km) model run for the period, December 1978 to July 1992, for the Bear River Basin in Utah, Idaho, and Wyoming."

Progress Summary. We have run over 1100 days of our 10 km run. Running from Dec/5/1978 through 1981 and now into 1982.

Changes to the Local RegCM2 model system to date include:

Changes to nest program  
Modify Lake Model for Great Salt Lake  
Modifications to the Precipitation Auto-Conversion term  
Modifications to the Radiation routine  
Changes in horizontal diffusion  
Breakup of runs into sectional chunks.

### **2.3.4 ZOOM, the Specific Point Weather Generator**

Distributed ZOOM:

Slope / Aspect Angle inclusion in ZOOM:

This section has not yet been completed.

### 2.3.4.1 ZOOM/WEPP

Goal: "Incorporation of the WEPP soils model and the USU snow model into ZOOM."

Progress Summary. In order to integrate the functions of the higher level climate models with the Water Erosion Prediction Project (WEPP) model more closely some of the functions of ZOOM/BATS are being incorporated into the WEPP model. This method will allow better modeling of the effects of vegetation on microclimate. In particular, evapotranspiration from forests is strongly dependent on soil moisture conditions, and a better representation of temperatures below the canopy is made possible by incorporating the more detailed hydrology of the WEPP model. Growth of shrub, grass, and herb species is strongly affected by this modification to the climate, and their growth and survival strongly affect erosion.

The flowcharts in Figures 2.1 - 2.3 show conceptually the layout of ZOOM/WEPP. The flowchart has been split into three pages, corresponding to three sets of processes. On figure 2.1, processes are predominantly atmospheric and are modified primarily by the topography. The second figure (figure 2.2) includes mostly surface hydrology routines that are strongly affected by site topography and vegetation and figure 2.3 covers the interaction between soil water, microclimate, and plant growth. The plant growth routines are based on the BGC models.

The initial section of ZOOM/WEPP adjusts modeled weather at the nodes to the horizontal coordinates of the site of interest. When looking at historical weather (1978-1992), the program will interpolate the data using the 4 corner points around the site. Because the stochastic generator does not handle spatial correlations between nearby nodes, a weather sequence for the nearest node to the site is used when stochastic sequences are used. Pseudo-historic weather sequences (from RegCM2 10 km model runs) are recorded hourly at the nodes, and stochastic weather is generated on daily time steps.

The next phase adjusts the weather variables (precip, temperature, humidity, wind) for elevation and aspect. This section will use modules already in ZOOM/BATS and incorporate the rain/snow partitioning and drift components of the snowmelt model described in a later chapter. An interception model will be designed to account for the effect of forest vegetation on precipitation.

The surface hydrology routines start with the snowmelt model. Either snowpack outflow or direct rainfall is then converted to a hyetograph for the infiltration and overland flow model of WEPP.

Once infiltration is determined, the plant physiology model determines how much water is released through evapotranspiration through each layer of the canopy. This is used to determine temperature for lower layers in the canopy and how much soil water is withdrawn. The temperature is important in describing plant growth in the lower layers and the probability of plant regeneration. The remainder of the model is plant growth and erosion calculations based on the weather and hydrology calculations.

Several components of ZOOM/WEPP have been completed. The basic hydrology components of WEPP have been coded and validated. The snowmelt model has been incorporated into the WEPP model and is awaiting verification. Validation of the snowmelt model is described in a chapter 3. The adjustments for elevation and aspect are already in ZOOM/BATS and will be incorporated into WEPP/ZOOM.

Remaining development tasks are the plant physiology, reproduction, and growth models and a new interception routine. Most of the plant physiology and growth routines will be based on a variant of BGC. The reproduction model is in development. Existing interception routines (ie. BATS) are oversimplified, and a new one will be designed for this model.



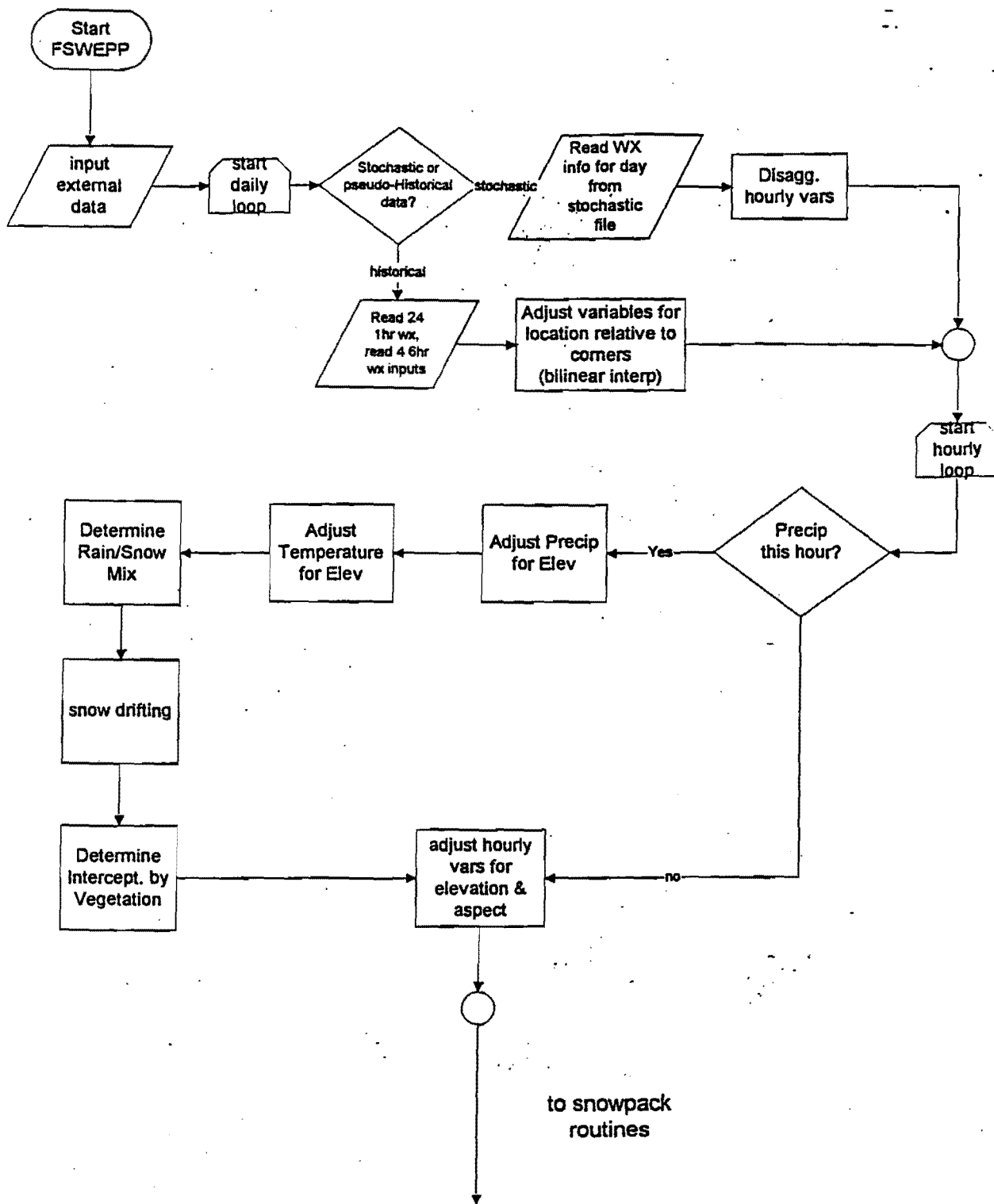


Figure 2.1. Atmospheric Processes.

ZOOM / WEPP  
(SNOWPACK & IRS)

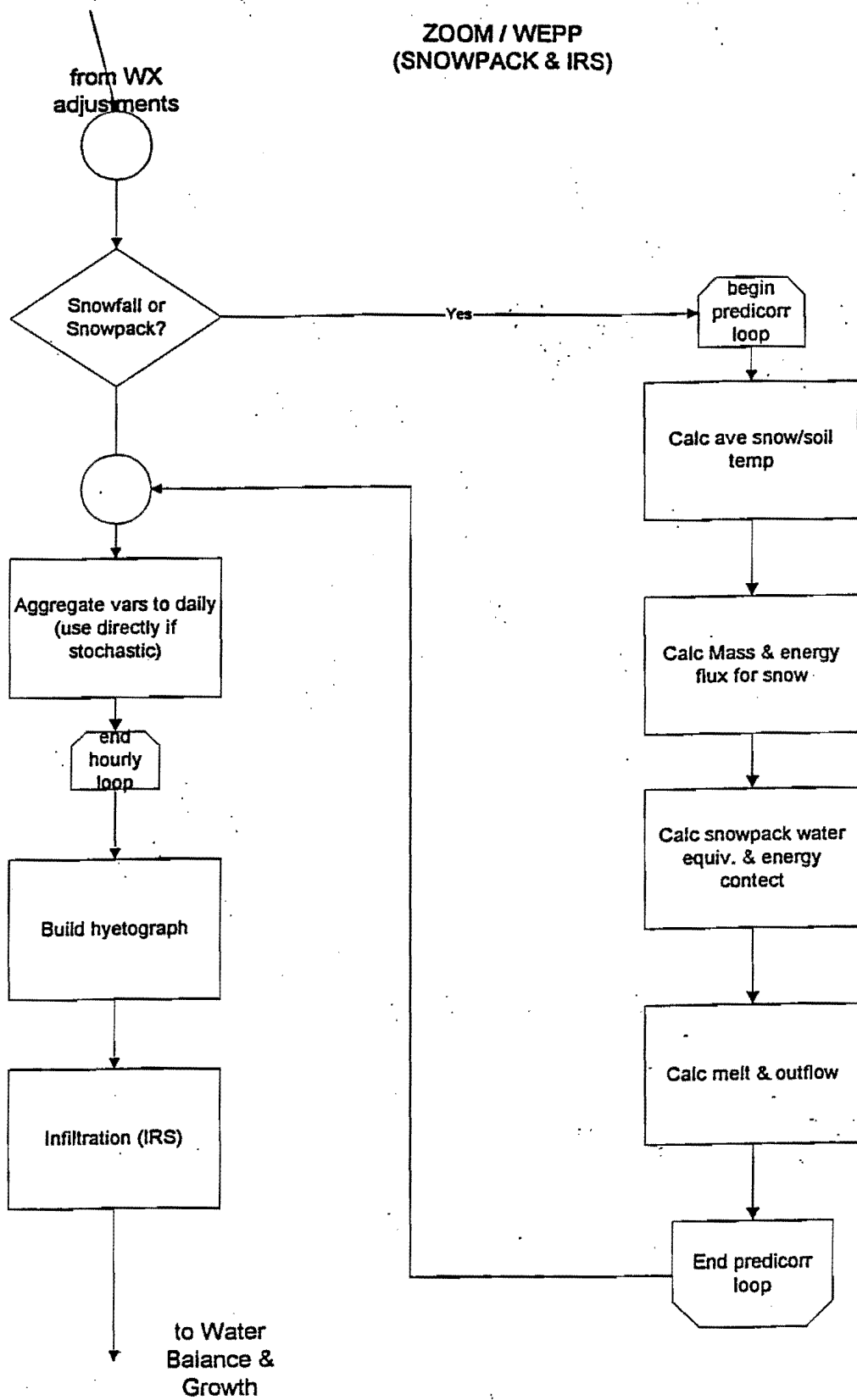


Figure 2.2 Surface Hydrology Processes.

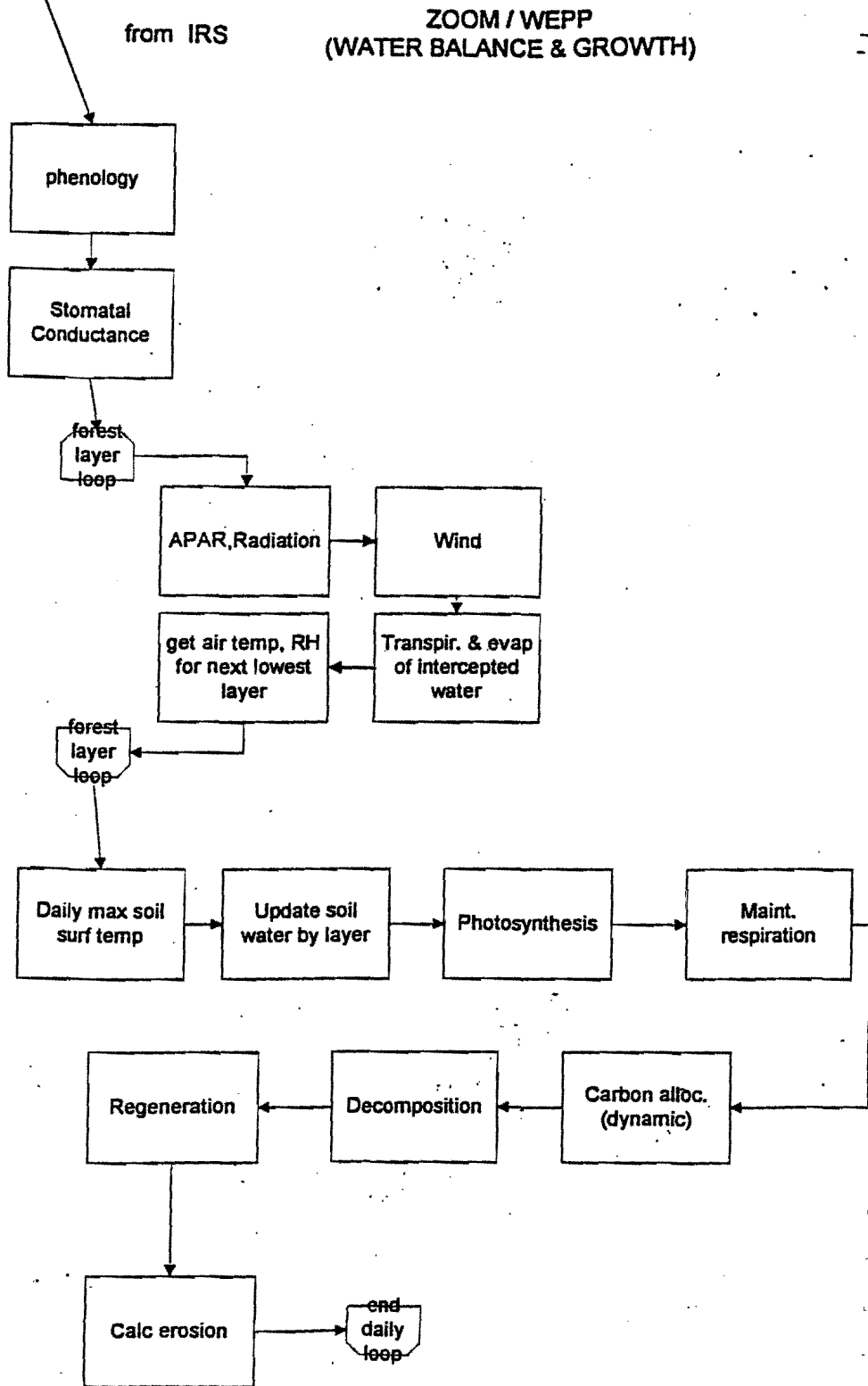


Figure 2.3 Interaction Between Soil, Water, Microclimate, and Plant Growth.

### **2.3.4.2 ZOOM/BATS Snow and USU Snow models**

Goal: "Intercomparison of the performance of the BATS soils and snow submodels with the WEPP soil model and the USU snow model to determine if the models significantly affect the calculation of local climate variables."

Progress Summary. For completion of the other tasks and to give highest emphasis on getting the highest quality regional and local scale model outputs we did not complete this task. We did get an older version of the USU snowmodel and have done comparisons of the computations to the BATS model.

### **2.3.5 MCLIGEN Model Validation**

Principal Component Analysis (PCA analysis),

Daily point COOP/SNOTEL comparisons of Tmin, Tmax, and Precip.

Comparisons to gridded station data.

These efforts were delayed by the available funding, and have been transferred to the next reporting period.

#### **2.3.5.1 Regional Scale Climate Model Evaluation (50 km Western U.S.)**

Goal: "Evaluate the regional model output through comparisons with data from NOAA Coop Stations, SNOTEL, RAWS, NWS rawinsonde data, and satellite data. (Funding for this task was not received as part of the first increment of FY 94 funding.)"

#### **2.3.5.2 Local Scale Climate Model Evaluation (10 km Bear River Basin)**

Goal: "Evaluate the local model output through comparisons with data from NOAA Coop Stations, SNOTEL, RAWS, NWS rawinsonde data, and satellite data. (Funding for this task was not received as part of the first increment of FY 94 funding.)"

#### **2.3.5.3 ZOOM Model Evaluation**

Goal: "Comparison of the BATS soil and snow submodels, the WEPP soils model, and the USU snow model, each driven by ZOOM climate data, with the snow course, SNOTEL, RAWS, and satellite data for all seasons."

## **2.4 WORK PLAN FOR FY 1995**

The following tasks are formulated for the next reporting period:

- C-1. Finish the regional RegCM2 (50 km) model run for the period, December 1978 to July 1992, for the Western United States.
- C-2. Finish the local RegCM2 (10 km) model run for the period, December 1978 to July 1992, for the Bear River Basin in Utah, Idaho, and Wyoming.

- C-3. Evaluate the regional and local model output through comparisons with data from NOAA Coop Stations, SNOTEL, RAWS, and NWS rawinsonde data.
- C-4. Incorporate WEPP/BATS adjustments into WEPP/ZOOM and complete the vegetation component of WEPP/ZOOM.
- C-5. Validate the WEPP/ZOOM model through comparisons with NOAA/Coop, SNOTEL, and RAWS stations.

# CHAPTER 3

## SNOWPACK MODELING

### 3.1 OBJECTIVE FOR THIS REPORTING PERIOD

The objective of the snowmelt modeling work was to develop and test a transportable energy balance snowmelt model to provide snowmelt inputs for the calculation of erosion in WEPP. The model is to be driven by climate generator inputs.

### 3.2 TASKS

Model development was essentially complete and is described in previous reports. The task for this phase of the work was:

1. To assist with the incorporation of the snow model into ZOOM/WEPP.

### 3.3 ACCOMPLISHMENTS AND PROBLEMS

The major accomplishment of this period was a thorough revision of the code to resolve some bugs, numerical instabilities, and to improve the structure. Following this revision the code was transmitted to the USDA Forest Service, Moscow, ID, and assistance was given in understanding and incorporation of the code into the ZOOM/WEPP components being developed there. The model was also tested further against data from Upper Sheep Creek, within the Reynolds Creek experimental watershed, Boise, ID, and against data collected in Cache Valley, near Logan, UT. These tests were reported in the following publications which are included as Appendices 4A and 4B (Tarboton et al., 1995, Tarboton, 1994) to this report. Based on this work, the following deficiencies and shortcomings need to be noted:

1. The parameterization of atmospheric instability based on Richardsons number was found to be deficient because the values obtained from typical input data frequently resulted in a Richardsons number greater in magnitude than the range for which the parameterizations were valid. An option to turn off the effect of atmospheric instability was included in the restructured code and until this is resolved we recommend that the instability parameterization be turned off.
2. The model has a tendency to underestimate the snow energy content early in the season prior to melt. This was discovered from the tests against data collected in Cache Valley where we had temperature measurements within the snowpack so we could check this aspect of the model. The calibration of adjustable parameters within the model apparently offsets this deficiency, but if possible it should still be corrected so that the model is as physically close to correct as possible. Melt predictions are still relatively good.
3. The representation of vegetation through a vegetation density factor F is still relatively primitive and has only had minimal testing against data at the CSSL forested site. We do not have other data on snow accumulation and melt rates, together with the necessary variables to drive the model under canopies of varying vegetation density.

### 3.4 WORK PLAN FOR FY 1995

Our objective for this period is to resolve the deficiencies noted above to the extent possible and verify the operation of the model within WEPP/ZOOM. More specific tasks have been highlighted as follows:

- C-1. Test the model against additional data. Additional data that are available includes: (i) data collected by Dennis Harr in the Pacific Northwest; (ii) data from the GLEES experimental watershed; and (iii) Snotel data. These tests should be performed in cooperation with USFS staff at Moscow, ID, using the model within an experimental version of WEPP. This will serve as a verification of the incorporation of the model into WEPP and provide a standard interface for data input and reporting of results. The utility and feasibility of these tests depend on the availability of the input variables necessary to drive the model. For some tests (Snotel in particular), the model will need to be driven by climate model inputs. This will serve as a check on the climate model, as well as on snow model components.
- C-2. Investigate the energy content underestimation problem. Diagnose the model runs using the Cache Valley data to understand why the energy content is under-represented and correct any deficiencies discovered. Re-calibrate if necessary.
- C-3. Research the parameterization and calculation of turbulent flux transfers to see if there is a reasonable way to account for atmospheric stability/instability while remaining within the ranges of values for which the parameterization has been developed.

Tasks 2 and 3 are open ended and, given the funding for this component of the research, substantial improvements cannot be guaranteed.

#### References

- Tarboton, D.G. (1994). Measurement and modeling of snow energy balance and sublimation from snow. In: Proceedings, International Snow Science Workshop, Snowbird, Utah, October 31 to November 2. Utah Water Research Laboratory Working Paper No. WP-94-HWR-DGT/002, pp.
- Tarboton et al. (1995).

# CHAPTER 4

## STOCHASTIC MODELING AND PARAMETER REGIONALIZATION

### 4.1 OBJECTIVE FOR THIS REPORTING PERIOD

The objective established for the current reported period was to improve the multivariate nonparametric generation scheme to facilitate simulation from the conditional density formed by conditioning on a key variable, precipitation. Also to explore other alternatives.

### 4.2 TASKS

The following tasks were established for the current reporting period:

Task IV-1: Improve the multivariate nonparametric generation scheme to facilitate simulation from conditional density formed by conditioning on the precipitation.

Task IV-2: Explore other alternatives to multivariate generation of weather variables.

### 4.3 ACCOMPLISHMENTS

1. The required improvements were made to the multivariate nonparametric generation scheme. The improved method samples the original data with replacement while smoothing the empirical conditional distribution function. The technique can be thought of as a smoothed conditional Bootstrap and is equivalent to simulation from a kernel density estimate of the multivariate conditional probability density function. This improves on the classical Bootstrap technique by generating values that have not occurred exactly in the original sample and by alleviating the reproduction of fine spurious details in the data.

Precipitation is generated from the nonparametric wet/dry spell model as described in Lall et al. (1995). A vector of other variables (solar radiation, maximum temperature, minimum temperature, average dew point temperature and average wind speed) is then simulated by conditioning on the vector of these variables on the preceding day and the precipitation amount on the day of interest. An application of the resampling scheme with 30 years of daily weather data at Salt Lake City, Utah, USA is provided. This method has been sent for publication as Rajagopalan et al. (1995) and is presented in Appendix 4A.

2. In the course of exploring other alternatives to nonparametric multivariate generation of weather variables, a nearest neighbor bootstrap method was developed. The motivation to this approach comes from a desire to preserve the dependence structure of the time series while bootstrapping (resampling it with replacement). This method is data driven, and is highly parsimonious. The method follows to the work of Lall et al. (1995) where they demonstrate this idea by applying it to resampling monthly streamflow by Lall et al. (1995). The development of the nearest neighbor approach to simulating weather variables has been sent for publication as Rajagopalan et al. (1995), and is presented in Appendix 4B.



3. In the course of improving the nonparametric wet/dry spell model for modeling daily precipitation, that was presented in Bowles et al. (1992), a new estimator for discrete probability distributions was developed and has been published as Rajagopalan and Lall (1995), besides incorporating this in the wet/dry spell model, that have been sent for publication as Lall et al. (1995) and Rajagopalan et al. (1995). The new estimator for discrete probability distributions is presented in Appendix 4C.
4. The wet/dry spell model as we know involves the breaking up of the year into seasons determined a priori. The seasons vary from place to place, as a result the same seasonal break of the year may not hold across all sites. Nonparametric techniques for studying the seasonal patterns in precipitation that could provide objective guidelines to breaking the year into meaningful seasons, were identified and were applied to precipitation data from numerous sites, along a meridional transect in western U.S. This is presented in Rajagopalan and Lall (1995) and also in Appendix 4D.
5. In the wake of significant change in seasonality of precipitation, an alternate representation to the wet/dry spell model for daily precipitation, that obviates the need for breaking the year into seasons was motivated. This alternate representation was developed as a nonhomogeneous Markov model that used the discrete kernel estimator developed in Rajagopalan and Lall (1995). The nonhomogeneous Markov model is presented in Rajagopalan et al. (1995) and also in Appendix 4E.

#### **4.4 WORK PLAN FOR FY 1995**

##### References

Lall, et al (1995)  
Rajagopalan et al. (1995)  
Rajagopalan and Lall (1995)

# APPENDIX 2A

## DOCUMENTATION OF THE MOUNTAIN CLIMATE GENERATOR INPUT FILES

## INTRODUCTION

Welcome to a new world -- the world of Mountain Climate Generator (MCLIGEN)! MCLIGEN is unique in that it links any interested Forest Service researcher, agent, or personnel, to state of the art meteorological data in way never before accomplished! It allows the Forester to access historical, modeled psuedo-historic, or even modeled doubled CO2 for any point he chooses within the set of files he has. Such data has always existed but never before at such high resolution and so easily accessed by personnel within the forest service. This is the vision and the reason for MCLIGEN now lets get into the details of the methodology of use.

ZOOM is the interface from the 10 km high resolution climate model output data to the point and basin surface hydrologic models (WEPP (point model), BATS (point or grid model) and/or CVHM (Basin model)). The files used for this interface are called local area (LA) files, because they only consist of a small portion of the 10 km model run. In general they should be approximately 60 km x 60 km, this way an entire 15 year series can be loaded onto a single 2 Gigabyte disk. Were the entire area to be loaded a 50 GigaByte disk would be needed. Also it was felt that in general a researcher would only be working with one small area at a time and the excess information stored for the other areas was a waste of disk storage. When a researcher needed a new area he could load the other set of files that included the area of interest. This allows a longer time-series to be loaded onto a single disk rather than either having to only load part of the files at a time or flip through a stack of CD's each time ZOOM/WEPP is run. If a larger area is actually needed the same format of files can be used just as well for a larger area as for a smaller one, but only a much shorter time-series can be loaded at a time. The other advantages of these files over the regular output files of the Penn State / National Center of Atmospheric Research's (NCAR) Mesoscale Model version 4 (MM4) Regional Climate Model number 2 (RegCM2) is: data is limited to only that needed for our work, it's in single rather than double-precision (saving disk storage), conversions for wind direction position etc. are already completed. Another problem taken care of with the LA files is elimination of the lateral boundary data. The outermost data points of the run are almost purely driven by the input Boundary Conditions to the run and are thus not useful. Because, the LA files include only the inner portions of the data this part of the data is already stripped out thus saving disk space and preventing Forest Service personnel from having to remember or determine these limits. Furthermore the LA files only use data stored from the regular RegCM2 output so new LA files can be added after a RegCM2 run is made. Thus, the researchers at Utah State University / Space Dynamics Lab. (USU/SDL) have already performed the difficult and obscure conversions needed to be done on the data, as well as eliminating the boundary data that is not relevant. This way the researchers at the Forest Service can concentrate on their work rather than being required to understand obscure conversions on the data.

The researchers at USU/SDL will need help from the Forest Service in deciding where the boundaries for each LA file should lie. There should be a larger amount of overlapping and redundancy to take care of as many of the possible variations that Forest Service researcher's will need. But, as the LA files can be constructed from the regular output files of RegCM2, new LA files can be added at any time.

## NOMENCLATURE

LA files are addressed as LA##SURF??, LA##BATS??, LA##ELEV??, and LA##RAD?? files. The "##" corresponds to a 2 digit number to identify which area this file corresponds too, the "??" corresponds to the last two digits of the year the file stores information on. LA files are designed to store an entire year worth of data in each file. "SURF" is the surface information file, BATS is the BATS surface information file, "ELEV" is the altitude profile information file, and "RAD" is the radiation/ explicit moisture information file.

The files are organized so that the standard driving variables are in one file. The other files contain different data useful in utilizing the advanced options of ZOOM. This way the disk storage is kept to a minimum. One researcher may want/need one advanced option but not another, thus he can load the files needed and avoid the unneeded files. Because the disk storage requirements large, it is important to provide flexibility in the creation of ZOOM input files. By nature ZOOM needs a long time series, but not necessarily a large area, so we created files that contain a complete year of information for a limited area. At the same time it was realized that at a later time some more advanced options may require certain new files thought to be unnecessary at first. For example, the dynamic option of ZOOM using the ELEV LA file to get Lapse rate. Or using the RAD LA file to get rainwater at higher elevations. Also as more experience is built up with ZOOM we will find which options are the most important for making it run the most efficiently with the best results. Following we explain the different file types, explaining if the file is a main driver file or an advanced option file, the variables stored, the number and type of vertical layers in the file, and how often the data is stored.

### File Types:

- |      |   |
|------|---|
| BATS | The BATS information file is the main driver of ZOOM. It contains the standard driving information for ZOOM. Containing the driving temperature, pressure, humidity, radiation, and rainfall data. It also contains other surface data either to be used for initial conditions or for model comparison for possible adjustment to the driving variables for more advanced options of ZOOM. Surface data is stored hourly.  |
| ELEV | This file is an advanced option file for the dynamic option of ZOOM. It contains the temperature, N-S, E-W winds, and humidity for each model sigma layer.(the model has 20 layers, but we will most likely restrict this file to the bottom 5 to 10 layers). ELEV data is stored every 6 hours.  |
| SURF | This file stores some of the surface information used in conjunction with the ELEV file. It contains ground pressure, total precipitation, and convective precipitation. SURF data is stored every 6 hours.   |
| RAD  | This file is an advanced option file for ZOOM. Currently the only data stored in the file is the fractional cloud cover. When the RegCM2 model is run with the explicit moisture scheme (which it will be at the 10 km resolution) the file will also store the cloud water and rain water. It stores data for each model sigma layer (the model has 20 layers, but we will most likely restrict this file to the bottom 5 to 10 layers). RAD data is stored every 6 hours. |

## MODEL NOTES

It must be kept in mind at all times when using this data that it is modeled data not real data. It does use a state of the art physically comprehensive model in the same class of model as used by the National Weather Service for forecasting, but a model nonetheless. It does use boundary conditions that are derived from observational data, and the boundary conditions are updated every 12 hours, and these boundary conditions do have the greatest impact in driving the model. But, no matter the 10 km run, the real observational data is almost always several 100's of Km away from your point of interest. Also the data contained within the model represents that for the surface slope, aspect angle, height, and land-use types that the model represents. In general this is the average of each of these over the 10 km grid square of interest. Because this average is done on a 10 km grid square basis, each of these can be significantly different than at any specific point within the grid square. For example the height could be as much as 1000 m higher or lower than the actual height at that point (For the extreme example of a very high and narrow peak or canyon). The model sets the surface slope and aspect angle to zero so this could be off by up to 45 degrees. Also the land-use type could be as wildly diverging as desert to evergreen forest or visa-versa in the most extreme example of an elevation difference that puts a forest next to a desert area. Remember, that the model data more or less represents what the variables would be, as the average of each 10 km grid square, if the topography, land-use, slope and aspect angles were those that the model uses.

This is why it is important to run ZOOM beneath the 10 km data-set to begin to account for some of the local effects of different: elevation, slope, aspect, and land-use. Some local effects can not be accounted for very well -- such as a canyon breeze, or fog or frost pockets. Also, as there is not a two way interaction between the driving variables and the surface variables, some information will indeed be lost. For example, the model may be at a high enough temperature to melt all the snow in the 10 km grid square, but at the zoomed point, say 500 m higher, the snow may be still quite deep. This would significantly cool the driving air temperature, but because this two way interaction is gone the driving air temperature will heat the snow more than it should.

But in comparisons of the model with averaged observed data at 50 km resolution we expect the following. Daily and hourly information can be significantly off, mostly due to the model inaccurately predicting the speed, track, and extent of storms. The model may have the storm moving faster or slower than the storms actual rate and arrive at the right position as much as a day or two before or after the actual storm did. The position may also be off by tens to hundreds of Km, and the extent may be off also. All of this means that in the short term and over a small area the model is not very accurate. But, if you average over a longer time period and a larger area the model produces a reasonable climate. The monthly averaged temperature is likely to be within 2-5 degrees of observed with most of this difference accounted for by the ZOOM adjustments. Precipitation, however, is much more difficult to model. Wintertime precipitation is likely to be within 4% to 50% of actual precipitation, depending on the region. Summertime precipitation can be very bad however, with a strong tendency for over prediction. If you just measure precipitation yes or no and compare this to observations a good correspondence will be shown, but the amounts themselves may be off, by up to 3X actual precipitation. Numerical point storm events are likely and will need to be filtered out by ZOOM. Also for this reason fudge factors may need to be built into ZOOM to compensate for the models over prediction of summertime precipitation. Also orographic effects within a 10 km grid square can effect the distribution of precipitation within a 10 km grid square. A simple parameterization of this effect will be accounted for in ZOOM. Radiation will be reasonable but completely dependent on the clouds and storms predicted by the model. Thus comparisons here can not be on amounts, but comparing cloudy days with modeled days that are cloudy.

## Lambert Conformal Grid Projection

The model is done on a grid that is not even with latitude longitude, but even on a Lambert conformal conic grid projection. In other words, because of the curvature of the earth the grid squares are not quite 10 Km in each direction they are distorted in the manner prescribed by a Lambert Conformal conic projection with standard parallels at 30 degrees North latitude and 60 degrees North latitude. Because, of this distortion the model stores both the latitude and longitude at the center of each model grid square, thus the actual distance on the globe can be calculated for each 10 km grid square to the next or converted to a different projection. However, even though the projection is Lambert Conformal conic, the data has been converted so that all vectors are in standard directional coordinates (N-S, E-W).

## RegCM2 Model Physics Packages

RegCM2 encompasses seven major packages to model the physics of the atmosphere. The seven areas are: The basic atmospheric model itself, the time integration scheme, the cumulus cloud parameterization, the radiation model, the surface physics model (BATS), the atmospheric planetary boundary layer (APBL) model, and the lake model.

**Atmospheric Model** RegCM2 uses the standard Navier-Stokes equations of: continuity, momentum, and Thermodynamic equations to model the basic atmosphere. In addition to this water vapor is tracked in the model and the latent heat of condensation is accounted for. Because, of the complexities of ice-phase physics the latent heat of fusion of water is not tracked. Although at mid-latitudes these processes can be important as most precipitating clouds will likely have at least a portion below freezing. An explicit option can be chosen that tracks cloud water, rain water, and water vapor separately. This provides a more accurate scheme for tracking water especially at high-resolution.

For speed and accuracy RegCM2 uses a Arakwara-B grid. This means that some of the data is horizontally offset by half a grid cell from the other data. The horizontal winds (u and v) are on the dot grid and the rest of the data are on the cross grid points. Similar to the following representation:

Figure 1. (x-y grid)

```
  .  .  .  .  .  .  .  .
  X  X  X  X  X  X  X  X
  .  .  .  .  .  .  .  .
  X  X  X  X  X  X  X  X
  .  .  .  .  .  .  .  .
  X  X  X  X  X  X  X  X
  .  .  .  .  .  .  .  .
```

Note that as in the example there is one less row and column for the cross grid points than for the dot grid points.

## Time integration Scheme

Like most fluid dynamics models, RegCM2 becomes numerically unstable if the time step is not sufficiently small for the size of the horizontal (and vertical) grid steps. This is mainly due to any types of waves that propagate through the medium (such as gravity, sound, or Rossby waves). The time

step must be small enough that the fastest of these waves can only step forward at most one grid square in each time step.

The basic model uses a leap frog approach for advancing forward the differential equation. This has the advantage of using centered differences for calculating derivatives which increases the accuracy of the method. The model calculates the derivatives at the current time-step and then uses them to advance the previous time step values forward two time-steps. The high energy terms are then dropped out with a Asselin filter.

Now, to speed up the computations even further we use a split-explicit time integration scheme to allow for larger time-steps to be used. The fastest waves are gravity waves that are quasi-linear and they only effect a small portion of the mass of the atmosphere. So we use linear theory to correct the solution for these fastest waves. After the normal time step we correct the solution by adding in the correction that results when integrating the linear solution for a time step a fourth the size of the regular time-step, and then again at half the size of the regular time-step. In this way the regular time-step can be about 3X larger than it would have to be otherwise. And even though the split-explicit scheme is running at time-steps a fourth the size of the regular time-step, because it's solution is linear it can be calculated in a small fraction of the time for the regular time-steps.

#### Sub-grid scale Cumulus Cloud Parameterization Model

A very large percentage of rainfall especially in the summer comes from convective type (Cumulus) clouds which are on the scale of 1-20 km in horizontal extent. Because, RegCM2 attempts to model with grid sizes from 10-100 km, these Cumulus clouds must be parameterized as a sub-grid scale process that occurs within a model grid square. RegCM2 has two one-dimensional Cumulus parameterization schemes: a Kuo-Anthes type approach and a Grell scheme that has flexibility in the closure type used (Arakwara-Shubert or Fritsch-Chappell type closure schemes). The Cumulus cloud scheme will test each layer to see if it is convectively unstable (if you lift a parcel of air to the next layer and the layer is stable or rising -- the layer is convectively unstable). If a layer is convectively unstable it then checks if any water vapor will condense from the rising air parcel. If not it goes on, but if it does -- it has found the base of a Cumulus cloud. At this point there are further checks and the cloud model parameterizes the properties of the cloud itself as well as how it's creation will effect the large-scale parameters. This is the point where the different closure schemes differ.

#### Radiation Model

RegCM2 uses a one-dimensional column model to track radiation. It divides the solar spectrum into 12 bands from 0.2 to 5 micrometers. Seven UV bands, one visible (0.35-0.7 micrometers), and four IR bands. Outside this range the solar spectrum is insignificant, but because of absorption and re-emission by the atmosphere a long wave energy flux is present also. The model includes scattering and absorption by the surface, clouds, water vapor, Ozone, Carbon dioxide and molecular Oxygen. The model tracks the flow of water vapor. Ozone is set by a constant profile. And the mass mixing ratios for Carbon Dioxide and molecular Oxygen are assumed constant. Direct and diffuse radiation is differentiated all the way down to the surface. The Long Wave flux incident from the atmosphere is the total

flux over all wavelengths that is emitted from the atmospheric sources, but as it represents a black-body at about 280 K its peak is around 20 micrometers and is only significant from about 5-30 micrometers.

### Surface Physics Model

Modeling the surface physics and hydrology is important in running an atmospheric model that spans more than a week of simulation time. The surface layer model used by RegCM2 is the Biosphere Atmosphere Transfer Scheme (BATS) model developed by Richard Dickenson and Ann Henderson-Sellers. This model keeps track of the soil moisture in a 3-layer hydrology model. It tracks the soil temperature, surface skin temperature, temperature of the foliage and temperature of the air in the canopy. It tracks the snow depth and areal coverage, and the vegetation growth. As well as determining the drag the vegetation has on the atmosphere.

### Planetary Boundary Layer Model

A significant amount of the energy of the atmosphere near the surface is put into driving small eddies on the scale of 1 cm to 10 meter. Because these eddy's are too small to be resolved by the atmospheric model directly they are parameterized by the stability conditions of the grid cell. These eddys tend to dominate within the Atmospheric Planetary Boundary Layer (APBL), which is the region of the atmosphere from the surface up to 500-1500 meters up -- depending on stability conditions. The model parameterization used by RegCM2 is the Holtslag APBL.

### Lake Model

RegCM2 also has a fresh or salt water columnar lake model. Ocean temperatures are set using the CAC monthly averaged 2.5 degree Sea Surface Temperature produced by NMC. It uses observations from ships, buoys, and satellite. But, to get the surface temperatures over lakes the Hostetler Lake model is used. This lake model takes each model grid point over the lake and models it as a column at 1 meter vertical steps. Attenuation of radiation with depth is accounted for and mixing due to eddy diffusion is parameterized.

### Description the LA file Header

General: All the LA files share the same header type. The header gives the dimension of the file arrays, the limits for this LA file, the sigma levels, the latitude, longitude, surface height, and BATS land-use type for the center of each grid square stored in the file. It also contains a character string description of each of the variables stored in the particular file, a general title and specific LA file type description. And a listing of the valid dates and times the file stores information on. Remember the model is on a Lambert Conformal Conic grid projection so the grid points are evenly spaced in this projection, but distorted in latitude and longitude. So there is not an even spacing in latitude or longitude of the grid points. Because of this the value for both latitude and longitude is given for each and every point stored in the file.

### List of variables:

TITLE(1) This is a 100 character title giving the model version and short description of it.



TITLE(2) This is a 100 character title giving a description of this file type. (Radiation file, Surface file etc.)

IY This is the dimension of the entire model run grid in the "Y" "N-S" direction. The model is on a Lambert Conformal grid projection so moving in the Y grid direction does not directly correspond to moving in N-S, except near the grid center.

JX This is the dimension of the entire model run grid in the "X" "E-W" direction. The model is on a Lambert Conformal grid projection so moving in the X grid direction does not directly correspond to moving in E-W, except near the grid center.

KZ This is the dimension of the number of sigma layers stored in the file. It can be up to the number of layers stored in run. Files with just surface information will just have KZ=1. Because, MCLIGEN does not need the whole atmospheric profile up to the tropopause, but maybe just the bottom 1000 m or only about 10 layers will likely be stored.

IY0 This is the starting index in the "Y" direction for this local area file. LA files only have data stored in the "Y" direction between IY0 and IYF.

IYF This is the ending index in the "Y" direction for this local area file. LA files only have data stored in the "Y" direction between IY0 and IYF.

JX0 This is the starting index in the "X" direction for this local area file. LA files only have data stored in the "X" direction between JX0 and JXF.

JXF This is the ending index in the "X" direction for this local area file. LA files only have data stored in the "X" direction between JX0 and JXF.

NTYPES This is the number of different data-types stored in the LA file.

DATYPE This is a character string description of each data type stored in the file.

NCH This is the number of characters stored in the directory name.

DIR This is the directory of where the file came from.

BEGDAY This is the beginning day for the whole run. In the YYJJJHH format, where YY=last two digits of the year the file is storing, JJJ=the julian day, and HH = hour in Universal time (UT). (Ie. 7833500 = year 1978 julian day 335 (or Dec/1) 00:00 hours UT).

IDATEBEG This is the beginning day that is stored in the file. In the YYJJJHH format, where YY=last two digits of the year the file is storing, JJJ=the julian day, and HH = UT in hours. (Ie. 8200100 = year 1982 julian day 1 (or Jan/1) 00:00 hours UT).

IDATEEND This is the last day that is stored in the file. In the YYJJJHH format, where YY=last two digits of the year the file is storing, JJJ=the julian day, and HH = UT in hours. (Ie. 8236523 = year 1982 julian day 365 (or Dec/31) 23:00 hours UT).

DAYINC	This is the number of days that are stored in each RegCM2 output file. It isn't really needed for LA files but is given anyway.
TOPO	This is the array of elevation heights in meters for the center of each model grid square.
XLAT	This is the array of latitudes in degrees for the center of each model cross-point grid square.
XLON	This is the array of longitudes in degrees for the center of each model cross-point grid square.
XLUSE	This is the array of BATS landuse types that are used for each model grid square. There are 18 BATS land-use types classified as follows:

Table I

- 1 = Crop
- 2 = Short Grass
- 3 = Evergreen Needle forest
- 4 = Deciduous Needle forest
- 5 = Deciduous Broad-leaf forest
- 6 = Evergreen Broad-leaf forest
- 7 = Tall grass
- 8 = Desert
- 9 = Tundra
- 10 = Irrigated Crop
- 11 = Semi-desert
- 12 = Glacier
- 13 = Swamp
- 14 = Lake
- 15 = Ocean
- 16 = Evergreen shrub
- 17 = Deciduous shrub
- 18 = Mixed forest (decid. and evergreen mix)

Because of the lack of soil texture and color information, the land-use type is also used to estimate the soil texture and color class. BATS has 12 soil texture types and 8 soil color types. Texture class 1 corresponds to sand, class 6 is loam and class 12 is clay. Soil color class 1 is light and 8 is dark.

Table II

Landuse type = Texture class, color [veg]

- 1 = Tex 6, Color 5 [Crop]
- 2 = Tex 6, Color 3 [Short Grass]
- 3 = Tex 6, Color 4 [Everg. Needle]
- 4 = Tex 6, Color 4 [Decid. Needle]
- 5 = Tex 7, Color 4 [Decid. Broad]
- 6 = Tex 8, Color 4 [Ever. Broad]
- 7 = Tex 6, Color 4 [Tall grass]
- 8 = Tex 3, Color 1 [Desert]
- 9 = Tex 6, Color 3 [Tundra]

- 10 = Tex 6, Color 3 [Irr. Crop]
- 11 = Tex 5, Color 2 [Semi-desert]
- 12 = Tex 12, Color 1 [Glacier]
- 13 = Tex 6, Color 5 [Swamp]
- 14 = Tex 6, Color 5 [Lake]
- 15 = Tex 6, Color 5 [Ocean]
- 16 = Tex 6, Color 4 [Ever. shrub]
- 17 = Tex 5, Color 3 [Decid. shrub]
- 18 = Tex 6, Color 4 [Mix forest]

The porosity of the soil (The fraction of the soil volume that is void and can hold water) is associated with its texture class the 12 classes are assigned the following values for porosity:

Table III

Texture class (type)	= Porosity (unitless (volume/volume))
1 (Sand)	= .33
2	= .36
3	= .39
4	= .42
5	= .45
6 (Loam)	= .48
7	= .51
8	= .54
9	= .57
10	= .60
11	= .63
12 (Clay)	= .66

#### SIGMA

If KZ is larger than one than the (KZ+1) full sigma levels are stored. Sigma is a terrain following coordinate related to the pressure. Defined as:

P top is the pressure at the model top. The standard P top is 80 mb. Data is stored on the half sigma level, layer 1 is average of full sigma height 1 and 2, layer 2 is average of full sigma height 2 and 3, etc. The 21 "standard" full sigma levels are:

1.0, 0.995, 0.987, 0.977, 0.96, 0.945, 0.925, 0.89, 0.84, 0.79, 0.71, 0.62, 0.53, 0.44, 0.35, 0.27, 0.19, 0.12, 0.07, 0.03, 0.00

Sigma = 1.0 is the surface and sigma = 0.0 is at the model top (80 mb). So the 20 "standard" half sigma levels where the data is, are:

0.9975, 0.991, 0.982, 0.9685, 0.9525, 0.935, 0.9075, 0.865, 0.815, 0.75, 0.665, 0.575, 0.485, 0.395, 0.31, 0.23, 0.155, 0.095, 0.05, 0.015

This corresponds to about 25 m above the surface to about 20 km above the surface at the highest layers. To get the actual heights you need to integrate the hydrostatic equation, which has temperature, pressure and humidity as input. Using the standard atmosphere we can give approximate ranges for the sigma levels as follows:

Table IV. Half step Sigma layers for 10 km run:

	Sigma	Pressure (mb)	Height (m)
1	0.9975	1010	15 - 21
2	0.991	1004	53 - 77
3	0.982	996	106 - 150
4	0.9685	983	190 - 270
5	0.9525	969	280 - 410
6	0.935	953	390 - 570
7	0.9075	927	560 - 820
8	0.865	887	830 - 1200
9	0.815	841	1200 - 1700
10	0.75	780	1600 - 2300
11	0.665	700	2200 - 3200
12	0.575	617	3000 - 4300
13	0.485	533	3800 - 5500
14	0.395	449	4700 - 6900
15	0.31	370	5800 - 8400
16	0.23	295	6900 - 10000
17	0.155	225	8200 - 12000
18	0.095	169	9500 - 14000
19	0.05	127	11000 - 16000
20	0.015	94	12000 - 19000

#### Description of the Output data section

**General:** For each time output: the time, date, and all data is stored. For each LA file there is a header and then the output data section is repeated for each time period stored until the end of file. Different file-types store data at different frequencies. But all store a years worth of data in each file.

#### List of Variables:

**XTIMEC**      The time in minutes since the start of the original run.

**IDATEX**      The date in YYJJHH format. Where YY=Last two digits of the year, JJ=Julian day, HH=Hour (Ie. 8236506 = year 1982, julian day 365 (Dec/31) and hour (Universal time) = 06:00 UT.

SURFDATA The surface (or elevation) data array stored from IY0-IYF, JX0-JXF, and 1 to KZ, and data-types 1-NTYPES.

#### Description of the BATS LA file data-types

General: The BATS LA files correspond to the main driving variables for ZOOM, plus some variables added for initialization or comparison.

#### List of Variables:

T ground This is the surface (snow or soil) skin temperature calculated using the "Force restore" method.

Input to this calculation is the radiation balance (R), and a sub-soil temperature calculated at a depth of:

Where  $K_{soil}$  is the thermal conductivity of the soil (and/or snow layer) ( $m^2/sec$ ) and  $TAU$  is the number of seconds in a day. Thermal conductivity of the soil depends both on the soil type and the soil moisture content. So this depth actually changes from location to location and changes in time as the soil moisture content changes. But, basically it is the depth at which diurnal variation is damped out and only longer time-scale variance is seen (around 10 cm).

Ground temperature is included as a comparison for the model and for initialization.

T air This is the temperature of the lowest model sigma layer (usually  $\sigma = 0.9975$  or about 25 m above the surface). So it represents the temperature of the air averaged over an entire 10 km grid square between the surface and the next closest full sigma level (usually 0.995). To find the representative height of this layer you must integrate the hydrostatic equation to get the height of the layer (See the section on sigma for this expression).

Air temperature is a driving parameter for zoom.

EW-Wind This is the wind in the east direction (West is negative) for the lowest model sigma layer. To get these values the model data had to interpolate the values from the dot grid points (a grid with points half-way between four neighboring cross grid points) to cross grid points, take into account curvature of the earth from the projection and then rotate the vectors such that the E-W component is given rather than the model "u" component in the "X" direction of the map projection.

The winds are a driving parameter for zoom.

NS-Wind	<p>This is the wind in the north direction (South is negative) for the lowest model sigma layer. To get these values the model data had to interpolate the values from the dot grid points (a grid with points half-way between four neighboring cross grid points) to cross grid points, take into account curvature of the earth from the projection and then rotate the vectors such that the N-S component is given rather than the model "v" component in the "Y" direction of the map projection.</p>
Total Soil Water	<p>This is the soil water in the total soil column down to a 3 m depth. (Note the total soil layer includes the root-zone soil layer as it's highest portion.) To get the percent of saturation that this represents take the soil porosity for this soil texture class (given by the land-use type) and multiply it by the depth of this layer (3 m).</p> <p>This is added for comparison and/or initialization.</p>
Snow	<p>This is snow water equivalent for the precipitation that has fallen as snow and hasn't melted as yet. BATS sets all precipitation that falls at a temperature below 1 degree Celsius as snow and everything above that mark as rain.</p> <p>This is added for comparison and/or initialization.</p>
Precipitation	<p>This is total precipitation accumulated since the beginning of the model run in cm. It includes both the grid scale resolvable precipitation as well as the parameterized sub-grid scale resolvable.</p> <p>This is a driving parameter for zoom.</p>
Conv. Precip.	<p>This is convective precipitation accumulated since the beginning of the model run in cm. It only includes the parameterized sub-grid scale resolvable portion of the precipitation. Normally it is generated by the Fritsch and Chappell type closure scheme from the Grell Cumulus cloud sub-grid scale non-resolvable convective cumulus cloud parameterization scheme.</p> <p>This is a driving parameter for zoom.</p>
Upper Soil Water	<p>This is the amount of water in the soil in mm in the upper soil layer (10 cm depth). To get the percent of saturation that this represents take the soil porosity for this soil texture class (given by the land-use type) and multiply it by the depth of this layer (10 cm).</p> <p>This is added for initialization and/or comparison.</p>
Visible SW radiation	<p>This is visible radiation in <math>W/m^2</math> of the radiation reaching the surface in the band between 0.35 and 0.7 micrometers. Including direct and diffuse radiation.</p> <p>This is a driving variable for zoom.</p>

SW incident This is total solar radiation in  $W/m^{**2}$  of the radiation reaching the surface in the band between 0.2 and 2.0 micrometers. Including direct and diffuse radiation.

This is a driving variable for zoom.

net SW This is net solar radiation in  $W/m^{**2}$  of the radiation absorbed by the surface in the band between 0.2 and 2.0 micrometers. Including direct and diffuse radiation.

This is added for comparison.

Humidity ground This is the mixing ratio of mass of water vapor to mass of dry air (Kg/Kg) for the surface layer.

This is added for comparison in zoom.

Humidity air This is the mixing ratio of mass of water vapor to mass of dry air (Kg/Kg) for the lowest model sigma layer (usually 0.9975 or about 20 m above the surface).

This is a driving variable for zoom.

Net LW radiation This is incident longwave radiation minus the longwave radiation released from the surface (proportional to  $T_g^{**4}$ ).

This is added for comparison in zoom.

Incident LW radiation Incident LW is the horizontal longwave radiation incident on the surface. It is the radiation at all wavelengths that is emitted from that atmosphere.

This is a driving variable in zoom.

Pressure ground This is the pressure at the model surface in Pascals.

This is a driving variable in zoom.

Diffuse visible rad. This is the diffuse component of the visible radiation between 0.35 and 0.7 micrometers.

This is a driving variable for zoom.

Root-zone soil water This is the amount of water in the soil column down to the root-zone soil level in mm. The depth of this layer depends on the vegetation class but varies between 1 to 2 meters (note the root-zone soil layer includes the upper layer soil layer as it's upper portion).

Table V.

1 = Crop	[1 m]
2 = Short Grass	[1 m]
3 = Evergreen Needle	[1.5 m]
4 = Deciduous Needle	[1.5 m]
5 = Deciduous Broad	[2 m]
6 = Evergreen Broad	[1.5 m]
7 = Tall grass	[1 m]
8 = Desert	[1 m]
9 = Tundra	[1 m]
10 = Irrigated Crop	[1 m]
11 = Semi-desert	[1 m]
12 = Glacier	[1 m]
13 = Swamp	[1 m]
14 = Lake	[1 m]
15 = Ocean	[1 m]
16 = Evergreen shrub	[1 m]
17 = Deciduous shrub	[1 m]
18 = Mixed forest	[2 m]

Note: Although a value is given for Lake and Ocean the model doesn't really use these values.

To get the percent of saturation that this represents take the soil porosity for this soil texture class (given by the land-use type) and multiply it by the depth of this layer.

This is added for initialization and/or comparison in zoom.

#### Description of the ELEV LA file data-types

General: This file contains the driving variables for zoom under the dynamic option. Primarily it is used to find the slope with elevation that each of the driving variables have. After finding the slope each variable can more accurately be adjusted for elevation. This file is only stored every 6 hours so the slope must be considered constant or varying only linearly during this time period.

#### List of Variables:

T air This is the temperature of the each respective model sigma layer (See the section on sigma to get the sigma levels). So it represents the temperature of the air averaged over an entire 10 km grid square between the nearest sigma layer.

EW-Wind This is the wind in the east direction (West is negative) for the respective model sigma layer. To get these values the model data had to interpolate the values from the dot grid points (a grid with points half-way between four neighboring cross grid points) to cross grid points, take into account



curvature of the earth from the projection and then rotate the vectors such that the E-W component is given rather than the model "u" component in the "X" direction of the map projection.

**NS-Wind** This is the wind in the north direction (South is negative) for the respective model sigma layer. To get these values the model data had to interpolate the values from the dot grid points (a grid with points half-way between four neighboring cross grid points) to cross grid points, take into account curvature of the earth from the projection and then rotate the vectors such that the N-S component is given rather than the model "v" component in the "Y" direction of the map projection.

**Humidity air** This is the mixing ratio of mass of water vapor to mass of dry air (Kg/Kg) for the respective model sigma layer.

#### Description of the SURF LA file data-types

**General:** This is the surface data that goes with the elevation data. It is included only to get the values that correspond with the ELEV files. All of this information is repeated in the BATS files but at a hourly basis rather than every 6 hours.

#### List of Variables:

**Pressure ground** This is the pressure at the model surface in centibar.

**Precipitation** This is total precipitation accumulated since the beginning of the model run in cm. It includes both the grid scale resolvable precipitation as well as the parameterized sub-grid scale resolvable.

**Conv. Precip.** This is convective precipitation accumulated since the beginning of the model run in cm. It only includes the parameterized sub-grid scale resolvable portion of the precipitation. Normally it is generated by the Fritsch and Chappell type closure scheme from the Grell Cumulus cloud sub-grid scale non-resolvable convective cumulus cloud parameterization scheme.

#### Description of the RAD LA file data-types

**General:** The RAD files keep track of some specialized information not normally used in Zoom, but useful for some advanced options. Cloud water and rain water is only included for the explicit option in RegCM2, which is the normal case for 10 km runs.

#### List of Variables:

**Cloud Cover** This is the fractional cloud cover at each sigma layer ranging from 0.0, for no clouds, to 1.0, for 100% coverage of the cloud over the grid cell. The bottom 3 sigma layers are not allowed to have clouds. For the cumulus parameterization cloud cover is assumed to be 100% for 10 km runs.

**Cloud Water** This is the mass mixing ratio of water that has condensed into cloud form. Over a grid cell without clouds it is zero and within a cloud will be higher. The units are Kg/Kg.

**Rain Water** This is the mass mixing ratio of water that is now falling as precipitation. Because, a significant portion of the rain may evaporate as it falls the rain water is useful in adjusting the precipitation with elevation. The units are Kg/Kg.

# **APPENDIX 2B**

**13 YEAR MESOSCALE MODELING STUDIES VIA A MULTI-PROCESSING  
VERSION OF A REGIONAL CLIMATE MODEL**

## ABSTRACT

Utah State University (USU) and the Interdisciplinary Climate Studies (ICS) group at the National Center of Atmospheric Research (NCAR) have teamed up. Using NCAR's RegCM2 (Regional Climate Model version 2) model which includes a surface physics and hydrology model (BATS -- Biosphere Atmosphere Transfer Scheme) we are producing a multi-year mesoscale climate sequence. These are historical 13 year simulations (December 1978 to April 1992) from nested climate models at three resolutions, T42 (the observational analysis of the European Center of Medium Range Weather Forecasting (2.81 x 2.81 degrees), 50 km and 10 km respectively. Without access to a super-computer we ported RegCM2 to our IBM-6000 and Indigo-2 work stations. However, running RegCM2 at 10 km resolution over a 400 km square grid for a 13 year sequence proved to be too CPU intensive, taking over 2 years of CPU time to process.

Therefore, a parallel version of RegCM2 was developed by us to cut the CPU time by the number of available processors. A preliminary version has cut the time for 8 processors by four times, allowing a 13 year run to be done in 12 months. The current version we are developing will half the CPU time once again.

Our method of parallelization was designed for capability on both a shared-memory multiple-CPU (such as our 8-CPU SGI Onyx System) or a distributed system of SGI Indigo's on an FDDI fiber optic ring operating up to 10 nodes at once. However, we found the model to be too tightly coupled to take advantage of a distributed system as it requires passing too high of an amount of information each time-step.

Timing information for the model and an overview of the method of parallelization will be presented by us, along with some of the hydrologic outputs of the model. We will also discuss the impact of a longer time series of mesoscale model output.

### Serial RegCM2

A 15 year simulation of MM4 at the 10 km level will take nearly 26.5 months of CPU time. For a 50 km run with the Kuo scheme it is 11 months of CPU.

Synopsis of MM4 code:

```
MM4 main program:
  initialization
  do while ( xtime < timax )
    if( input time ) Input BC's
  tend ----- Model Dynamics.
    splitf ---- Split-explicit time      integration.
    Track days/date/time etc.
    if( output time) Output
  end do
  output /closing etc.
```

The TEND subroutine and it's children takes over 95% of the CPU time. So it is most natural to work with it.

TEND Subroutine (Model Dynamics):

Initialization

Decouple first 4 "j" (longitude) slices put in temporary arrays.

do j = 2, jx-1

Calculate pressure, temperature, humidity and wind tendencies (time derivatives), horizontal and vertical diffusion, pressure gradient, horizontal and vertical advection (a gradient term multiplied by a wind term).

Calculate forecast values for next time step.

HOLTBL --- Holtslag PBL routine. (Which calls both BATS (every 360 simulation seconds) and RAD every 30 simulation minutes).

Decouple next "j" slice.

Copy forecast values into current arrays, and current values into previous arrays.

Juggle the temporary indices around. To allow for next "j" slice.

end do

Copy the last j slices from forecast array into current array and current array into previous array.

bdyval ----- Boundary values.

nconvp ----- Non-convective

precipitation.

conmas ----- Mass conservation check.

solar1 ----- Compute solar zenith angle.

The diffusion and horizontal advection routines are 4th order so they need j levels + or - 2 levels away from j level solving for. Holtslag PBL needs to have the previous j-slice calculated before doing the current j-slice and it needs to have the drag calculated from the BATS routine before it can calculate the PBL. Other routines just need current j level.

### Parallel RegCM2

Because, the j-loop in TEND takes up most of the CPU (92% of the total computation) and as it provides a natural place to parallelize the code, we have worked at distributing this loop over each CPU. First machine solves j=1-5, next 6-10, etc. For initialization we must calculate the geopotential on the previous j-slice and the vertical wind. Also for the Holtslag PBL we must first calculate the tendencies for temperature and water (vapor, cloud, and rain), and then run BATS to get the drag for the j-2 j slice. Then we calculate the temperature and water tendencies for the j - 1 j-slice and calculate the PBL for the j-1 j-slice. At this point we can begin the j loop for the given node. We save the results for these previous 2 j-slices so that we don't have to recalculate them for the previous node.

Because, this method does not parallelize the entire code it will not increase the speed directly proportional to the number of nodes but as given by Amadahl's Law:

$$F = \frac{1}{\frac{\text{percent}}{\text{\#nodes}} + (100 - \text{percent})}$$

So for our case where the percent parallelized is 92% and the number of nodes is 4 we get 3.2 as our speedup. Now in actual runs the speed-up was 2.2. We believe the discrepancy is due partially to the additional time taken for parallelization, having CPU's wait for resources, and problems with load balancing.

Table of Run Times for 15-Year Climate Run

Resolution (km)	Grid Size	Machine	Run-Time (months)
10	36x37x20	Cray YMP	3.5
10	36x37x20	SGI indigo	30
10	36x37x20	IBM 6000	25
10	36x37x20	SGI Indigo-2	25
10	36x37x20	HP-750	50
10	36x37x20	SGI Risc-8800	15
10	36x37x20	SGI Onyx-1 node	25
10	36x37x20	SGI Onyx-2 nodes	16
10	36x37x20	SGI Onyx-4 nodes	11
10	36x37x20	SGI Onyx-8 nodes	8

### Fully Parallel RegCM2

The fully parallel version of RegCM2 builds upon the previous version that only does the TEND j-loop in parallel. Making the rest of the TEND subroutine parallel is not too difficult. Most of the routines are straight forward and the only difference is restricting the j loops. However, the SPLITF routine is more difficult. SPLITF adjusts the output of TEND using linear theory.

### Distributed RegCM2

Distributed RegCM2 takes the fully parallel RegCM2 and distributes the work on totally separate machines on a FDDI network. The problem here is that each node has to have information on the other nodes before proceeding. So rather than doing computations we do data transfer. Because, of the complexities with the Holtslag PBL, the explicit moisture scheme, and the surface physics model (BATS) 1.7 MB of data has to be transferred for a 4 node system (for a 36x37x20 grid) each time step. As the number of nodes increases this amount increases and the timestep itself decreases. This means that the percent of time spent transferring data increases exponentially with the number of nodes. Thus, eventually curtailing any speedup by adding more CPU's. We show in our example estimates for a 4 node case. A four node case spends 15% of it's time transferring data, and a 6 node case spends 30% transferring data. This is about the practical limit unless you can allow unlimited traffic on your network. For some of the simpler options of RegCM2 this percentage is much smaller. So the distributed system is useful for a small number of nodes up to 4 or 5, but with access to shared memory systems the time is much better spent.

Distributed System transfer rates:

Because, the CPU times for each time step are so short it is not necessarily cost effective to transfer the data (approx. 3.22 MB) for each time step. The theoretical transfer rate is 10 MB/sec. However, tests show that realistically a much lower rate is expected. The results are shown below.

### TRANSFER RATE TEST RESULTS

MBytes	MB/sec	MBytes	MB/sec	MBytes	MB/sec
.00002	.005	.8	5.03	1.8	4.68
.001	.2	.9	5.26	1.9	4.66
.01	2.0	1.0	5.29	2.0	4.52
.1	3.70	1.1	5.26	3.0	4.39
.2	3.57	1.2	5.17	5.0	4.38
.3	3.30	1.3	5.14	--	--
.4	6.35	1.4	5.02	--	--
.5	11.11	1.5	4.89	--	--
.6	6.98	1.6	4.72	--	--
.7	4.96	1.7	4.72	--	--

Load Balance:

In tests with MM4 the load is fairly evenly divided among each j-slice. The CPU's that contain the boundary data will be unbalanced by up to about 10%, as they have additional calculations for the BC's while at the same time aren't calling BATS or CCM2 Radiation. But this is an acceptable level, this means that CPU's will only spinning for less than a second. Changing the way j-s are divided among the CPU's does not improve the load balance. Outside the boundaries the load is balanced to much less than 1%.

The biggest area of concern with the load balance is the jx that you choose for the run. You must choose a jx so that:

$$(jx-3) / \# \text{ nodes}$$

is an even number. Where jx is the number of grid points in the j (longitude) direction, and # nodes is the number of compute nodes (threads) that you want to use.

Distributed Parallel version [4 nodes] (2.7 X speedup)

Routine	CPU (sec) (1 pass)	Load Balance	Fraction called	Seconds per step	Percent
TEDJLOOP	6.412/4	1.2	1	1.923	74.5%
BATS	0.725/4	1.1	1/10	0.020	0.75%
RAD	5.0/4	1.1	1/50	0.028	1%
Input	0.093	1.0	1/600	0.0002	neg
Input transfer	1.1 MB /5 MB/sec	1.0	1/600	0.0004	neg
Output	5.713	1.0	1/600	0.0095	0.5%
Out transfer	3.2 MB /4MB/sec	1.0	1/600	0.0013	neg
BATS output	0.093	1.0	1/100	0.0009	neg
BATS output transfer	0.11 MB /3MB/sec	1.0	1/100	0.0004	neg
splitf	0.329/4	1.5	1	0.1235	4.75%
transfer	1.6MB /4MB/sec	1.0	1	0.4	15.5%
Other TEND	0.088/4	1.5	1	0.033	1.25%
Overhead	0.046	1.0	1	0.046	1.75%
Total CPU Time per loop	2.586 [7.07 serial]				100.0%



## Conclusions

We have developed a preliminary version of parallel RegCM2 that opens up the capability to modeling 15 year climate sequences with a high resolution atmospheric model with doable run times on small machines. A fully parallel version of RegCM2 is also in the works that will decrease the run times even further. Having longer high resolution climate sequences will allow better regional and local scale assessments of GCM's and GCM doubled CO2 runs than currently available. We believe this is imperitative to the understanding of the impacts of global climate change for the local decision makers trying to plan what impact a changing climate may have on their situation. We also believe that a model such as RegCM2 that runs in parallel on small machines is one of the most cost-effective ways of obtaining high resolution climate sequences.

# **APPENDIX 2C**

**A VIDEO DISPLAY SYSTEM FOR QUALITY CONTROL OF HYDRO-  
METEOROLOGICAL MODEL DATA**

The huge quantity of data generated by GCMs and regional hydro-meteorological models makes it difficult to evaluate the detailed features of the output of these models. Often, the primary variables of these runs are only evaluated as monthly, seasonal or annual means, while the secondary variables are seldom evaluated. To allow us to evaluate the port of the RegCM2 model to workstations, we developed a video display system for the data. The display system provided color scaled images for each of the model outputs at each time step. A clock symbol shows the passage of time as each data image is placed onto video tape. The tape can be played back to show the temporal and spatial range of variation of each variable.

The Visual Images Data Display System (VIDDS) is controlled by a graphical control system that allows the choice of variables, the resolution, and speed of display. It also provides a choice of the type of data display. Data can be displayed on a terrain contour or a map overlay type background. The angle and orientation of the terrain painted data can be changed as desired.

This poster will show a video demonstration of the display of 50 km data from RegCM2 in the western US from our historical comparison for the period 1979-1981. The display also demonstrates the display on a nested 10 km resolution area in northern Utah. The model shows the range of temperature, precipitation, soil moisture, vegetative cover, and evaporation rates for the models domain. While the air temperatures and precipitation rates are reasonable, the display shows that, for particular time periods, the leaf temperatures and evaporation rates of the model in the complex terrain are not reasonable. A technical description of the VIDDS is provided.

### **Technical Description**

VIDDS uses a Graphical User Interface (GUI) of "C" X-11 Motif code produced by Builder-X to set up the files parameters to be plotted etc. The GUI then spawns the graphics plotting program which uses "C" and Silicon Graphics Inc. Graphics Language code (SGI GL) to plot the surface in 3D as a triangular mesh interpolating between grid points with Lambert shading. The colors are picked using a color look-up table defined by the user or picked using the minimum and maximum of the data array.

### **Conclusion**

With the recent advances in graphical workstations we can now plot a time series of visual images of atmospheric model data in 3D and in real time. This method of viewing the model data can allow the researcher to understand details hidden by other methods. This method also allows a researcher to view all the data output by the model in a reasonable amount of time. This method can bring out deficiency's and model errors not noticed by other methods. In our case it has helped us to catch several errors in the model and better understand the dynamics of the model and the system.

# **APPENDIX 2D**

**A NESTED MODEL CHAIN BETWEEN GCM SCALE AND RIVER FLOW**

A NESTED MODEL CHAIN BETWEEN GCM SCALE AND RIVER FLOW:  
A Testbed For Vegetation, Erosion, and Water Yield Scaling Studies\*.

G.E. Bingham\*\*, D.S. Bowles, E. Kluzek, A.S. Limaye and J.P. Riley

Utah State University  
Logan, UT. USA

1. INTRODUCTION

The interpretation of GCM indicated climate change (Mearns et al., 1990) at regional and local levels is complex due to subgrid scale effects that are ignored by the larger scale models. This is particularly evident in mountainous terrain, such as the western U.S., where whole mountain ranges are smoothed over at GCM scales. When attempting to apply GCM climatic data to vegetation and hydrologic studies, the short term variations in temperature and precipitation are often lost in climatic averaging. It is these fluctuations that provide the stresses that limit vegetative development and cause extreme events. In this paper, we describe the development of a modeling system and database designed to allow us to study the accuracy of nested regional model predictions and the scaling effects found in the models being developed to examine these issues. A validation data base is also being developed that provides the basis for multi-year model to historical data intercomparisons in the Rocky Mountains. The climate modeling portion of this effort is a joint project with the Interdisciplinary Climate Systems Section (ICS-CGD), National Center for Atmospheric Research (NCAR), Boulder CO. The system provides a physically based robust modeling system to link global, decade scale climate inputs to vegetation change and river basin flow studies.

2. THE MODEL SYSTEM

An overview of the Mountain Climate - Hydrometeorology Model System (MCHS) is shown in Figure 1. The atmospheric portion of the system has its foundation in two well documented models, (RegCM2 Giorgi, et, 1993) and BATS (Dickinson et al., 1986). RegCM2 is nested to run at two resolutions, 50 and 10

Km. BATS provides surface inputs to RegCM2 at both

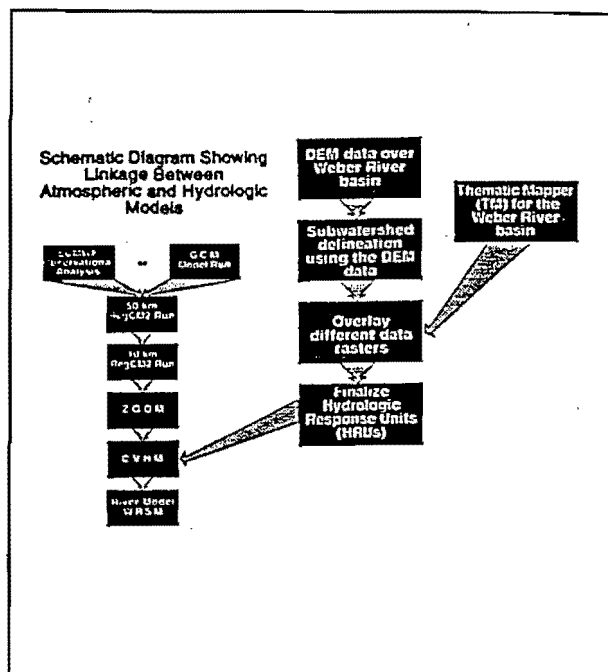


Figure 1. Overview of MCHS.

resolutions. Input to the current version of RegCM2 is either the T42 (2.81°) output of NCAR's CCM2 GCM or NCAR's archive of ECMWF data that has been formatted to match the CCM2 output. Data inputs to the model and the output file structure of the 50 Km and 10 Km model runs are shown in Figures 2 and 3. Model output is stored on an on-line mass storage device at the Space Dynamics Laboratory at USU. Super computer time to run nested model studies like those described here is expensive and difficult to acquire. We have modified the NCAR provided RegCM2 model to run on RISC workstations, but runtimes on even the latest versions limit detailed studies. A fully parallel version of RegCM2 for multi-processor workstations has been completed (Kluzek et al., 1994). This version reduces the run time requirements significantly. A one year 10 Km run (36x37x20 grid points) on an 8 processor Silicon Graphics Onyx requires about 12 days. The outputs of the nested RegCM2 model provide all of the radiation, wind, temperature, precipitation and soil moisture data required for detailed hydrometeorological and vegetation studies.

\*This research was supported in part by funds provided by the Intermountain Research Station, Forest Service, U.S. Department of Agriculture.

\*\*Corresponding author address: Gail E. Bingham, Space Dynamics Laboratory, Utah State University, 1695 North Research Park Way, Logan, Utah 84321-1942

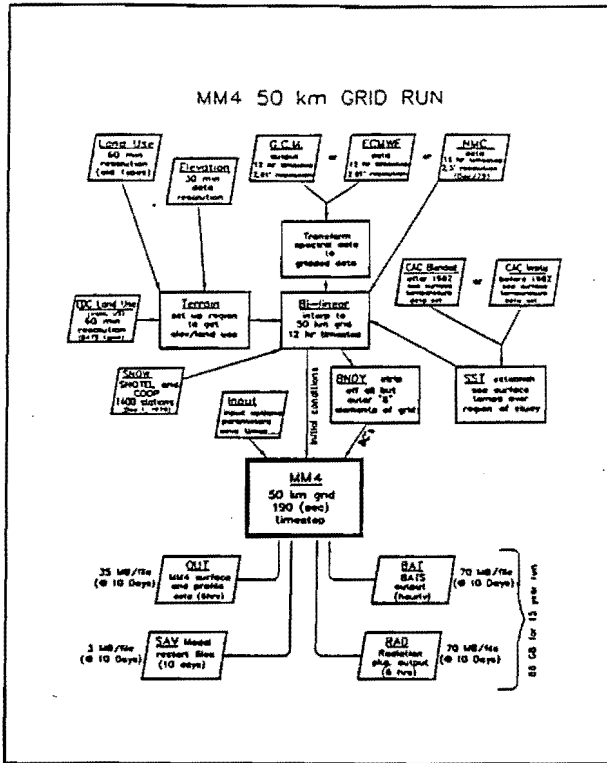


Figure 2. Output file structure of 50 Km.

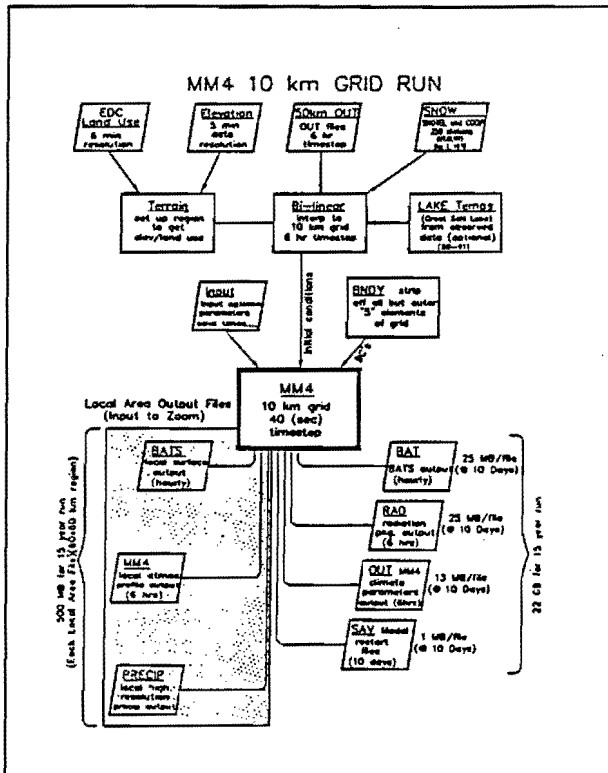


Figure 3. Output file structure of 10 Km.

of the MCHS and the vegetation and hydrologic models is a USU developed model (Zoom). Zoom interpolates the RegCM2 data from the "smooth" surfaces of the 10 Km model to the topographical scale desired for the next layer of models. Two versions of Zoom are now available. The first option provides single location time series data such as would be provided by a complex weather station. Our current single station option is used to drive a single point version of the BATS model. Some of the basic BATS surface types have been replaced with conditions more suitable to western mountain surface types. A second single point model is being developed in cooperation with the U.S. Forest Service to provide a version of the WEPP soil erosion model (NSERL, 1991) for use in complex terrain. The second Zoom option provides horizontally distributed output temporal data sequences to drive HRU or grid based models.

Our hydrologic model (CVHM - Sikka, et al 1993) is a modified version of the PRMS model, (Leavesley, 1983). The structure of CVHM is shown in Figure 4. CVHM operates at daily time steps using daily minimum and maximum temperatures, radiation and precipitation. Model outputs are surface and ground water flows, evaporation and transpiration, soil moisture at two levels, and leaf and air temperatures. The model uses a parameterized soil moisture leaf conductance submodel to calculate transpiration as a function of soil moisture, radiation, humidity and temperature. When used in multiple year simulations, CVHM has the ability to annually adjust leaf area index based on cumulative transpiration.

CVHM is an Hydrologic Response Unit (HRU) based model, requiring the input of HRU boundaries, vegetation type and leaf area index (LAI). This information is developed from DEM and LandSat TM data. This process is time consuming and somewhat subjective. Optimization of this process, both in scale and procedure, is one of our initial research thrusts. Currently, water shed definition is accomplished using the MIPS analysis system. Vegetation type and leaf area are derived using the MIPS system and LandSat TM data. The TM data is sorted into vegetation classes using an automated, fuzzy classification scheme developed by Gunderson, et al., (1992). The DEM and TM data sets are then merged using a rule based system to provide the HRUs that allow CVHM to be applied across the region being studied. Runoff, both surface and subsurface, is collected into a river basin model which provides the routing and hydrograph calculations.

### 3. DATABASE DEVELOPMENT - An Example

Figure 5 shows an example application of MCHS. We are currently applying the MCHS to the western U.S. for the period 1979-1993. This is the period for which ECMWF global observation data exist. Using observed upper air data to constrain RegCM2 should allow the output data to be compared with recorded

The interface between the atmospheric portion

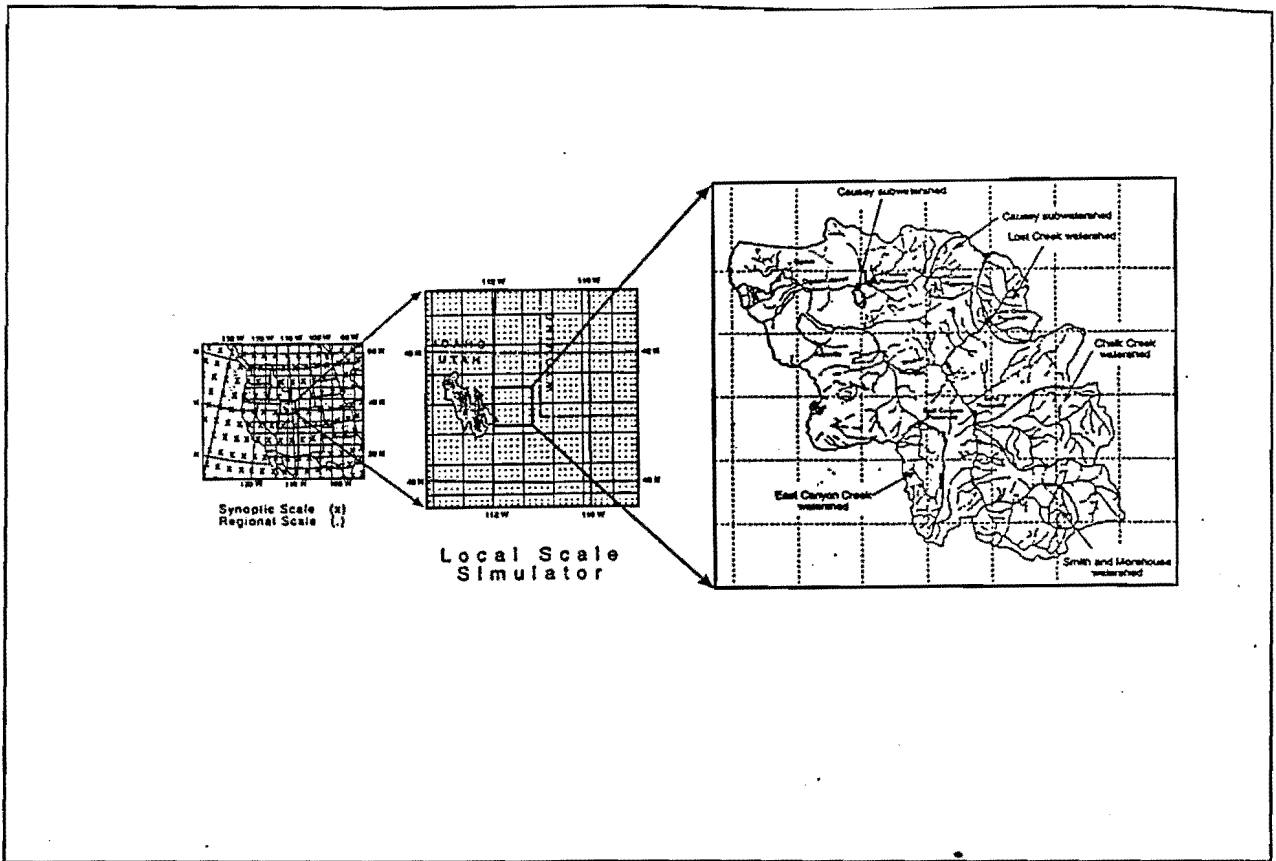


Figure 4. An example application of MCHS.

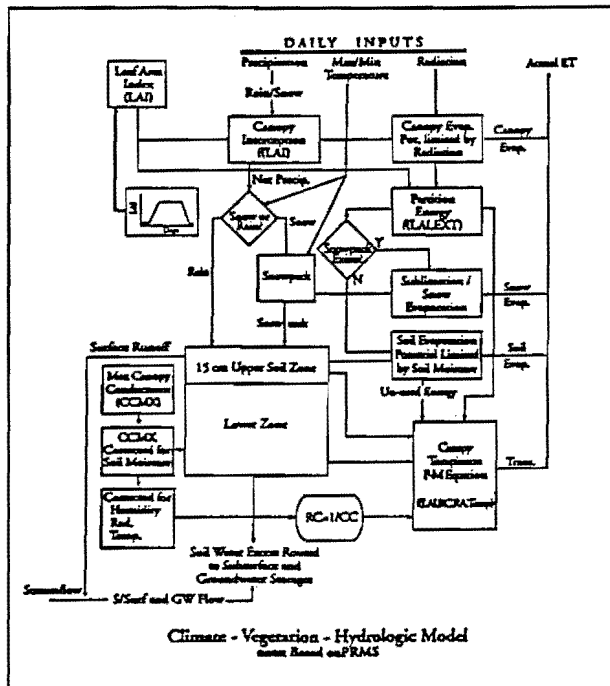


Figure 5. The structure of CVHM.

surface climate data. The 50 Km resolution regional model domain includes all of the western U.S. Two local (10 Km) regions are also being modeled for this period. The first area to be completed covers the Bear River drainage in Wyoming, Idaho and Utah; the second will center on northern Idaho. Hydrologic and vegetative response studies are being conducted on the Weber River basin in north-central Utah. The Weber is a sub watershed in the Bear River drainage. Our intent is to eventually expand the study to include the whole Bear River basin.

To examine the fidelity of the RegCM2 model output, a historic climate data base (which includes the NOAA Coop, RAWS and Snotel data for the region) has been collected and quality controlled. This data set is currently being gridded and adjusted for terrain effects (Jensen, 1994). Model and climate data intercomparisons are planned to take two forms. Zoom model outputs are being developed for the location of each of the existing weather stations. In addition a gridded data set based on historical observations is being prepared to allow direct comparison with the gridded model data. This two way comparison will allow us to test the assumptions in the model to local data

adjustment procedures.

During the development, CVHM and the Weber River System Model were tested using some simplified but fairly standard climate change scenarios (Sikka et al., 1994). The hydrographs for these conditions are shown in Figure 6. Base data for the tests were scaled historical climate data sets. The conditions include +/- 10% changes in precipitation coupled with 4° and 6° temperature increases. All of the test simulations show significant decreases in flow and in the timing of the peak. If actually experienced, these changes would have significant impact on the population and agriculture which has developed along the Wastach Mountains in Utah. Before these kinds of assessments work their way into planning documents, more robust studies using much more defensible data inputs need to be conducted. MCHS was designed to provide the foundation for these types of studies.

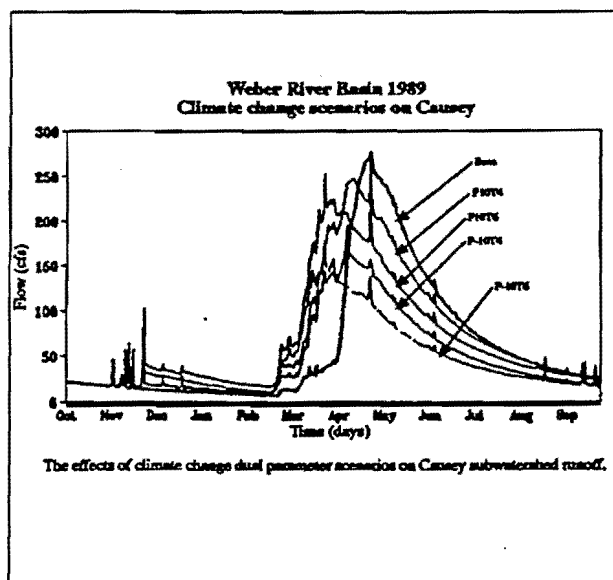


Figure 6. Hydrographs for CVHM and Weber River System Model output under simplified climate change scenarios.

#### 4. CONCLUSIONS

The Mountain Climate Hydrometeorological modeling System was designed to provide a testbed to study the scaling effects and models being proposed to scale GCM climate change data for regional and local studies. Developing the climate model and vegetation data sets for this testbed has been a major effort. Now that this effort is almost complete, detailed studies of both the models and the scaling assumptions for a wide range of subjects can be easily studied. The current round of watershed scale surface studies now underway at USU cover easily historically based 14 year period, December 1978 to April 1993. We hope to expand our efforts to include studies of vegetation response, erosion and

hydrologic responses using existing single point and hydrologic response unit based models. Model outputs are being compared with a combined point data set based on NOAA COOP and USDA/SCS SNOTEL stations as well as SSM/I and GOES image data. We anticipate that many additional cooperative studies will benefit from this extensive modeling effort.

#### 5. REFERENCES

- Dickenson, R.E., A. Henderson-Sellers, and P.J. Kennedy, 1993: Biosphere-Atmosphere Transfer Scheme (BATS) version 1E as coupled to the NCAR Community Climate Model. NCAR Tech. NCAR/TN-387+STR.
- Giorgi, F., M.R. Marinucci, G.T. Bates, and G.D. Canio, 1993: Monthly Weather Review. Vol. 121:2814-2832.
- Gunderson, R.W., 1983: An adaptive FCV clustering algorithm. International Journal on Man-machine Studies. Vol. 19:97-104.
- Leavesley, G.H., R.W. Lichty, B.M. Troutman, and L.G. Saindon, 1983: Precipitation Runoff Modeling System-users manual. USGS Water Resources Invest. Rep. 83-4239, 207.
- Meams, L.O., S.H. Schneider, S.L. Thompson, and L.R. McDaniel, 1990: Analysis of Climate Variability in General Circulation Models: Comparison with observations and Changes in Variability in 2xCO<sub>2</sub> Experiments. Journal of Geophysical Research, Vol. 95, November 20, 1990, No. D2, pages 20, 469-20, 490.
- Sikka, A.K., 1993: A Hydrologic Model for Studying the Influence of Climate Change on Evapotranspiration and Water Yield. Unpublished PhD dissertation. Utah State University, Logan, Utah. 250 p.
- NSREL Report No.6, 1991. Water Erosion Prediction Project - Hillslope Profile Model: Documentation Corrections and Additions. National Soil Erosion Research Laboratory, USDA-Agricultural Research Service, West Lafayette, Indiana.



# **APPENDIX 3A**

## **MEASUREMENTS AND MODELING OF SNOW ENERGY BALANCE AND SUBLIMATION FROM SNOW**

Measurements and Modeling of  
Snow Energy Balance and  
Sublimation from Snow

David G. Tarboton

*Presented at International Snow Science  
Workshop, Snowbird, Utah, November 1,  
1994*

Working Paper WP-94-HWR-DGT/002

November 1994

# MEASUREMENTS AND MODELING OF SNOW ENERGY BALANCE AND SUBLIMATION FROM SNOW

David G. Tarboton  
Utah Water Research Laboratory,  
Utah State University,  
Logan, Utah 84322-8200

Telephone: 801-797-3172; Fax: 801-797-3663; email: dtarb@cc.usu.edu

## Abstract

Snow melt runoff is an important factor in runoff generation for most Utah rivers and a large contributor to Utah's water supply and periodically flooding. The melting of snow is driven by fluxes of energy into the snow during warm periods. These consist of radiant energy from the sun and atmosphere, sensible and latent heat transfers due to turbulent energy exchanges at the snow surface and a relatively small ground flux from below. The turbulent energy exchanges are also responsible for sublimation from the snow surface, particularly in arid environments, and result in a loss of snow water equivalent available for melt. The cooling of the snowpack resulting from sublimation also delays the formation of melt runoff. This paper describes measurements and mathematical modeling done to quantify the sublimation from snow. Measurements were made at the Utah State University drainage and evapotranspiration research farm. I attempted to measure sublimation directly using weighing lysimeters. Energy balance components were measured, by measuring incoming and reflected radiation, wind, temperature and humidity gradients.

An energy balance snowmelt model was tested against these measurements. The model uses a lumped representation of the snowpack with two state variables, namely, water equivalent and energy content relative to a reference state of water in the solid phase at 0°C. This energy content is used to determine snowpack average temperature or liquid fraction. The model is driven by inputs of air temperature, precipitation, wind speed, humidity and solar radiation. The model uses physically based calculations of radiative, sensible, latent and advective heat exchanges. An equilibrium parameterization of snow surface temperature accounts for differences between snow surface temperature and average snowpack temperature without having to introduce additional state variables. This is achieved by incorporating the snow surface thermal conductance, which with respect to heat flux is equivalent to stomatal and aerodynamic conductances used to calculate evapotranspiration from vegetation. Melt outflow is a function of the liquid fraction, using Darcy's law. This allows the model to account for continued melt outflow even when the energy balance is negative.

The purpose of the measurements presented here was to test the sublimation and turbulent exchange parameterizations in the model. However the weighing lysimeters used to measure sublimation suffered from temperature sensitive oscillations that mask short term sublimation measurements. I have therefore used the measured data to test the models capability to represent the overall seasonal accumulation and ablation of snow.

## Description of Experiment

The experiment reported here was conducted at the USU drainage and evapotranspiration research farm in Cache Valley. Instrumentation in place is designed for the study of evapotranspiration from agricultural lands, but for this study was utilized for the study of winter snow cover. The instrumentation consisted of two  $1\text{ m}^2$  weighing lysimeters and meteorological and energy balance equipment. The weighing lysimeters are  $1\text{ x }1\text{ x }1\text{ m}$  metal boxes embedded flush with the surface and filled with soil, vegetated with grass similar to the surrounding agricultural field. Load cells (underneath in the case of one lysimeter and at the corners for the other) record the weight of soil, grass, soil moisture and snow over the  $1\text{ m}^2$  area. Meltwater infiltrates into the lysimeter so does not result in a weight change. Changes in weight are due only to addition or removal of mass from the surface, which in the case of snow can be due to precipitation, condensation, sublimation and wind drifting.

Meteorological and energy balance instrumentation used is listed in table 1.

Table 1: Meteorological Instrumentation

---

2 Net Radiometers (Fritchen type Q6 and Q4) installed 1m above the snow surface.
2 Lycor pyranometers that record solar radiation. One was pointed down to estimate albedo.
1 Eppley pyranometer to record incident solar radiation.
2 Everest Interscience model 4000 Infrared surface temperature sensors.
4 Anemometers at heights 0.6, 0.9, 1.4 and 2.4 m above the ground surface.
4 Viasala temperature and relative humidity sensors at height 0.58, 0.90, 1.44, 2.57 m above the ground surface.
2 REBS Ground heat flux plates
Thermocouple ladder. This consisted of 14 copper/constantine thermocouples at the following levels: -0.075, -0.025, 0, 0.05, 0.125, 0.2, 0.275, 0.35, 0.425, 0.5, 0.575, 0.65, 0.725, 0.8 m, from the ground surface. The first two thermocouples were buried and the third placed on the ground. The remainder were suspended on fishing line strung between two upright posts.
Heated (unshielded) tipping bucket rain/snow gage.
Wind direction sensor

---

Two campbell scientific 21X dataloggers powered by a deep cycle 12 volt battery charged by a solar panel were used to take measurement readings every minute and record 30 minute averages for output.

The dataloggers were downloaded during biweekly visits at which time the sensors were also inspected and cleared of snow and grime buildup. During these visits, snow depth and water

equivalent was measured at eight locations using an Adirondack snow tube sampler. To guard against the danger of bridging in the snow between snow over the lysimeters and surrounding snow which would distort the weights and inferred sublimation a plastic batten and saw was used to saw the snow between the lysimeter and surrounding. This was done from a ladder supported between two trestles over the lysimeter so as not to disturb the snow on or near the lysimeter. This procedure was only partly successful as we did notice some abrupt changes in lysimeter weight that coincided with the sawing. We also found that the lysimeter weight measurements had a diurnal temperature sensitivity that precluded using them for short term sublimation measurements. They still provide an overall measurement of snow accumulation.

The USU drainage and irrigation experimental farm is located in Cache Valley near Logan, Utah, USA (41.6° N, 111.6° W, 1350m elevation). The weather station and instrumentation are in a small fenced enclosure at the center of a large open field. There are no obstructions to wind in any direction for at least 500m. Cache valley is a flat bottomed enclosed valley surrounded by mountains that reach elevations of 3000m. During winter periods of settled weather strong temperature inversions accompanied by very cold (-20 °C) nighttime temperatures and night and morning fog develop. Unsettled stormy periods serve to break the inversion. During the period of this experiment the ground was snow covered from November 20, 1992 to March 22, 1993. Air temperatures ranged from -23 °C to 16 °C and there was 190 mm of precipitation (mostly snow, but some rain). The snow accumulated to a maximum depth of 0.5 m with maximum water equivalent of 0.14 m. Table 2 gives a chronology of the events and measurements. The instrumentation was only fully functional for the latter half of the winter, which will be the focus of the analysis.

Table 2. Chronology.

<u>From</u>	<u>To</u>	<u>Day</u>	<u>Event</u>
11/20/92		-41	First snowfall 6 mm.
11/20/92	1/13/93	-41 to 13	Several snowstorms resulting in an accumulation of 86 mm of water equivalent and depth of 400 mm.
1/13/93		13	<i>Supplementary equipment (thermocouple ladder and air temperature and humidity profile) is finally functional.</i>
1/17/93		17 to 19	<i>Datalogger battery failure, some data lost.</i>
1/18/93	1/25/93	18 to 25	Period of unsettled weather (12 mm precipitation).
1/25/93		25	<i>Heated precipitation gage and downward pointing pyranometer installed and functional.</i>
1/26/93	2/8/93	26 to 39	Inversion and fog.
2/8/93	2/25/93	39 to 56	Period of unsettled weather (45 mm precipitation).
2/26/93	3/9/93	57 to 68	Inversion and fog.
3/10/93	3/11/93	69 to 70	Rain and snow (20 mm precipitation). Highest water equivalent accumulation of 139 mm was recorded just prior to this event which initiated melt.
3/11/93	3/14/93	70 to 73	Clear warm weather. Melt continues.
3/15/93	3/16/93	74 to 75	Light rain (2 mm).
3/17/93	3/18/93	76 to 77	Heavy rain (18.5 mm) that caused considerable snowmelt.
3/19/93	3/22/93	78 to 81	Remaining snow melted rapidly.

## Energy Balance Snowmelt Model

The energy balance model used (Chowdhury et al., 1992; Bowles et al., 1992; Bowles et al., 1994; Tarboton et al., 1995) was developed for purposes of erosion prediction and water balance modeling. The snowpack is characterized by two primary state variables, water equivalent,  $W$  [m], and energy content,  $U$ , [ $\text{kJ}/\text{m}^2$ ]. The state variable, energy content  $U$ , is defined relative to a reference state of water at  $0^\circ\text{C}$  in the ice (solid) phase.  $U$  greater than zero means the snowpack (if any) is isothermal with some liquid content and  $U$  less than zero can be used to calculate the snowpack average temperature,  $T$ , [ $^\circ\text{C}$ ]. Energy content is defined as the energy content of the snowpack plus a top layer of soil with depth  $D_e$  [m]. This provides a simple buffering against numerical instabilities when the snowpack is shallow, as well as simple approximations of frozen ground and melting of snow falling on warm ground. We discuss below the choice of  $D_e$  and the role it plays in the model.

The model is designed to be driven by inputs of air temperature,  $T_a$  [ $^\circ\text{C}$ ]; wind speed,  $V$  [m/s]; relative humidity,  $RH$ ; precipitation,  $P$  [m/hr]; incoming solar  $Q_{si}$  and longwave  $Q_{li}$  radiation [ $\text{kJ}/\text{m}^2/\text{hr}$ ]; and ground heat flux  $Q_g$  [ $\text{kJ}/\text{m}^2/\text{hr}$ ] (taken as 0 when not known) at each time step. When incoming solar radiation is not available it is estimated as extra terrestrial radiation (from sun angle) times an atmospheric transmission factor,  $Tr$ , estimated from the daily temperature range using the procedure given by Bristow and Campbell (1984). When incoming longwave radiation is not available it is estimated based on air temperature, the Stefan-Boltzman equation and a parameterization of air emissivity due to Satterlund (1979) adjusted for cloudiness using  $Tr$ .

Given the state variables  $U$  and  $W$ , their evolution in time is determined by solving energy and mass balance equations.

$$\frac{dU}{dt} = Q_{sn} + Q_{li} + Q_p + Q_g - Q_{le} + Q_h + Q_e - Q_m \quad (1)$$

$$\frac{dW}{dt} = P_r + P_s - M_r - E \quad (2)$$

In the energy balance equation terms are (all in  $\text{kJ}/\text{m}^2/\text{hr}$ ):  $Q_{sn}$ , net shortwave radiation;  $Q_{li}$ , incoming longwave radiation;  $Q_p$ , advected heat from precipitation;  $Q_g$ , ground heat flux;  $Q_{le}$ , outgoing longwave radiation;  $Q_h$ , sensible heat flux;  $Q_e$ , latent heat flux due to sublimation/condensation;  $Q_m$ , advected heat removed by meltwater. In the mass balance equation (all in m/hr of water equivalent) terms are:  $P_r$ , rainfall rate;  $P_s$  snowfall rate;  $M_r$ , meltwater outflow from the snowpack;  $E$ , sublimation from the snowpack. Many of these fluxes depend functionally on the state and input driving variables. We elaborate on the parameterization of these functional dependencies below. Equations (1) and (2) form a coupled set of first order, nonlinear

ordinary differential equations. They can be summarized in vector notation as:

$$\frac{d\mathbf{X}}{dt} = \mathbf{F}(\mathbf{X}, \text{driving variables}) \quad (3)$$

where  $\mathbf{X} = (U, W)$  is a state vector describing the snowpack. With  $\mathbf{X}$  specified initially, this is an initial value problem. A large variety of numerical techniques are available for solution of initial value problems of this form. Here we have adopted a Euler predictor-corrector approach (Gerald, 1978).

$$\mathbf{X}' = \mathbf{X}_1 + \Delta t \mathbf{F}(\mathbf{X}_1, \text{driving variables}) \quad (4)$$

$$\mathbf{X}_{i+1} = \mathbf{X}_i + \Delta t \frac{\mathbf{F}(\mathbf{X}_i, \text{driving variables}) + \mathbf{F}(\mathbf{X}', \text{driving variables})}{2} \quad (5)$$

where  $\Delta t$  is the time step,  $\mathbf{X}_i$  refers to the state at time  $t_i$  and  $\mathbf{X}_{i+1}$  refers to the state at time  $t_{i+1} = t_i + \Delta t$ . This is a second order finite difference approximation, with global error proportional to  $\Delta t^2$  (Gerald, 1978, p257).

### Parameterization

**Depth averaged temperature - T:** The snow and interacting soil layer average temperatures are obtained from the energy content and water equivalent, relative to 0°C ice phase.

$$\text{If } U < 0 \quad T = U / (\rho_w W C_s + \rho_g D_e C_g) \quad \text{All solid phase} \quad (6)$$

$$\text{If } 0 < U < \rho_w W h_f \quad T = 0^\circ\text{C}. \quad \text{Solid and liquid mixture} \quad (7)$$

$$\text{If } U > \rho_w W h_f \quad T = \frac{U - \rho_w W h_f}{\rho_g D_e C_g + \rho_w W C_w} \quad \text{All liquid} \quad (8)$$

In the above the heat required to melt all the snow water equivalent is  $\rho_w W h_f$  [kJ] where  $h_f$  is the heat of fusion [333.5 kJ kg<sup>-1</sup>] and  $U$  in relation to this determines the solid-liquid phase mixtures. The heat capacity of the snow is  $\rho_w W C_s$  [kJ/°C] where  $\rho_w$  is the density of water [1000 kg m<sup>-3</sup>] and  $C_s$  the specific heat of ice [2.09 kJ kg<sup>-1</sup> °C<sup>-1</sup>]. The heat capacity of the soil layer is  $\rho_g D_e C_g$  [kJ/°C] where  $\rho_g$  is the soil density [ $\approx 1700$  kg m<sup>-3</sup>] and  $C_g$  the specific heat of soil [ $\approx 2.1$  kJ kg<sup>-1</sup> °C<sup>-1</sup>]. These together determine  $T$  when  $U < 0$ . In practice, unless we allow ponded water (which we don't)  $W$  will always be 0 in (8). The heat capacity of liquid water,  $\rho_w W C_w$ , where  $C_w$  is the specific heat of water [4.18 kJ kg<sup>-1</sup> °C<sup>-1</sup>], is however retained in (8) for numerical

consistency during time steps when the snowpack completely melts.

Heat flow in snow and soil is governed by Laplace's equation. The depth of penetration of changes in surface temperature can be evaluated from the expression (Rosenberg, 1974):

$$\frac{R_z}{R_s} = \exp\left(-z \left(\frac{\pi}{\alpha P}\right)^{\frac{1}{2}}\right) \quad (9)$$

where  $R_s$  is the range of temperature oscillation at the surface,  $R_z$  the range of temperature oscillation at depth  $z$ ,  $P$  the period of oscillation, and  $\alpha$  the thermal conductivity. For soil  $\alpha$  is typically in the range 0.004 to 0.006 cm<sup>2</sup>/s. Figure 1 shows  $R_z/R_s$  versus  $z$  for  $\alpha = 0.005$  cm<sup>2</sup>/s for various periods. This shows that for oscillations less than one week the effect at 40 cm is damped to less than 30% and even for monthly oscillations is still damped 50% at 40 cm depth. This suggests using  $D_e = 40$  cm in our model. Rosenberg (1974) also suggests this as an effective depth. The state variable  $U$  represents energy content above this level. The ground heat flux represents heat transport at this depth and is therefore a long term average. Diurnal oscillating ground heat fluxes above this depth are absorbed into  $U$ , the energy stored in the snow and soil above depth  $D_e$ .

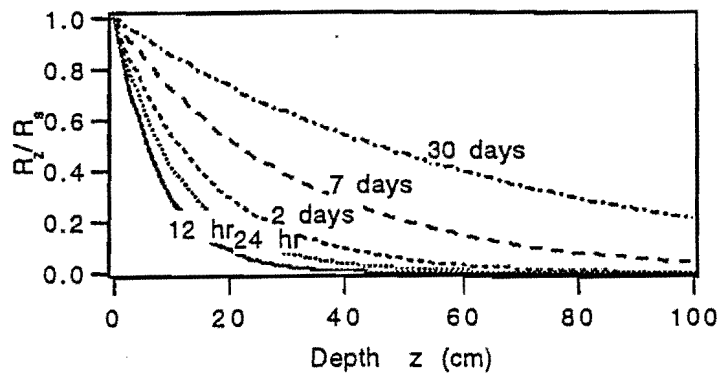


Figure 1. Depth of penetration of temperature fluctuations into soil with  $\alpha = 0.005$  cm<sup>2</sup>/s.

**Net Shortwave Radiation,  $Q_{sn}$ :** This is calculated as

$$Q_{sn} = Q_{si} (1-A) \quad (10)$$

where Albedo,  $A$ , is calculated based on the age of the snow surface using a parameterization described by Dickinson et al. (1993). For shallow snowpacks (depth less than 0.1 m) the albedo is interpolated between the bare ground value (0.25) and snow value.



**Outgoing Longwave Radiation,  $Q_{le}$ :** Snow is essentially a black body, with emissivity  $\epsilon_s \approx 0.99$ . Outgoing radiation is

$$Q_{le} = \epsilon_s \sigma \left( T_s^{abs} \right)^4 \quad (11)$$

where  $\sigma$  is the Stefan Boltzmann constant [ $2.07 \times 10^{-7} \text{ kJ m}^{-2} \text{ hr}^{-1} \text{ K}^{-4}$ ] and the superscript "abs" in  $T_s^{abs}$  indicates that this is absolute temperature [K].

**Snow fall accumulation and heat with precipitation:** Measured precipitation rate,  $P$ , is partitioned into rain,  $P_r$ , and snow,  $P_s$ , (both in terms of water equivalent depth) using the following rule based on air temperature,  $T_a$ , (U.S. Army Corps of Engineers, 1956)

$$\begin{aligned} P_r &= P & T_a &\geq T_r = 3 \text{ }^\circ\text{C} \\ P_r &= P(T_a - T_b)/(T_r - T_b) & T_b &< T_a < T_r \\ P_r &= 0 & T_a &\leq T_b = -1 \text{ }^\circ\text{C} \\ P_s &= P - P_r \end{aligned} \quad (12)$$

where  $T_r$  is a threshold air temperature above which all precipitation is rain and  $T_b$  a threshold air temperature below which all precipitation is snow.

The temperature of rain is taken as the greater of the air temperature and freezing point and the temperature of snow the lesser of air temperature and freezing point. The advected heat is the energy required to convert this precipitation to the reference state ( $0^\circ\text{C}$  ice phase).

$$Q_p = P_s C_s \rho_w \min(T_a, 0 \text{ }^\circ\text{C}) + P_r \left[ h_f \rho_w + C_w \rho_w \max(T_a, 0 \text{ }^\circ\text{C}) \right] \quad (13)$$

**Turbulent fluxes,  $Q_h$ ,  $Q_e$ ,  $E$ :** Sensible and latent heat fluxes between the snow surface and air above are modeled using the concept of flux proportional to temperature and vapor pressure gradients with constants of proportionality, the so called turbulent transfer coefficients or diffusivity a function of windspeed and surface roughness. Considering a unit volume of air, the heat content is  $\rho_a C_p T_a$  and the vapor content  $\rho_a q$ , where  $\rho_a$  is air density (determined from atmospheric pressure and temperature),  $C_p$  air specific heat capacity [ $1.005 \text{ kJ kg}^{-1} \text{ }^\circ\text{C}^{-1}$ ], and  $q$  specific humidity [kg water vapor per kg air]. Heat transport towards the surface,  $Q_h$  [ $\text{kJ/m}^2/\text{hr}$ ] is given by:

$$Q_h = K_h \rho_a C_p (T_a - T_s) \quad (14)$$

where  $K_h$  is heat conductance [ $\text{m/hr}$ ] and  $T_s$  is the snow surface temperature. Vapor transport

away from the surface (sublimation),  $M_e$  [kg/hr] is:

$$M_e = K_e \rho_a (q_s - q) \quad (15)$$

where  $q_s$  is the surface specific humidity and  $K_e$  the vapor conductance [m/hr].

By comparison with the usual expressions for turbulent transfer in a logarithmic boundary layer profile (Male and Gray, 1981; Anderson, 1976; Brutsaert, 1982; Calder, 1990) for neutral condition, one obtains the following expression:

$$K_h = K_e = \frac{k^2 V}{[\ln(z/z_0)]^2} = K \quad (16)$$

where  $V$  is wind speed [m/hr] at height  $z$  [m];  $z_0$  is roughness height at which the logarithmic boundary layer profile predicts zero velocity [m]; and  $k$  is von Karman's constant [0.4]. The subscript  $n$  denotes that these are conductances in neutral conditions. Recognizing that the latent heat flux towards the snow is:

$$Q_e = -h_v M_e \quad (17)$$

and using the relationship between specific humidity and vapor pressure and the ideal gas law one obtains:

$$Q_e = K_e \frac{h_v 0.622}{R_d T_a^{\text{abs}}} (e_a - e_s(T_s)) \quad (18)$$

where  $e_s$  is the vapor pressure at the snow surface snow, assumed saturated at  $T_s$ , and calculated using a polynomial approximation (Lowe, 1977);  $e_a$  is air vapor pressure,  $R_d$  is the dry gas constant [287 J kg<sup>-1</sup> K<sup>-1</sup>] and  $h_v$  the latent heat of sublimation [2834 kJ/kg]. The water equivalent depth of sublimation is:

$$E = -\frac{Q_e}{\rho_w h_v} \quad (19)$$

When there is a temperature gradient near the surface, buoyancy effects may enhance or dampen the turbulent transfers. This can be quantified in terms of the Richardson number or Monin-Obukhov length. We had hoped that the lysimeter measurements made here would have provided

data to allow us to determine the effect of stability on snow sublimation. However since that did not work out the results presented here use neutral buoyancy.

**Snow Surface Temperature,  $T_s$ :** Since snow is a relatively good insulator,  $T_s$  is in general different from  $T$ . This is accounted for using an equilibrium approach that balances energy fluxes at the snow surface. Heat conduction into the snow is calculated using the temperature gradient and thermal diffusivity of snow, approximated by:

$$Q = \kappa \rho_s C_s (T_s - T)/Z_e = K_s \rho_s C_s (T_s - T) \quad (20)$$

where  $\kappa$  is snow thermal diffusivity [ $m^2 \text{ hr}^{-1}$ ] and  $Z_e$  [m] an effective depth over which this thermal gradient acts. The ratio  $\kappa/Z_e$  is denoted  $K_s$  and termed snow surface conductance analogous to the heat and vapor conductances. A value of  $K_s$  is obtained by assuming a depth,  $Z_e$  equal to the depth of penetration of a diurnal temperature fluctuation calculated from equation (9) (Rosenberg, 1974).  $Z_e$  is chosen so that  $R_z/R_s$  is small. In fact  $K_s$  is used as a tuning parameter, with this calculation used to define a reasonable range. Then assuming equilibrium at the surface, the surface energy balance gives.

$$Q = Q_{sn} + Q_{li} + Q_h(T_s) + Q_e(T_s) + Q_p - Q_{le}(T_s) \quad (21)$$

where the dependence of  $Q_h$ ,  $Q_e$ , and  $Q_{le}$  on  $T_s$  is through equations (14), (18) and (11).

Analogous to the derivation of the Penman equation for evaporation the functions of  $T_s$  in this energy balance equation are linearized about a reference temperature,  $T^*$  and the equation is solved for  $T_s$ :

$$T_s^{abs} = \frac{Q_{sn} + Q_{li} + Q_p + K T_a^{abs} \rho_a C_p - 0.622 K h_v \rho_a (e_s(T^*) - e_a - T^* \Delta) / P_a + 3 \epsilon_s \sigma T^{*abs4} + \rho_s C_s T^{abs} K_s}{\rho_s C_s K_s + K \rho_a C_p + 0.622 \Delta K h_v \rho_a / P_a + 4 \epsilon_s \sigma T^{*abs3}} \quad (22)$$

where  $\Delta = de_s/dT$ . This equation is used in an iterative procedure with an initial estimate  $T^* = T_a$ , in each iteration replacing  $T^*$  by the latest  $T_s$ . The procedure converges to a final  $T_s$  which if less than freezing is used to calculate surface energy fluxes. If the final  $T_s$  is greater than freezing it means that the energy input to the snow surface cannot be balanced by thermal conduction into the snow. Surface melt will occur and the infiltration of meltwater will account for the energy difference and  $T_s$  is then set to  $0^\circ\text{C}$ .

**Meltwater Outflux,  $M_r$  and  $Q_m$ :** The energy content state variable  $U$  determines the liquid content of the snowpack. This, together with Darcy's law for flow through porous media, is used to determine the outflow rate.

$$M_r = K_{sat} S^*{}^3 \quad (23)$$

where  $K_{sat}$  is the snow saturated hydraulic conductivity [ $=160 \text{ m hr}^{-1}$ ] and  $S^*$  is the relative saturation in excess of water retained by capillary forces. This expression is based on Male and Gray (1981 p400 eqn 9.45).  $S^*$  is given by:

$$S^* = \frac{\text{liquid water volume} - \text{capillary retention}}{\text{pore volume} - \text{capillary retention}} = \left( \frac{L_f}{1 - L_f} - L_c \right) / \left( \frac{\rho_w}{\rho_s} - \frac{\rho_w}{\rho_i} - L_c \right) \quad (24)$$

where  $L_f = U / (\rho_w h_f W)$  denotes the mass fraction of total snowpack (liquid and ice) that is liquid,  $L_c$  [0.05] the capillary retention as a fraction of the solid matrix water equivalent, and  $\rho_i$  the density of ice [ $917 \text{ kg m}^{-3}$ ].

This melt outflow is assumed to be at  $0^\circ\text{C}$  so the heat advected with it, relative to the solid reference state is:

$$Q_m = \rho_w h_f M_r \quad (25)$$

### Model parameters

Apart from known physical constants and readily estimable quantities the model has adjustable parameters listed in Table 3. The values used were taken from previous work with the model calibrated against data collected at the Central Sierra Snow Laboratory. These results therefore present an independent check of the model in a different setting.

Table 3. Adjustable parameter values

Parameter	Notation	Value
Surface aerodynamic roughness	$z_0$	0.002 m
Surface conductance	$K_s$	0.015 m/hr
Snow density	$\rho_s$	$450 \text{ kg m}^{-3}$
Saturated hydraulic conductivity	$K_{sat}$	160 m/hr
Capillary retention fraction	$L_c$	0.05

## Results and Discussion

Figure 2 gives the measured lysimeter weights, measured snow water equivalent and accumulated precipitation. The measured snow water equivalent values shown are the average from the 8 snow core measurements made each visit. The individual water equivalent measurements usually varied within a range of 10 to 20% from this average. This shows general agreement between weight accumulation on the lysimeters, snow accumulation and precipitation. Figure 3 compares model and measured snow water equivalent for the model run from day 26 to the end of melt. Two model runs are shown, one with the model driven by measured net radiation and the other with the model driven by incoming solar radiation. The first run bypasses the albedo and outgoing longwave radiation calculations so serves only to test the models sensible and latent heat flux components. The second run is a more realistic check on overall model performance. For both runs the model was initialized with the measured day 26 water equivalent of 0.104 m and energy content based on the average temperature of thermocouples in the snow and soil. This energy content was,  $-1136 \text{ kJ/m}^2$ . These results show that the model does reasonably well at representing snow accumulation and melt. The second model run, with solar radiation as the primary energy input, was used for the remainder of the comparisons in this paper.

Figure 4. shows modeled and measured snow (and soil) energy content. The measured energy content was estimated from the measured water equivalent (linearly interpolated between measurements) and snow and soil temperatures averaged from the thermocouple ladder measurements. There is obviously a large discrepancy between modeled and measured energy content early on, and given this it is surprising how well the model does at representing other aspects of the snow accumulation and melt processes. The lowest energy content on day 39 would predict an average snow and soil temperature of  $-14 \text{ }^\circ\text{C}$ . This is well below the observed snow temperatures shown on figure 5. These discrepancies indicate that the model loses too much energy during cold periods, suggesting that the snow surface conductance may be too large. It also indicates that temperature fluctuations do not penetrate to the full interacting soil layer depth,  $D_e$  [0.4 m] suggesting that perhaps  $D_e$  should be reduced. After day 70 (March 20) the model energy content is above zero due to the liquid water content of the snow. This is the melt period. The measured energy, estimated from thermocouple measurements of snow and soil temperatures, does not account for liquid water in the snow.

Figures 6a-f present detailed results for the period from January 26 to February 7 (day 26 to 38) during which there was a strong temperature inversion and no measurable precipitation, although there was condensation and accumulation of frozen fog. During this period the snow depth was 0.4 m. The sensor heights are given with respect to the ground so the lowest vapor pressure and temperature sensors were only 0.2 m above the snow surface. The lysimeters (only lysimeter 2 is shown in fig 6a, but lysimeter 1 was similar) recorded a diurnal oscillation in weight that is I believe an effect of the cold temperatures on the electronics or load cell system. The oscillations which correlate well with air temperature amount to 2 mm of water equivalent. Based on net radiation measurements the net radiation could only supply energy to sublimate a maximum

of 0.6 mm/day (if all energy goes to sublimation) in this period. The oscillations therefore mask any sublimation signal and preclude the use of these lysimeter measurements for the study of short term sublimation. Figure 6b shows the model water equivalent on an expanded scale where you can see that it does go through a very small diurnal oscillation (up to 0.1 mm/day) with nighttime condensation and daytime sublimation. This oscillation is out of phase with the vapor pressure measurements which increase during the day then drop at night. This suggests a recycling process where the snow surface layer is sublimated during the day then redeposited during nighttime cooling. There is a net accumulation from day 32 to day 33 when the vapor pressures (figure 6d) are high. Then on day 34 there is a period of relatively strong wind (figure 6c) and low vapor pressure (figure 6d) that results in a relatively large modeled sublimation and drop in water equivalent (figure 6b). Gradients in vapor pressure (the difference between the lines on figure 6d) coincide with modeled condensation and sublimation periods (figure 6b). Figure 6f compares model and measured infrared snow surface temperatures. This indicates that the equilibrium procedure for calculation of snow surface temperature works reasonably well.

Detailed results for the melt period (March 19, day 69 to March 23, day 82) are shown in figures 7a-h. The onset of melt was triggered by the 20 mm of precipitation, rain and snow mix on day 69 and 70. Following the precipitation strong winds and low humidity (vapor pressure, figure 7g) induces sublimation in the model over days 71 and 72 (figure 7h). There is some suggestion of a downward trend (implying sublimation) in the lysimeter trace on figure 7a. With this sublimation and cooler air temperatures there is minimal melt modeled on days 71 and 72. Freezing of the snow surface is well modeled as indicated by the model and measured snow surface temperatures (figure 7f). Warmer weather and higher humidity from day 73 on are characterized by positive sensible heat (higher temperatures at the upper sensor, fig 7e) and condensation (higher vapor pressure at the higher sensor, fig 7g) which both add energy to the snowpack, which consequently melts rapidly. The horizontal dashed line on figure 7g is 6.1 mb, the saturation vapor pressure of water over ice at freezing point. Vapor pressures higher than this imply a downward vapor pressure gradient which will result in condensation. Rain on day 76 makes melting even more rapid. Figure 7a indicates that over the whole season, according to the model, net sublimation was only a small fraction (the difference between the dashed lines) of the snow mass. This was due to the persistent inversions and high humidity associated with valley fog.

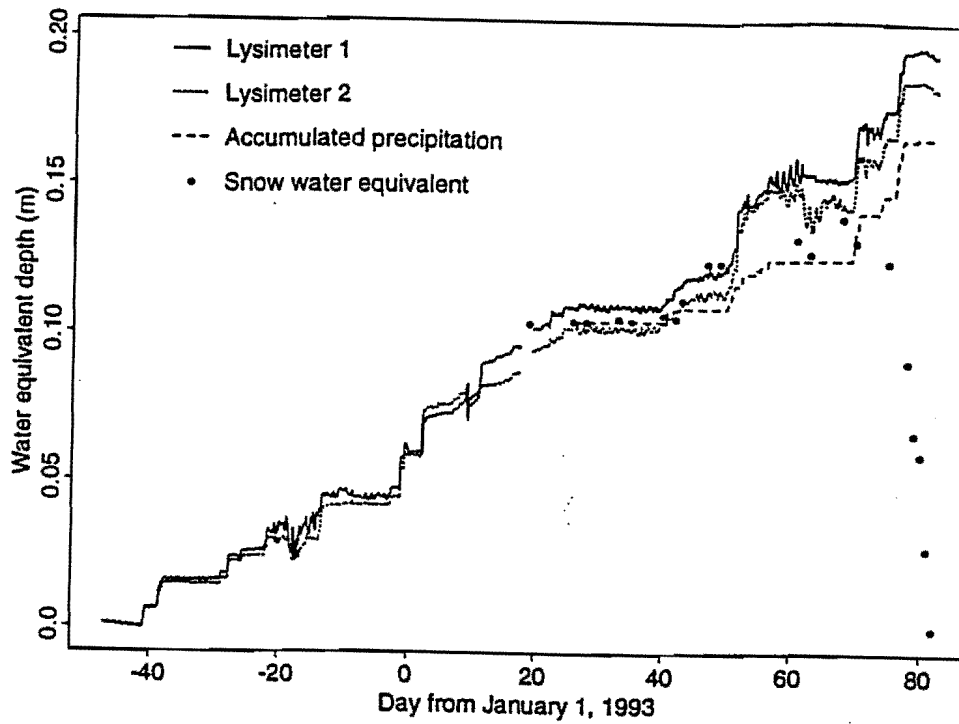


Figure 2. Overall snow accumulation and ablation measurements.

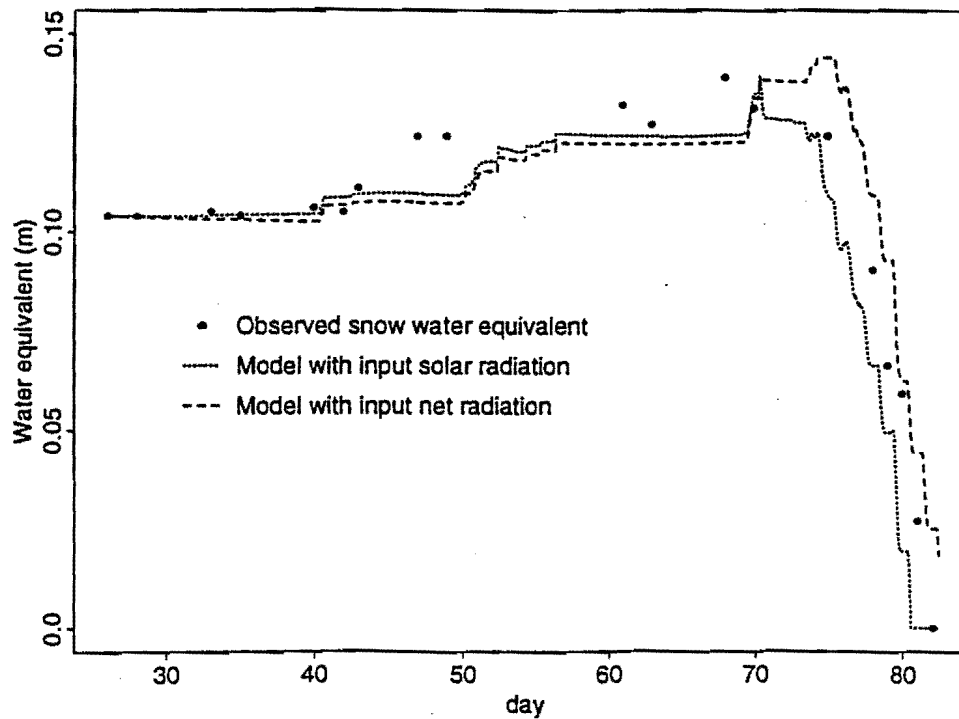


Figure 3. Comparison of observed and modeled snow water equivalent.

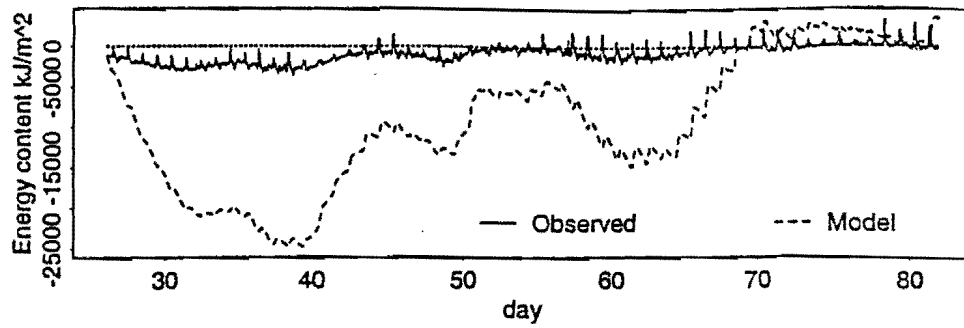


Figure 4. Comparison of measured and modeled energy content of the snow and top 0.4 m of soil.

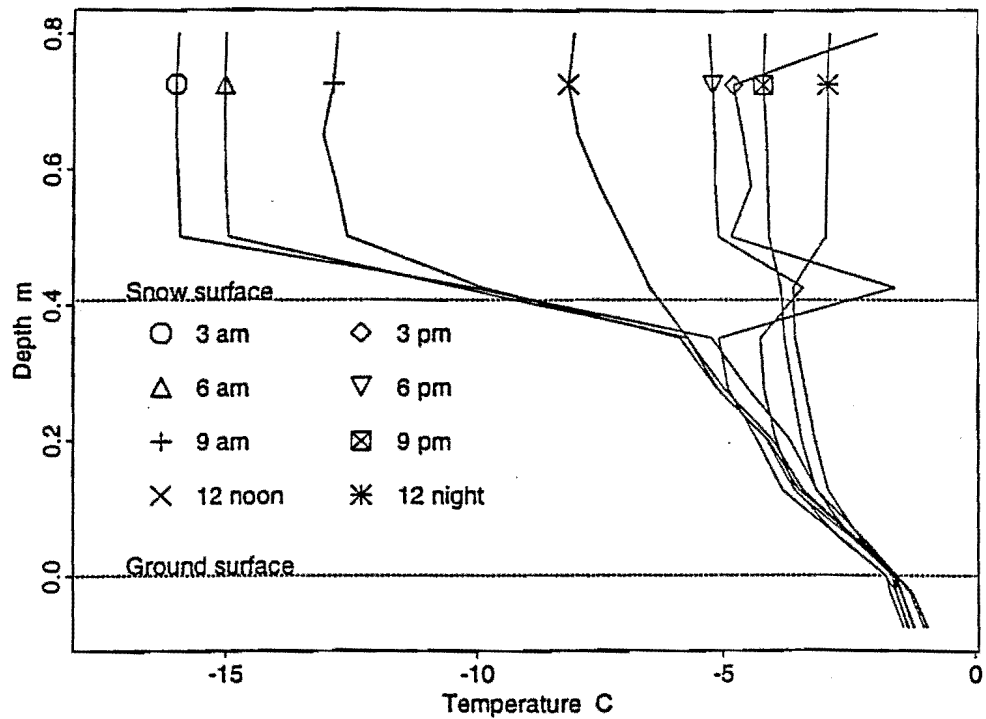
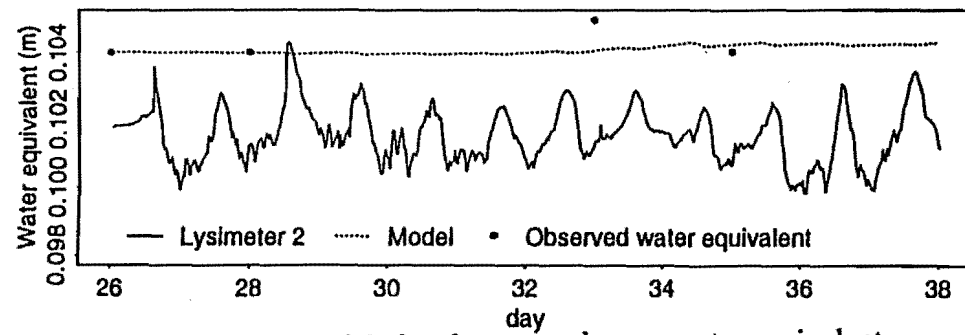
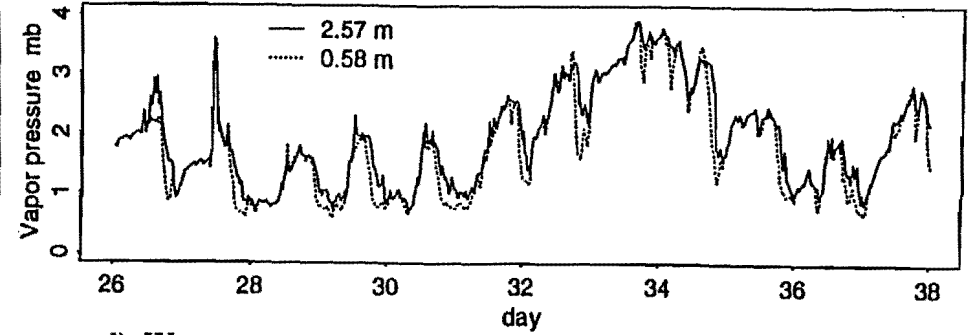


Figure 5. Measured soil and snow temperatures on February 8, 1993 (day 39).

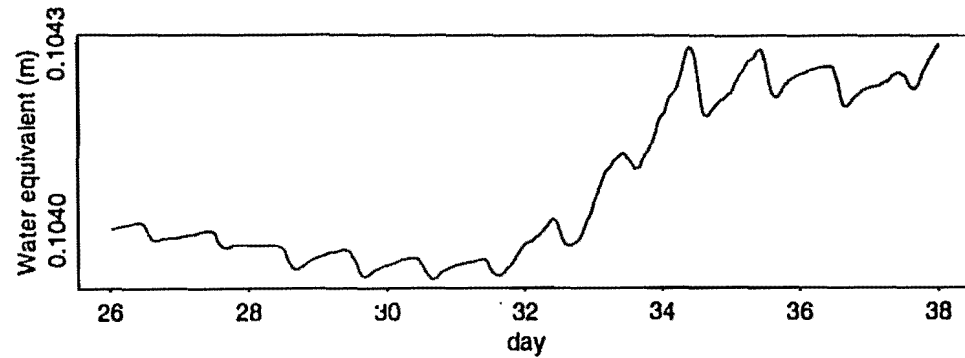




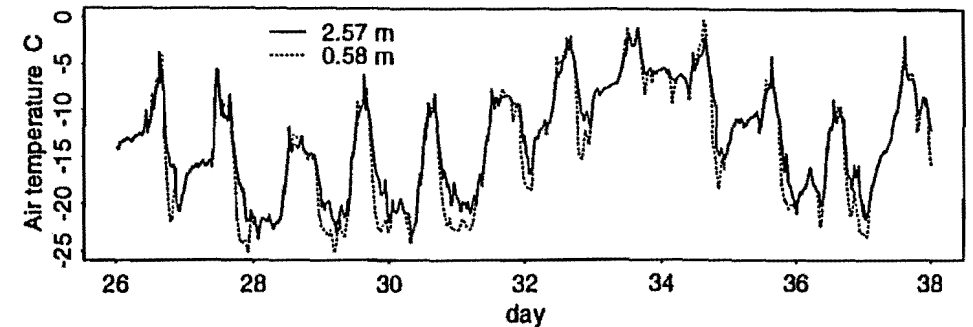
a) Lysimeter, modeled and measured snow water equivalent.



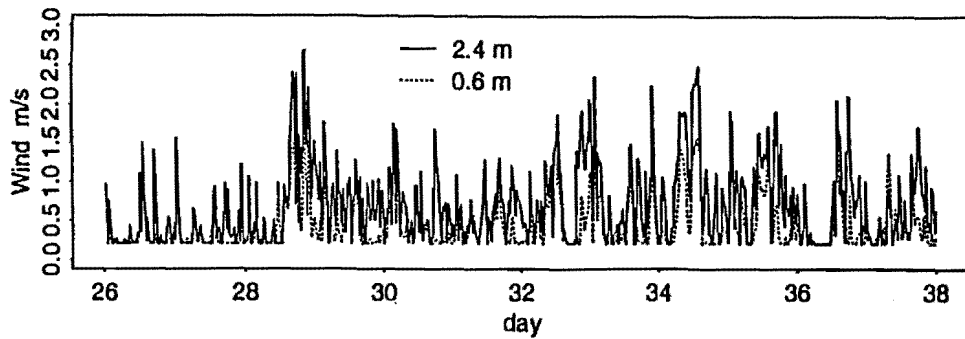
d) Water vapor pressure at upper and lower levels.



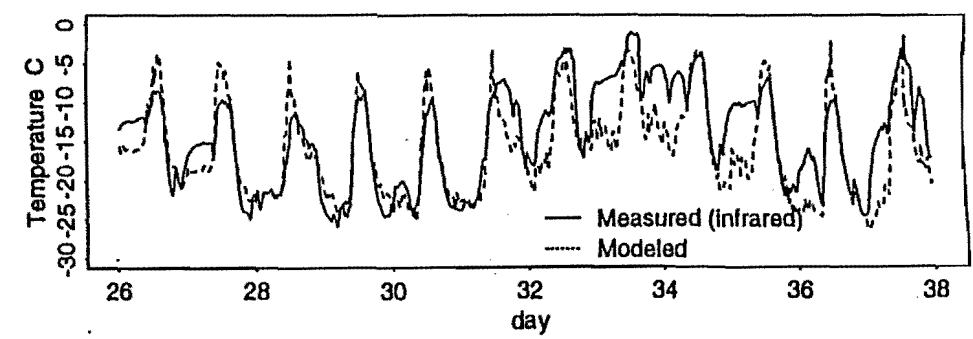
b) Expanded scale model snow water equivalent.



e) Air temperature at upper and lower levels.

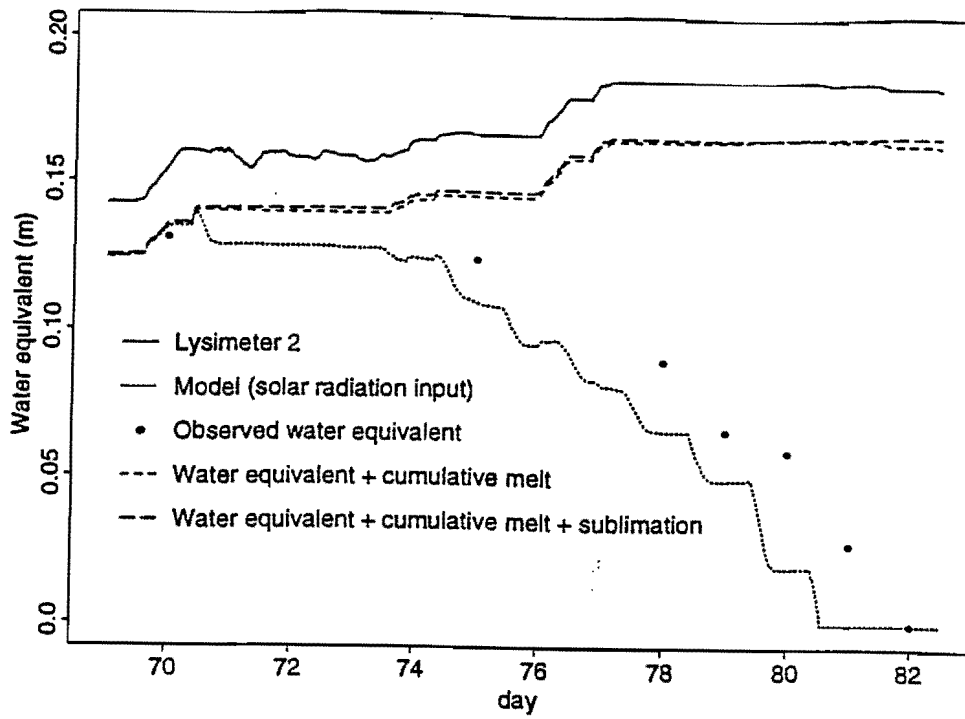


c) Wind velocity at upper and lower levels.

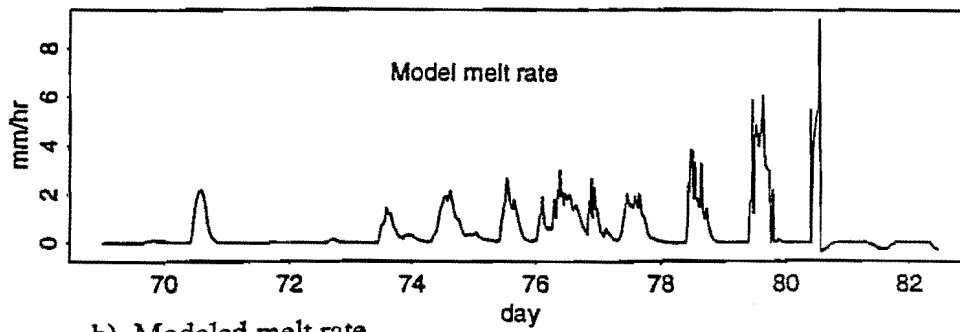


f) Snow surface temperature.

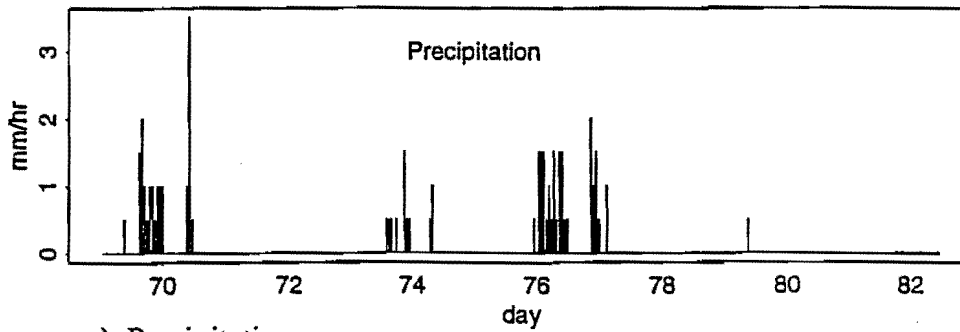
Figure 6. Detailed results for January 16, 1993 to February 7, 1993 (Days 26 to 38).



a) Lysimeter, modeled and measured snow water equivalent, accumulated melt and sublimation.

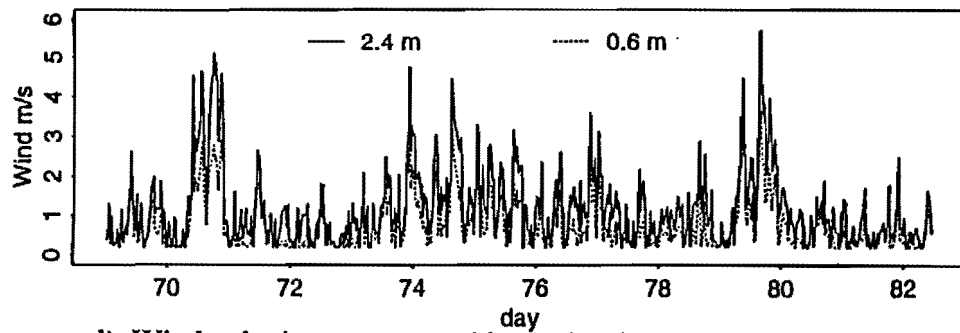


b) Modeled melt rate.

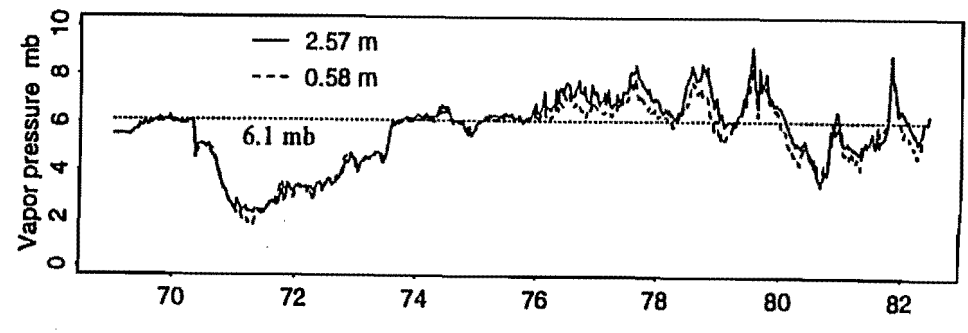


c) Precipitation.

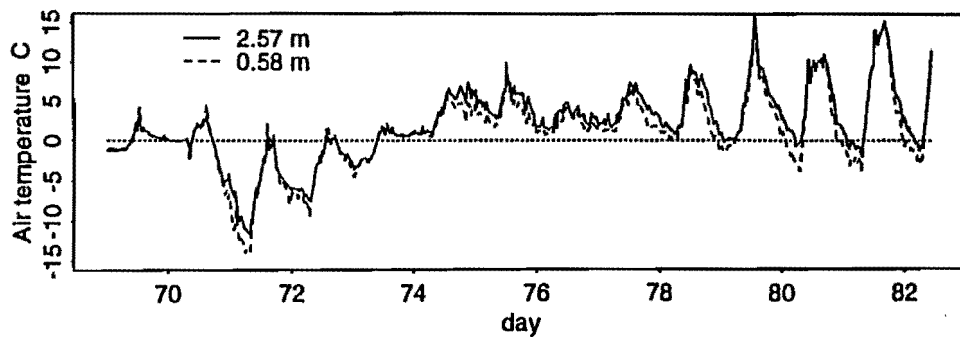
Figure 7. Detailed results for melt period, March 9, 1993 to March 23, 1993 (Days 69 to 82).



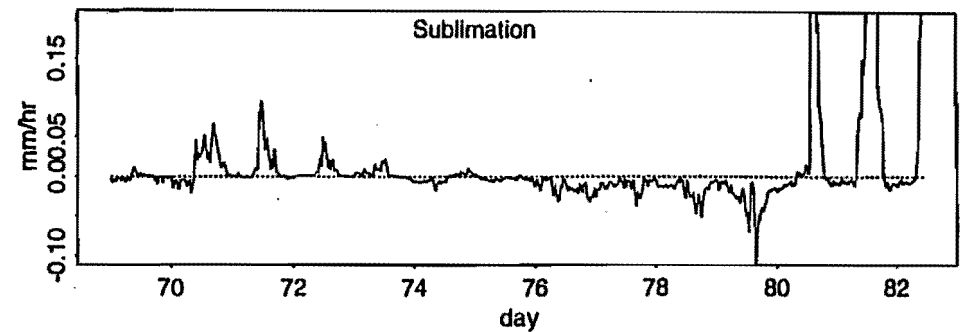
d) Wind velocity at upper and lower levels



g) Water vapor pressure at upper and lower levels.

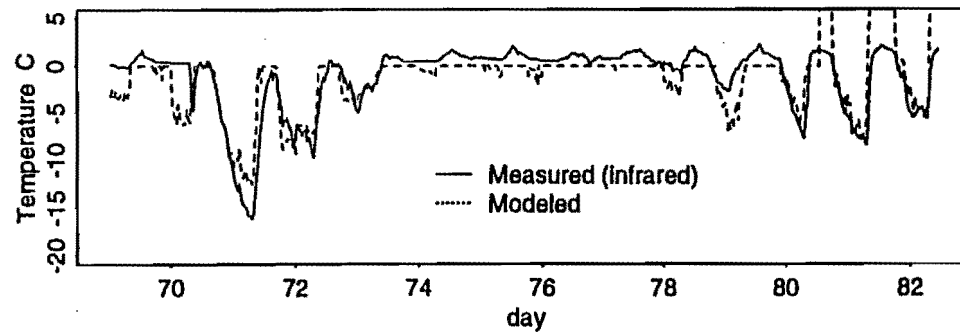


e) Air temperature at upper and lower levels.



h) Modeled sublimation.

17



f) Snow surface temperature.

Figure 7. Detailed results for melt period, March 9, 1993 to March 23, 1993 (Days 69 to 82), continued.

## **Conclusions**

An experiment to quantify the sublimation and energy balance of snow was conducted the winter of 1992/93 at the Utah State University drainage and evapotranspiration research farm near Logan, Utah, USA. The experiment was not altogether successful in that large temperature dependent oscillations in the weight recorded by the lysimeters precluded the measurement of sublimation. However the meteorological variables measured were used to test an energy balance snowmelt model. Comparisons against measured snow water equivalent and measured snow surface temperatures indicate satisfactory performance of the model in representing these aspects of the snow accumulation, energy and melt processes. Deficiencies in the models representation of the snow energy content were found and will need to be addressed in future work. Future work with this data set could also attempt to remove the temperature dependence from the lysimeter measurements and obtain estimates of measured sublimation. There is also the information necessary to quantify heat flux, somewhat tenuously, as the residual from net radiation, ground heat flux and changes in energy content of the snow. This could then be compared to temperature gradients and modeled heat flux based on wind. It may also be possible to obtain useful information and learn something about the turbulent transfers of sensible and latent heat fluxes from the analysis of gradient information. This will however be difficult as the air temperature and humidity differences measured were small and approach the resolution limit of the sensors. Overall the improvement of our understanding of turbulent processes over snow will require more study and more precise measurements.

## **Acknowledgements**

This work was supported by a Utah State University faculty research grant and the USDA Forest Service joint venture agreement INT-92660-RJVA. Thank you Richard Allen for access to the USU drainage and evapotranspiration research farm instrumentation, data and tremendous assistance in instrumentation setup and data interpretation. Thank you Arijit Chattopadhyay for your efforts as field assistant.

## **References**

- Anderson, E. A., (1976), "A Point Energy and Mass Balance Model of a Snow Cover," NOAA Technical report NWS 19, U.S. Department of Commerce.
- Bowles, D. S., G. E. Bingham, U. Lall, D. G. Tarboton, E. Malek, B. Rajagopalan, T. Chowdhury and E. Kluzek, (1992), "Development of Mountain Climate Generator and Snowpack Model for Erosion Predictions in the Western United States using WEPP," Phase IV research completion report submitted to the U.S.D.A. Forest Service Intermountain research station under

joint venture agreement INT-92660-RJVA, Utah Water Research Laboratory.

Bowles, D. S., G. E. Bingham, U. Lall, D. G. Tarboton, B. Rajagopalan, T. Chowdhury and E. Kluzek, (1994), "Development of Mountain Climate Generator and Snowpack Model for Erosion Predictions in the Western United States using WEPP," Research Completion Report for the funding period January 1, 1993 to September 30, 1993 of Phase IV, submitted to the U.S.D.A. Forest Service Intermountain research station under joint venture agreement INT-92660-RJVA, Utah Water Research Laboratory.

Bristow, K. L. and G. S. Campbell, (1984), "On the Relationship Between Incoming Solar Radiation and the Daily Maximum and Minimum Temperature," Agricultural and Forest Meteorology, 31: 159-166.

Brutsaert, W., (1982), Evaporation into the Atmosphere, Kluwer Academic Publishers, 299 p.

Calder, I. R., (1990), Evaporation in the Uplands, John Wiley & Sons, Chichester, 148 p.

Chowdhury, T. G., D. G. Tarboton and D. S. Bowles, (1992), "An Energy Balance Snowmelt Model for Erosion Prediction," Eos Transactions AGU, 73(43): Fall Meeting Suppl., 174.

Dickinson, R. E., A. Henderson-Sellers and P. J. Kennedy, (1993), "Biosphere-Atmosphere Transfer Scheme (BATS) Version 1e as Coupled to the NCAR Community Climate Model," NCAR/TN-387+STR, National Center for Atmospheric Research.

Gerald, C. F., (1978), Applied Numerical Analysis, 2nd Edition, Addison Wesley, Reading, Massachusetts, 518 p.

Lowe, P. R., (1977), "An Approximating Polynomial for the Computation of Saturation Vapour Pressure," Journal of Applied Meteorology, 16: 100-103.

Male, D. H. and D. M. Gray, (1981), "Snowcover Ablation and Runoff," Chapter 9 in Handbook of Snow. Principles. Processes. Management and Use, Edited by D. M. Gray and D. H. Male, Pergammon Press, p.360-436.

Rosenberg, N. J., (1974), Microclimate The Biological Environment, John Wiley & Sons, Inc., 315 p.

Satterlund, D. R., (1979), "An Improved Equation for Estimating Long-wave Radiation From the Atmosphere," Water Resources Research, 15: 1643-1650.

Tarboton, D. G., T. G. Chowdhury and T. H. Jackson, (1995), "A Spatially Distributed Energy Balance Snowmelt Model," Paper in preparation for presentation at IAHS symposium, July 3-14, Boulder Colorado.

U.S. Army Corps of Engineers, (1956), "Snow Hydrology, Summary report of the Snow Investigations," , U.S. Army Corps of Engineers, North Pacific Division, Portland, Oregon.

# APPENDIX 3B

A SPATIALLY DISTRIBUTED ENERGY BALANCE SNOWMELT MODEL

# A Spatially Distributed Energy Balance Snowmelt Model

David G. Tarboton  
Tanveer G. Chowdhury  
Thomas H. Jackson

*Submitted for publication in IAHS  
Proceedings of Symposium on  
Biogeochemistry of Seasonally Snow  
Covered Catchments, June 3-14, 1995,  
Boulder, CO.*

Working Paper WP-94-HWR-DGT/003

November 1994



## **A Spatially Distributed Energy Balance Snowmelt Model**

DAVID G. TARBOTON, TANVEER G. CHOWDHURY

Utah Water Research Laboratory, Utah State University, Logan, Utah 84322-8200, USA

THOMAS H. JACKSON

Turner Collie & Braden, Houston, Texas, USA

**Abstract** This paper describes an energy balance snowmelt model developed for the prediction of rapid snowmelt rates responsible for soil erosion and water input to a distributed water balance model. The model uses a lumped representation of the snowpack with two state variables, namely, water equivalent and energy content relative to a reference state of water in the ice phase at 0°C. This energy content is used to determine snowpack average temperature or liquid fraction. This representation of the snowpack is used in a distributed version of the model with each of these state variables modeled at each point on a rectangular grid corresponding to a digital elevation model. Inputs are air temperature, precipitation, wind speed, humidity and radiation at hourly time steps. The model uses physically based calculations of radiative, sensible, latent and advective heat exchanges. An equilibrium parameterization of snow surface temperature accounts for differences between snow surface temperature and average snowpack temperature without having to introduce additional state variables. Melt outflow is a function of the liquid fraction, using Darcy's law. This allows the model to account for continued outflow even when the energy balance is negative. A detailed description of the model is given together with results of tests of individual components and the complete model against data collected at the Central Sierra Snow Laboratory, California; Reynolds Creek Experimental Watershed, Boise Idaho; and at the Utah State University drainage research farm, Logan Utah. The testing includes comparisons against melt outflow collected in lysimeters and melt collectors, surface snow temperatures collected using infrared temperature sensors and depth and water equivalent measured using snow core samplers.

## **INTRODUCTION**

Snowmelt is a significant surface water input of importance to many aspects of hydrology including water supply, erosion and flood control. Snowmelt is driven primarily by energy

exchanges at the snow-air interface. The model described here was developed initially to predict the rapid melt rates responsible for erosion. It has also been used to provide the spatially distributed surface water input in a water balance study. In developing a new snowmelt model our goal was to incorporate ideas from the many existing models and parameterize the processes involved in as simple, yet physically correct a manner as possible. We hoped to develop a simple, parsimonious, physically based model that could be driven by readily available inputs and applied anywhere with no (or minimal) calibration. The striving for simplicity led us to parameterize a snowpack in terms of lumped (depth averaged) state variables so as to avoid having to model the complex processes that occur within a snowpack. We have still however attempted to capture important physical differences between bulk (depth averaged) properties and the surface properties that are important for surface energy exchanges. Due to space limitations a detailed literature review is not given. We have relied heavily on an understanding of snowmelt processes gleaned from Gray and Male (1981) and the descriptions of existing models (Anderson, 1973; 1976; Morris, 1982; Leavesley *et al.*, 1983). In what follows we first give a detailed description of the model. We then describe the data sets we used to test the model and show results comparing model calculations to observations.

## MODEL DESCRIPTION

The snowpack is characterized by state variables, water equivalent,  $W$  [m], energy content,  $U$ , [ $\text{kJ}/\text{m}^2$ ] and the age of the snow surface which is only used for albedo calculations. These are, we believe, sufficient to characterize the snowpack for the surface water inputs of interest. The state variable, energy content  $U$ , is defined relative to a reference state of water at  $0^\circ\text{C}$  in the ice (solid) phase.  $U$  greater than zero means the snowpack (if any) is isothermal with some liquid content and  $U$  less than zero can be used to calculate the snowpack average temperature,  $T$ , [ $^\circ\text{C}$ ]. Energy content is defined as the energy content of the snowpack plus a top layer of soil with depth  $D_e$  [m]. We discuss below the choice of  $D_e$  and the role it plays in the model.

The model is designed to be driven by inputs of air temperature,  $T_a$  [ $^\circ\text{C}$ ]; wind speed,  $V$  [m/s]; relative humidity,  $RH$ ; precipitation,  $P$  [m/hr]; incoming solar  $Q_{si}$  and longwave  $Q_{li}$  radiation [ $\text{kJ}/\text{m}^2/\text{hr}$ ]; and ground heat flux  $Q_g$  [ $\text{kJ}/\text{m}^2/\text{hr}$ ] (taken as 0 when not known) at each time step. Time steps of 0.5, 1 and 6 hours have been used in data comparisons here. When incoming solar radiation is not available it is estimated as an extra terrestrial radiation (from sun angle and

solar constant) times an atmospheric transmission factor,  $Tr$ , estimated from the daily temperature range using the procedure given by Bristow and Campbell (1984). When incoming longwave radiation is not available it is estimated based on air temperature, the Stefan-Boltzman equation and a parameterization of air emissivity due to Satterlund (1979), adjusted for cloudiness using  $Tr$ .

Given the state variables  $U$  and  $W$ , their evolution in time is determined by solving energy and mass balance equations.

$$\frac{dU}{dt} = Q_{sn} + Q_{li} + Q_p + Q_g - Q_{le} + Q_h + Q_e - Q_m \quad (1)$$

$$\frac{dW}{dt} = P_r + P_s - M_r - E \quad (2)$$

In the energy balance equation terms are (all in  $\text{kJ/m}^2/\text{hr}$ ):  $Q_{sn}$ , net shortwave radiation;  $Q_{li}$ , incoming longwave radiation;  $Q_p$ , advected heat from precipitation;  $Q_g$ , ground heat flux;  $Q_{le}$ , outgoing longwave radiation;  $Q_h$ , sensible heat flux;  $Q_e$ , latent heat flux due to sublimation/condensation;  $Q_m$ , advected heat removed by meltwater. In the mass balance equation (all in  $\text{m/hr}$  of water equivalent) terms are:  $P_r$ , rainfall rate;  $P_s$  snowfall rate;  $M_r$ , meltwater outflow from the snowpack;  $E$ , sublimation from the snowpack. Many of these fluxes depend functionally on the state and input driving variables. We elaborate on the parameterization of these functional dependencies below. Equations (1) and (2) form a coupled set of first order, nonlinear ordinary differential equations. They can be summarized in vector notation as:

$$\frac{d\mathbf{X}}{dt} = \mathbf{F}(\mathbf{X}, \text{driving variables}) \quad (3)$$

where  $\mathbf{X} = (U, W)$  is a state vector describing the snowpack. With  $\mathbf{X}$  specified initially, this is an initial value problem. A large variety of numerical techniques are available for solution of initial value problems of this form. Here we have adopted a Euler predictor-corrector approach (Gerald, 1978).

$$\mathbf{X}' = \mathbf{X}_i + \Delta t \mathbf{F}(\mathbf{X}_i, \text{driving variables}) \quad (4)$$

$$\mathbf{X}_{i+1} = \mathbf{X}_i + \Delta t \frac{\mathbf{F}(\mathbf{X}_i, \text{driving variables}) + \mathbf{F}(\mathbf{X}', \text{driving variables})}{2} \quad (5)$$

where  $\Delta t$  is the time step,  $\underline{X}_i$  refers to the state at time  $t_i$  and  $\underline{X}_{i+1}$  refers to the state at time  $t_{i+1}=t_i+\Delta t$ . This is a second order finite difference approximation, with global error proportional to  $\Delta t^2$  (Gerald, 1978, p257). Numerical instabilities sometimes occur under melting conditions when the snowpack is shallow due to the nonlinear nature of the melt outflow parameterization. To deal with this we compare  $\underline{X}_{i+1}$  to  $\underline{X}'$  and if they differ by more than a specified tolerance (0.025 m for W and 2000 kJ/m<sup>2</sup> for U) iterate up to four times setting  $\underline{X}'$  to  $\underline{X}_{i+1}$  then recalculating  $\underline{X}_{i+1}$  at each iteration. If convergence is still not achieved we take the solution that would keep the liquid fraction of the snow constant. Following I describe how each of the processes involved in equations (1) and (2) are parameterized.

### Depth averaged temperature - T

The snow and interacting soil layer average temperatures are obtained from the energy content and water equivalent, relative to 0°C ice phase.

$$\text{If } U < 0 \quad T = U / (\rho_w W C_s + \rho_g D_e C_g) \quad \text{All solid phase} \quad (6)$$

$$\text{If } 0 < U < \rho_w W h_f \quad T = 0^\circ\text{C}. \quad \text{Solid and liquid mixture} \quad (7)$$

$$\text{If } U > \rho_w W h_f \quad T = \frac{U - \rho_w W h_f}{\rho_g D_e C_g + \rho_w W C_w} \quad \text{All liquid} \quad (8)$$

In the above the heat required to melt all the snow water equivalent is  $\rho_w W h_f$  [kJ] where  $h_f$  is the heat of fusion [333.5 kJ kg<sup>-1</sup>] and U in relation to this determines the solid-liquid phase mixtures. The heat capacity of the snow is  $\rho_w W C_s$  [kJ/°C] where  $\rho_w$  is the density of water [1000 kg m<sup>-3</sup>] and  $C_s$  the specific heat of ice [2.09 kJ kg<sup>-1</sup> °C<sup>-1</sup>]. The heat capacity of the soil layer is  $\rho_g D_e C_g$  [kJ/°C] where  $\rho_g$  is the soil density and  $C_g$  the specific heat of soil. These together determine the T when  $U < 0$ . The heat capacity of liquid water,  $\rho_w W C_w$ , where  $C_w$  is the specific heat of water [4.18 kJ kg<sup>-1</sup> °C<sup>-1</sup>], is included in (8) for numerical consistency during time steps when the snowpack completely melts.

The parameter  $D_e$  is intended to quantify the depth of soil that interacts thermally with the snowpack. Heat flow in snow and soil is governed by Laplace's equation. The depth of penetration of changes in surface temperature can be evaluated from the expression (Rosenberg, 1974):

$$\frac{R_z}{R_s} = \exp\left(-z \left(\frac{\pi}{\alpha P}\right)^{\frac{1}{2}}\right) \quad (9)$$

where  $R_s$  is the range of temperature oscillation at the surface,  $R_z$  the range of temperature oscillation at depth  $z$ ,  $P$  the period of oscillation, and  $\alpha$  the thermal conductivity. For soil  $\alpha$  is typically in the range 0.004 to 0.006 cm<sup>2</sup>/s. Fig. 1 shows  $R_z/R_s$  versus  $z$  for  $\alpha = 0.005$  cm<sup>2</sup>/s for various periods. This shows that for oscillations less than one week the effect at 40 cm is damped to less than 30% and even for monthly oscillations is still damped 50% at 40 cm depth. This suggests using  $D_e = 40$  cm in our model since the time scale of interest is the seasonal accumulation then melting of snow. The state variable  $U$  represents energy content above this level. The ground heat flux represents heat transport at this depth and is therefore a long term average. Oscillating, high frequency, ground heat fluxes above this depth are absorbed into  $U$ , the energy stored in the snow and soil above depth  $D_e$ . This procedure provides a simple approximation of the effects of frozen ground, or snow falling on warm ground.

## Radiation

Net shortwave radiation is calculated as

$$Q_{sn} = Q_{si} (1-A) \quad (10)$$

where Albedo,  $A$ , is calculated based on the age of the snow surface using a parameterization due to Dickinson et al. (1993). The age of the snow surface is retained as a state variable, and is updated each time step, dependent on snow surface temperature and snowfall. When the snowpack is shallow (depth  $z < h = 0.1$  m) the albedo is taken as  $r A_{bg} + (1-r) A$  where  $r = (1-z/h)e^{-z/2h}$ . This interpolates between the snow albedo,  $A$ , and bare ground albedo,  $A_{bg}$ , with the exponential term approximating the exponential extinction of radiation penetration of snow.

Outgoing longwave radiation is

$$Q_{le} = \epsilon_s \sigma (T_s^{abs})^4 \quad (11)$$

where  $\epsilon_s$  is emissivity,  $\sigma$  the Stefan Boltzmann constant [ $2.07 \times 10^{-7} \text{ kJ m}^{-2} \text{ hr}^{-1} \text{ K}^{-4}$ ] and the superscript "abs" in  $T_s^{abs}$  indicates that this is absolute temperature [K].

### Snow fall accumulation and heat with precipitation

Measured precipitation rate,  $P$ , is partitioned into rain,  $P_r$ , and snow,  $P_s$ , (both in terms of water equivalent depth) using the following rule based on air temperature,  $T_a$ , (U.S. Army Corps of Engineers, 1956)

$$\begin{aligned} P_r &= P & T_a &\geq T_r = 3 \text{ }^\circ\text{C} \\ P_r &= P(T_a - T_b)/(T_r - T_b) & T_b &< T_a < T_r \\ P_r &= 0 & T_a &\leq T_b = -1 \text{ }^\circ\text{C} \\ P_s &= (P - P_r) F \end{aligned} \quad (12)$$

where  $T_r$  is a threshold air temperature above which all precipitation is rain and  $T_b$  a threshold air temperature below which all precipitation is snow. The accumulation of snow is sometimes subject to considerable wind redistribution with drifts forming on lee slopes. We account for this in the model through a snow drift factor,  $F$ , dependent on location. Ideally  $F$  needs to be related to topography. In the application to Reynolds Creek,  $F$  was estimated by calibrating the snow water equivalents obtained from the snow model (with  $F = 1$ ) at each cell,  $W_m$ , against the observed values,  $W_o$ . The discrepancy between observations and predictions over an interval between measurements is attributed to drifting and suggests  $F = 1 + (W_o - W_m)/P_s$  where  $P_s$  is the gage snowfall (calculated from  $P$  with  $F = 1$ ) during the interval. Values of  $F$  less than one correspond to locations of depletion or wind scour. This approach models drifting which actually occurs after snowfall as concurrent with snowfall. The calibration of  $F$  assumes that the snowmelt model correctly accounts for all other processes (melt, sublimation, condensation etc.) affecting the accumulation and ablation of snow water equivalent. Further details are given in Jackson (1994).

The temperature of rain is taken as the greater of the air temperature and freezing point and the temperature of snow the lesser of air temperature and freezing point. The advected heat is the energy required to convert this precipitation to the reference state ( $0^\circ\text{C}$  ice phase).

$$Q_p = P_s C_s \rho_w \min(T_a, 0 \text{ }^\circ\text{C}) + P_r [h_f \rho_w + C_w \rho_w \max(T_a, 0 \text{ }^\circ\text{C})] \quad (13)$$

### Turbulent fluxes, $Q_h$ , $Q_e$ , $E$

Sensible and latent heat fluxes between the snow surface and air above are modeled using the concept of flux proportional to temperature and vapor pressure gradients with constants of proportionality, the so called turbulent transfer coefficients or diffusivity a function of windspeed and surface roughness. Considering a unit volume of air, the heat content is  $\rho_a C_p T_a$  and the vapor content  $\rho_a q$ , where  $\rho_a$  is air density (determined from atmospheric pressure and temperature),  $C_p$  air specific heat capacity [ $1.005 \text{ kJ kg}^{-1} \text{ }^\circ\text{C}^{-1}$ ], and  $q$  specific humidity [kg water vapor per kg air]. Heat transport towards the surface,  $Q_h$  [ $\text{kJ/m}^2/\text{hr}$ ] is given by:

$$Q_h = K_h \rho_a C_p (T_a - T_s) \quad (14)$$

where  $K_h$  is heat conductance [m/hr] and  $T_s$  is the snow surface temperature. Vapor transport away from the surface (sublimation),  $M_e$  [kg/hr] is:

$$M_e = K_e \rho_a (q_s - q) \quad (15)$$

where  $q_s$  is the surface specific humidity and  $K_e$  the vapor conductance [m/hr].

By comparison with the usual expressions for turbulent transfer in a logarithmic boundary layer profile (Male and Gray, 1981; Anderson, 1976; Brutsaert, 1982) for neutral condition, one obtains the following expression:

$$K_h = K_e = \frac{k^2 V}{[\ln(z/z_0)]^2} = K \quad (16)$$

where  $V$  is wind speed [m/hr] at height  $z$  [m];  $z_0$  is roughness height at which the logarithmic boundary layer profile predicts zero velocity [m]; and  $k$  is von Karman's constant [0.4]. Recognizing that the latent heat flux towards the snow is:

$$Q_e = -h_v M_e \quad (17)$$

and using the relationship between specific humidity and vapor pressure and the ideal gas law, one obtains:

$$Q_e = K_e \frac{h_v 0.622}{R_d T_a^{\text{abs}}} (e_a - e_s(T_s)) \quad (18)$$

where  $e_s$  is the vapor pressure at the snow surface snow, assumed saturated at  $T_s$ , and calculated using a polynomial approximation (Lowe, 1977);  $e_a$  is air vapor pressure,  $R_d$  is the dry gas constant [ $287 \text{ J kg}^{-1} \text{ K}^{-1}$ ] and  $h_v$  the latent heat of sublimation [ $2834 \text{ kJ/kg}$ ]. The water equivalent depth of sublimation is:

$$E = -\frac{Q_e}{\rho_w h_v} \quad (19)$$

When there is a temperature gradient near the surface, buoyancy effects may enhance or dampen the turbulent transfers. This can be quantified in terms of the Richardson number or Monin-Obukhov length. Adjustments to the neutral transfer coefficients to account for these effects exist and were tried based on the temperature difference between the air and snow surface. However we found that it was quite common that large temperature differences and low wind speeds resulted in unreasonable correction factors, beyond the range for which they had been developed, so for the purposes of the results presented here we have used neutral transfer coefficients.

### Snow Surface Temperature, $T_s$

Since snow is a relatively good insulator,  $T_s$  is in general different from  $T$ . This is accounted for using an equilibrium approach that balances energy fluxes at the snow surface. Heat conduction into the snow is calculated using the temperature gradient and thermal diffusivity of snow, approximated by:



$$Q = \kappa \rho_s C_s (T_s - T)/Z_e = K_s \rho_s C_s (T_s - T) \quad (20)$$

where  $\kappa$  is snow thermal diffusivity [ $m^2 \text{ hr}^{-1}$ ] and  $Z_e$  [m] an effective depth over which this thermal gradient acts. The ratio  $\kappa/Z_e$  is denoted  $K_s$  and termed snow surface conductance analogous to the heat and vapor conductances. A value of  $K_s$  is obtained by assuming a depth,  $Z_e$  equal to the depth of penetration of a diurnal temperature fluctuation calculated from equation (9) (Rosenberg, 1974).  $Z_e$  should be chosen so that  $R_z/R_s$  is small. Here  $K_s$  is used as a tuning parameter, with this calculation used to define a reasonable range. Then assuming equilibrium at the surface, the surface energy balance gives.

$$Q = Q_{sn} + Q_{li} + Q_h(T_s) + Q_e(T_s) + Q_p - Q_{le}(T_s) \quad (21)$$

where the dependence of  $Q_h$ ,  $Q_e$ , and  $Q_{le}$  on  $T_s$  is through equations (14), (18) and (11).

Analogous to the derivation of the Penman equation for evaporation the functions of  $T_s$  in this energy balance equation are linearized about a reference temperature,  $T^*$ , and the equation is solved for  $T_s$ :

$$T_s^{abs} = \frac{Q_{sn} + Q_{li} + Q_p + K T_a^{abs} \rho_a C_p - 0.622 K h_v \rho_a \left( e_s(T^*) - e_a - T^* \Delta \right) / P_a + 3 \epsilon_s \sigma T^{*abs4} + \rho_s C_s T^{abs} K_s}{\rho_s C_s K_s + K \rho_a C_p + 0.622 \Delta K h_v \rho_a / P_a + 4 \epsilon_s \sigma T^{*abs3}} \quad (22)$$

where  $\Delta = de_s/dT$ . This equation is used in an iterative procedure with an initial estimate  $T^* = T_a$ , in each iteration replacing  $T^*$  by the latest  $T_s$ . The procedure converges to a final  $T_s$  which if less than freezing is used to calculate surface energy fluxes. If the final  $T_s$  is greater than freezing it means that the energy input to the snow surface cannot be balanced by thermal conduction into the snow. Surface melt will occur and the infiltration of meltwater will account for the energy difference and  $T_s$  is then set to  $0^\circ\text{C}$ .

### Meltwater Outflux, $M_r$ and $Q_m$

The energy content state variable  $U$  determines the liquid content of the snowpack. This, together

with Darcy's law for flow through porous media, is used to determine the outflow rate.

$$M_r = K_{sat} S^*{}^3 \quad (23)$$

where  $K_{sat}$  is the snow saturated hydraulic conductivity and  $S^*$  is the relative saturation in excess of water retained by capillary forces. This expression is based on Male and Gray (1981 p400 eqn 9.45).  $S^*$  is given by:

$$S^* = \frac{\text{liquid water volume} - \text{capillary retention}}{\text{pore volume} - \text{capillary retention}} = \left( \frac{L_f}{1 - L_f} - L_c \right) / \left( \frac{\rho_w}{\rho_s} - \frac{\rho_w}{\rho_i} - L_c \right) \quad (24)$$

where  $L_f = U / (\rho_w h_f W)$  denotes the mass fraction of total snowpack (liquid and ice) that is liquid,  $L_c$  [0.05] the capillary retention as a fraction of the solid matrix water equivalent, and  $\rho_i$  the density of ice [917 kg m<sup>-3</sup>]. This melt outflow is assumed to be at 0°C so the heat advected with it, relative to the solid reference state is:

$$Q_m = \rho_w h_f M_r \quad (25)$$

## Forest Cover

The presence of vegetation, especially forests, significantly influences energy exchanges at the snow surface. A forest canopy reduces windspeed, thus reducing sensible and latent heat transfers. It also affects the radiation exchanges. The penetration of radiation through vegetation has been widely studied (Sellers *et al.*, 1986; Verstraete, 1987a; 1987b; Verstraete *et al.*, 1990; Dickinson *et al.*, 1993), and models developed that discretize the canopy into layers treating the energy balance of each layer separately (Bonan, 1991). Here we avoid these complexities and adopt a pragmatic parameterization modeled after the representation of snowmelt used by the WEPP winter routines (Young *et al.*, 1989; Hendrick *et al.*, 1971). Forest cover is parameterized by the canopy density parameter  $F_c$ , representing the canopy closure fraction (between 0 and 1). Windspeed, and therefore the corresponding heat and vapor fluxes are reduced by a factor  $(1 - 0.8F_c)$ . Radiative fluxes  $Q_{SN}$ ,  $Q_{li}$  and  $Q_{le}$  in equation (1) are reduced by a factor  $(1 - F_c)$ . Adjustments are also made to the radiation terms in the calculation of snow surface

temperature (equation 22).

## **DATA**

In this paper data collected at the Central Sierra Snow Laboratory (CSSL); Utah State University drainage and evapotranspiration research farm and Reynolds Creek Experimental Watershed are used to calibrate and test the model.

### **Central Sierra Snow Laboratory**

The CSSL located 1 km east of Soda Springs, California, measures and archives comprehensive data relevant to snow. It is at latitude 39°19'N and at elevation 1100m. We obtained the meteorological and snow observation data for the winter of 1985 - 1986. The meteorological data is reported each hour and consists of temperature, radiation, humidity, precipitation, and wind measurements at two levels in a 40 x 50 m clearing and in a mixed conifer fir forest with 95% forest cover. Only data from the clearing are used here. Snow depths and water equivalent are measured daily (except on weekends) and eight lysimeters record melt outflow each hour. We used the temperature, precipitation, radiation (incoming solar and net), humidity and wind measurements to drive our model and compared model output to measurements of snow water equivalent, melt outflow and snow surface temperature (infrared sensor).

### **USU drainage and evapotranspiration research farm**

An experiment to measure snow energy balance and sublimation from snow the winter of 1992 - 1993 is described more fully by Tarboton , 1994 #1669]. Data from this work included measurements of snow water equivalent, snow surface temperature and the meteorological variables necessary to drive our model.

### **Reynolds Creek Experimental Watershed**

Upper Sheep creek is a 26 ha catchment within the semi-arid Reynolds creek experimental watershed. Snowmelt is the main hydrologic input and its areal distribution is heavily influenced

by wind induced drifting. Detailed descriptions of the various features of the area are given in Flerchinger et al. (1992) and references therein. Snow water equivalent measurements are made biweekly (as weather permits) on a 30.48 m (100 ft) grid over the watershed. A digital elevation model (DEM) was constructed from a 1:1200 map with 0.61 m (2 ft) contour interval developed from low level aerial photography. The DEM grid was constructed to coincide with the grid used for field measurements and provided slope and aspect inputs to the model radiation calculations. Fig. 2 shows the topography and grid over Upper Sheep creek together with locations of the instrumentation. Data from the winters of 1985 - 1986 and 1992 - 1993 were used in this study to test the model running in a distributed mode at each grid cell. Snow melt outputs were used as hydrologic inputs for a water balance study (Jackson, 1994; Tarboton *et al.*, 1995).

## RESULTS

The model was calibrated against the CSSL data for the winter 1985 - 1986. The energy balance and overall accumulation and ablation of the snowpack is governed primarily by surface energy exchange processes. The adjustable parameters involved in these are  $z_0$  and  $K_s$ , which were adjusted to obtain a match between water equivalent, modeled and observed (shown in Fig. 3), and snow surface temperatures, modeled and observed (Fig. 4) with the model driven by the measured net radiation input. We then used measured incoming solar radiation to drive the model and found that the melt is delayed (Fig. 3). Discrepancies were analyzed and attributed to differences in daytime net radiation, primarily affected by albedo. The albedo parameterization (Dickinson *et al.*, 1993) has parameters  $A_{VO} = 0.95$  and  $A_{NIR} = 0.65$  which represent the albedo of new snow in the visible and infrared ranges.  $A_{VO}$  was reduced to 0.85 to match the daytime net radiation when compared to measured CSSL 1985 - 1986 data (Fig. 5). The resulting snow water equivalent comparison (Fig. 3) indicates that some early season melt is not modeled resulting in slight over accumulation, but the main melt is well modeled. In all results except the line indicated on Fig. 3,  $A_{VO} = 0.85$  was used. Melt outflow rate was compared to the average from the eight melt lysimeters, with  $K_{sat}$  adjusted to get a good fit. Results are shown in Fig. 6.

Table 1 lists the adjustable parameters that were calibrated against the CSSL data. Table 2 lists the remaining model parameters which were held fixed at their nominal values. The model was tested against the data from Reynolds Creek and USU drainage and evapotranspiration research farm without further adjustment of parameters. The Reynolds Creek study applied the

model to each 30.48 x 30.48 m grid cell over Upper Sheep creek (Fig. 2). The drift factor to adjust snow input was estimated from the observed grided snow data for 1985-1986 (Jackson, 1994). Fig. 7 shows the drift factors and Fig. 8 compares measured and modeled spatial distribution of snow about halfway through the snowmelt phase in 1992-1993. Due to space limitations not all of the comparisons are shown. They indicate that the model correctly represents the spatial accumulation and melt patterns. Fig. 9 compares measured and modeled snow water equivalent at the USU drainage and evapotranspiration research farm.

## CONCLUSIONS

The tests described have shown that this simple, depth averaged, mass and energy balance snowmelt model is able to capture the essential physics of the snow accumulation and melt processes and provide distributed hydrologic inputs. Using parameter values calibrated against CSSL data the model performed well when tested at other locations. This suggests that the model is transportable and parameter values listed may be acceptable for wider application. However further testing against additional data is necessary. In particular we need to test the parameterization of forest cover and further evaluate the parameterization of albedo and the effect of atmospheric stability on turbulent fluxes.

The model is available electronically from David Tarboton (dtarb@cc.usu.edu).

**Acknowledgements** Thank you Bruce McGurk for access to the CSSL data, and Keith Cooley and the USDA ARS Northwest Watershed Research Center staff for access to and collaboration in Reynolds Creek. This work was funded in part by the US Department of the Interior, Geological Survey, under USGS Grant No. 14-08-0001-G2110, and the US Department of Agriculture, Forest Service joint venture agreement INT-92660-RJVA. The views and conclusions are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the U.S. Government.

## REFERENCES

Anderson, E. A. (1973) National Weather Service River Forecast System-Snow Accumulation and Ablation Model. NOAA Technical Memorandum NWS HYDRO-17, U.S. Dept of Commerce.

- Anderson, E. A. (1976) A Point Energy and Mass Balance Model of a Snow Cover. NOAA Technical report NWS 19, U.S. Department of Commerce.
- Bonan, G. B. (1991) A Biophysical Surface Energy Budget Analysis of Soil Temperature in the Boreal Forests of Interior Alaska. *Water Resources Research*. 27(5): 767-781.
- Bristow, K. L. and Campbell, G. S. (1984) On the Relationship Between Incoming Solar Radiation and the Daily Maximum and Minimum Temperature. *Agricultural and Forest Meteorology*. 31: 159-166.
- Brutsaert, W. (1982) *Evaporation into the Atmosphere*, Kluwer Academic Publishers.
- Dickinson, R. E., Henderson-Sellers, A. and Kennedy, P. J. (1993) Biosphere-Atmosphere Transfer Scheme (BATS) Version 1e as Coupled to the NCAR Community Climate Model. NCAR/TN-387+STR, National Center for Atmospheric Research.
- Flerchinger, G. N., Cooley, K. R. and Ralston, D. R. (1992) Groundwater Response to Snowmelt in a Mountainous Watershed. *Journal of Hydrology*. 133: 293-311.
- Gerald, C. F. (1978) *Applied Numerical Analysis*, 2nd Edition, Addison Wesley, Reading, Massachusetts.
- Gray, D. M. and Male, D. H. ed. (1981) *Handbook of Snow, Principles, processes, management & use*. Pergamon Press.
- Hendrick, R. L., Filgate, B. D. and Adams, W. M. (1971) Application of Environmental Analysis to Watershed Snowmelt. *Journal of Applied Meteorology*. (10): 418-429.
- Jackson, T. H. R. (1994) A Spatially Distributed Snowmelt-Driven Hydrologic Model applied to the Upper Sheep Creek Watershed. Ph.D Thesis, Civil and Environmental Engineering, Utah State University.
- Leavesley, G. H., Lichty, R. W., Troutman, B. M. and Saindon, L. G. (1983) Precipitation-runoff modeling system--Users manual:. Water resources Investigations Report 83-4238, U.S. Geological Survey.
- Lowe, P. R. (1977) An Approximating Polynomial for the Computation of Saturation Vapour Pressure. *Journal of Applied Meteorology*. 16: 100-103.
- Male, D. H. and Gray, D. M. (1981) Snowcover Ablation and Runoff. Chapter 9 in *Handbook of Snow, Principles, Processes, Management and Use*, Edited by D. M. Gray and D. H. Male, Pergammon Press, p.360-436.
- Morris, E. M. (1982) Sensitivity of the European Hydrological System snow models. *Hydrological Aspects of Alpine and High Mountain Areas*, Proceedings of the Exeter

- Symposium, IAHS Publ no 138, 221-231.
- Rosenberg, N. J. (1974) *Microclimate The Biological Environment*, John Wiley & Sons, Inc.
- Satterlund, D. R. (1979) An Improved Equation for Estimating Long-wave Radiation From the Atmosphere. *Water Resources Research*. 15: 1643-1650.
- Sellers, P. J., Mintz, Y., Sud, Y. C. and Dalcher, A. (1986) A simple biosphere model (SiB) for use with general circulation models. *Journal of the Atmospheric Sciences*. 43(6): 505-531.
- Tarboton, D. G., Jackson, T. H., Liu, J. Z., Neale, C. M. U., Cooley, K. R. and McDonnell, J. J. (1995) A Grid Based Distributed Hydrologic Model: Testing Against Data from Reynolds Creek Experimental Watershed. Preprint submitted for presentation at AMS Conference on Hydrology, 15-20 January, Dallas, Texas.
- U.S. Army Corps of Engineers (1956) Snow Hydrology, Summary report of the Snow Investigations. , U.S. Army Corps of Engineers, North Pacific Division, Portland, Oregon.
- Verstraete, M. M. (1987a) Radiation Transfer in Plant Canopies: Scattering of Solar Radiation and Canopy Reflectance. *Journal of Geophysical Research*. 93(D8): 9483-9494.
- Verstraete, M. M. (1987b) Radiation Transfer in Plant Canopies: Transmission of Direct Solar Radiation and the Role of Leaf Orientation. *Journal of Geophysical Research*. 92(D9): 10985-10995.
- Verstraete, M. M., Pinty, B. and Dickinson, R. E. (1990) A physical model of the bidirectional reflectance of vegetation canopies, 1. Theory. *Journal of Geophysical Research*. 95(D8): 11755-11765.
- Young, R. A., Benoit, G. R. and Onstad, C. A. (1989) Snowmelt and frozen soil. Chapter 3 in *USDA Water Erosion Prediction Project, Hillslope profile model documentation*, Edited by L. J. Lane and M. A. Nearing, NSERL Report #2, USDA-ARS National Soil Erosion Research Laboratory, West Lafayette, Indiana, 47907.

## Figure Captions

Figure 1. Depth of penetration of temperature fluctuations into soil with  $\alpha = 0.005 \text{ cm}^2/\text{s}$ .

Figure 2. Upper Sheep Creek topography and instrumentation.

Figure 3. Comparison between observed and modeled snow water equivalent, CSSL.

Figure 4. Comparison between observed and modeled snow surface temperatures, CSSL.

Figure 5. Comparison between observed and modeled net radiation, CSSL.

Figure 6. Comparison between observed and modeled melt outflow rate, CSSL.

Figure 7. Drift factor from Jackson (1994). Contours at 0.5, 0.9, 1.5, 2.5, 4 and 6.

Figure 8. Observed and modeled spatial distribution of snow at Upper Sheep creek, April 8, 1993.

Figure 9. Observed and modeled snow water equivalent, USU research farm.



Table 1. Adjustable parameter recommended values.

Parameter	Notation	Calibrated Value
Surface aerodynamic roughness	$z_0$	0.005 m
Surface conductance	$K_s$	0.02 m/hr
Saturated hydraulic conductivity	$K_{sat}$	20 m/hr
New snow visible albedo	$A_{v0}$	0.85

Table 2. Snowmelt model fixed parameters.

Parameter	Notation	Reference Value
Ground Heat Capacity	$C_g$	2.09 kJ kg <sup>-1</sup> °C <sup>-1</sup>
Density of Soil Layer	$\rho_g$	1700 kg m <sup>-3</sup>
Snow density	$\rho_s$	450 kg m <sup>-3</sup>
Capillary retention fraction	$L_c$	0.05
Emissivity of Snow	$\epsilon_s$	0.99
Temperature above which precipitation is rain	$T_r$	3°C
Temperature below which precipitation is snow	$T_s$	-1°C
Wind/Air temperature measurement height	$z$	2 m
Soil Effective Depth	$D_e$	0.4 m
Bare ground albedo	$A_{bg}$	0.25
Albedo extinction depth	$h$	0.1 m

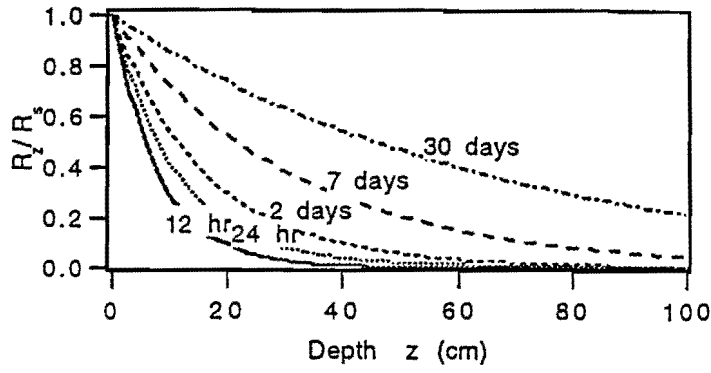


Figure 1.

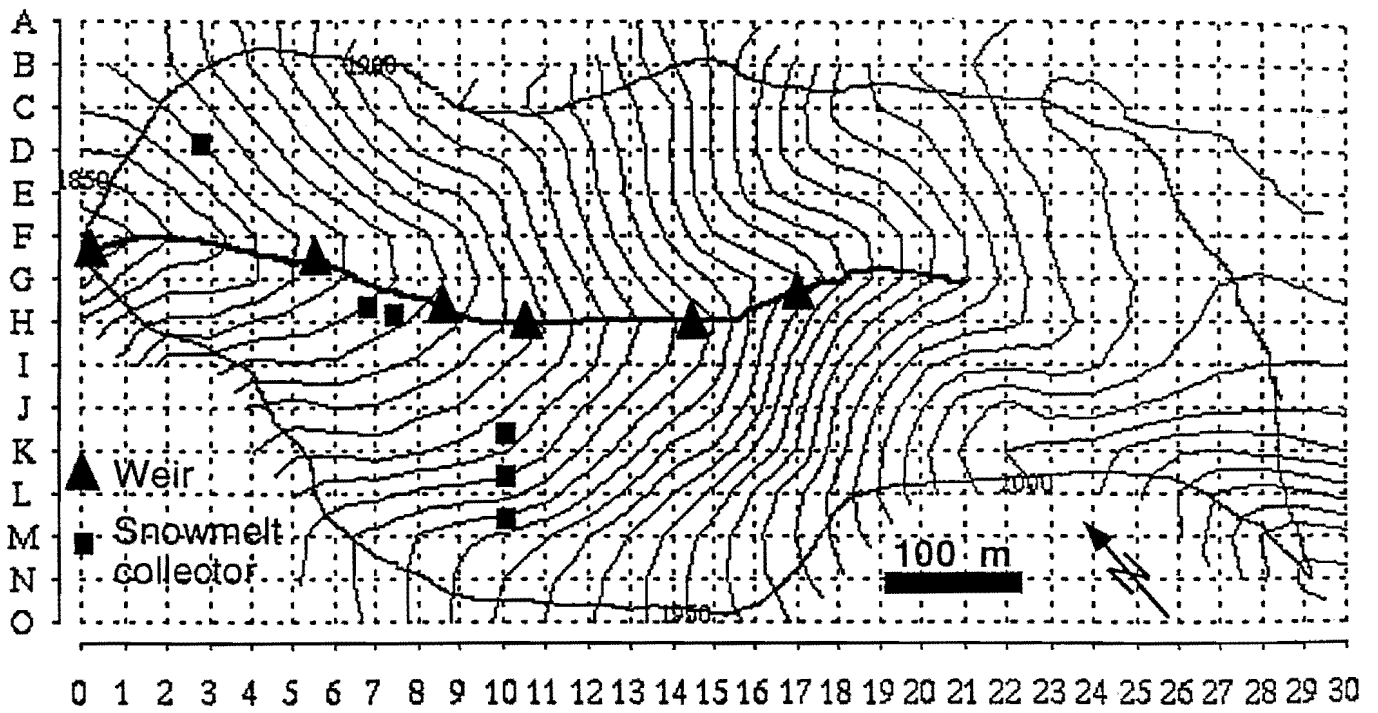


Figure 2.

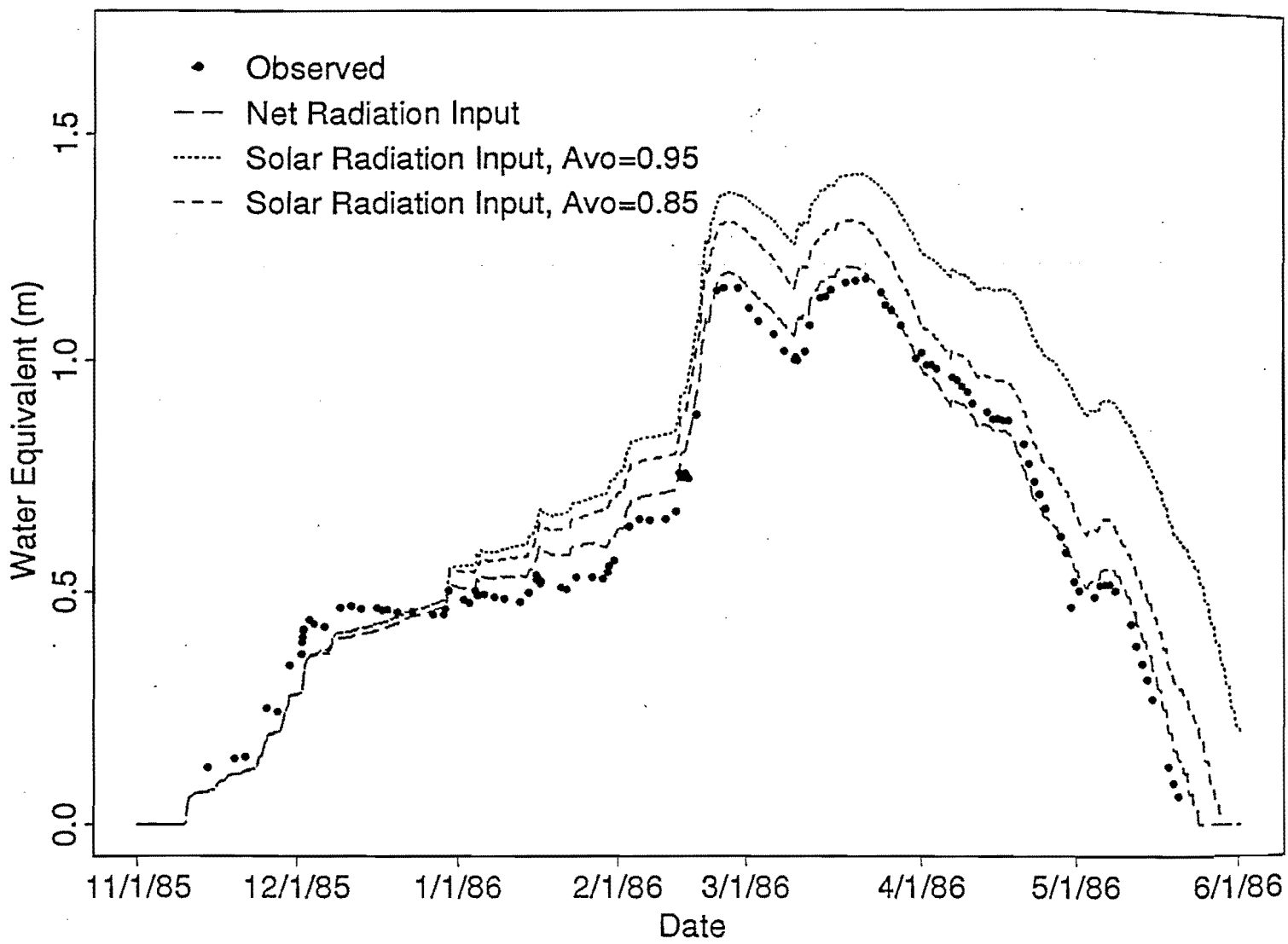


Figure 3.

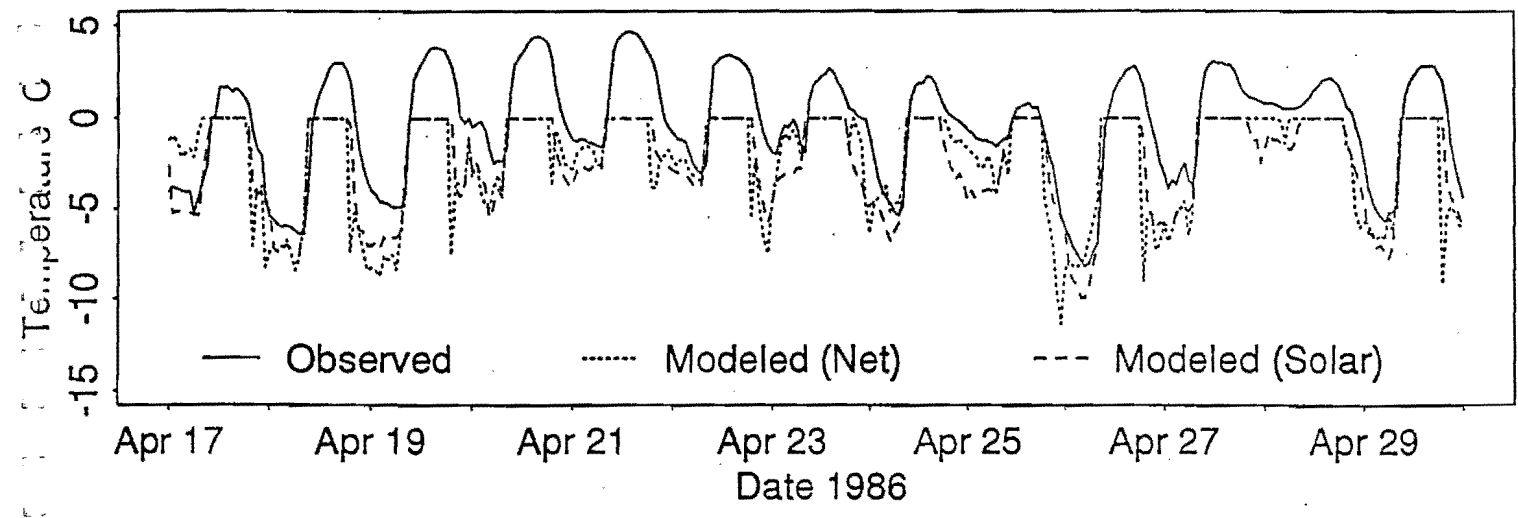
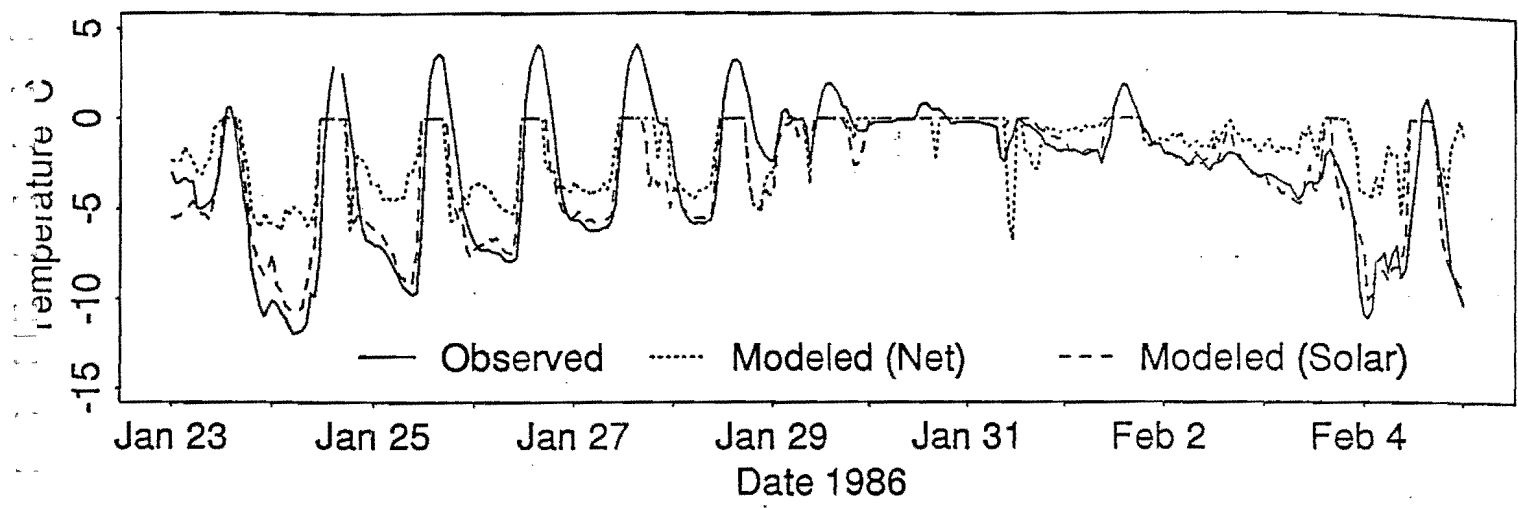


Figure 4.

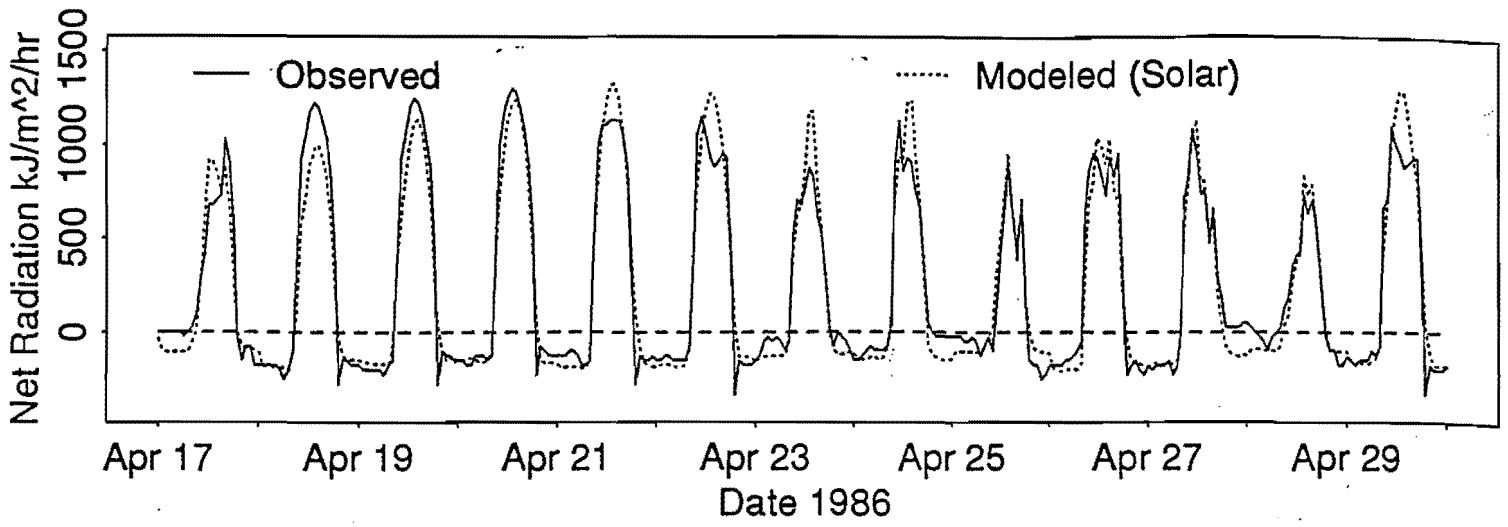


Figure 5.

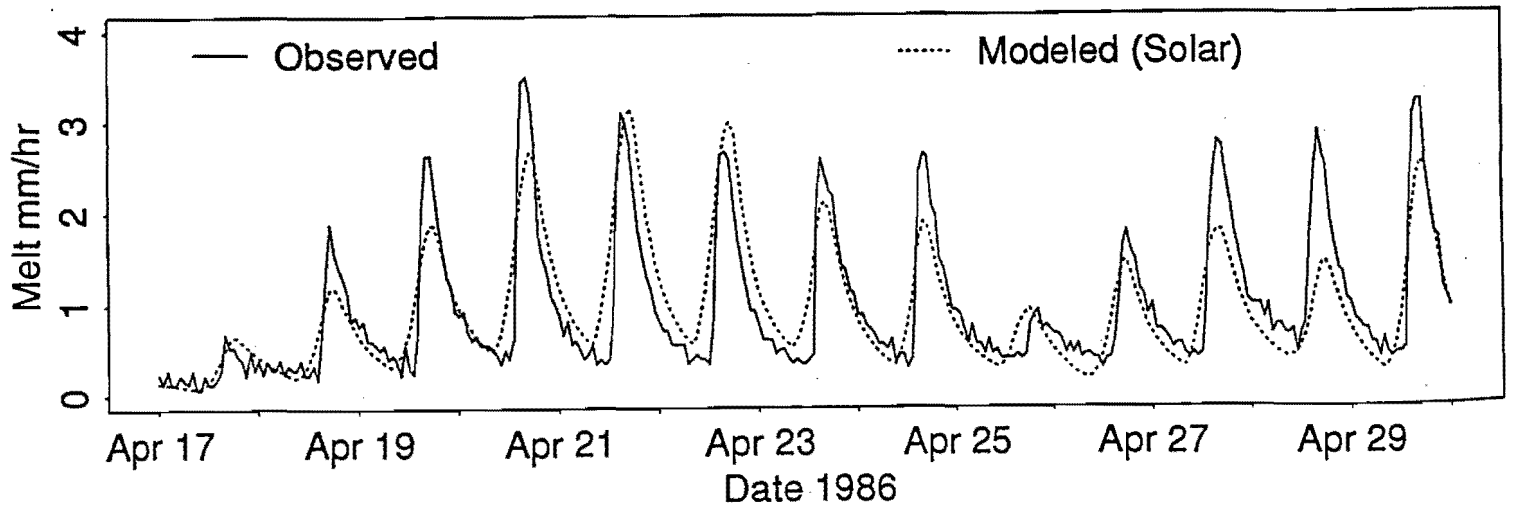


Figure 6.

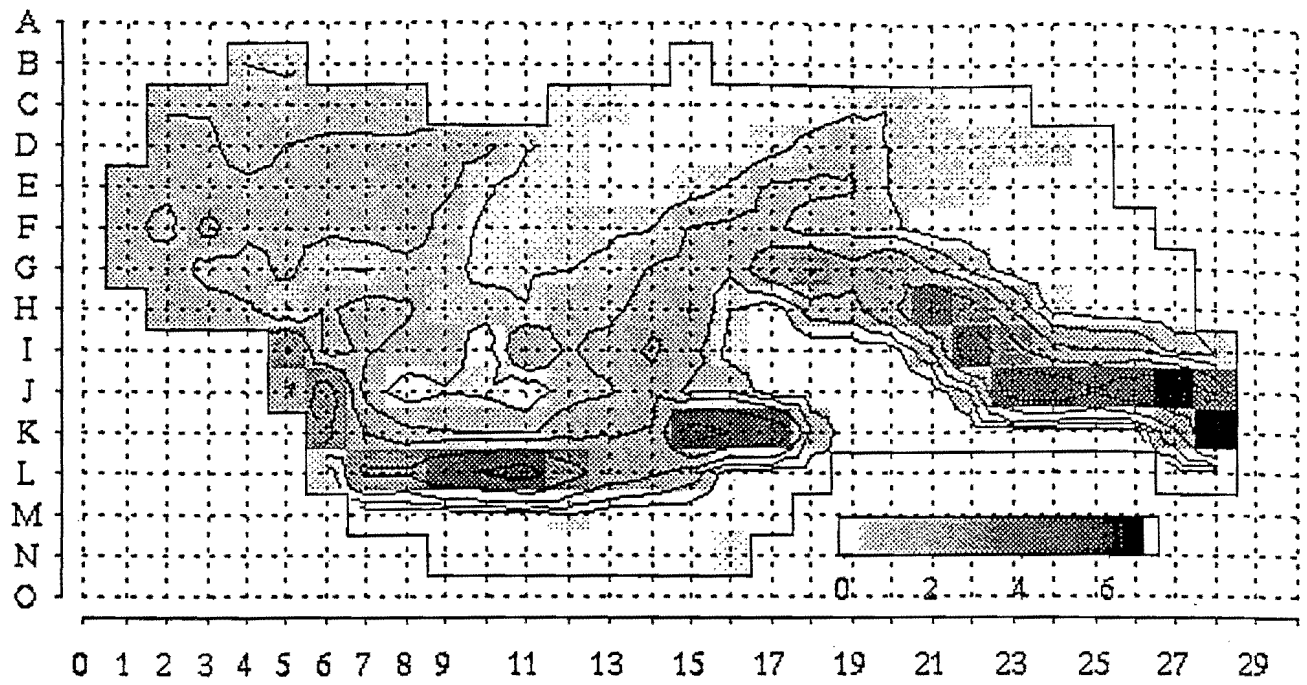


Figure 7.

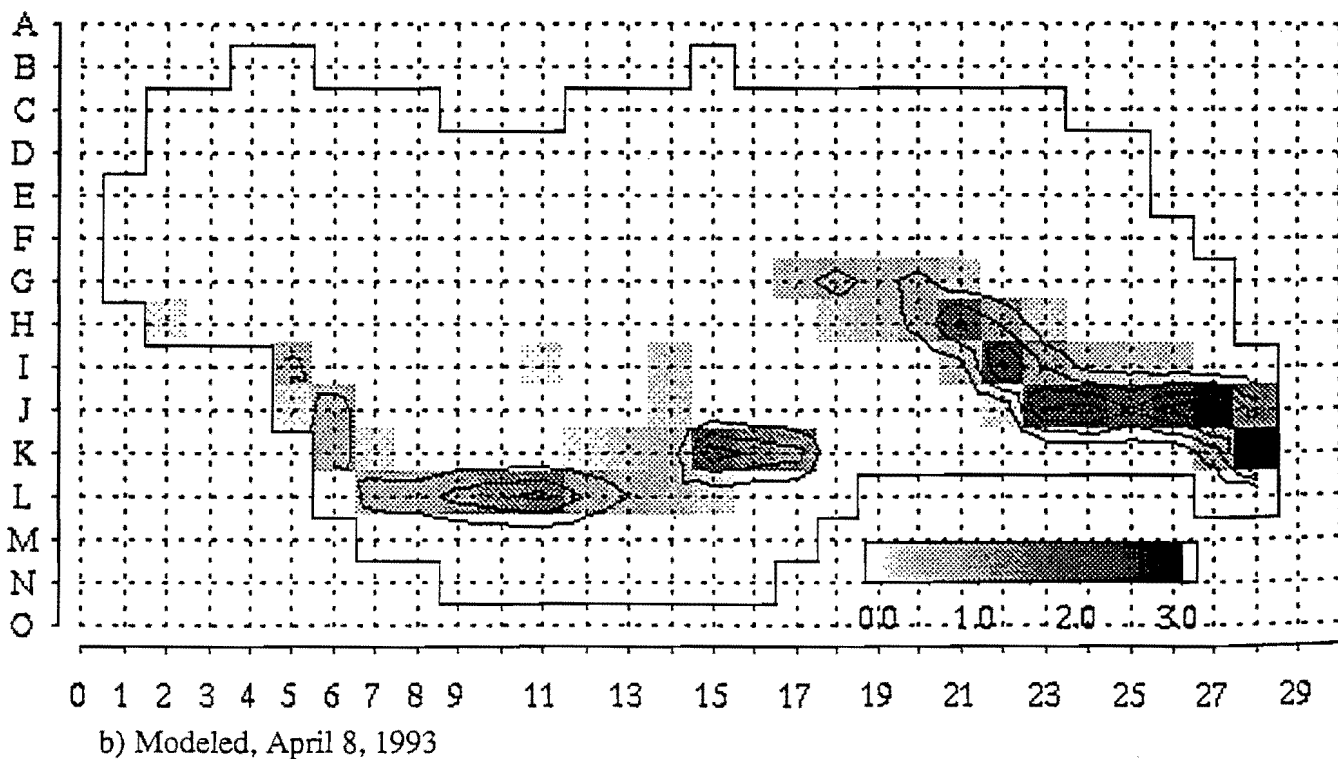
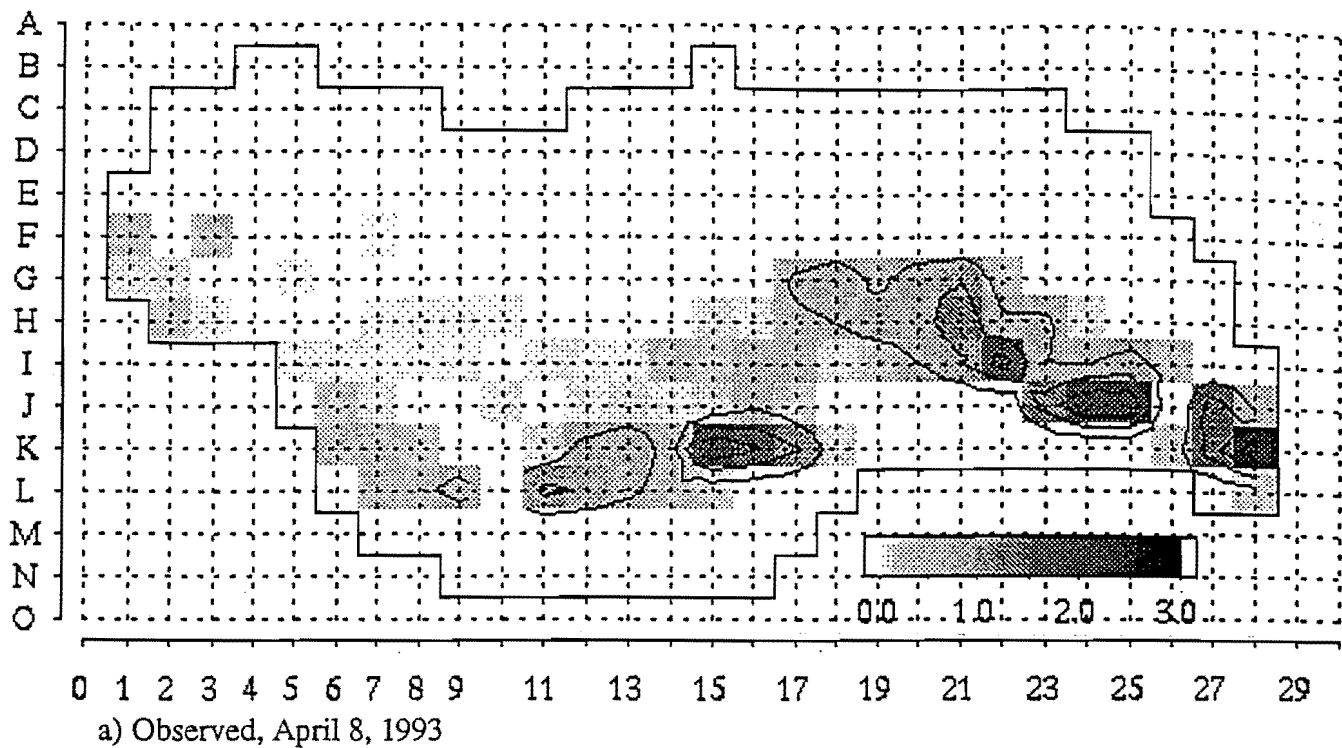


Figure 8.

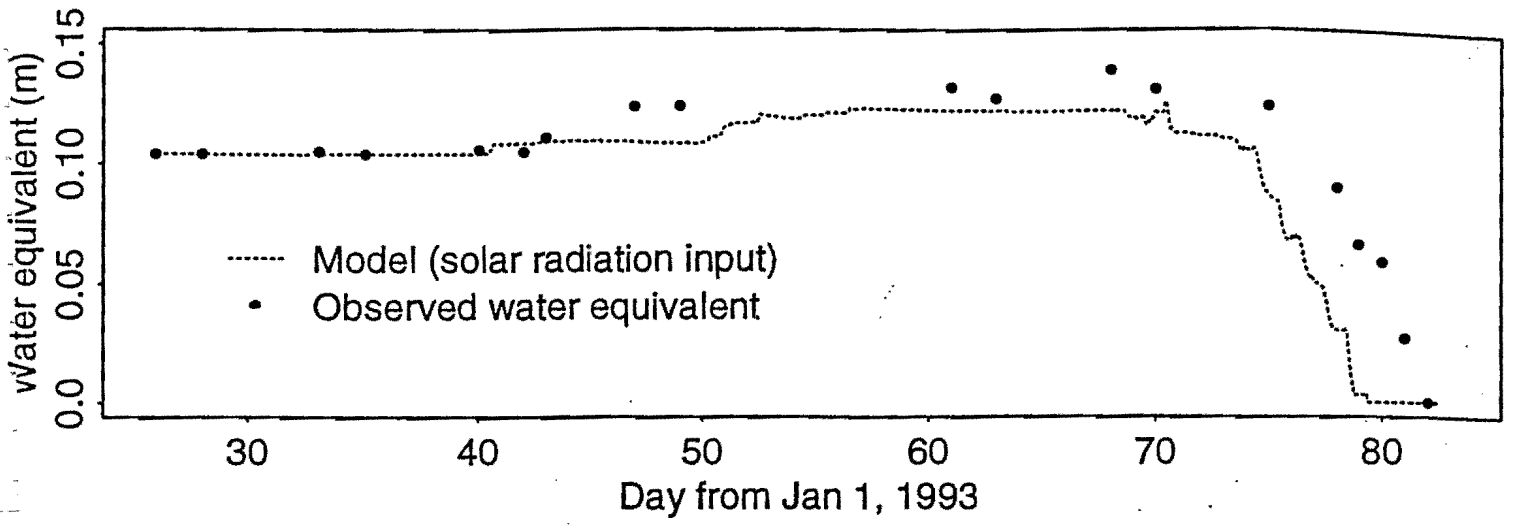


Figure 9.



# APPENDIX 4A

## Multivariate Nonparametric Resampling Scheme for Generation of Daily Weather Variables

Balaji Rajagopalan, Upmanu Lall, David G. Tarboton and David S. Bowles

### Abstract

A nonparametric resampling technique for generating daily weather variables at a site is presented. The method samples the original data with replacement while smoothing the empirical conditional distribution function. The technique can be thought of as a smoothed conditional Bootstrap and is equivalent to simulation from a kernel density estimate of the multivariate conditional probability density function. This improves on the classical Bootstrap technique by generating values that have not occurred exactly in the original sample and by alleviating the reproduction of fine spurious details in the data. Precipitation is generated from the nonparametric wet/dry spell model as described in Lall et. al. (1995). A vector of other variables (solar radiation, maximum temperature, minimum temperature, average dew point temperature and average wind speed) is then simulated by conditioning on the vector of these variables on the preceding day and the precipitation amount on the day of interest. An application of the resampling scheme with 30 years of daily weather data at Salt Lake City, Utah, USA is provided.

# 1. INTRODUCTION

Daily weather variations influence agricultural and engineering management decisions. Crop yields and hydrological processes such as runoff and erosion are very sensitive to weather. Recognizing the inherent variability in climate, it is often necessary to assess management scenarios for a number of likely input sequences. Stochastic models are consequently useful for simulating weather scenarios. Such models need to simulate sequences that are representative of the data. While there is a substantial literature for rainfall simulation and for other variables one at a time, only a few "multivariate" models have been developed.

In this paper we develop and exemplify nonparametric procedures for resampling a vector of daily weather variables, such that selected lag 0 and lag 1 dependence characteristics are preserved. Dependence is defined in terms of joint or conditional probabilities, rather than correlation.

This work is an off shoot of the ongoing Water Erosion Prediction Project (WEPP) of the United States Department of Agriculture (USDA). WEPP, is a key model for soil and forest conservation studies. WEPP, includes a Climate Generator (CLIGEN) and the work presented here intends to improve it. Hillslope erosion is driven largely by precipitation and a suite of other weather variables. Hence, the main objective is to generate weather sequences which will be used by WEPP to estimate hillslope erosion. In this study, we chose a set of five daily variables (Solar Radiation (SRAD), Maximum temperature (TMX), Minimum temperature (TMN), Avg. Wind speed (WSPD) and Avg. Dew point temperature (DPT) in addition to Precipitation (P), that are of interest for erosion prediction. Most of these weather variables are sensitive to precipitation. Solar radiation, dew point temperature, maximum temperature and minimum temperature are more likely to be below normal on rainy days than on dry days, while the wind speed may be above normal on rainy days than on dry days. Consequently precipitation is chosen as the driving variable of the models developed so far. Typically (see Jones et al. 1972, Nicks and Harp 1980, Richardson 1980), daily precipitation is generated independently and the other variables are generated by conditioning on precipitation events (i.e. whether a day is wet or dry).

Throughout this paper we denote the historical time series of the five weather variables chosen above as  $[z]_{mkj}$  ( $m=1,\dots,NY$ ,  $k = 1,\dots,366$ ,  $j=1,\dots,NV$ ), where  $NY$  is the number of years of record,  $NV(=5)$  is the number of variables considered (SRAD, TMX, TMN, DPT and WSPD). Further, define  $[\bar{Z}]_{kj}$  and  $[STD]_{kj}$  as the corresponding mean and standard deviation vector for each calendar day  $k$  ( $k=1,\dots,366$ ) of each variable  $j$  ( $j=1,\dots,5$ ). The historical time series of the precipitation is denoted as  $[P]_{mk}$ .

We now discuss key attributes of some strategies for resampling or synthesizing vectors of these variables.

## *1.1 Resampling Approaches*

Multivariate stochastic simulation of weather variables has not been studied as extensively as streamflow or precipitation. Two broad approaches that are possible are:

1. Parametric
2. Nonparametric - Bootstrap (Raw, Conditional and Smoothed)

### *1.1.1 Parametric*

The parametric approach is the traditional method (see Jones et al., 1972, Bruhn et al. 1980, Nicks and Harp 1980, Lane and Nearing 1989 and Richardson 1980) for stochastic daily weather simulations. Figure 1 summarizes the general structure of the parametric approaches. The general strategy is to generate precipitation independently and the other variables conditioned on the status of precipitation (i.e. rain or no rain on the day). The other variables are generated from either independently fitted statistical distributions to each of the variables and separately for each of the two precipitation states (i.e. rain, no rain), or independently or jointly fitted auto regressive models of order 1 (AR-1) to the variables.

Usually the year is divided into periods (seasons) and moments (mean standard deviation etc.) are calculated for each variable for each period for each precipitation state. The moments are used to fit statistical distributions or models. Dividing the year into various periods assumes homogeneity within each period and offers a treatment of seasonality. Jones et al. (1972), Bruhn et al. (1980), Nicks and Harp (1980) and CLIGEN (Lane and Nearing, 1989) divide the year into 14 day and one month periods respectively in their works. Richardson (1980) adopted a method, wherein the means and standard deviations of each periods and each precipitation state are smoothed using Fourier series. The smoothed daily values of the means and standard deviations are subsequently used for deseasonalization.

Daily Precipitation is typically generated from a fitted first order Markov Chain for precipitation occurrence and by sampling from the distribution (such as Gamma, Exponential, Truncated Normal etc.) fitted for the daily amounts for each period.

One approach to generate the other variables is to fit distributions independently for each variable for each period and for each precipitation state. Here, the simulations are made under the assumption that each variable is independent and identically distributed (i.i.d). This approach and its variants are used by Jones et al. (1972), Bruhn et al. (1980) and CLIGEN (Lane and Nearing, 1989). In CLIGEN each variable is assumed to be an independent Gaussian variable for each month, with parameters dependant on the precipitation state transition (e.g. wet to wet, dry to wet etc.). This

approach does not consider the dependence between the variables nor the serial dependence for each variable. Only the dependence on the precipitation state or the precipitation transition is considered.

Serial dependence was incorporated by Nicks and Harp (1980) who fit Auto Regressive models of order one (AR-1) independently to each variable for each period. Consideration of dependence across variables is added by Richardson (1981) who used a Multivariate Auto Regressive model of order one (MAR-1). When the cross dependence terms are neglected in MAR-1, it reduces to an AR-1 process. These AR models suffer from the drawback of assuming the data to be normally distributed. As a result only linear dependence can be reproduced. In practice data may not be normally distributed. Transformation of the data to be multivariate normal may be difficult and may lead to biased statistics upon back transforming to the original space.

The parametric approaches discussed have four main drawbacks, which are (i) Choice of a model (a statistical distribution or the order of an AR or MAR model) is often subjective and rarely formally tested on a site by site basis (ii) Reliance on an implicit Gaussian framework (e.g. AR or MAR) which preserves only linear dependence and is not appropriate for bounded variables (iii) The fitted models have limited portability in the sense that procedures/distributions used at one site may not be best at other sites. The last point is important where an agency wishes to prescribe a uniform procedure over its domain.

### *1.1.2 Nonparametric*

Nonparametric techniques do not require pre-selected distributions or models to be fit to data. The Bootstrap (or Raw Bootstrap) is a nonparametric technique introduced by Efron (1979). It is often used for constructing a confidence region, attaching a standard error to an estimate, carrying out a test of a hypothesis, or estimating the sampling distribution of some statistic. Historical data is resampled with replacement. Since it is the same data, the simulations by construction have the same distributional properties as that of the historical data. Since each resampled observation is drawn independently, serial dependence is not preserved. Serial dependence can be accommodated by using the '*block-resampling scheme*' (a Conditional Bootstrap) developed by Kunsch (1989) and Liu and Singh (1992). Here a block of 'k' observations is resampled as opposed to a single observation in the Bootstrap. Serial dependence is preserved within, but not across a block. The block length 'k' determines the order of the serial dependence that can be preserved.

A property of the Bootstrap technique is that the simulated samples will only have values that have occurred in the historical data and consequently the simulations are restricted to the historical set of values. Silverman (1986, p. 142) points out that this behaviour may reproduce spurious fine structure in the original data. This is not a desirable feature while applying the technique to simulation of daily weather variables, where we may wish to have simulated values that have not been observed

in the historical data and may be also beyond the maximum/minimum of the observed data. This problem can be alleviated by using 'Smoothed Bootstrap'.

In the Smoothed Bootstrap (Silverman 1987, p 144), each observation  $y_i$  ( $i=1,\dots,n$ ) is considered to be representative of a region  $(y_i-h, y_i+h)$  around it. The extent of this region  $h$  is called the bandwidth and is determined from the data. Intuitively, it is desirable to resample such that the maximum weightage is given to the observation  $y_i$  and weights decrease when moving towards  $y_i-h$  or  $y_i+h$ . This is accomplished by having a weight function centered at each observation. The weight function is usually chosen to be a valid probability density function, such as the Gaussian  $(N(0,1))$ . The simulation proceeds by picking an observation  $y_i$  with replacement from  $\{y_1,\dots,y_n\}$  and then generating a value from  $N(y_i, h)$  with  $h$  specified. Formally, the Smoothed Bootstrap is equivalent to resampling from a kernel density estimate (k.d.e). Kernel density estimation is a nonparametric procedure described in section 2.3.

In this paper, we develop a Smoothed Conditional Bootstrap that considers multivariate and serial dependence amongst the variables of interest. Hereafter, we refer to the scheme presented as the NP model. We first provide the motivation and main ideas of the model. The simulation algorithm is outlined next. The utility of the model is then illustrated through application to daily weather data at Salt Lake City, Utah, USA.

## 2. MAIN IDEAS OF THE NP MODEL

Our goal is to develop an approach that is driven directly by the observed data with reasonable assumptions, is easy to implement, is readily transferable from site to site and captures the relative frequencies of the data in a natural manner. We do this by defining the appropriate probability densities that we need to resample from and then discuss their estimation.

### 2.1 Overview of the NP model

A conceptual flow chart of the model is shown in Figure 2. The historical data of the other weather variables is standardized as  $[x]_{lkj} = ([z]_{lkj} - [\bar{Z}]_k) / [STD]_{kj}$  where  $l, k$  and  $j$  are the same as defined in section 1. This removes the seasonality present in each variable. Precipitation for day 't' ( $P_t$ ) is generated from the wet/dry spell model as described in Lall et al. (1995) briefly summarized in section 2.2. However, the user can generate the daily precipitation from his favourite model.

In the NP model the year is divided into four periods or seasons (for the Salt Lake City example these are Season 1 (Jan-mar), Season 2 (Apr-Jun), Season 3 (Jul-Sep), Season 4 (Oct-Dec)). Simulations for days in any particular period are made using the historical data of that period. Subsequently, the comparison between the simulations and the historical data are also made the same

scale. One could choose different periods (e.g. monthly, weekly etc.). We chose the above four periods so as to be consistent with the wet/dry spell model (Lall et al., 1993) for daily precipitation.

The aim of the model is to capture the day-to-day dependence present between the variables. The standardized vector of variables  $\mathbf{x}_t$  for any day 't' is simulated from the multivariate conditional p.d.f  $f(\mathbf{x}_t | \mathbf{V}_t = \mathbf{V}^*)$ . Where,  $\mathbf{x}_t$  = standardized vector of [SRAD, TMX, TMN, WSPD, DPT]<sub>t</sub> that is to be generated for day t,  $P_t$  is the generated precipitation for day t from the wet/dry spell model;  $\mathbf{x}_{t-1}$  = standardized vector of [SRAD, TMX, TMN, WSPD, DPT]<sub>t-1</sub> already generated for day t-1,  $\mathbf{V}^* = [\mathbf{x}_{t-1}, P_t]$  is the conditioning vector, 'd (=5)' is the number of variables to be generated, 'd' (=6) is the number of conditioning variables and  $dg=d+d'$ .

The conditional density  $f(\mathbf{x}_t | \mathbf{V}_t = \mathbf{V}^*)$  is defined as,

$$f(\mathbf{x}_t | \mathbf{V}_t = \mathbf{V}^*) = \frac{f(\mathbf{x}_t, \mathbf{V})}{\int f(\mathbf{x}_t, \mathbf{V}) d\mathbf{x}_t} = \frac{f(\mathbf{x}_t, \mathbf{V} = \mathbf{V}^*)}{f_V(\mathbf{V} = \mathbf{V}^*)} \quad (1)$$

where  $f_V(\mathbf{V} = \mathbf{V}^*)$  is the marginal density of  $\mathbf{V}$ , evaluated at a current vector  $\mathbf{V}^*$ . The standardized sequences  $\mathbf{x}_t$  are transformed to  $\mathbf{z}_t = \mathbf{x}_t * [\mathbf{STD}]_k + [\mathbf{\bar{Z}}]_k$ , where k is the calendar day associated with day 't'. Thus, the key idea here is the estimation of this conditional probability density function from the historical data using nonparametric density estimators (kernel estimators) and subsequently simulating or bootstrapping from it. The mechanism of kernel density estimation is described in section 2.3, and the algorithm for simulation from a conditional p.d.f (as in Equation 1) using kernel density estimators is developed and outlined in section 3.

## 2.2 Precipitation Model

The seasonal wet/dry spell model for daily precipitation described fully in Lall et al (1995) has three random variables - wet spell length,  $L_w$  days, dry spell length,  $L_d$  days, and wet day precipitation amount, P inches. The periods(seasons) are as defined in the previous section. Variables wsp and dsp are defined through the set of integers between 1 and the season length, and P is defined as a continuous, positive random variable. A mixed set of discrete and continuous random variables is thus considered. Successive wet day's precipitation amount is taken to be independent and the precipitation is independent of the wet spell length ( $L_w$ ). Correlation statistics computed for the data sets analyzed supported these assumptions.

The p.d.f.'s of wet day precipitation amount  $f(P)$  and the probability mass functions (p.m.f.'s) of wet spell length  $f(L_w)$  and dry spell length  $f(L_d)$  are estimated for each season using kernel density estimators.

A dry spell is first generated using  $f(L_d)$ . Then a wet spell is generated using  $f(L_w)$ . Precipitation for each of the ' $L_w$ ' wet days is then generated from  $f(P)$ . The process is repeated with the generation of another dry spell. If a season boundary is crossed, the p.d.f.'s used for generation are switched to those for the new season. This procedure continues until a synthetic sequence of the desired length has been generated. The p.d.fs  $f(L_w)$ ,  $f(L_d)$  and  $f(P)$  are estimated using kernel density estimators as described in Lall et al. (1993) and Rajagopalan et al. (1995). At this point the kernel density estimation is generically described and the estimators used in this work are outlined below.

### 2.3 Kernel Density Estimation

The kernel density estimator generalizes the frequency histogram as an estimator of the p.d.f. While the histogram is capable of showing some features of the data, it has several drawbacks. It is difficult to manipulate analytically, it is not easy to visualize for multivariate situations, and it allows for no extrapolation beyond the data. The indicated frequency distribution is sensitive to the class width, as well as the origin of each class. Silverman (1986, p.9-11) illustrates these problems graphically. One can improve the histogram by centering rectangular boxes at each observation (to gain independence from choice of origin). A kernel density estimator, introduced by Rosenblatt (1956), is formed by centering a smooth kernel function at each observation. Kernel density estimators for univariate continuous variables, univariate discrete variables and multivariate continuous variables are now defined.

#### 2.3.1 Univariate Continuous Variables

We stated earlier that the Smoothed Bootstrap is equivalent to sampling from a kernel density estimate. The kernel density estimator for a continuous variable (such as the wet day precipitation  $P$ ) is defined as

$$f(P) = \sum_{i=1}^n \frac{1}{nh} K\left(\frac{P-P_i}{h}\right) \quad (2)$$

where  $K(\cdot)$  is a kernel function centered on the observation  $P_i$ , and can be any valid probability density function and  $h$  is a bandwidth. The bandwidth  $h$  controls the amount of smoothing of the data in the density estimate. An estimator with constant bandwidth  $h$  (like in Equation 2) is called a fixed kernel estimator. Commonly used kernels are:

Gaussian Kernel 
$$K(t) = (2\pi)^{-1/2} e^{-t^2/2} \quad (3a)$$

$$\text{Epanechnikov Kernel} \quad K(t) = 0.75 (1 - t^2) \quad |t| \leq 1 \quad (3b)$$

$$\text{Bisquare Kernel} \quad K(t) = (15/16) (1 - t^2)^2 \quad |t| \leq 1 \quad (3c)$$

An evaluation of  $K(\cdot)$  represents the weight given to the observation  $P_i$  that is based on distance between  $P$ , and  $P_i$ . One can see from Equation 2, that the kernel estimator is a convolution estimator that forms a local weighted average of the relative frequency of observations in the neighborhood of the point of estimate. The kernel function,  $K(\cdot)$  prescribes the relative weights,  $h$  prescribes the range of data values over which the average is computed. This is illustrated in Figure 3.

The p.d.f of wet day precipitation  $f(P)$  is obtained by applying, a kernel density estimator to log transformed data. Note that most of the data of wet day precipitation is concentrated near the lower boundary (i.e. 0.), as a result the p.d.f estimates using the kernel estimators are highly biased due to the boundary problem. The log transformation on such heavily skewed data alleviates the boundary problem. The resulting estimator is given as:

$$f(P) = \frac{1}{P} f(\log(P)) = \sum_{i=1}^n \frac{1}{nhP} K\left(\frac{\ln(P) - \ln(P_i)}{h}\right) \quad (4)$$

The Epanechnikov kernel is used and the bandwidth  $h$  is chosen for the log transformed data using the recursive approach of Sheather and Jones (1991) to minimize the Mean Integrated Square Error (MISE) of estimate of  $f(\log(P))$ .

Note that no assumptions regarding the parent density of  $P$  have been made thus far. We need to specify only the bandwidth  $h$  and the kernel function. Silverman (1986) points out that the kernel density estimator is more sensitive to the choice of the bandwidth than to that of the kernel.

### 2.3.2 Univariate Discrete Variables

In this section, we present procedures for the estimation of the univariate probability mass functions for discrete variables (such as wet spell lengths  $w$ , dry spell lengths  $d$ ). We recommend the Discrete Kernel (DK) estimator developed in Rajagopalan and Lall (1995). The DK estimator for the p.m.f.  $\hat{f}(L)$ , where  $L$  is either  $w$  or  $d$ , and  $n$  is the corresponding sample size is given as:

$$\hat{f}(L) = \sum_{j=1}^{L_{\max}} K_d\left(\frac{L-j}{h}\right) \tilde{\alpha}_j \quad (5)$$



where  $\tilde{\alpha}_j$  is the sample relative frequency ( $n_j/n$ ) of spell length  $j$ ,  $n_j$  is the number of spells of length  $j$ ,  $L_{\max}$  is the maximum observed spell length (note that  $\sum_{j=1}^{L_{\max}} \tilde{\alpha}_j = 1$ ),  $K_d(\cdot)$  is a discrete kernel function, and  $L$ ,  $j$  and  $h$  are positive integers. The kernel function  $K_d(\cdot)$  is given as:

$$K_d(t) = at_j^2 + b \quad \text{for } |t| \leq 1 \quad (6)$$

The expressions for  $a$  and  $b$  for the interior of the domain,  $L > h+1$  and the boundary region  $L < h$  are developed in Rajagopalan and Lall (1995).

The bandwidth  $h$  is estimated by minimizing a Least Squares Cross Validation (LSCV) function given as,

$$\text{LSCV}(h) = \sum_{j=1}^{L_{\max}} (\hat{f}(j))^2 - 2 \sum_{j=1}^{L_{\max}} \hat{f}_{-j}(j) \tilde{\alpha}_j \quad (7)$$

where,  $\hat{f}_{-j}(j)$  is the estimate of the p.m.f of spell length  $j$ , formed by dropping all the spells of length  $j$  from the data. This method has been shown by Hall and Titterington (1987) to automatically adapt the estimator to an extreme range of sparseness types. Monte Carlo results showing the effectiveness of the DK estimator with bandwidth selected by LSCV are presented in Rajagopalan and Lall (1995).

### 2.3.3 Multivariate Continuous Variables

Extending the idea of the kernel density estimator for univariate continuous variables, a kernel density estimate of the multivariate p.d.f of a vector  $\mathbf{y}$  is defined as (Silverman, 1986, p. 76-78):

$$f(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n K(\mathbf{u}) \quad (8)$$

where  $\mathbf{u} = \frac{(\mathbf{y} - \mathbf{y}_i)^T \mathbf{S}^{-1} (\mathbf{y} - \mathbf{y}_i)}{h^2}$ , and  $K(\mathbf{u})$  is a multivariate kernel function.  $\mathbf{y} = [y_1, y_2, \dots, y_d]^T$  denotes the  $d$  dimensional random vector whose density is being estimated with,  $\mathbf{y}_i = [y_{1i}, y_{2i}, \dots, y_{di}]^T$   $i = 1$  to  $n$  the sample values of  $\mathbf{y}$ ;  $n$  is the number of sample vectors;  $h$  the kernel bandwidth and  $\mathbf{S}$  the covariance matrix. Here we use a Gaussian kernel function given as,

$$K(u) = \frac{1}{(2\pi)^{d/2} \det(S)^{1/2} h^d} \exp(-u/2) \quad (9)$$

Just as in the univariate case described in section 2.3.1,  $K(u)$  represents the weight given to an observation  $y_i$  that is based on distance between  $y$ , and  $y_i$ . The distance used here is the Euclidean distance modified to recognize the covariance of the  $y$ . It can be seen that the estimator in Equation 8, is similar to the univariate estimator in Equation 2 since it is a local weighted average of the relative frequency of observations in the neighborhood of the point of estimate. Here too the kernel function,  $K(\cdot)$  prescribes the relative weights,  $h$  prescribes the range of data values over which the average is computed and the covariance  $S$  provides the orientation.

Here we choose the bandwidth as the one that minimizes mean integrated square error in  $f(y)$  if the underlying distribution is assumed to be multivariate Gaussian. Silverman (1986, p 86-87) gives an appropriate  $h$  to use for a multivariate Gaussian p.d.f. using the Gaussian kernel as,

$$h = \{(4/(2d+1))^{1/(d+4)}\} n^{-1/(d+4)} \quad (10)$$

Here  $n$  is the number of observations and  $d$  is the dimension. Note that  $h \rightarrow 0$  as  $n \rightarrow \infty$  so that the kernel density estimate is consistent. However, as the dimension  $d$  increases  $h$  also increases. This is because in higher dimensions large regions of high density may be completely devoid of observations in a sample of moderate size. The bandwidth in such a situation has to be bigger to cover large regions. The above choice of bandwidth, is optimal for p.d.fs that are near Gaussian and is an adequate choice for many cases (Silverman, 1986, p 45-48). Cross Validation or Plug in methods could be used here to choose  $h$  as in the wet/dry spell model. However, this increases the computational burden substantially. Recall that the parametric approaches often assume a Gaussian distribution. In a Bayesian context using this bandwidth can be thought of as developing a posterior kernel density estimate with a Gaussian prior. The resulting tail behaviour and degree of smoothing supplied will be consistent with an underlying Gaussian p.d.f, with some adaption to local features.

An attractive feature of kernel estimators of the p.d.f is that they are local (use only a neighborhood around the point of estimate) and hence are not overly effected by outliers. Since they make no prior assumptions of the underlying probability density function, they are data driven and robust and are portable across sites/data sets.

In the bootstrap context we have a region that each observation  $y_i$  represents. The orientation and shape of the region is given by the scaling factor  $hS$  and the kernel function  $K(u)$ . Resampling from the kernel density estimate entails picking a point  $y_i$  uniformly in  $[y_1, \dots, y_n]$  and then simulating from the kernel  $K(u)$ , i.e.  $N(y_i, h^2S)$ . We extend this approach formally for simulation from a

multivariate conditional p.d.f in the following section. For details on kernel density estimation refer to Silverman (1986) and Scott (1992).

### 3. Kernel Density Estimation of Multivariate Conditional p.d.f

For the simulation of interest here an estimate of the conditional p.d.f  $f(\mathbf{x}_t | \mathbf{V}_t = \mathbf{V}^*)$  is needed. The strategy used here is similar to the one used by Tarboton et al (1993) for streamflow simulation. Applying the estimator in Equation 8 to the conditional p.d.f in Equation 1 with sample vectors  $\mathbf{x}_i = [\mathbf{x}_t, \mathbf{x}_{t-1}, P_t]_i$  denoted as  $[\mathbf{x}_i, \mathbf{V}_i]$  we get:

$$f(\mathbf{x}_t | \mathbf{V}_t = \mathbf{V}^*) = \frac{1}{nh^d} \frac{1}{f_v(\mathbf{V}^*)} \sum_{i=1}^n \frac{1}{\det(S)^{1/2}} K\left(\frac{[\mathbf{x}_t - \mathbf{x}_i; (\mathbf{V}^* - \mathbf{V}_i)^T] S^{-1} \begin{bmatrix} \mathbf{x}_t - \mathbf{x}_i \\ \mathbf{V}^* - \mathbf{V}_i \end{bmatrix}}{h^2}\right) \quad (11)$$

where  $S$  is the  $dg$  by  $dg$  covariance matrix of the vector  $(\mathbf{x}_i, \mathbf{V}_i)$  estimated from historical data. Let the matrix  $S$  be partitioned as,

$$S = \begin{bmatrix} S_x & S_{xv}^T \\ S_{xv} & S_v \end{bmatrix} \quad (12)$$

where  $S_x$  is the  $d$  by  $d$  covariance matrix of  $\mathbf{x}$ ,  $S_v$  is the  $d'$  by  $d'$  covariance matrix of  $\mathbf{V}$  and  $S_{xv}$  the  $d$  by  $d'$  cross covariance between  $\mathbf{x}$  and  $\mathbf{V}$ . Using the Gaussian kernel function (i.e. Equation 9) Equation (11) can be reduced to a weighted sum of Gaussian functions,

$$f(\mathbf{x}_t | \mathbf{V}_t = \mathbf{V}^*) = \sum_{i=1}^n w_i N(\mathbf{b}_i, \mathbf{c}_i) \quad (13)$$

where,

$$w_i = w'_i / \sum_{i=1}^n w'_i, \quad w'_i = \exp(-a_i/2); \quad a_i = \frac{([\mathbf{V}^* - \mathbf{V}_i]^T [S_v]^{-1} [\mathbf{V}^* - \mathbf{V}_i])}{h^2}, \quad (14)$$

$$\mathbf{b}_i = \mathbf{x}_i + ([\mathbf{V}^* - \mathbf{V}_i]^T [S_v]^{-1} [S_{xv}]); \quad \mathbf{c} = h^2 (S_x - S_{xv}^T S_v^{-1} S_{xv}) \quad (15)$$

Note that  $\sum_{i=1}^n w_i = 1$

From Equation (13) we see that the conditional p.d.f reduces to a weighted sum of Gaussian functions. It can be thought of as a slice through a multivariate density function, estimated as a weighted sum of slices with the same orientation through the kernels placed on each observation.

Simulation from the conditional p.d.f can be achieved by picking a point  $\mathbf{x}_i$  with probability  $w_i$ , then sampling from  $N(\mathbf{b}_i, \mathbf{c})$ .

### 3.1 NP Simulation algorithm

The simulation proceeds as:

- I. Simulate Precipitation for all the days of the year from the wet/dry spell model.
- II. Estimate the model parameters (bandwidth  $h$  and the covariance matrix  $S$ ) from the data for each season.
- III. At the start of each period of interest initialize  $t=0$ ,  $\mathbf{x}_t = [\mathbf{0}]$
- IV. Generate  $\mathbf{x}_t$  sequentially (day by day) from  $f(\mathbf{x}_t | \mathbf{V}_t)$ , where the conditioning vector  $\mathbf{V}_t$  consists of the previous day's vector  $\mathbf{x}_{t-1}$  and the current day's generated precipitation  $P_t$  (i.e.  $\mathbf{V}_t = [\mathbf{x}_{t-1}, P_t]$ ) as:
  1. Estimate weights ( $w_i$ ) associated with each data point ( $\mathbf{x}_i$ ) ( Equation 14)
  2. Resample an index  $i$  using  $w_i$  ( $i = 1, \dots, n$ ) as probabilities.
  3. Estimate the conditional mean ( $\mathbf{b}_i$ ) and conditional variance ( $\mathbf{c}_i$ ) using the picked point  $\mathbf{x}_i$  and  $\mathbf{V}_t$  (Equation 15)
  4. Generate vector  $\mathbf{x}_t = \mathbf{b}_i + \boldsymbol{\varepsilon}$  where  $\boldsymbol{\varepsilon}$  is from a multivariate normal distribution with mean  $\mathbf{0}$  and variance  $\mathbf{c}$  (following Devroye, 1986, p. 565)
  5. Recover  $\mathbf{z}_t$  as  $\mathbf{z}_t = \mathbf{x}_t * [\mathbf{STD}]_k + [\bar{\mathbf{X}}]_k$  where  $k$  is the calendar day corresponding to day  $t$ .
- V. At the start of a new simulation go to III.

## 4. MODEL APPLICATION AND PERFORMANCE MEASURES

To demonstrate the utility of the resampling model for generation of daily weather variables, the model was applied to daily weather data from the station Salt Lake City in Utah. Thirty years of daily weather data was available from the period 1961-1991. Salt Lake City is at 40°46' N latitude, 111° 58' W longitude and at an elevation of 1288 m. Most of the precipitation comes in the form of winter snow. Rainfall occurs mainly in Spring, with some in Fall.

We shall first outline the experimental design and then some measures of performance used to judge the utility of the model.

### *4.1 Experiment design*

Our purpose here is to test the utility of the NP generation scheme. The main steps involved in accomplishing this are:

1. Daily precipitation is generated from the wet/dry spell model.
2. The other variables are generated following the simulation algorithm described in section 3.1
3. Twenty five synthetic records of thirty years each (i.e. the historical record length) are simulated using the NP model.
4. The statistics of interest, described below are computed for each simulated record, by each period and compared to statistics of the historical record using boxplots.

### *4.2 Performance measures*

The following statistics were considered to be of interest in comparing the historical record and the NP simulated record of other weather variables.

Moments:

1. Mean of each variable for each season.
2. Standard deviation of each variable for each season.
3. Skew of each variable for each season.
4. Co-efficient of variation of each variable for each season.

Relative Frequencies:

5. 25% quantile of each variable for each season.
6. 75% quantile of each variable for each season.

Dependence:

7. Cross correlation on any given day between the variables for each season.
8. Lag-1 daily Cross correlation between the variables for each season.

9. Lag-1 daily correlation of each variable for each season.

## 5. RESULTS

The statistics of interest calculated from the simulations are compared with those for the historical record using boxplots. A box in the boxplots (e.g. Figure 4) indicates the interquartile range of the statistic computed from twenty five simulations, the line in the middle of the box indicates the median simulated value. The solid lines correspond to the statistic of the historical record. The boxplots show the range of variation in the statistics from the simulations and also show the capability of the simulations to reproduce historical statistics.

Figures 4 through 8 show the boxplots of moments and relative frequency measures of Solar Radiation, Maximum Temperature, Minimum Temperature, Average Wind Speed and Average Dew Point Temperature respectively. It can be seen that the historical values of mean, and the quantiles are well reproduced, while standard deviation, coefficient of skew and coefficient of variation are not quite well reproduced. This is to be expected as the kernel methods inflate the variance by a factor equal to  $(1+h^2)$  (see Silverman 1986, p. 143) which in turn effects the skew and the coefficient of variation. It may be desirable to have to have a slight increase in the variance of the simulations as compared to that of the historical.

Illustrative statistics of wet spell lengths, dry spell lengths and wet day precipitation for the simulations from the wet/dry spell model are also estimated and are shown in Figures 9,10 and 11 respectively. Figure 9 shows the boxplots of average wet spell length, standard deviation of wet spell length, fraction of wet days and length of longest wet spell length for each season. Figure 10 shows the boxplots of these statistics of the dry spell length. Figure 11 shows the boxplots of average wet day precipitation, standard deviation of wet day precipitation, percentage of yearly precipitation in each season. The boxplots in Figures 9, 10 and 11 show that the historical statistics are reproduced well by the simulations.

Figures 12 and 13 show the boxplots of the lag-0 cross correlation and lag-1 cross correlation between the variables. Figure 14 shows the lag-1 auto correlation of each variable for each of the four seasons. The correlations from the simulations and the historical correlations seem to be different in a number of cases.

One reason for this mismatch of the correlations is that the precipitation is supplied externally from the wet/dry spell model. As a result the covariance between  $x_{t-1}$  and  $P_t$  need not correspond to

that of the historical covariance between between them. This introduces a bias in the conditioning plane from which  $\mathbf{x}_t$  is generated and results in a mismatch of the correlations. One way to get around this problem is to generate the precipitation also in the multivariate model, i.e. simulate  $\mathbf{x}_t$  from  $f(\mathbf{x}_t | \mathbf{x}_{t-1})$  where both  $\mathbf{x}_t$  and  $\mathbf{x}_{t-1}$  are of dimension 6. This should reproduce the correlations statistics. However, negative values for precipitation may then be simulated. Since most of the precipitation is concentrated near 0., simulating precipitation also along with the other variables may lead to oversmoothing of the mode of the precipitation density.

Another reason, could be that there are two different process which are the one with zero precipitation and one without. Since there are two different processes with different correlation structure, the combined correlations need not match. A multivariate autoregressive model of lag 1 with precipitation supplied exogenously (like in the NP model described here), that can be thought of as a counterpart in the parametric frame work, will also suffer from the correlation mismatch.

## 6. SUMMARY AND CONCLUSIONS

A multivariate nonparametric model NP that aims at capturing dependence upto lag-1 was presented and illustrated. The simulations are made from the conditional p.d.f estimated from the data using kernel density estimators. The kernel estimators being local average estimators of the target function, have the advantage of readily admitting arbitrary probability densities without requiring that they be hypothesized or formally identified. Broader dependence structures can be consequently considered. The need to choose/justify and fit the best p.d.f is side stepped.

The bandwidth is the key parameter in the NP model, as it determines the degree of smoothness that will be imparted to the p.d.f. The larger the bandwidth the smoother the p.d.f and vice-versa. Choosing  $h$  automatically (using cross-validation or other approaches Scott 1992) from the data would be more appropriate than the choice used here. However, the additional variance in the choice of  $h$  induced by such an estimation process may detract from its use where the primary purpose is to resample the data. Bandwidth selection methods are undergoing continuous improvement. We expect to implement more formal selection procedures in due course. One could also use a local covariance matrix estimated at each data point using a few neighbors of that point (i.e.  $S_i$  instead of  $S$  in Equation 8). Tarboton et al (1993) use this method for streamflow simulation.

Another problem with simulations is the boundary effect. For the variables that are bounded (e.g. Solar Radiation and Precipitation), values that violate the bounds could be generated. Typically these are censored to the bound. This may introduce a bias in the simulations. Procedures to better address this problem in the kernel framework are described in Muller (1989) and Lall et al (1993).

We chose to apply the NP model on a seasonal time scale, because the precipitation model that was used to drive the NP model is a seasonal model. However, we checked the results of the seasonal NP model at monthly time scale, and found the performance to be similar (results are not presented here).

The NP model developed here underscores our growing conviction that nonparametric techniques have an important role to play in improving the synthesis of hydrologic time series. They can capture dependence structure present in the data, without imposing arbitrary distributional assumptions, and produce synthetic sequences that are statistically similar to the historic sequence. The idea of resampling the data with appropriate perturbation of each value while maintaining selected dependence characteristics (or data sequencing) is easy to accept as a practical matter. A Markovian interpretation of the NP model described here is apparent upon thinking about the manner in which the 1-step transition process works. The value to be simulated at the next time step can be thought of as a transition to any of the states within a bandwidth from the state of the current time. The conditional p.d.f can be viewed as approximation to the transition probabilities. Thus, the NP model can be seen as a 1-step Markov model with the transitions estimated nonparametrically.

We are working on improving the multivariate, nonparametric resampling scheme using nearest neighbor and similar methods.

#### ACKNOWLEDGEMENTS

Partial support of this work by the U.S. Forest Service under contract notes, INT-915550-RJVA and INT-92660-RJVA, Amend #1 is acknowledged.



## REFERENCES

- Bruhn, J.A., W.E., Fry, and G.W., Fick, Simulation of daily weather data using theoretical probability distributions, *Journal of Applied Meteorology*, 19(9), 1029-1036, 1980.
- Devroye, L., *Non-Uniform Random Variate Generation*. Springer-Verlag, New-York, 1986.
- Dong Jianping and J. Simonoff, The construction and properties of boundary kernels for sparse multinomials, *Journal of Computational and Graphical Statistics*, 3, 1-10, 1994.
- Efron, B., Bootstrap methods: Another look at the Jackknife, *Annals of Statistics*, 7, 1-26, 1979.
- Hall, P. and D.M. Titterton., On smoothing sparse multinomial data, *Australian Journal of Statistics*, 29(1), 19-37, 1987.
- Jones, W., Rex, R.C. and D.E. Threadgill, A simulated environmental model of temperature, evaporation, rainfall, and soil moisture, *Transactions of the ASAE*, 366-372, 1972.
- Kunsch, H.R., The Jackknife and the Bootstrap for General Stationary Observations, *Annals of Statistics*, 17, 1217-1241, 1989.
- Lall, U., Rajagopalan, B. and Tarboton, D.G. A nonparametric wet/dry spell model for resampling daily precipitation, *Submitted to Water Resour. Res.*, 1995.
- Lane, L.J. and M.A. Nearing, USDA - *Water Erosion Prediction Project: Hillslope Profile Model Documentation*, NSERL Report No.2, National Soil Erosion Research Laboratory, USDA-Agricultural Research Service, W. Lafayette, Indiana 47907, 1989.
- Liu R.Y. and K. Singh, *Using iid Bootstrap Inference for some Non-iid Models*, Preprint. Department of Statistics, Rutgers University, 1988.
- Nicks, A.D. and J.F. Harp, J.F., Stochastic generation of temperature and solar radiation data. *Journal of Hydrology*, 48, 1-7, 1980.
- Rajagopalan, B. U. Lall and D.G. Tarboton, Evaluation of kernel density estimation methods for daily precipitation resampling, *Submitted to Water Resour. Res.*, 1995.
- Richardson, C.W., Stochastic simulation of daily precipitation, temperature and solar radiation. *Water Resources Research*, 17(1), 182-190, 1981.
- Rosenblatt, M., Remarks on some nonparametric estimates of a density function. *Annals of Mathematical Statistics*, 27, 832-837, 1956.
- Scott, D.W., *Multivariate Density Estimation*. John Wiley and Sons, INC, New York, 1992.
- Sheather, S.J. and M.C. Jones, A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society*, B, 53, 683-690, 1991.
- Silverman, B.W., *Density estimation for statistics and data analysis*. Chapman and Hall, New York, 1986.

Tarboton, D.G., A. Sharma and U. Lall, The use of non-parametric probability distributions in streamflow modeling, *Proc. of the Sixth South African National Hydrological Symp., Pietermaritzburg, South Africa*, Ed. by S.W.K. Eds: S.A. Lorentz & M.C. Dent, 315-327, 1993.

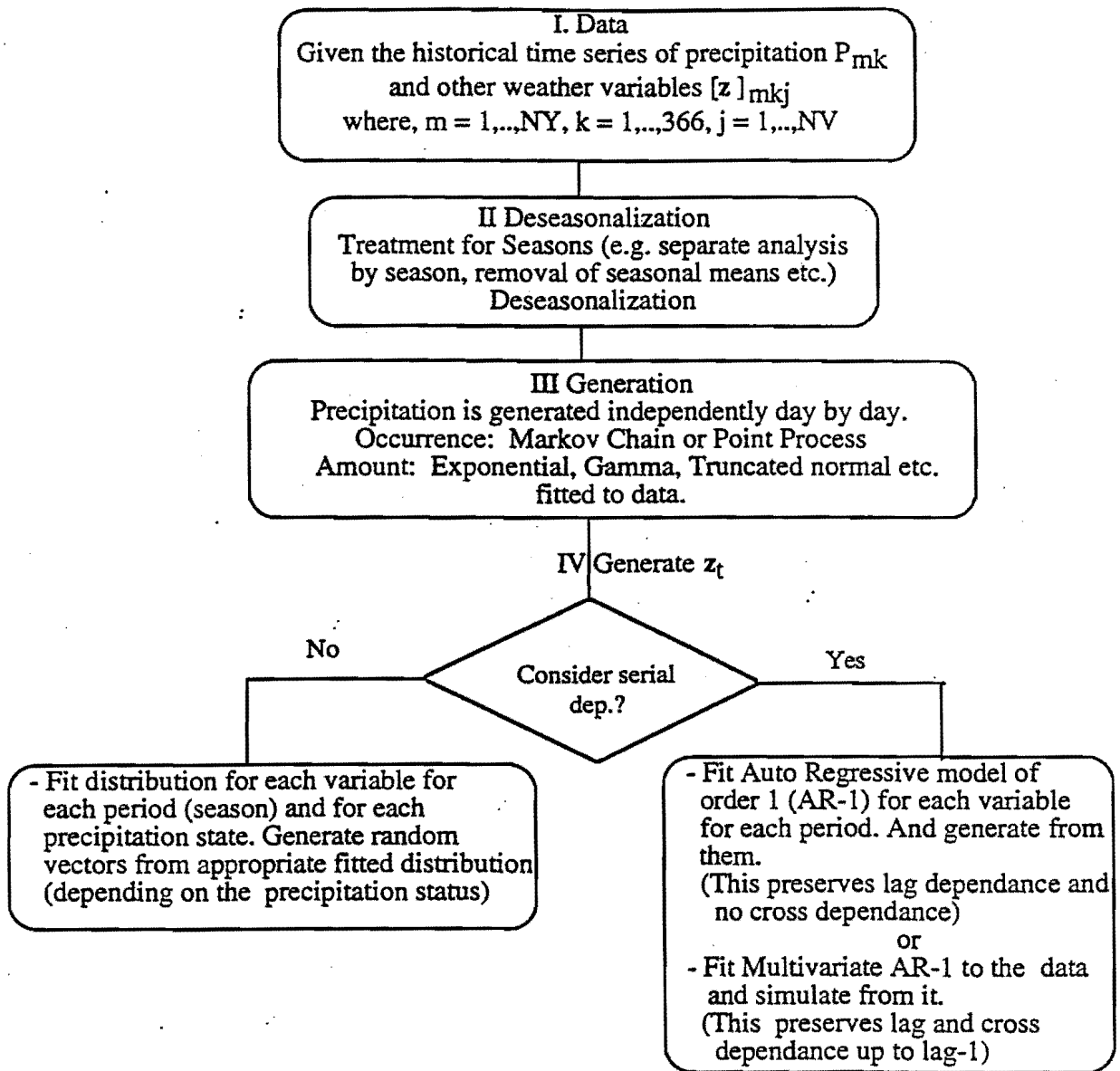


Figure 1: General Structure of Parametric Approaches.

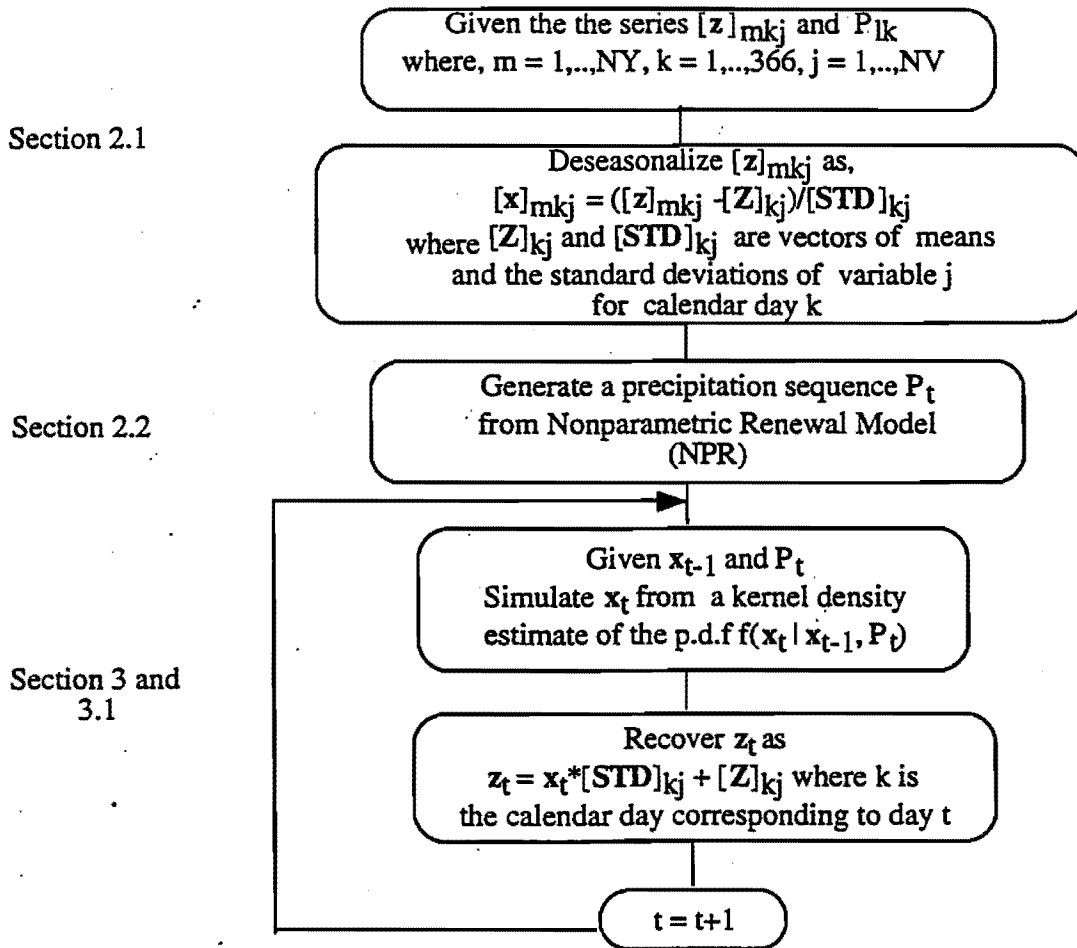
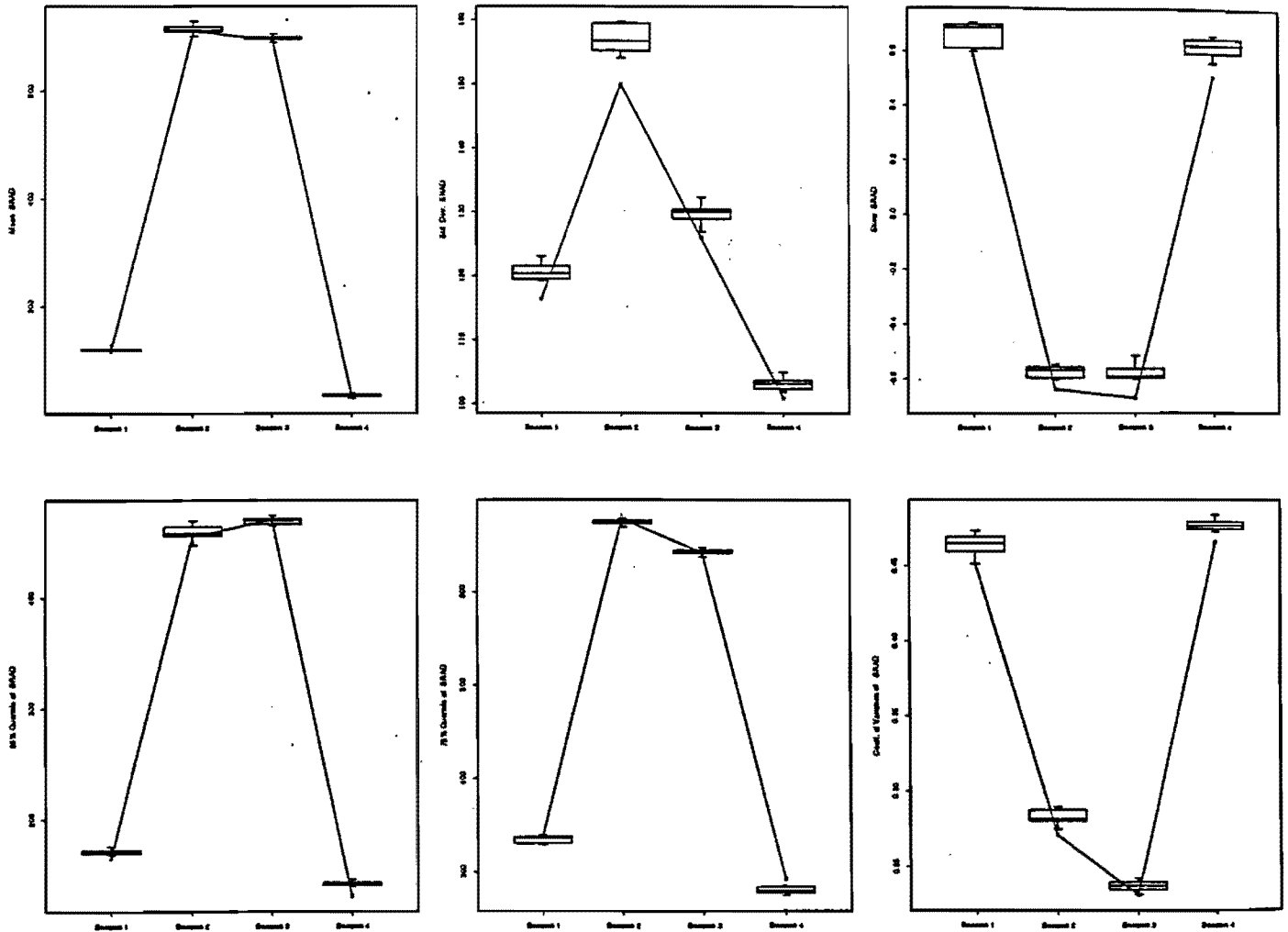
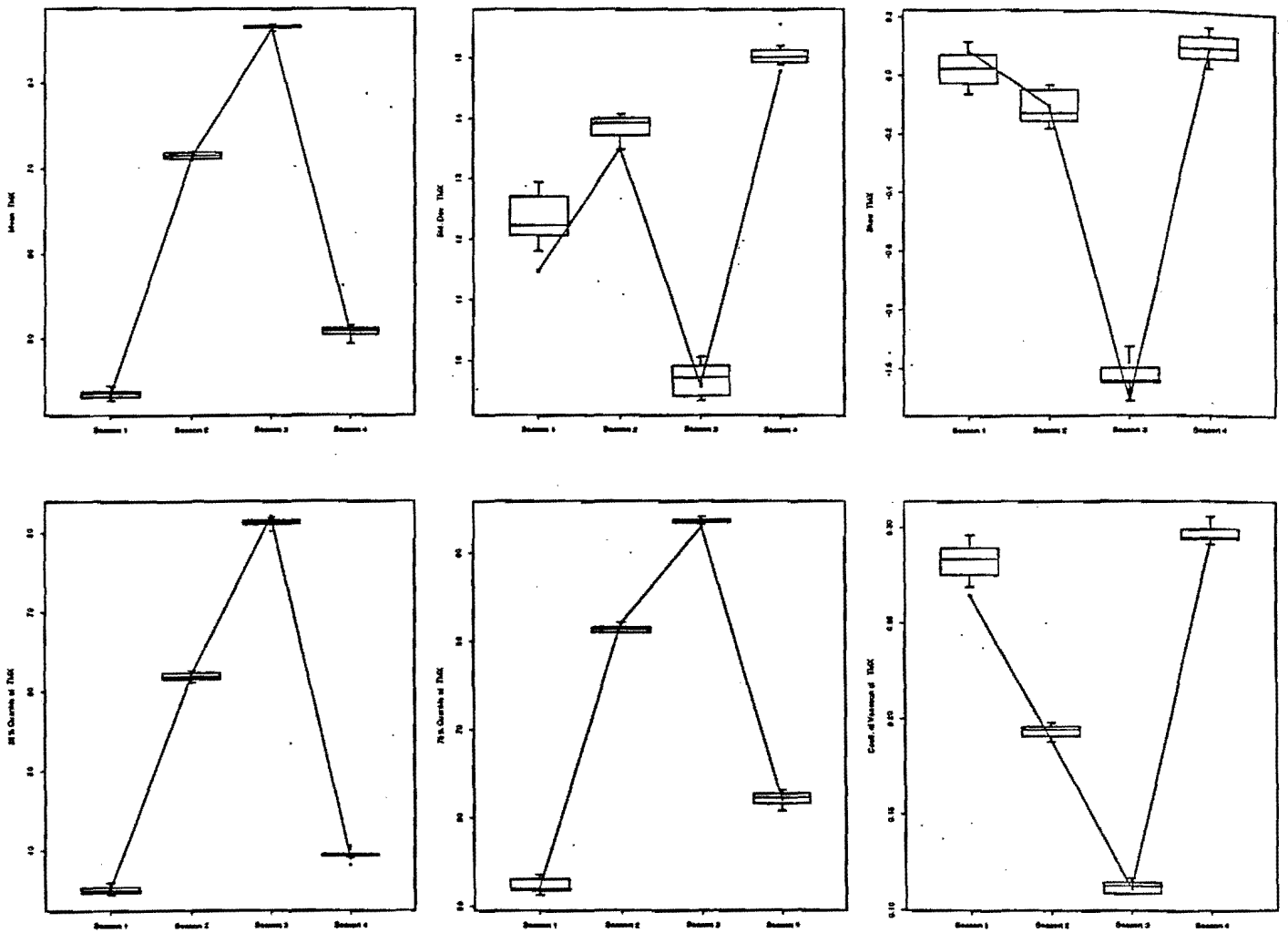


Figure 2: Overview of Development of the NP Model.

**Figure 3:** Example of kernel density estimation using 5 data points with Gaussian Kernel,  $h = 0.5$ .



**Figure 4:** Boxplots of mean, standard deviation, co-efficient, 25% quantile, 75% quantile and co-efficient of variation of SRAD for the four seasons.



**Figure 5:** Boxplots of mean, standard deviation, co-efficient, 25% quantile, 75% quantile and co-efficient of variation of TMX for the four seasons.

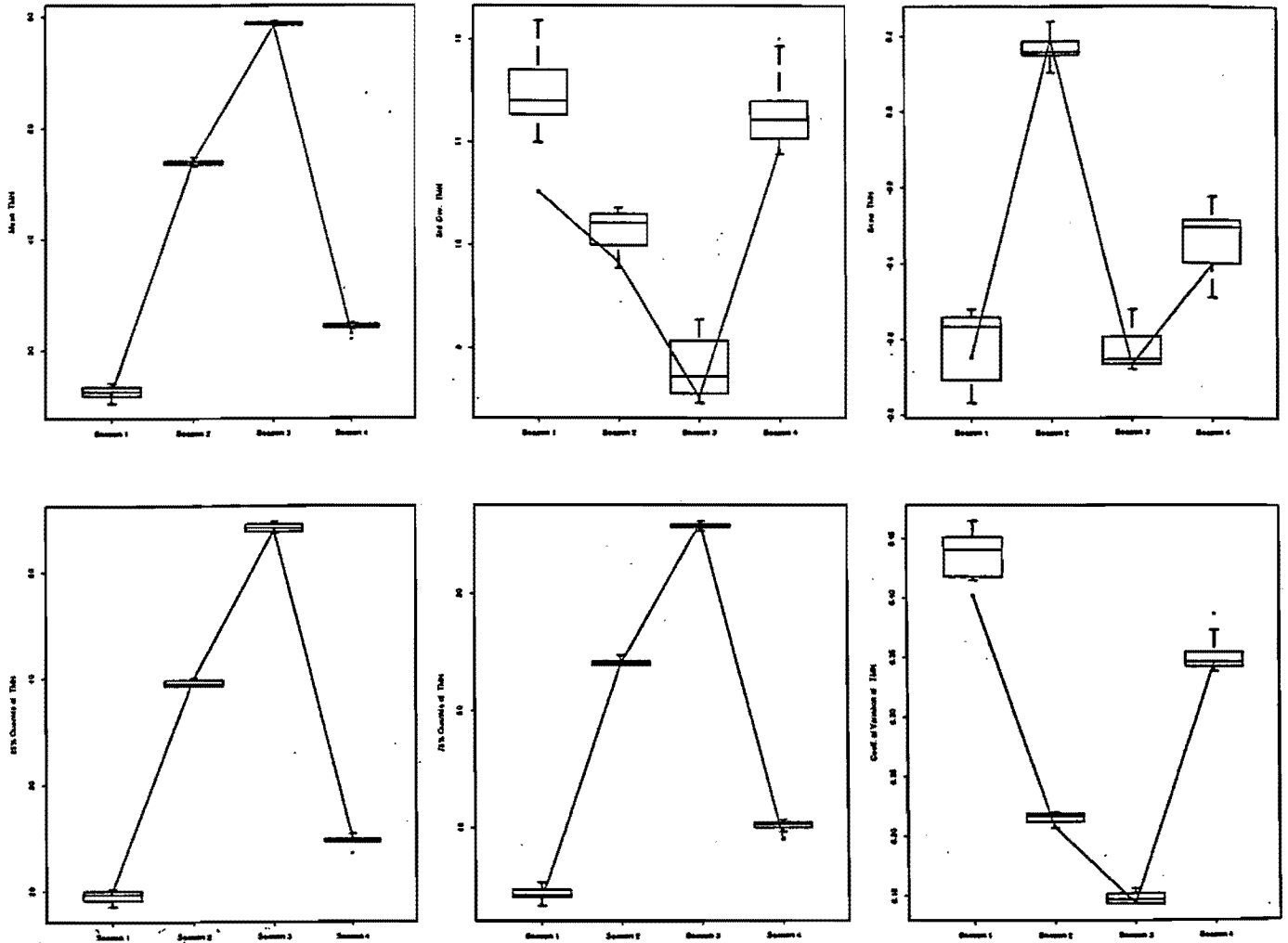
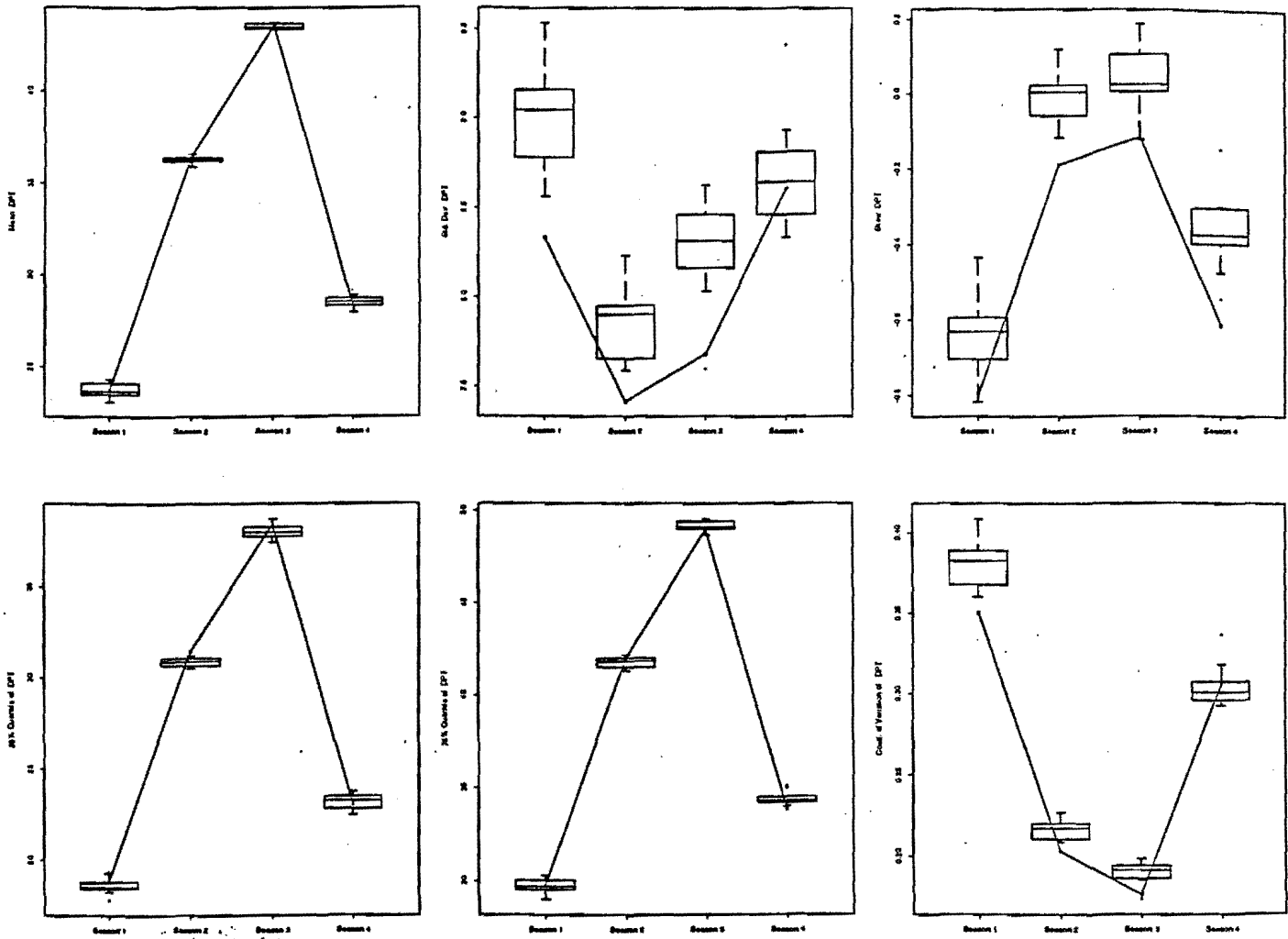
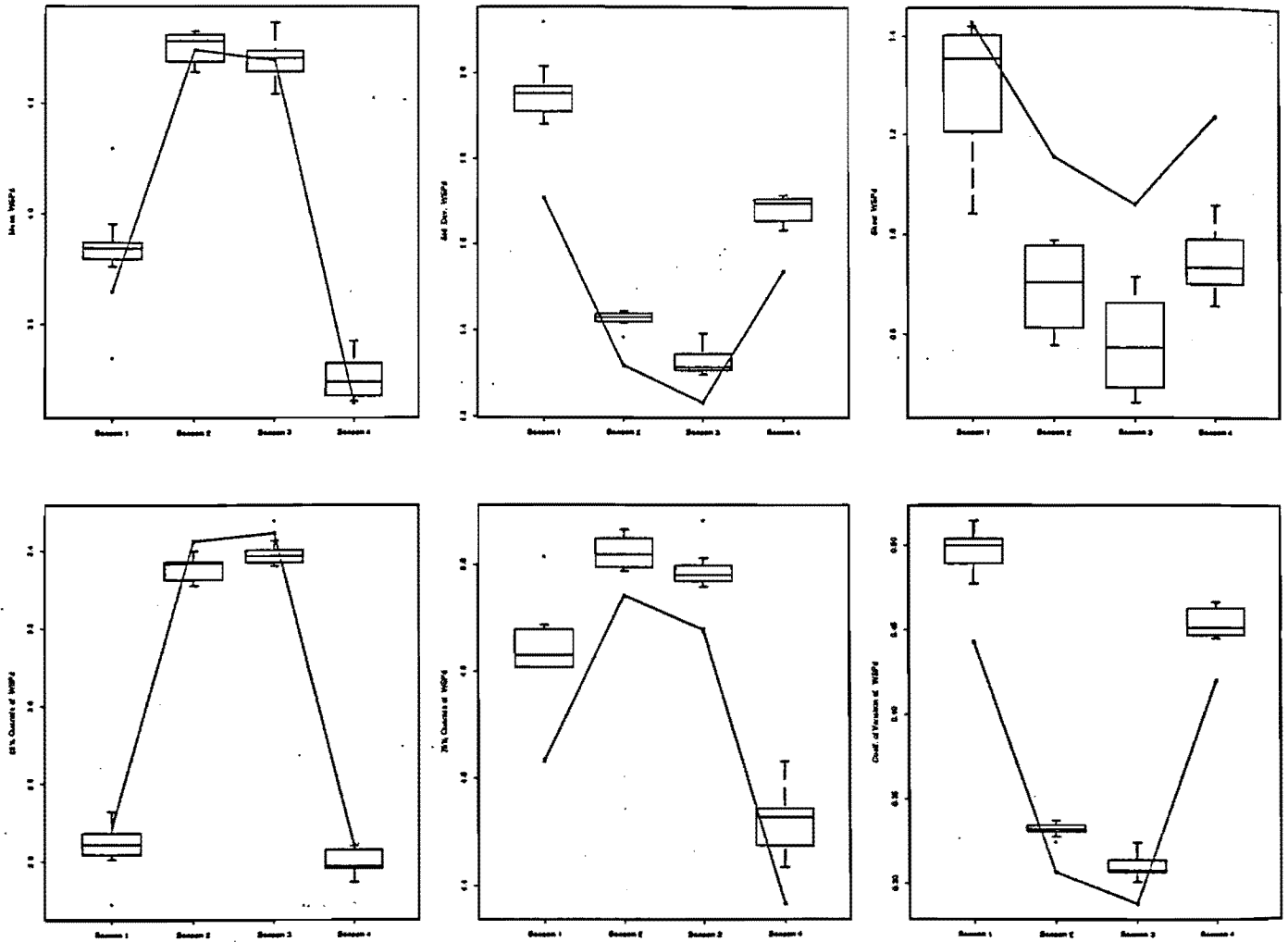


Figure 6: Boxplots of mean, standard deviation, co-efficient, 25% quantile, 75% quantile and co-efficient of variation of TMN for the four seasons.





**Figure 7:** Boxplots of mean, standard deviation, co-efficient, 25% quantile, 75% quantile and co-efficient of variation of DPT for the four seasons.



**Figure 8:** Boxplots of mean, standard deviation, co-efficient, 25% quantile, 75% quantile and co-efficient of variation of WSPD for the four seasons.

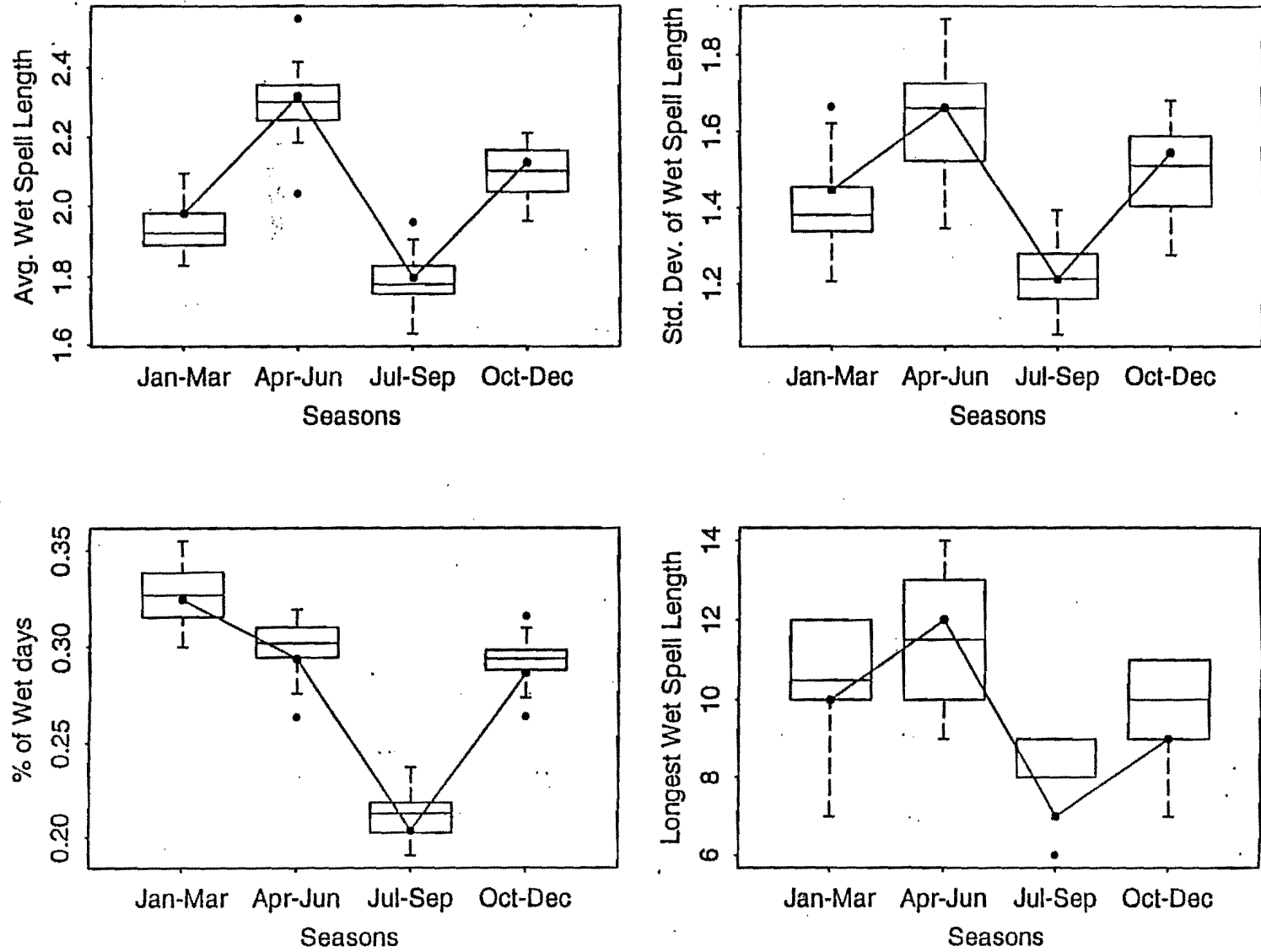


Figure 9: Boxplots of mean, standard deviation, fraction of wet days and longest wet spell length in each season, for simulations from the wet/dry spell model.

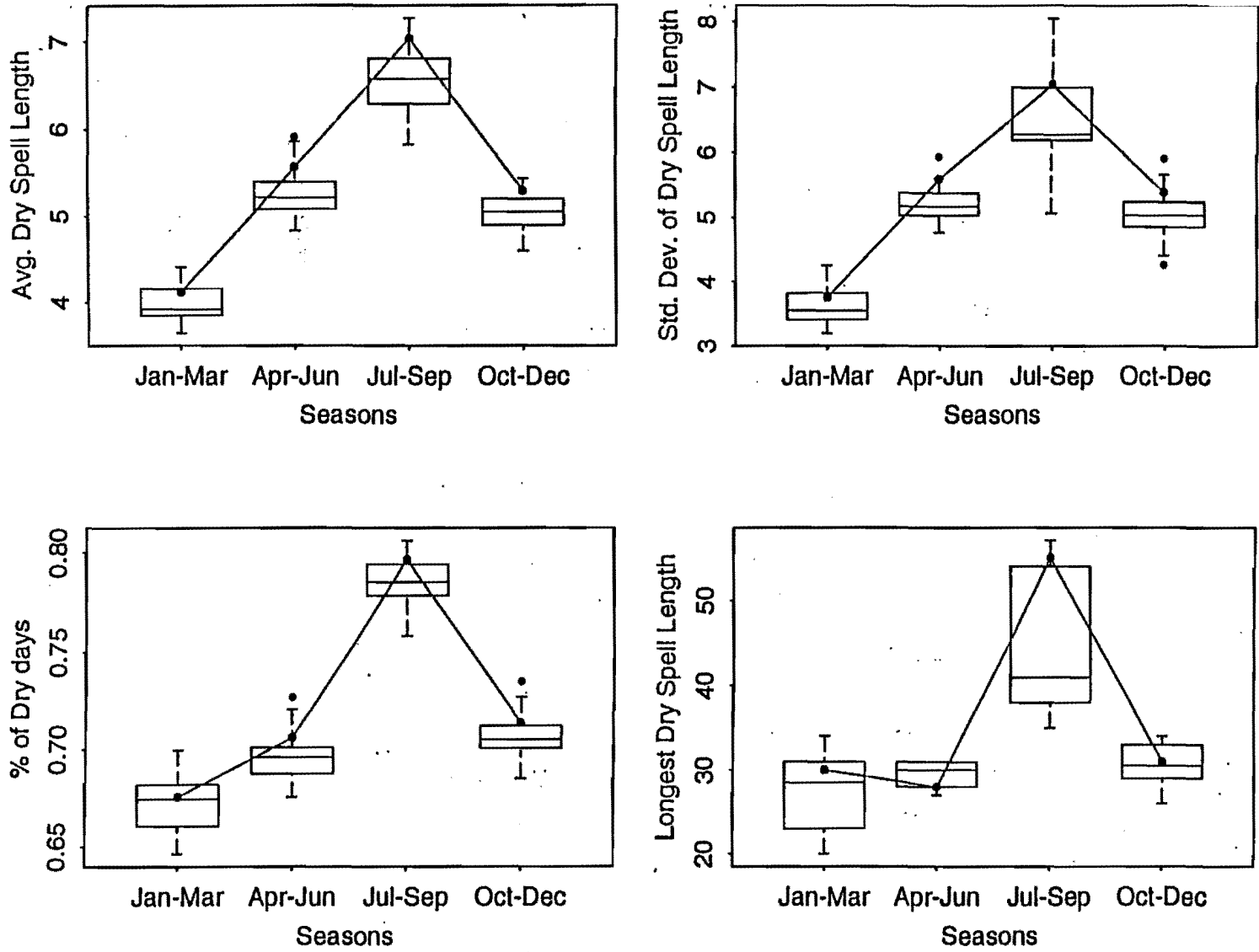


Figure 10: Boxplots of mean, standard deviation, fraction of dry days and longest dry spell length in each season, for simulations from the wet/dry spell model.

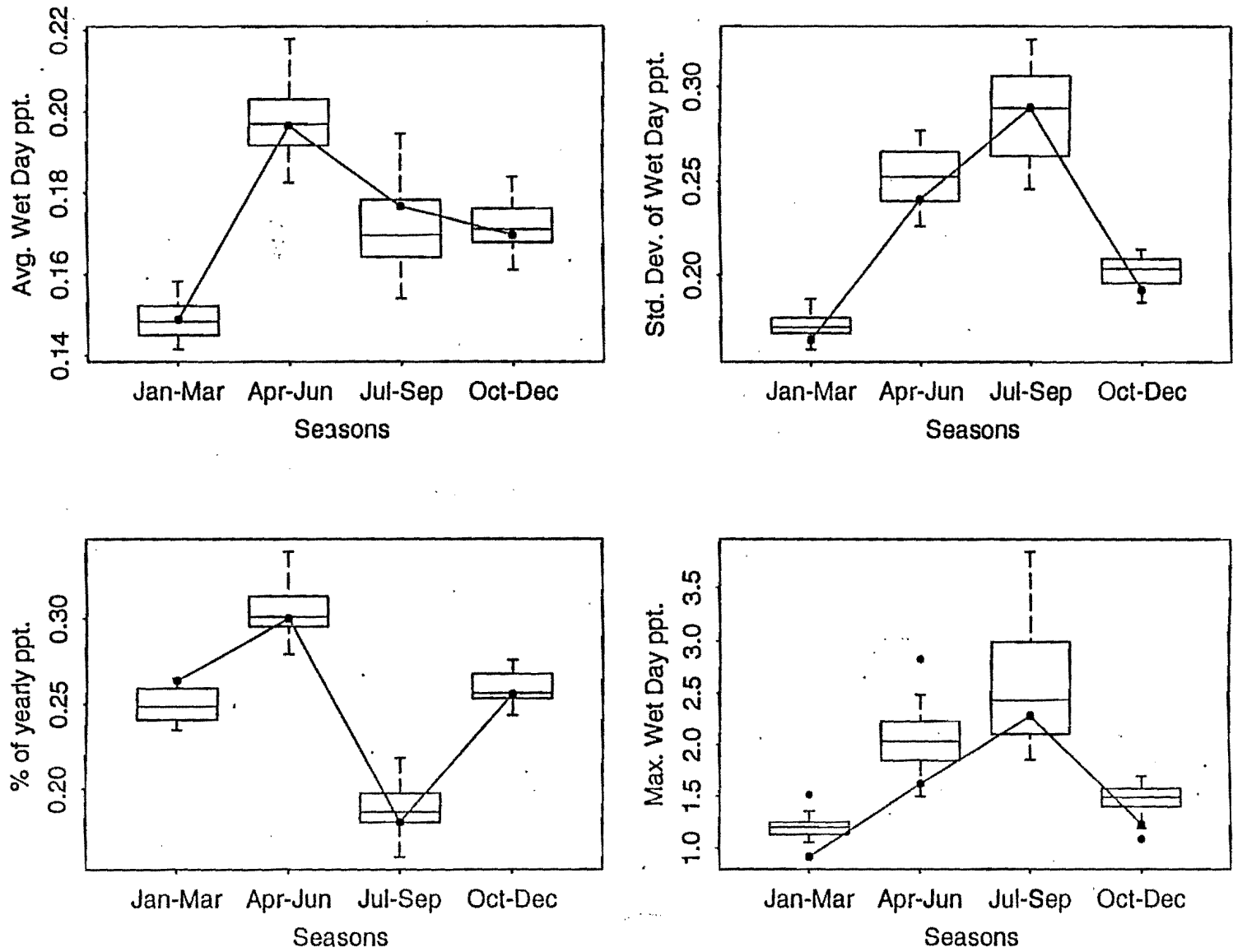


Figure 11: Boxplots of mean, standard deviation, fraction of yearly wet day precipitation and maximum wet day precipitation in each season, for simulations from the wet/dry spell model.

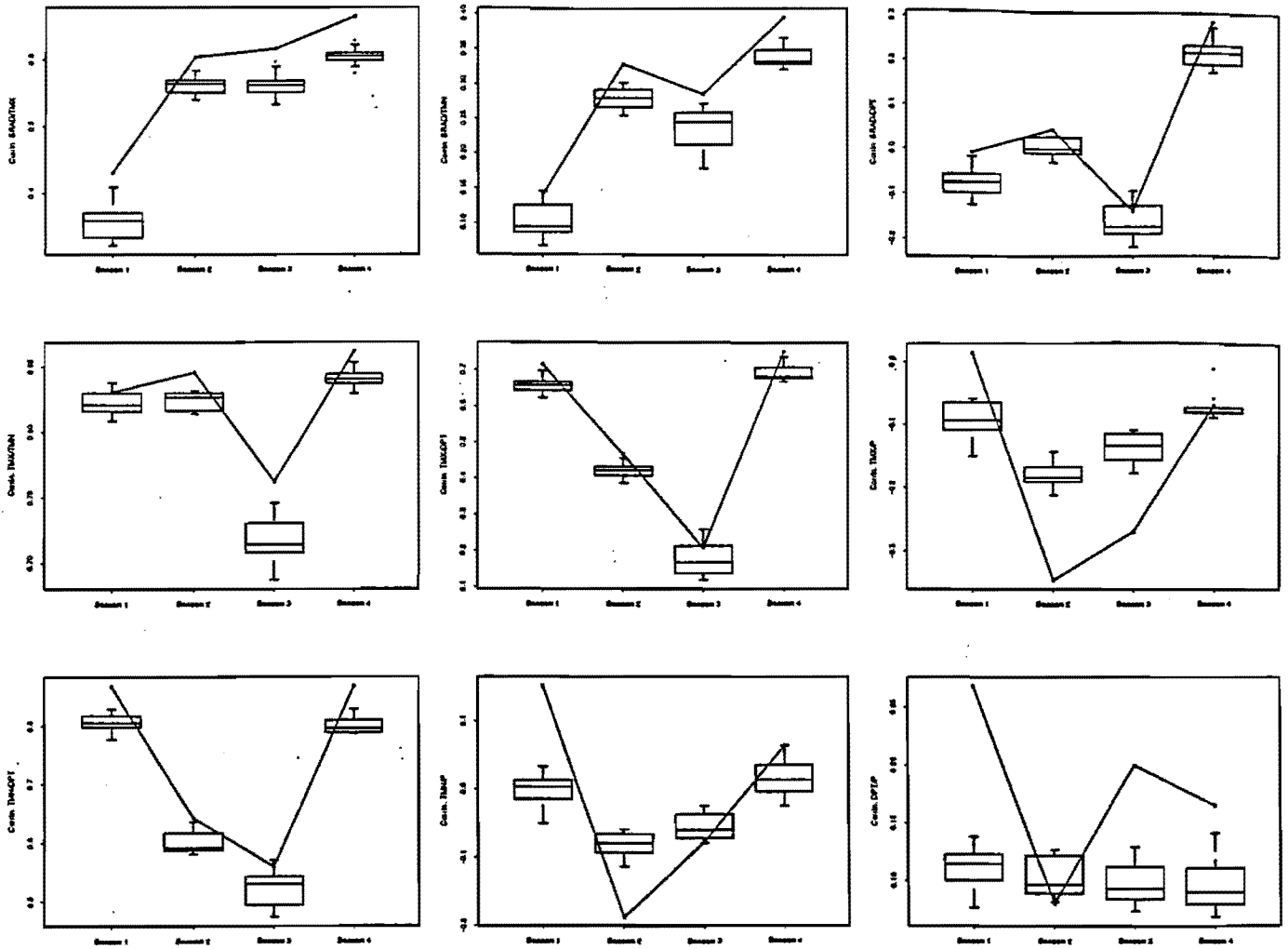
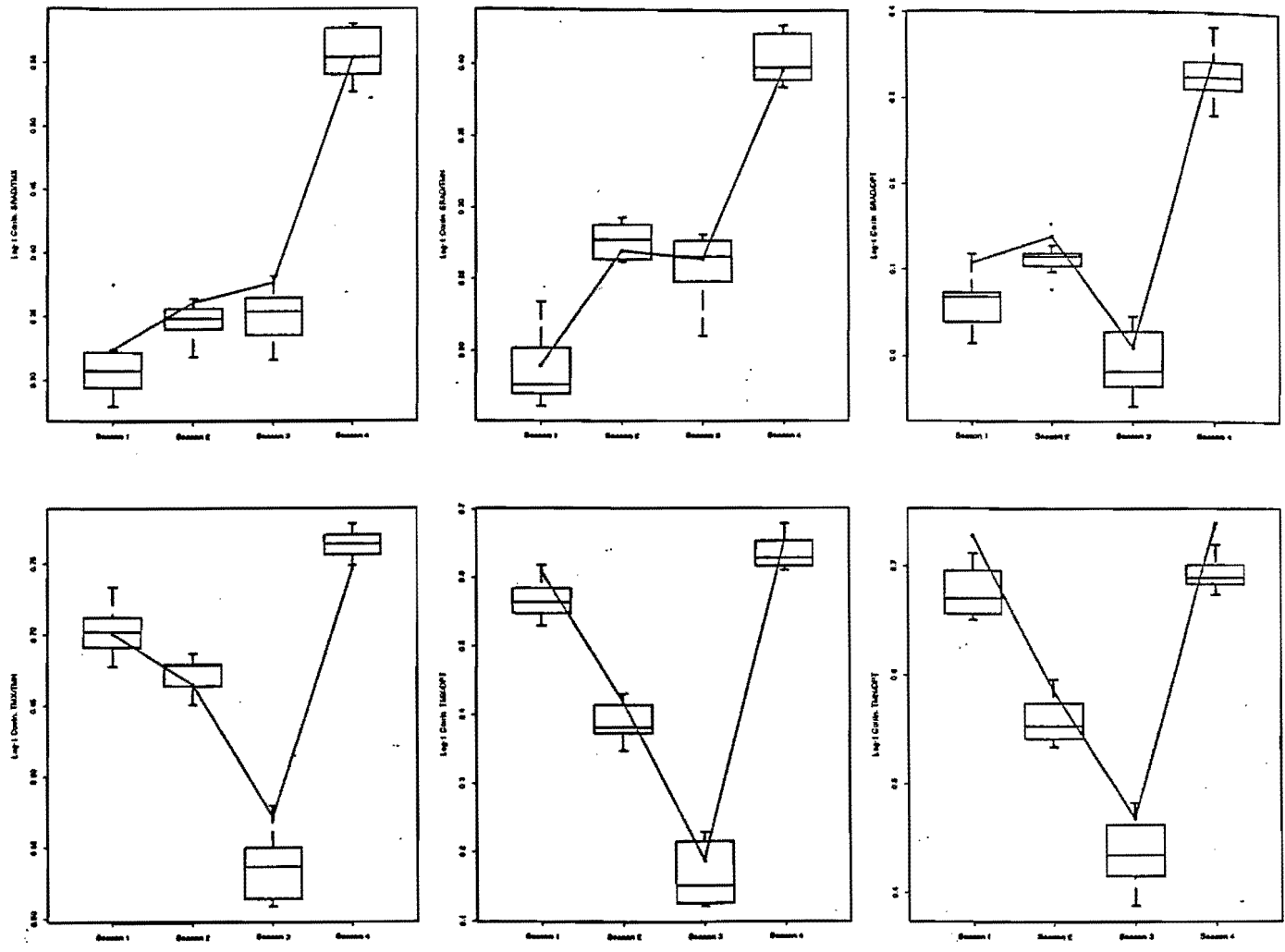
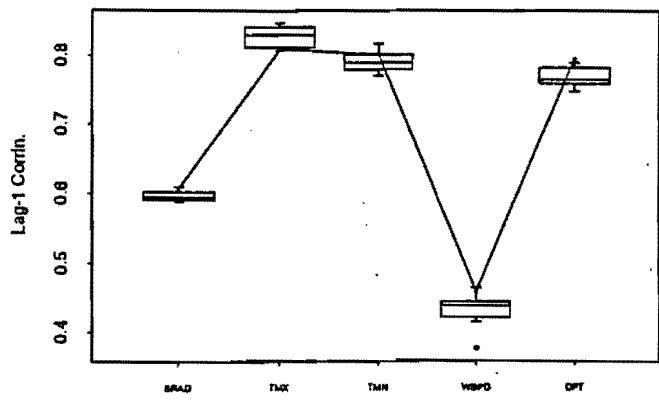


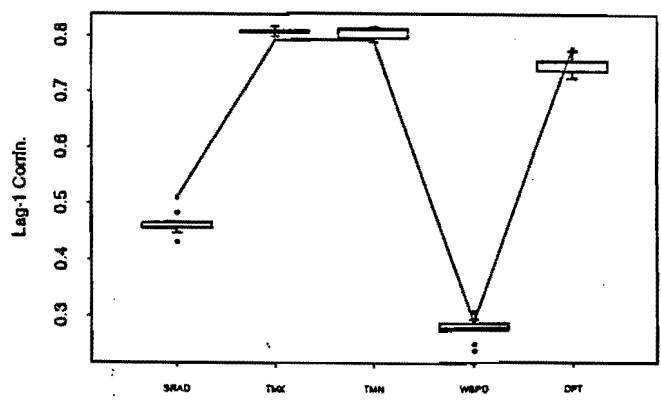
Figure 12: Boxplots of Lag-0 cross correlation between SRAD, TMX, TMN, DPT and P for the four seasons.



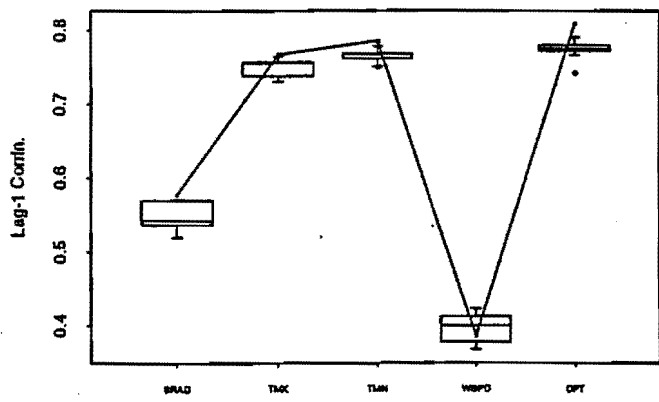
**Figure 13: Boxplots of Lag-1 cross correlation between SRAD, TMX, TMN and DPT for the four seasons.**



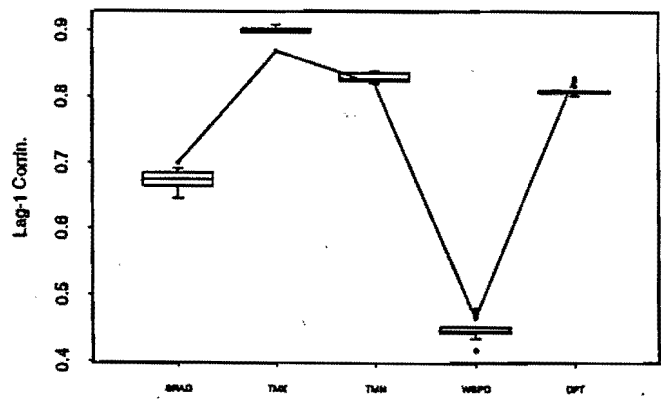
Season 1



Season 2



Season 3



Season 4

Figure 14: Boxplots of Lag-1 Auto Correlation of SRAD, TMX, TMN and DPT for the four seasons.



# APPENDIX 4B

## A Kernel Estimator for Discrete Distributions

Balaji Rajagopalan and Upmanu Lall  
Utah Water Research Laboratory  
Utah State University, Logan, UT - 84322-8200

### Abstract

We present a discrete kernel estimator appropriate for estimating probability mass functions (p.m.f.'s) for integer data. Discrete kernel functions analogous to the Beta functions used as kernels in the continuous case are derived for the interior and for the boundary of the domain. An integer bandwidth is considered. Cross validation is used for bandwidth selection. The estimator was motivated by the need to characterize processes (e.g., mixtures of geometric distributions) with long tailed distributions with high mass near the origin, and integer arguments of the random variable. Monte Carlo comparisons with the Hall and Titterington [8] (HT) estimator are offered. An application for estimating the p.m.f.'s of wet and dry spell lengths for a nonparametric renewal model of daily rainfall is also presented. Other possible methods for obtaining discrete weight sequences are also presented.

## 1. BACKGROUND

The problem of nonparametric smoothing of the empirical discrete p.m.f (or multinomial cell proportions) has been of interest in recent years. However, it has not been studied as intensively as nonparametric density estimation, its counterpart in the continuous case. Hall and Titterington [8] mention that smoothing can be beneficial when there are many cells with small or zero frequencies, i.e the data are sparse. Here we consider that we have a sample  $x_1, \dots, x_n$  for  $n$  multinomial trials with possible outcomes  $1, 2, \dots, k_{\max} \in V$  with probabilities of occurrence  $p_1, \dots, p_{k_{\max}}$  that are unknown. Estimates  $\hat{p}_i$  of the probabilities  $p_i$  may be obtained as sample relative frequencies ( $\hat{p}_i = n_i/n$ ) or cell proportions, or by smoothing the  $\hat{p}_i$ . In the latter case we presume that  $V$  is an ordered set and that "distance" between its members is definable through a standard Lebesgue measure. We consider cases where the set  $V$  may be bounded or unbounded, and focus on developing an appropriate smoother for the sample relative frequencies that properly deals with the discrete nature of the process.

Our practical interest lay in developing a discrete, nonparametric p.m.f for data on the length (in days) of dry or wet spells of rainfall. The shortest spell considered is 1 day. In general, the longest possible spell is not known a priori. Data suggests long right tailed distributions for dry spell length that may correspond to a mixture of geometric p.m.f.'s (see Rajagopalan et. al. [10]).

The concept of smoothing in the context of multinomial cell probability estimation was introduced by Good [6 and 7]. This was later studied and improved by Fienberg and Holland [5], Stone [13], Titterington [14], Titterington [15], Aitchison and Aitken [1], Titterington and Bowman [16] among others. Bishop et al [2] show that these estimators are often better than the cell proportion estimate under squared error loss. Hall and Titterington [8] argue that  $\hat{p}_i$  may not be consistent in data sparse situations. The smoothing estimators developed by Wang and Van Ryzin [17], Simonoff [12] and Hall and Titterington [8] formed a starting point for our work.

The general form of smoothing estimators in this context is given by

$$\hat{p}_i = \sum_{j=-\infty}^{j=\infty} K(i,j,h) \tilde{p}_j \quad i,j \in I, \text{ the set of integers} \quad (1)$$

$K(i,j,h)$  is a weight function or kernel,  $\tilde{p}_j$  is the relative frequency of cell  $j$  and  $h$  is called the bandwidth or window width.

Wang and Van Ryzin [17] developed a class of estimators of the form (1), using a Geometric kernel (WV) ( $K(i,j,h) = 0.5h(1-h)^{|i-j|}$  if  $|i-j| \geq 1$ ;  $K(i,j,h) = (1-h)$  if  $i=j$  and  $h \in [0,1]$ ). The "drop off" of weights associated with the Geometric kernel is rapid. Wang and Van Ryzin [17] estimate  $h$  under an approximate (MSE) criterion formed by truncating the Geometric kernel beyond two cells. As a result, very little smoothing is obtained in most cases and not much may be gained for sparse data.

By imposing a smoothness constraint on the cell probabilities, Simonoff [12] obtained relative consistency results for an estimator based on a maximum penalised likelihood criterion (MPLE). In this approach, the estimates  $\hat{p}_i$  are solved by minimizing a penalized likelihood function defined as,

$$L = \sum_{i=1}^{k_u} n_i \log(\hat{p}_i) - \beta \sum_{i=1}^{k_u} \{\log(\hat{p}_i/\hat{p}_{i+1})\}^2$$

such that

$$\sum_{i=1}^{k_u} \hat{p}_i = 1 \quad (2)$$

$\beta \geq 0$ , is a smoothing parameter, and  $V : [1, k_u]$

The estimates from MPLE depend significantly on the extent of estimation required (i.e.,  $k_u$ ) beyond the maximum observed cell (i.e.,  $k_{\max}$ ). This is of concern, because we would prefer

a natural extension of the tail of the p.m.f by the method used, rather than a prior specification of its extent.

The estimator developed by Hall and Titterington [8] (here after referred to as HT) is given as,

$$\hat{p}_i = \sum_{j=-\infty}^{j=\infty} W(i,j,h) \tilde{p}_j \quad (3)$$

where  $W(i,j,h) = \frac{K((i-j)/h)}{s(h)}$ ,  $h > 1$  and  $s(h) = \sum_{j=-\infty}^{j=\infty} K(j/h)$ .  $K(\cdot)$  is any suitable continuous univariate kernel function, with compact support satisfying the conditions of positivity, integration to unity, symmetry, and finite variance which are,

$$(a) K(u) > 0; (b) \int K(u)du = 1; (c) \int uK(u)du = 0; (d) \int u^2K(u)du = k_2 \neq 0 \quad (4)$$

where  $(u = (i-j)/h)$ , and  $s(h)$  is a multiplicative factor required to normalize the continuous variable kernel function for use with discrete data, such that the desired conditions on  $W(\cdot)$  viz.

$$\sum_{j=-\infty}^{j=\infty} W(i,j,h) = 1 \quad \text{and} \quad \sum_{j=-\infty}^{j=\infty} j W(i,j,h) = 0$$

are satisfied. Hall and Titterington [8] proposed a cross-validatory procedure for selecting  $h$ . This was later studied by Dong and Simonoff [3] who extended this estimator to boundary kernels.

It is well known that kernel estimators suffer from increased bias in the boundary region (i.e.  $1 \leq i \leq h+1$  in our situation of interest). For the estimates of cells in the boundary there is a lack of full complement of observations on either side of the cell of estimate. As a result, the desired conditions on  $W(i,j,h)$  mentioned above will not be preserved. To correct this, special boundary kernels that satisfy the required conditions are used (Müller [9]). Müller [9] formally developed special boundary kernels in the continuous case. Dong and Simonoff [3] developed

boundary kernels (condition 4(a) is relaxed) that could be used in the HT estimator for the discrete case. We refer to the HT estimator with the boundary modification of Dong and Simonoff [3] as HT/DS.

We performed comparisons of these three estimators (viz. WV, MPLE and HT/DS) on data generated from long tailed distributions (see Rajagopalan et. al. [10]) and found HT/DS to be the best. Hence, we compare the relative performance of the estimator we develop later in this paper with HT/DS.

For finite samples, some disquieting aspects of the HT estimator become apparent. The non-integer bandwidth leads to an effective kernel that also varies with  $h$  in a manner quite different from that prescribed by (4). The effective integer support of  $W(i,j,h)$  is  $[(i-h^*), (i+h^*)]$ , where  $h^*$  is the closest integer greater than or equal to  $h$ . HT/DS kernels are defined as quadratics or other polynomials over  $[i-h, i+h]$ . Since this is not the effective integer support of the kernel the effective kernel over the space of integers is not the quadratic defined.

Alternatively, it is possible to develop a kernel that recognizes the data to be in integer space, has an integer bandwidth and satisfies all the required conditions in the integer space. This also obviates the need for normalization of the kernel weights as done in HT/DS. We explored this line of thought and, sought a direct, discrete analog of the continuous kernel density estimator.

The estimator is first presented. Bandwidth estimation is described next. Monte Carlo comparisons with HT/DS are then present. Comparisons with real data sets follow. Discussion of the new estimator and other possible discrete estimators conclude the paper.

## 2. THE DISCRETE KERNEL ESTIMATOR (DKE)

We define our estimator  $\hat{p}_i$  for cell  $i$  through a weighted linear combination of the sample relative frequencies,  $\tilde{p}_i$  as,

$$\hat{p}_i = \sum_{j=1}^{k_{\max}} K(t_j) \tilde{p}_j \quad (5)$$

where  $i, j$  and  $h$  are positive integers,  $t_j = (i-j)/h$ ,  $K(t)$  is a kernel function, and  $V : [1, \infty]$ . In the continuous case, Epanechnikov [4] showed that the MSE optimal kernel of second order, is the quadratic kernel (QK), also known as the Epanechnikov kernel. The general form of the QK is,

$$K(u) = au^2 + b \quad \text{for } |u| \leq 1 \quad (6)$$

In the continuous case,  $a = -0.75$ ,  $b = 0.75$ . Scott [11], p. 140, Equation 6.25 points out that this corresponds to a Beta density function, defined for  $t \in [-1, 1]$ . Other members of this class can be used if additional smoothness is desired.

Here, we chose a discrete quadratic (DQ) kernel of the form  $K(t_j) = at_j^2 + b$ , where  $t_j = (i-j)/h$ . The main focus then is to specify the constants  $a$  and  $b$  for the interior ( $i > h+1$ ) and the boundary region ( $1 \leq i \leq h+1$ ). The constants  $a$  and  $b$  are solved to satisfy : (A) the kernel function goes to zero for  $|i-j| \geq h$ , i.e  $K(t_j) = 0$  for  $|t_j| \geq 1$ , (B) sum of the weights is unity, i.e  $\sum_{j=i-h}^{j=i+h} K\left(\frac{i-j}{h}\right) = 1$  and (C) the first moment of the kernel function is zero, i.e  $\sum_{j=i-h}^{j=i+h} K\left(\frac{i-j}{h}\right)t_j = 0$ . Note that the above conditions are the discrete versions of the conditions given in Equation (3) for continuous variable kernels. One could choose higher order Beta kernels and derive results similar to these that follow for DQ.

For the interior region ( $i > h+1$ ) using conditions (A) and (B) gives Equations (7) and (8),

$$K(t_{i+h}) = K(t_{i-h}) = 0 \quad (7)$$

$$\sum_{j=i-h}^{j=i+h} (at_j^2 + b) = 1, \quad \text{where } t_j = (i-j)/h \quad (8)$$

Condition (C) is satisfied if  $a=-b$ . The coefficients  $a$  and  $b$  can now be expressed in terms of the bandwidth  $h$  as,

$$a = \frac{-3h}{(1-4h^2)} \quad \text{and} \quad b = \frac{3h}{(1-4h^2)} \quad (9)$$

For the boundary region ( $1 < i \leq h+1$ ) condition A is modified as,

$$K(t) = 0 \quad \text{for } t \leq -1 \quad \text{and } t \geq q \quad \text{where } q = (i-1)/h. \quad (10)$$

Applying conditions (B) and (C) we get Equations (11) and (12).

$$\sum_{j=1}^{j=i+h} (at_j^2 + b) = 1 \quad (11)$$

$$\sum_{j=1}^{j=i+h} t_j(at_j^2 + b) = 0 \quad (12)$$

Solving for  $a$  and  $b$  we get,

$$a = \frac{-D}{2h(h+i)} \times \frac{1}{\left(\frac{E}{4h^3} - \frac{CD}{12h^3(h+i)}\right)}, \quad b = \left[1 - \frac{aC}{6h^2}\right] \frac{1}{(h+i)} \quad (13)$$

where,

$$C = h(h-1)(2h-1) + (i-2)(i-1)(2i-3)$$

$$D = -h(h-1) + (i-2)(i-1)$$

$$E = -(h(h-1))^2 + ((i-2)(i-1))^2$$

From Equation (10) it can be seen that at the boundary (i.e.,  $i = 1$ ) the weight associated with the kernel is zero. This is not desirable because, for long tailed distributions defined on the interval  $[1, \infty)$  most of the mass is concentrated right at  $i=1$ . Clearly, using the boundary modification in Equation (13) for estimation of p.m.f at the boundary (i.e.,  $i=1$ ) will introduce a large bias in the estimate. Therefore, we need a further modification for estimation at  $i=1$ . By not enforcing the  $K(t) = 0$  at  $i = 1$ , we modify (A) to be

$$K(t) = 0 \text{ for } t \leq -1 \quad (14)$$

while Equation (11) and (12) remain the same. Solving Equations 14, 11 and 12 for a and b we get,

$$a = \frac{-D}{2h^2} \times \frac{1}{\left(\frac{E}{4h^3} - \frac{CD}{12h^4}\right)}, \quad b = \left[1 - \frac{aC}{6h^2}\right] \frac{1}{h} \quad (15)$$

where,

$$C = h(h-1)(2h-1)$$

$$D = -h(h-1)$$

$$E = -(h(h-1))^2$$

From Equations (9), (13) and (15) note that the kernel and hence, the estimator  $\hat{p}_i$  is expressed strictly in terms of the bandwidth  $h$ . An optimal choice of  $h$  then completes the definition of the estimator.

Three criterion often used for bandwidth estimation are (1) direct minimization of average mean square error (MSE) (2) Maximum likelihood cross validation (MLCV) and (3) Least squares



cross validation (LSCV). These could be optimized over a discrete set of  $h$  values.

We tested all the three methods and found LSCV to be the best. Hall and Titterington [8] and Dong and Simonoff [3] also argue in favour of LSCV. The bandwidth is selected by minimizing the LSCV function given as,

$$\text{LSCV}(h) = \sum_{i=1}^{k_{\max}} (\hat{p}_i)^2 - \frac{2}{n} \sum_{i=1}^{k_{\max}} \hat{p}_i n_i \quad (16)$$

where,  $\hat{p}_i$  is the estimate of the  $i^{\text{th}}$  cell, by dropping the  $i^{\text{th}}$  cell and  $n$ . In a related context, Hall and Titterington [8] also show that cross-validation automatically adapts the estimator to an extreme range of sparseness types. If the multinomial is only slightly sparse, cross-validation will produce an estimator which is virtually the same as the cell-proportion estimator. As sparseness increases, cross-validation will automatically supply more and more smoothing, to a degree which is asymptotically optimal.

An example application comparing DKE (with DQ kernel) to HT/DS with QK based kernels for four data sets is shown in Figures 1, 2, 3 and 4. The data in Figure 1 was sampled from a Geometric distribution (G1) defined as  $G(\pi=0.2)$ . The data in Figure 2 was sampled from a mixture of two Geometric distributions (G2) defined as  $(0.3G(\pi=0.9) + 0.7G(\pi=0.2))$ . The sample sizes for G1 and G2 are 250. Figure 3 shows the p.m.f estimates estimated for the mines data, analysed by Dong and Simonoff [3]. Figure 4 shows the estimated p.m.f from both estimators of dry spell length data, for season 3 (i.e. Jul - Sep) for the station Woodruff, in Utah. The sample size in this case was 539. All four figures indicate that both DKE and HT/DS perform comparably. As both the estimators are similar this is expected. We investigate through Monte Carlo simulations, the behaviour of these estimates for selected situations. The behaviour of the weight sequence from both the estimators are also probed. The results are discussed in the following section.

### 3. MONTE CARLO COMPARISONS

We present results from Monte Carlo simulations, comparing our estimator with the HT/DS estimator using QK. Data sets were generated from situations that may be of interest in our particular context (e.g., geometric distribution, with a considerable boundary region). We generated 500 realizations from the two populations G1 and G2. Sample sizes chosen were  $n = 50, 100, 200, 300, 500$ .

The statistical measures computed to assess the relative performance of DKE and HT/DS estimators are:

1. Average Sum of Squared Errors (ASSE)  $( \sum_{j=1}^{j=nsim} ( \sum_{i=1}^{i=k_u} (\hat{p}_{ij} - p_i)^2 ) / nsim )$  across all realizations for each sample size.
2. Sum of Squared Error (SSE<sub>j</sub>)  $( \sum_{i=1}^{i=k_u} (\hat{p}_{ij} - p_i)^2 )$  for each realization  $j = 1, \dots, nsim$
3. Average Sum of Absolute Error (ASAE)  $( \sum_{j=1}^{j=nsim} ( \sum_{i=1}^{i=k_u} \text{abs}(\hat{p}_{ij} - p_i) ) / nsim )$  across all realizations for each sample size.
4. Cell Root Mean Square Error (CRMSE)  $\{ \sum_{j=1}^{j=nsim} ((\hat{p}_{ij} - p_i)^2) / nsim \}^{0.5}$  across all realizations for each sample size and for each cell  $i = 1, \dots, k_u$
5. Fractional Cell Root Mean Square Error :  $FCRMSE_i = CRMSE_i / p_i$
6. Average Cell Bias (CBIAS<sub>i</sub>)  $\sum_{j=1}^{j=nsim} ((\hat{p}_{ij} - p_i) / nsim )$  across all realizations for each sample size and for each each cell  $i = 1, \dots, k_u$
7. Fractional Cell Bias:  $FCBIAS_i = CBIAS_i / p_i$

8. Coefficient of variation of bandwidth  $C_v = s/\bar{h}$  for each sample size. Where  $s$  and  $\bar{h}$  are the standard deviation and mean of the bandwidths obtained for all the  $n_{sim}$  realizations.

Note that we chose  $k_{ij}$  to be 30 in this case, and  $p_i$ 's are the true p.m.f 's obtained from the known underlying distributions from the samples were generated,  $n_{sim}$  is the number of simulations, in our case it is 500.

Table 1 shows the ASSE and ASAE for the two estimators for the two populations G1 and G2 considered. It can be observed from Table 1 and Figures 5 and 6 that the performance of the two estimators over these two measures is quite close. Figures 5 and 6 indicate that the ASSE appears to decrease with  $n$  at rates -1.03 and -0.86 for HT/DS and -0.85 and -0.9 for DKE, for G1 and G2 respectively. These rates are very similar, and are close to the rate  $n^{-1}$  as anticipated in Hall and Titterington's[8] Theorem 2.1. However, the SSE for HT/DS has a larger spread than DKE as can be seen from Figures 7 and 8 for G1 and G2 respectively for a sample size of 50. The results were generally similar for other sample sizes.

As mentioned earlier we are interested in the behaviour of these estimators at the boundary (left boundary) and in the tails. To assess this,  $CRMSE_i$  and  $FCRMSE_i$  for different sample sizes  $n$  were estimated. As an illustration we present the estimates of  $FCRMSE_i$  for sample sizes 50 and 500 for G1 in Figures 9a and 9b respectively. Figures 10a and 10b are corresponding figures for G2. These figures suggest that DKE performs better than HT/DS in the tail region for all sample sizes, more so for smaller sample sizes. The results for other sample sizes were intermediate.

From Figures 11 and 12 we see that part of the poorer performance of HT/DS in the tails is due to higher bias.

The MSE expression of the estimate  $\hat{p}_i$  as given by Wang and Van Ryzin [17] is,

$$E\left[\sum_{i=1}^{k_{max}} \{\hat{p}_i - p_i\}^2\right] = \sum_{i=1}^{k_{max}} \sum_{j=1}^{k_{max}} W^2(i,j,h)p_j/n - \sum_{i=1}^{k_{max}} \left\{ \sum_{j=1}^{k_{max}} W(i,j,h)p_j \right\}^2/n +$$

$$\sum_{i=1}^{k_{\max}} \left\{ \sum_{j=1}^{k_{\max}} W(i,j,h) p_j - p_i \right\}^2 \quad (17)$$

where  $p_j$  is the true p.m.f,  $W(i,j,h)$  is the weight function,  $h$  is the bandwidth and  $n$  is the sample size. For the the two populations considered viz.  $G1$  and  $G2$  we know the true p.m.f. Substituting this for  $p_j$  in the above equation, the optimal bandwidth can be determined for various sample sizes. These bandwidth values are then compared with the corresponding average bandwidths obtained from the simulations. These along with the coefficient of variance of bandwidth  $C_v$  are summarized in Table 2. It can be observed that  $C_v$  is smaller for DKE for all the sample sizes for  $G1$  and  $G2$ . Note that DKE smooths the Geometric distribution data ( $G1$ ) more than HT/DS, and smmoths the mixture data ( $G2$ ) less than HT/DS. Also the average bandwidths from DKE are close to the MSE optimal bandwidths. This suggests that the bandwidth from DKE is more stable than from HT/DS.

The behaviour of HT/DS in these simulations is interesting. There is a tendency to undersmooth relative to the optimal bandwidth. As a result the boundary bias decreases with  $n$ , while the tail bias may be high. The higher coefficient of variance of the HT/DS bandwidth suggests a higher degree of adaptation to sample attributes. However, this fails to consistently provide a lower bias on MSE than DKE.

The need to choose a bandwidth in the boundary region that is different from the interior has been recognized by several researchers (e.g. Müller [9]). Generally variation in  $h$  across the range of the data, and especially in the tails is needed. The selection of a "local" bandwidth considering boundary kernels and tail regions remains an area of research.

#### 4. OTHER POSSIBLE ESTIMATORS

Müller [9] shows how one can develop minimum variance kernels and kernels belonging to different smoothness classes for continuous variates. Extensions of these ideas to the discrete case is also feasible. Here we outline two such extensions.

A discrete, minimum variance (DMV), second order kernel can be developed as the solution to:

$$\text{Minimize } \sum_{j=q}^{i+h} w_j^2 \quad (18)$$

Subject to:

$$w_q = w_{i+h} = 0 \quad (19)$$

$$\sum_{j=q}^{i+h} w_j = 1 \quad (20)$$

$$\sum_{j=q}^{i+h} t_j w_j = 0 \quad (21)$$

where  $t_j = (i-j)/h$ ,  $i, j, h$  are integers, and  $q = \max(i-h, 1)$ , recognizes whether we are in the boundary region or the interior.

A smooth, discrete (DS $\mu$ ) kernel of smoothness  $\mu$  can be defined by solving the problem:

Minimize  $\sum_{j=q}^{i+h-\mu} (w_{j+\mu} - w_j)^2$  subject to the conditions (19) through (21) above. Solutions to the two problems defined above can be readily obtained by defining the associated Lagrangian problems and solving them for the weights  $w_j$  that define the kernel sequence over the appropriate span of integers.

The weight sequences resulting for DMV and DS1 ( $\mu=1$ ) for selected values of  $h$ , and  $i$  are compared with the DQ and HT/DS weight sequences in Table 3. In the interior, the HT/DS, DQ

and DS1 weight sequences coincide. This is to be expected since they all converge to the quadratic kernel. The DMV sequence degenerates to uniform weights as expected. An examination of the weight sequences in the boundary region shows that the DQ sequences stay closer to the DS1 sequences than the HT/DS ones. Thus if a computationally fast approximation to the DS1 sequences was desired in the boundary region, DQ would be preferred. Note that the DMV sequences in the boundary region are still generally closer to the DS1 than the HT/DS.

An interesting aspect of the HT/DS sequence is the adaptation of the weight sequence as  $h$  varies between two integers. We observe that the weight sequences at the intermediate  $h$  value are not strictly in between the weight sequences at the end points. While this may lead to a high degree of adaptability of the HT/DS procedure, it makes it rather difficult to assess its impact on the estimation procedure. The high coefficient of variation of the bandwidth selected by HT/DS may be related to the nature of the resulting weight sequence.

The boundary kernels developed by Dong and Simonoff [3] do not correspond to the ones presented by Müller [9] for the continuous case. It may be interesting to try the Müller [9] boundary kernels, possibly with a floating boundary value, directly with the HT procedure.

Computational considerations have restricted our Monte Carlo investigations thus far to DQ and HT/DS. The relative utility of DMV and DS may be investigated subsequently. Except in the boundary region, our limited investigations show that differences between the different kernels may not be large. Consequently, kernels that are easier to compute are expedient. In this respect the DQ kernels are to be preferred.

## 5. SUMMARY AND CONCLUSIONS

The estimator presented here was motivated by practical considerations. We offer this work in the hope that it will stimulate interest and theoretical development. We show that the discrete kernel procedure advocated can give results comparable to those from the HT/DS procedure. Computational advantages of the DKE procedure and the similarity of its properties to kernel sequences based on smoothness criteria were demonstrated. The relative stability of the bandwidth selection procedure and the DQ weight sequence also recommend it as an alternative to the HT/DS method.

We present only one special case (a quadratic kernel in the interior and in the boundary region). Clearly other similar higher order kernels can be derived. However, as is typical in the kernel smoothing literature, bandwidth selection is likely to be a more tenuous issue than kernel specification. The LSCV choice of  $h$  appears to perform quite satisfactorily for the test cases. Extensions to the multivariate case are being investigated.

## ACKNOWLEDGEMENTS

Partial support of this work by the U.S. Forest Service under contract notes, INT-915550-RJVA and INT-92660-RJVA, Amend #1 is acknowledged. The principal investigator of the project is D.S. Bowles. We are grateful to J.Simonoff and J. Dong for discussions, provision of relevant manuscripts and computer code for p.m.f estimation using HT/DS estimator.

## REFERENCES

- [1] J. Aitchison and C.G. Aitken, Multivariate binary discrimination by the kernel method, Biometrika, 63 (1976), 413-420.
- [2] Y.M. Bishop, S.E. Fienberg and P.W. Holland, Discrete multivariate analysis: Theory and Practice, MIT Press, Cambridge, Mass., (1975).
- [3] J. Dong and J.S. Simonoff, The construction and properties of boundary kernels for sparse multinomials, Journal of Computational and Graphical Statistics, 3, (1), (1994) 57-66.
- [4] V.A. Epanechnikov, Nonparametric estimations of a multivariate probability density, Theor. Probab. Appl., 14 (1969), 153-158.
- [5] S.E. Fienberg and P.W. Holland, Simultaneous estimation of multinomial cell probabilities, J. Amer. Statist. Assoc., 68 (1973), 683-691.
- [6] I.J. Good, The estimation of probabilities, MIT Press, Cambridge, Mass, 1965.
- [7] I.J. Good, A Bayesian significance test for multinomial distributions (with discussion), J.Roy Statist. Soc., Ser. B., 29 (1967), 399-431.
- [8] P. Hall, and D.M. Titterington, On smoothing sparse multinomial data, Australian Journal of Statistics 29 (1989), 19-37.
- [9] H.G. Müller, Smooth optimum kernel estimators near endpoints, Biometrika 78(3) (1991), 521-530.
- [10] B. Rajagopalan, U. Lall, D.G. Tarboton, Simulation of daily precipitation from a nonparametric renewal model. Working Paper WP-93-HWR-UL/003. In Utah Water Research Laboratory, Utah State University, Logan, UT, (1993).
- [11] D.W. Scott, Multivariate density estimation, Wiley series in probability and mathematical statistics, John Wiley and Sons, New York, (1992).



- [12] J.S. Simonoff, A penalty function approach to smoothing large sparse contingency tables, Ann. Statist., 11 (1983), 208-218.
- [13] M. Stone, Cross-validation and multinomial prediction, Biometrika, 61 (1974), 509-515.
- [14] D.M. Titterington, Updating a diagnostic system using unconfirmed cases, Appl. Statist., 25 (1976), 238-247.
- [15] D.M. Titterington, A comparative study of kernel-based density estimates for categorical data, Technometrics, 22 (1980), 259-268.
- [16] D.M. Titterington and A.W. Bowman, A comparative study of smoothing procedures for ordered categorical data, J. Statist. Comput. Simul. 21, (1985), 291-312.
- [17] M.C. Wang and J. Van Ryzin, A class of smooth estimators for discrete distributions, Biometrika 68(1) (1981), 301-309.

Table 1  
Comparison of ASSE and ASAE

	ASSE			ASAE		
	DKE	PAR	HT/DS	DKE	PAR	HT/DS
Samples generated from G1 (Geometric ( $\pi=0.2$ ))						
n = 50	0.0058	0.0008	0.0084	0.2032	0.0816	0.2737
n = 100	0.0032	0.0006	0.0038	0.1558	0.0599	0.1814
n = 200	0.0019	0.0003	0.0019	0.1183	0.4250	0.1264
n = 300	0.0013	0.0002	0.0012	0.1000	0.0323	0.0987
n = 500	0.0008	0.0000	0.0008	0.0780	0.0226	0.0797
Samples generated from G2 (0.7* Geometric ( $\pi=0.2$ )+0.3* Geometric ( $\pi=0.9$ ))						
n = 50	0.0080	---	0.0081	0.2300	---	0.2481
n = 100	0.0039	---	0.0038	0.1676	---	0.1638
n = 200	0.0021	---	0.0022	0.1261	---	0.1194
n = 300	0.0016	---	0.0016	0.1071	---	0.0978
n = 500	0.0010	---	0.0011	0.0855	---	0.0785

Note

PAR is the fitted parametric (in this case the fitted Geometric distribution)

Table 2  
Bandwidth statistics

	<u>Coefficient of Variation</u>		<u>Average Bandwidth</u>		<u>Optimal Bandwidth from MSE criteria</u>	
	DKE	HT/DS	DKE	HT/DS	DKE	HT/DS
<u>Sample from G1</u>						
n = 50	0.349	0.442	6.73	5.48	7.00	8.06
n = 100	0.305	0.401	6.13	4.97	6.00	8.06
n = 200	0.316	0.361	4.96	4.36	5.00	7.14
n = 300	0.290	0.314	4.51	4.21	4.00	6.25
n = 500	0.275	0.341	4.00	3.47	4.00	5.56
<u>Sample from G2</u>						
n = 50	0.309	0.291	2.844	3.067	3.00	4.10
n = 100	0.210	0.220	2.280	2.931	2.00	4.03
n = 200	0.007	0.213	2.020	2.902	2.00	4.03
n = 300	0.000	0.212	2.000	2.912	2.00	4.03
n = 500	0.000	0.214	2.000	2.844	2.00	4.03

Table 3  
Comparison of weight sequences

	h = 2	h = 2.5	h = 3
<u>Interior</u>			
DQ	0,.3,.4,.3,0	---	0,.14,.23,.26,.23,.14,0
HT/DS	0,.3,.4,.3,0	0,.11,.25,.29,.25,.11,0	0,.14,.23,.26,.23,.14,0
DMV	0,.33,.33,.33,0		0,.2,.2,.2,.2,0
DS1	0,.28,.44,.28,0		0,.14,.23,.26,.23,.14,0
<u>Boundary</u>			
i = 1			
DQ	1,0,0	---	.75,.5,-.25,0
HT/DS	0,1,0	0,1.7,-.7,0	0,.88,.12,0
i = 2			
DQ	0,1,0,0	---	0,.75,.5,-.25,0
HT/DS	0,.63,.37,0	0,.62,.45,-.07,0	0,.5,.4,.1,0
DMV	0,1,0,0		0,.83,.33,-.16,0
DS1	0,1,0,0		0,.8,.4,-.2,0
i = 3			
DQ		---	0,.3,.4,.3,0,0
HT/DS		0,.28,.35,.28,.08,0	0,.28,.32,.28,.12,0
DMV			0,.4,.3,.2,.1,0
DS1			0,.34,.37,.23,.06,0

Notes:

i is the point of estimate, on which the kernel is placed, h is the bandwidth.

DQ, DMV and DS1 do not admit non integer bandwidths.

The HT/DS weights correspond to a quadratic kernel, and admits non-integer h.

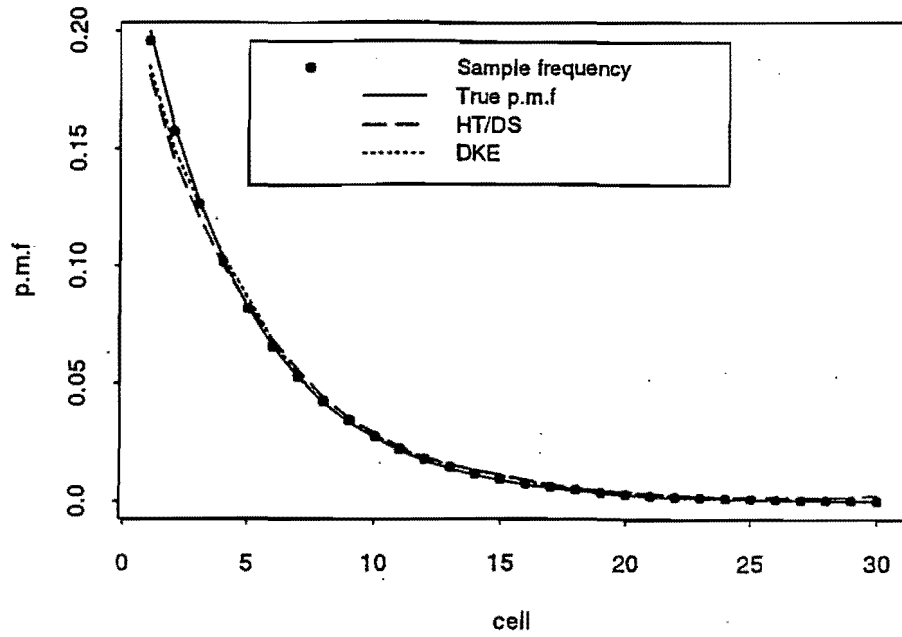


Figure 1. True p.m.f, estimated p.m.f from HT/DS and DKE of a sample of size 250, data generated from Geometric ( $\pi=0.2$ ), along with the sample frequency.

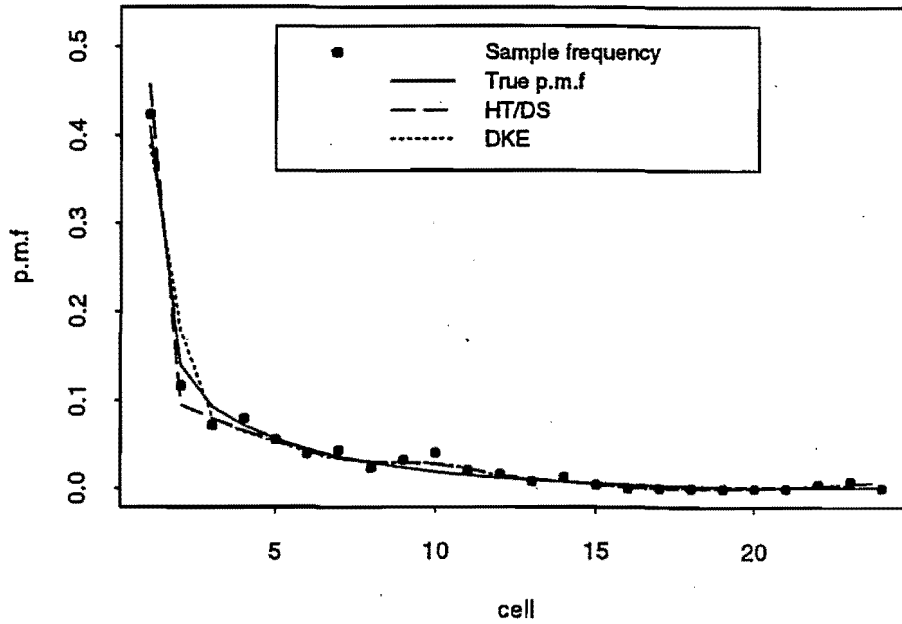


Figure 2. True p.m.f, estimated p.m.f from HT/DS and DKE of a sample of size 250 generated from  $0.7 * \text{Geometric} (\pi=0.2) + 0.3 * \text{Geometric} (\pi=0.9)$ , along with the sample frequency.

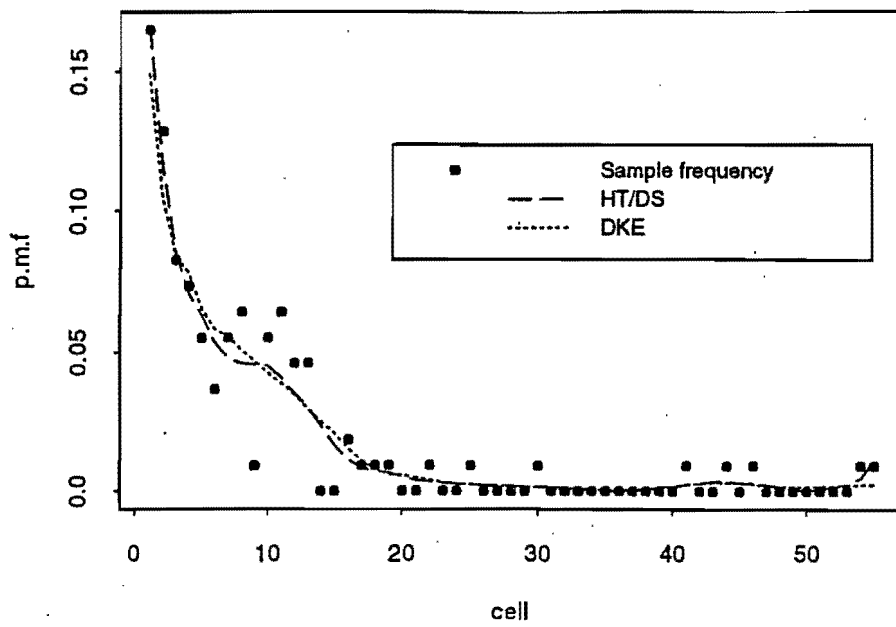


Figure 3. Estimated p.m.f from HT/DS and DKE of the mines data, along with sample frequency.

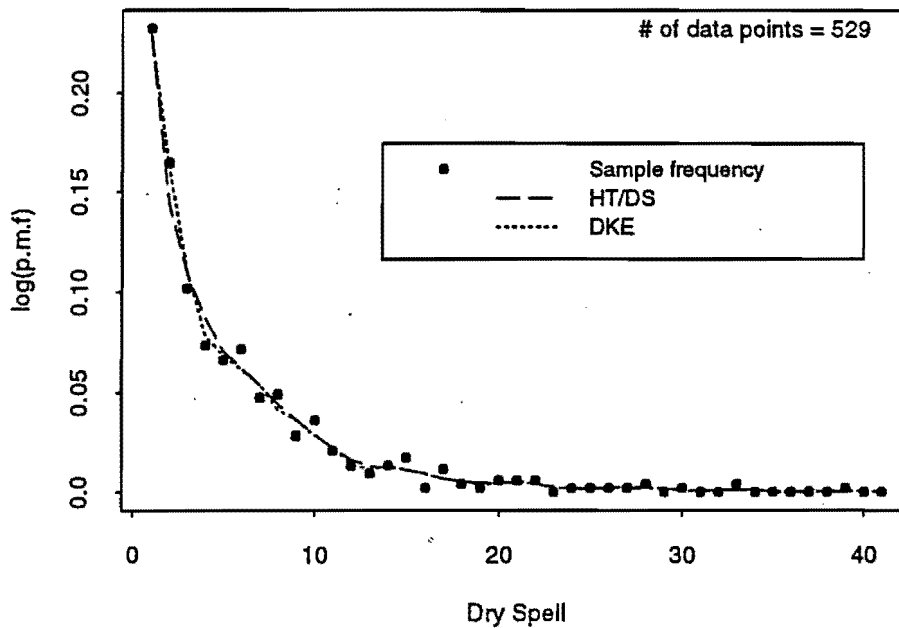


Figure 4. Estimated p.m.f from HT/DS and DKE of the dry spell length data of Woodruff, Utah, along with the sample frequency.

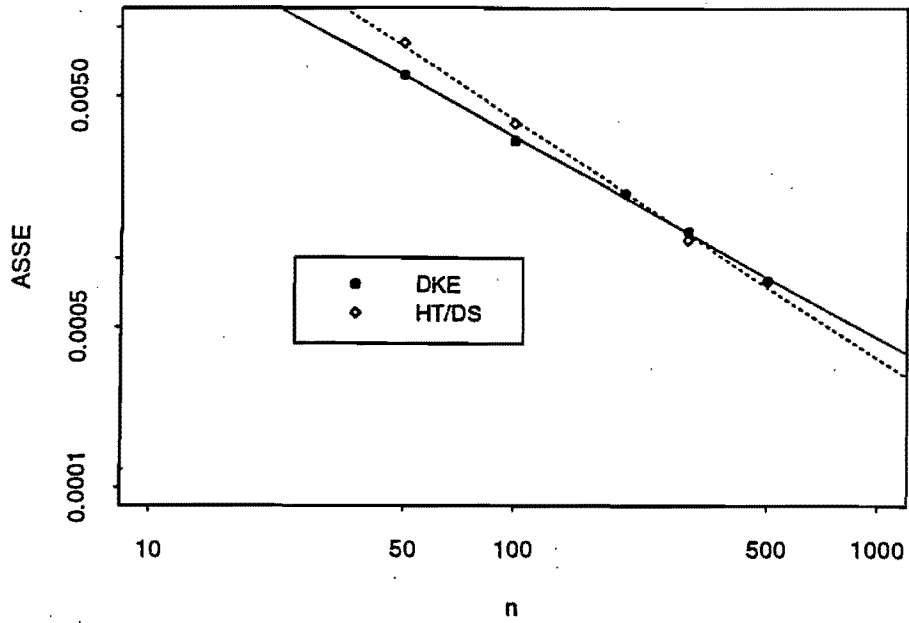


Figure 5. Log-Log plot of ASSE with sample size  $n$ , of samples generated from Geometric ( $\pi=0.2$ ) along with the fitted lines.

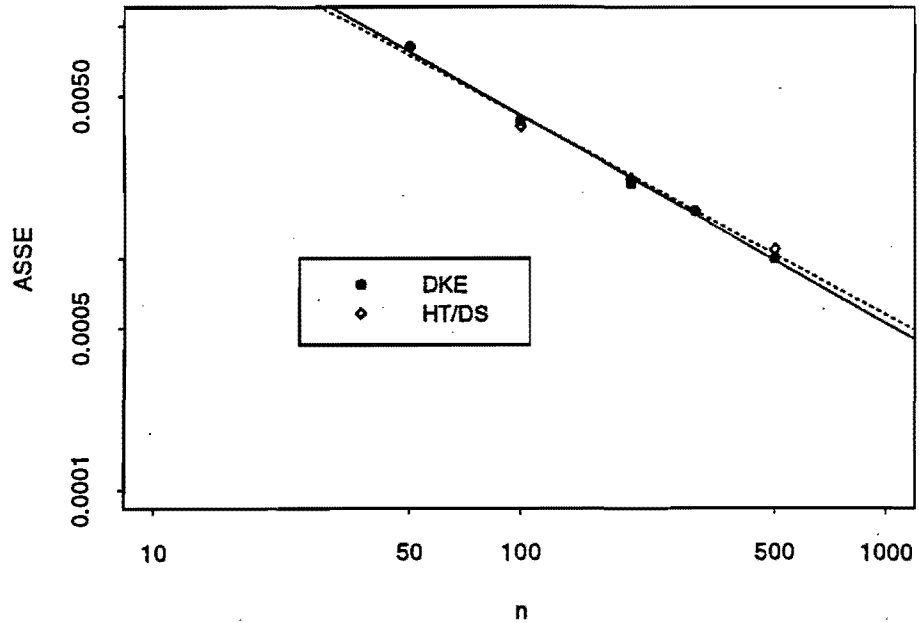


Figure 6. Log-Log plot of ASSE with sample size  $n$ , of samples generated from  $0.7 * \text{Geometric} (\pi=0.2) + 0.3 * \text{Geometric} (\pi=0.9)$  along with the fitted lines.

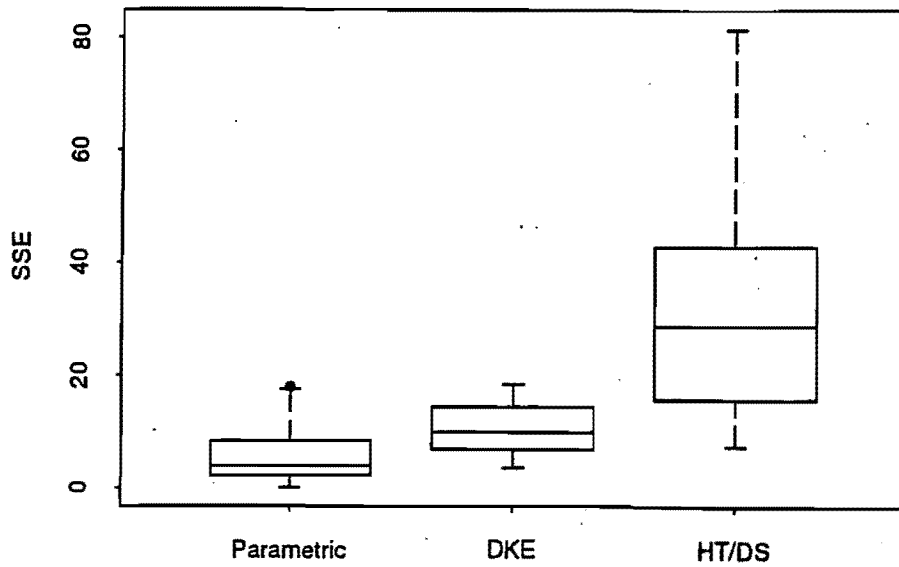


Figure 7. Boxplots of  $SSE_j$  from HT/DS, DKE and fitted Parametric distribution, of samples generated from Geometric ( $\pi=0.2$ ) of sample size 50.

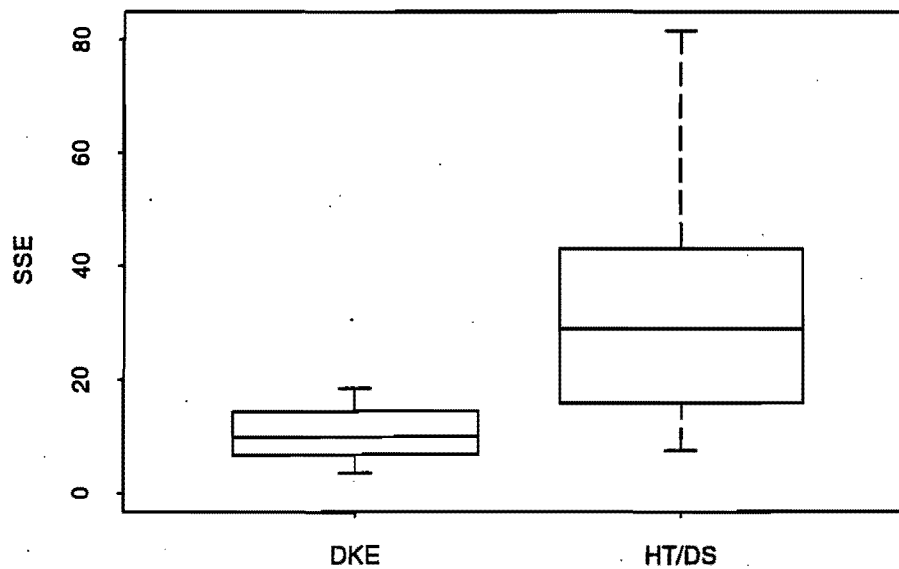


Figure 8. Boxplots of  $SSE_j$  from HT/DS and DKE of samples generated from  $0.7 \cdot \text{Geometric}(\pi=0.2) + 0.3 \cdot \text{Geometric}(\pi=0.9)$  of sample size 50.



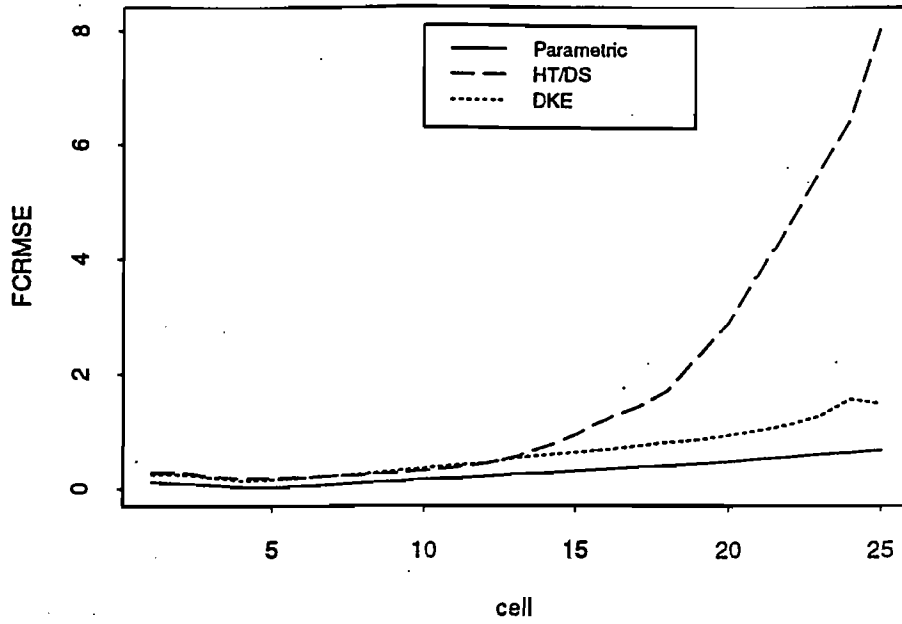


Figure 9(a).  $FCRMSE_i$  from HT/DS, DKE and fitted Parametric distribution, of samples generated from Geometric ( $\pi=0.2$ ) of sample size 50.

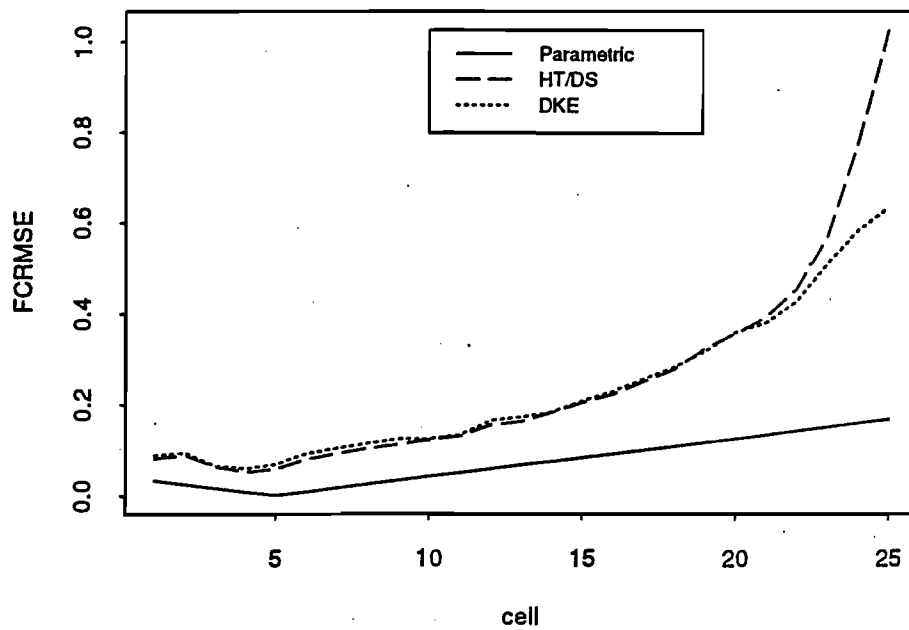


Figure 9(b).  $FCRMSE_i$  from HT/DS, DKE and fitted Parametric distribution, of samples generated from Geometric ( $\pi=0.2$ ) of sample size 500.

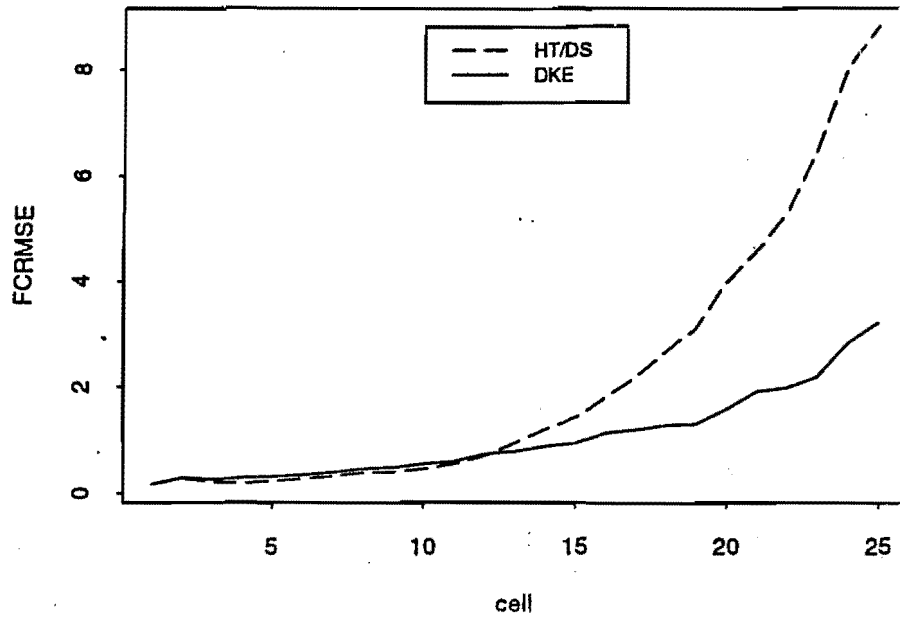


Figure 10(a).  $FCRMSE_i$  from HT/DS and DKE, of samples generated from  $0.7 * \text{Geometric} (\pi=0.2) + 0.3 * \text{Geometric} (\pi=0.9)$  of sample size 50.

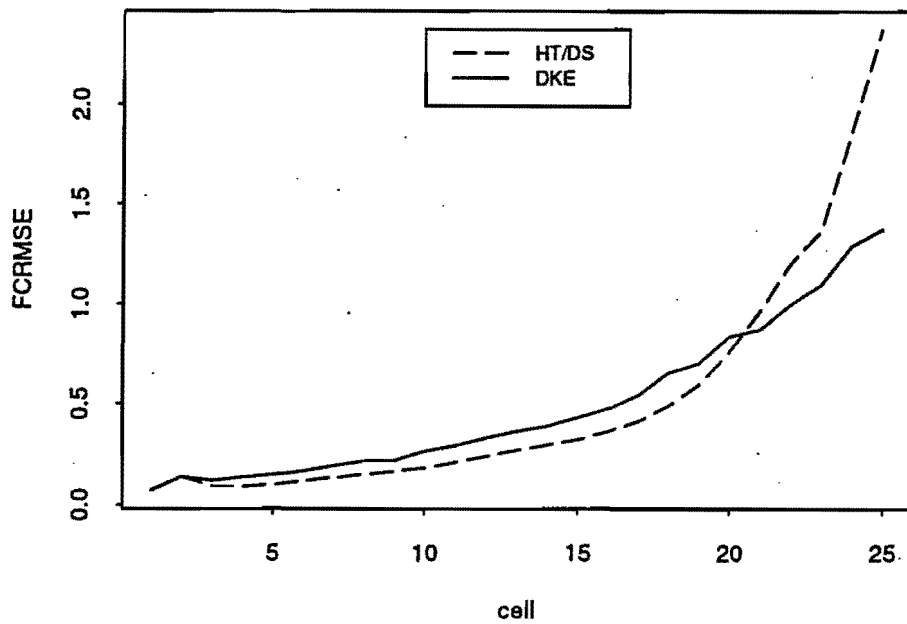


Figure 10(b).  $FCRMSE_i$  from HT/DS and DKE, of samples generated from  $0.7 * \text{Geometric} (\pi=0.2) + 0.3 * \text{Geometric} (\pi=0.9)$  of sample size 500.

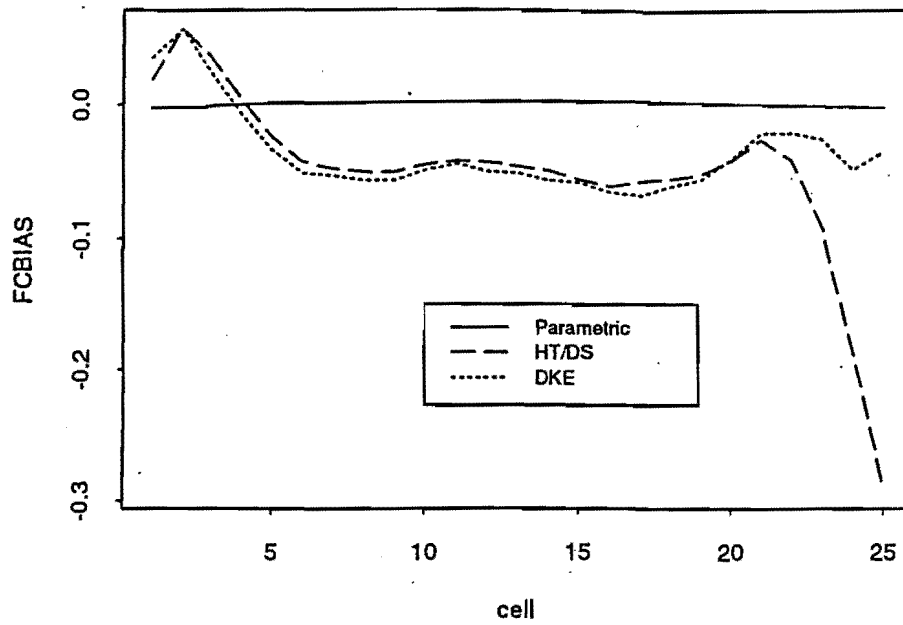


Figure 11. FCBIAS<sub>i</sub> from HT/DS and DKE, of samples generated from Geometric ( $\pi=0.2$ ) of sample size 500.

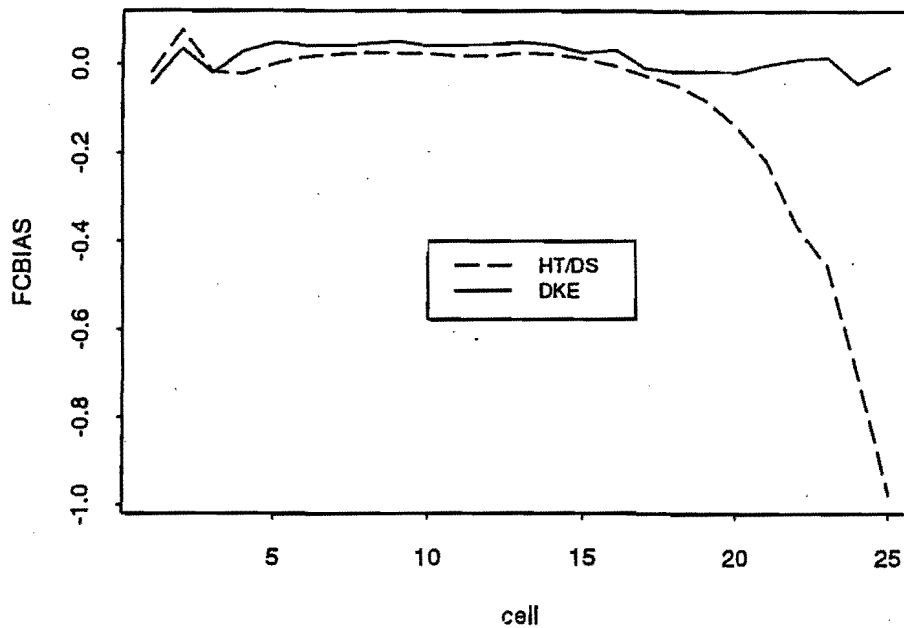


Figure 12. FCBIAS<sub>i</sub> from HT/DS, DKE of samples generated from  $0.7 \cdot \text{Geometric}(\pi=0.2) + 0.3 \cdot \text{Geometric}(\pi=0.9)$  of sample size 500.

# APPENDIX 4C

## Seasonality of Precipitation along a Meridian in the Western U.S

Balaji Rajagopalan and Upmanu Lall  
Utah Water Research Laboratory  
Utah State University, Logan, UT 84322-8200

### Abstract

We investigate seasonality of daily precipitation along a meridian in the Western U.S. using a nonparametric technique. The occurrence of daily precipitation is treated as a nonhomogeneous Poisson process and the time varying intensity function is estimated for every calendar day using a kernel estimator. The technique is fully data adaptive. We apply this technique to selected long record stations along a meridional transect spanning from Tuscon, AZ to Priest River ID. Differences in the seasonality of precipitation occurrence and magnitude are revealed as a function of latitude and topographic factors. A monotonic trend in the seasonality of precipitation over the length of record is also observed.

## 1. Introduction

Seasonality in hydroclimatic variables is usually related to the unequal heating of the earth's surface over the year, particularly as one moves to higher latitudes. Precipitation is an important hydrologic variable since it is a primary input into surface hydrologic models. The timing and duration of the "seasons" of high precipitation at a site is important since they indicate the form (rain or snow) of precipitation as well as the nature of the input "signal" for the surface hydrologic system.

Here we were interested in dynamically visualizing how the seasonality of rainfall varies by latitude along a transect in the western U.S. (approx. longitude  $112^{\circ}$  W). Long record precipitation stations which had essentially complete records were selected from latitude  $48^{\circ} 17'$  N to latitude  $32^{\circ} 15'$  N . We were interested in daily precipitation because of its use for agriculture, crop management and forest management. The attributes of interest considered are precipitation 'magnitude' and 'relative frequency of occurrence'.

Stochastic precipitation models as well as other hydrologic models often deal with the nonstationarity in precipitation and other climatic inputs by dividing the year into a number of seasons and then fitting model parameters independently for each season. The leading terms (one or two) of a Fourier series representation of the precipitation data are commonly used to identify seasonality, for time varying parameter description and for delineating seasons.

An attractive alternative to Fourier series methods is provided in this paper. We focus first on the rate of occurrence of precipitation as a function of calendar date (1 to 366) within the year. A kernel estimator is used to estimate the "rate" of rainfall occurrence of precipitation by calendar day, by "smoothing" a binary (1 or 0) indicator sequence that represents precipitation occurrence on a given day in the historical record. This rate is interpretable as the time varying rate parameter of a nonhomogeneous Poisson process. Variation in precipitation magnitude over a 90 day moving window is also investigated.

An interesting trend in seasonality is exhibited by the stations we analyzed. There appears to be a consistent shift in the seasons identified on the basis of precipitation rate. The calendar dates associated with the highest and the lowest precipitation rates for a given year appear to move forward each year of the record.

## 2. Methodology

Precipitation is an intermittent process. For understanding climatic variations it is often useful to consider adaptive representations that allow a smooth, continuous time interpretation of precipitation. The Poisson process has been used to describe rainfall occurrence as a point process (Waymire and Gupta 1981 a; Cox and Isham, 1980). In the stationary point process, the number of events (e.g., the events are occurrence of wet days)  $n(T)$  occurring in a duration  $T$  is a random variable with a Poisson distribution with mean  $\lambda T$ :

$$p(n(T) = k) = (\lambda T)^k e^{-(\lambda T)}/k! \quad k = 0,1,2 \quad (1)$$

where  $\lambda$  is called the rate or intensity parameter. Often, it is hard to distinguish between changing intensity of the process and event clustering. This situation can be addressed by explicitly allowing changing event intensity in the model and consequently, modeling the daily precipitation as a nonhomogeneous Poisson process (same as Equation 1 but with a time varying rate parameter  $\lambda$ , i.e.  $\lambda(\tau)$ ,  $\tau = 1, \dots, 366$ ) to capture the changing precipitation pattern over the year. Our thesis here is that this time varying rate parameter is a useful indicator of precipitation seasonality at a site.

Kernel intensity estimators (see Diggle, 1985; Solow, 1991) can be used to estimate  $\lambda(\tau)$  from the record, through an optimal, weighted moving average of the rate of rainfall occurrence over time. To form such an estimate, we need to define an appropriate weight function, a span over

which to average and a criteria for choosing the weight function and span in an optimal way. Our presentation here is informal and is restricted to a description of the estimation process used.

Daily precipitation data from about a dozen sites spread along Arizona, Utah and Idaho were used to estimate the intensity parameter for each day of the historical record. Table 1 summarizes the site and data information.

### 2.1 Estimation Procedure

We considered the estimation of  $\lambda(\tau)$ , for each calendar day  $\tau$  (1,2,...,366), for each year of record  $y$ . The average across years of the estimates of  $\lambda(\tau)$  provides a measure of the typical seasonality at the site.

The kernel estimator used for  $\lambda_y(\tau)$ , the rate on calendar day  $\tau$ , in year  $y$  is,

$$\widehat{\lambda}_y(\tau) = \frac{1}{h_y} \sum_{i=1}^{n_y} K\left(\frac{\tau - \tau_{i,y}}{h_y}\right) \quad (2)$$

In equation 2,  $\tau$  (1,2,...,366) is the calendar day on which the estimate is required,  $\tau_{i,y}$  is the index of a calendar day on which there was rain in year  $y$ ;  $K(\cdot)$  is a kernel function which is taken to be a positive function that integrates to unity, is symmetric and has finite variance;  $h_y$  is a bandwidth or "scale" parameter (for year  $y$ ) of the kernel function, that controls the smoothness of  $\widehat{\lambda}_y(\tau)$ .

The estimator in Equation (2) is very similar to a kernel density estimator (see Silverman, 1986; Scott, 1992). The choice of a kernel function is considered secondary (Silverman, 1986; Scott, 1992) to the choice of the bandwidth in terms of the Mean Square Error (MSE) of the resulting estimate  $\widehat{\lambda}_y(\tau)$ . Different kernels can be made equivalent in this sense through an appropriate choice of the bandwidth. Diggle and Marron (1988) show the equivalence between density and intensity (or rate) estimation and show that the same bandwidth is optimal in both

cases under a mean square error criterion. The “plug in” or recursive bandwidth estimator due to Sheather and Jones (1991), has worked the best in our tests for kernel density estimation (Rajagopalan et al., 1995). This procedure strives to minimize the average mean integrated square error in density estimation through a data driven estimate of the pointwise bias and variance of the estimate. We used this procedure to select the bandwidth  $h_y$ .

For this study we used the Epanechnikov kernel, given as:

$$K(x) = \frac{3}{4}(1-x^2)^2 \quad |x| \leq 1 \quad \text{where } x = \frac{\tau - \tau_i}{h_y} \quad (3)$$

Periodic boundaries are used for the estimation process by (a) recognizing that dates from the end of one year can be within a bandwidth  $h_y$  of dates in the beginning of the next year, and (b) using data from year (y-1) or (y+1) for estimates on days within such a bandwidth in year y.

The intensity parameter of the nonhomogeneous Poisson process is estimated for each calendar day ( $\tau = 1, \dots, 366$ ) of each year (1, ..., y) in the historical record using the estimator in Equation (2). Weighted average precipitation for each calendar day of each year in the historical record is also estimated using the Epanechnikov weight function with a bandwidth of 90 days.

#### 4. Results

The average rate across years and the average weighted precipitation for each calendar day, estimated as described above is plotted for all the twelve stations. The x-axis on all the figures is the calendar day (i.e. 1 to 366), where 1 corresponds to January 1 and 366 to December 31 respectively. In all these figures the solid line denotes the average daily rate, and the dotted lines indicate the average weighted precipitation. The following observations are offered from the figures.



1. The average daily rate and the average weighted precipitation fluctuate in about the same way at all the stations (see Figures 1a through 1l). Thus, the use of the rate to describe seasonality seems to be a useful notion.
2. Stations in the north of the meridional transect (namely, SNP, PRR, LAK, LOG, SIL, SNC, HEB and SPF ) have similar shape of the rate and precipitation curves as can be seen from Figures 1a,1b,1c,1d,1f,1g,1h and 1i. These stations seem to have higher than average values of the rate function around the first 70 to 100 days and the last 70 to 100 days of the year, with the exact number of days varying from station to station. A similar trend is seen in the precipitation.
3. The curves of rate and precipitation are similar for stations near the southern end of the meridional transect (namely, ALT, MIA and TUS) as seen from Figures 1j,1k,1l. These stations appear to have high rates during the middle 100 days of the year and increased rates during the first and last 30 to 60 days of the year. This is prominent at ALT, and is subdued in MIA and TUS. The "wet" seasons in the north appear to correspond to "dry" seasons in the South and vice versa. This observation corresponds to the largely zonal flow driven winter/spring precipitation in the north, as opposed to the largely convective summer precipitation in the South (Ropelewski and Halpert (1986,1987)).
4. Station WOD exhibits an interesting pattern (see Figure 1e). The rate appears to be high during day 70 to 130 of the year (i.e., in spring) and is low the rest of the time. WOD lies in a rain shadow region with respect to the large scale atmospheric flow and hence gets very little precipitation during the general wet period and gets all its precipitation during the spring time due to local orographic/convective effects. There are two periods with higher than average daily precipitation at this station. One that corresponds to the high rate (day 70 through 130) and another during day 190 to 290. Apparently this station can receive high convective rainfall in the summer/fall even though the number of rainy days is low then.

### *Seasonality trends over this century*

Schneider (1995) reports that D. J. Thomson found significant changes in the timing of seasons since around 1940 in the Northern hemisphere by analyzing the 1651-1991 Central England temperature record. The seasonality of temperature in the Northern Hemisphere is determined by radiative heating which peaks on June 22, and transport of heat from other parts of the globe. The peak temperature occurs later in the year as one moves to higher latitudes in the Northern hemisphere reflecting the delay in transport of heat. Thomson's thesis is that in an atmosphere enriched by Carbon Dioxide, heating and transport of heat are more efficient, and the advance in the seasons in the Northern hemisphere is evidence of global warming.

Consequently, it was of interest to examine changes in the seasonality of precipitation along our meridional transect, as reflected by the estimated rate and average weighted precipitation amounts. We estimate the average rate for the periods before and after 1950 ( a time approximately in the middle of the data sets) at four stations with long records, which are PRP, SAN, MIA and TUS, and plot them in Figures 2a,2b, 2c and 2d respectively. In these four figures the thick line is the average rate from the entire historical record, the dotted line is the average rate from the historical record before 1950 and the dashed line is the average rate from the historical record after 1950. The average rate curves for the periods before and after 1950 are shifted from the average rate curve estimated from the entire historical record. It can be seen that the average rate after 1950 is shifted to the left (i.e., the peaks and valleys are shifted left) relative to the average rate before 1950. Similar observations can be seen from the above analysis on the average weighted precipitation amounts, in Figures 3a,3b,3c and 3d at the four stations PRN,SAN,MIA and TUS respectively.

On observing these patterns in seasonality, we decided to analyze the records to see how this shift was occurring over time, i.e., is it a sudden or continuous trend. The calendar day in each year on which the estimated rate was maximum and the date on which it was a minimum were selected. The maximum (minimum) rate at PRR/TUS occur near the end (or beginning) of the

calendar year. Thus a change in seasonality could move this date across calendar year boundaries. It is easier to analyze the transition in the date of the maximum rate at PRR and the minimum rate at TUS if we change the year boundaries away from these dates. Consequently, the date associated with maximum rate at PRR and the minimum rate at TUS is computed on a calendar year that runs from July 1 to June 30, rather than Jan. 1 to December 31. The dates for the minimum rate at PRR and the maximum rate at TUS are computed using the standard calendar.

These dates are plotted for two stations PRR and TUS (the northern and the southern extremes of our data set), in Figures 4a and 4b for maximum rate and Figures 5a and 5b for minimum rates respectively. The line in these figures is a nonparametric smooth fitted by LOWESS (Cleveland, 1979). One can see that the date for both the maximum and minimum rates has a decreasing trend with year. The nonparametric Mann-Kendall test (Gilbert (1987)) for monotonic trend showed that these trends were significant (p-values in all cases were of the order of  $e^{-10}$ ). Robust estimates of the linear trend, the Sen slopes (see Gilbert (1987)) range from -0.33 to -1 days per year. We performed the above analysis with the average weighted precipitation and a similar behaviour was observed. Results are not presented for brevity. It is rather curious that the march of seasons as measured by the precipitation rate and also the average weighted precipitation is advancing at these sites at roughly a constant rate over the whole record.

## 5. Closure

The nonparametric methods presented here were shown to be useful for identifying seasonal variations in precipitation occurrence as a function of latitude and also for variations in seasonality across years. For the data sets analyzed, we remarkable differences were seen in the timing and duration of the precipitation seasons along the meridional transect selected west of the Rockies. An interesting trend in the seasonality across the sites was also identified. If this trend is related to global warming it has important implications for the form of precipitation in these areas,

and also for crop water requirements in the growing season. Further investigation of such trends and their relationship to atmospheric circulation is warranted.

### Acknowledgments

Partial support of this work by the U.S. Forest Service under contract notes, INT-915550-RJVA and INT-92660-RJVA, Amend #1 is acknowledged. The principal investigator of the project is D.S. Bowles.

## References

- Cleveland, W.S, Robust locally weighted regression and smoothing scatter plots, *J. Amer. Statist. Assoc.*, 74, 829-836, 1979.
- Cox, D.R. and V. Isham, *Point Processes*, Chapman and Hall, London, 1980.
- Diggle, P.J, A kernel method for smoothing point-process data, *Applied Statistics*, 34, 138-147.
- Diggle, P.J. and J.S. Marron, Equivalence of smoothing parameters selectors in density and intensity estimation, *J. Amer. Statist. Assoc.*, 83, 793-800.
- Gilbert, R.O, *Statistical methods for environmental pollution monitoring*, Van Nostrand Reinhold Company, New York, 1987.
- Rajagopalan, B. U. Lall and D.G. Tarboton, Simulation of Daily Precipitation from A Nonparametric Renewal Model, *UWRL Working Paper WP-93-HWR-UL/003*, Utah State University, 1993.
- Ropelewski, C.F. and M.S. Halpert, North American precipitation and temperature patterns associated with the El Niño/Southern Oscillation (ENSO), *Monthly Weather Review*, 114, 2352-2362, 1986.
- Ropelewski, C.F. and M.S. Halpert, Global and regional scale precipitation patterns associated with the El Niño/southern oscillation, *Monthly Weather Review*, 115, 1606-1626, 1987.
- D. Schneider, Global Warming Is Still a Hot Topic: Arrival of the seasons may show greenhouse effect, *Scientific American*, 272(2), 13, 1995.
- Scott, D.W., *Multivariate density estimation: Theory, Practice and Visualization*, Wiley series in probability and mathematical statistics, John Wiley and Sons, New York, 1992.
- Sheather, S.J. and M.C. Jones, A reliable data-based bandwidth selection method for kernel density estimation, *Journal of the Royal Statistical Society*, B. 53, 683-690, 1991.
- Silverman, B.W., *Density estimation for statistics and data analysis*, Chapman and Hall, New York, 1986.

Solow, A.R., The nonparametric analysis of point process data: the freezing history of lake konstanz, *Journal of Climate*, 4, 116-119, 1991.

Waymire, E. and V.K. Gupta, The mathematical structure of rainfall representations, 2, A review of the theory of point processes, *Water Resour. Res.*, 17(5), 1273-1286, 1981a.

TABLE 1 Data Sets Analyzed

	Latitude	Longitude	Elevation (ft.aboveMSL)	Record Length
Priest River, Idaho (PRR)	48° 21' N	116° 50' W	2380	1911-1992
Sandpoint, Idaho [SNP]	48° 17' N	116° 34' W	2100	1910-1992
Laketown, Utah [LAK]	41° 49' N	111° 19' W	5980	1948-1992
Logan, Utah [LOG]	41° 45' N	111° 48' W	4790	1928-1992
Woodruff, Utah [WOD]	41° 32' N	111° 09' W	6320	1948-1992
Silverlake, Utah [SIL]	40° 36' N	111° 35' W	8740	1948-1992
Snake Creek, Utah [SNC]	40° 33' N	111° 30' W	6010	1928-1992
Heber, Utah [HEB]	40° 30' N	111° 25' W	5630	1928-1992
Spanish Fork, Utah [SPF]	40° 05' N	111° 36' W	4720	1932-1992
Alton, Utah [ALT]	37° 26' N	112° 29' W	7040	1929-1992
Miami, Arizona [MIA]	33° 24' N	110° 53' W	3560	1914-1992
Tucson, Arizona [TUS]	32° 15' N	110° 57' W	2440	1901-1992

Note: All data sets were obtained from Earth Info, CD-ROM

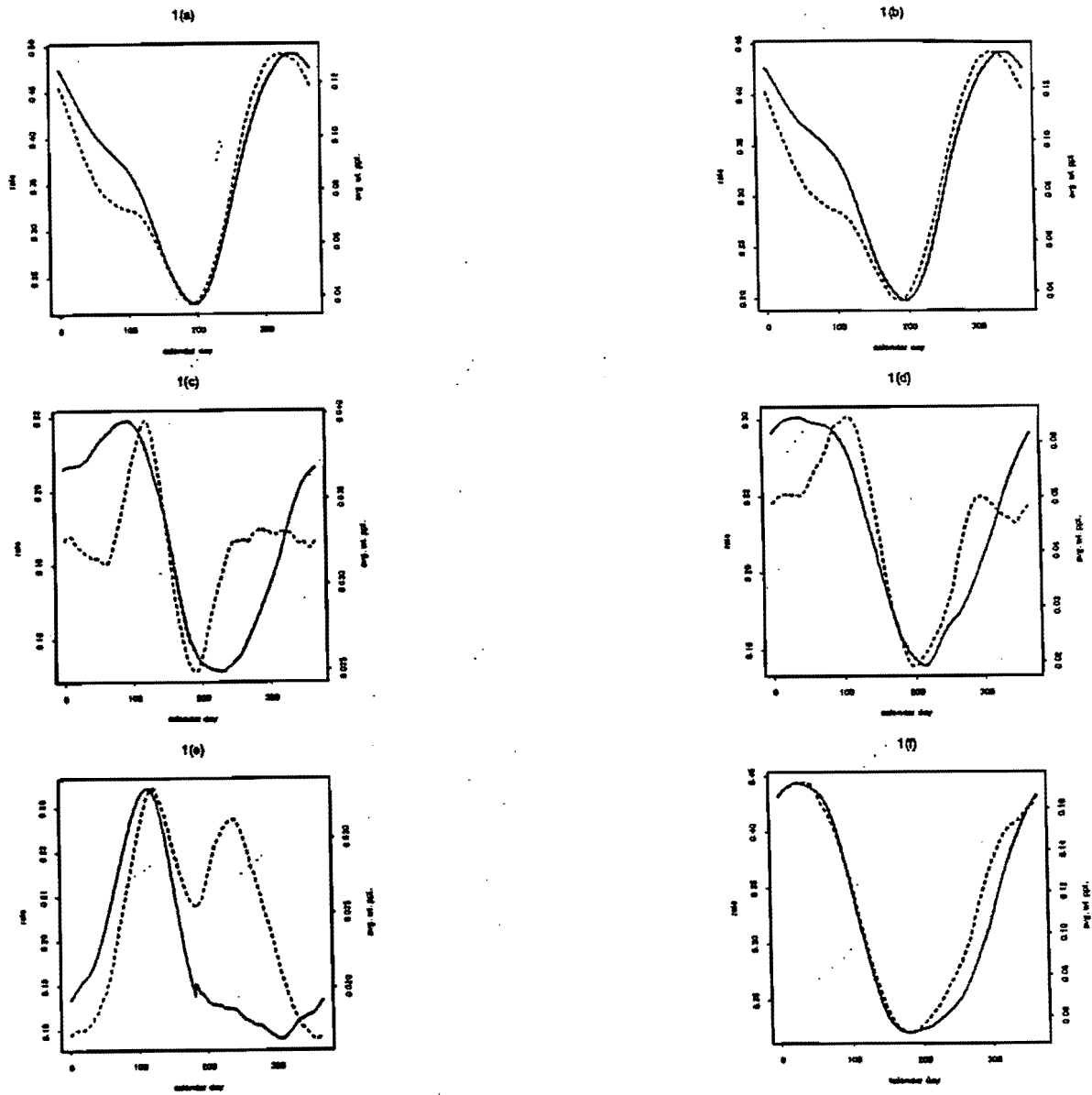


Figure 1. Average daily rate (solid line) and average weighted precipitation (dotted line) for each calendar day, at (a) at Priest River, ID, (b) at SandPoint, ID, (c) Laketown, UT, (d) Logan, UT, (e) Woodruff, UT, (f) Silverlake, UT, (g) Snake Creek, UT, (h) Heber, UT, (i) Spanishfork, UT, (j) Alton, UT, (k) Miami, AZ and (l) Tucson, AZ.



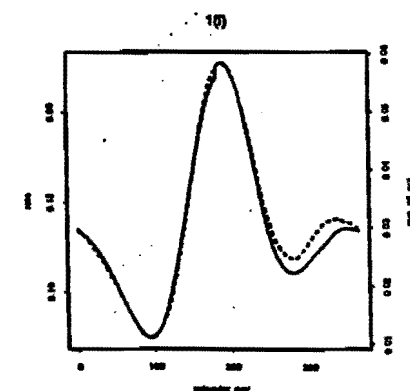
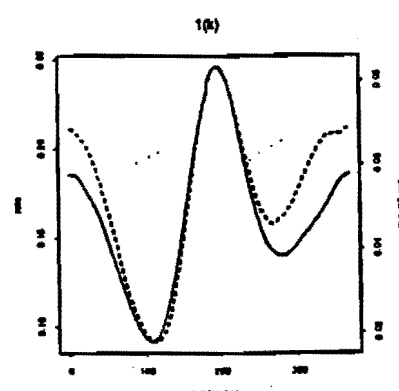
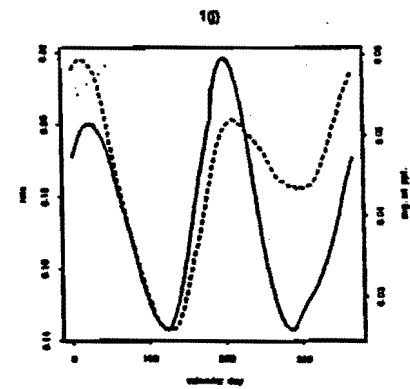
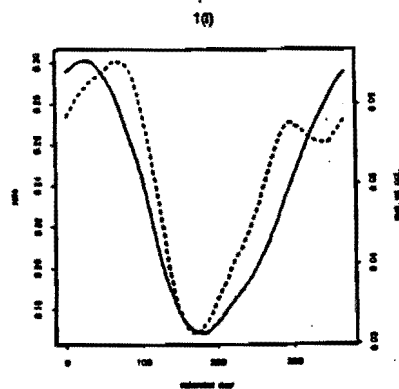
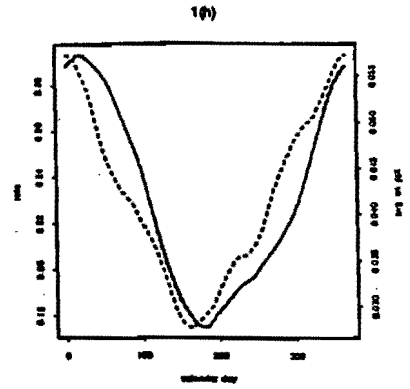
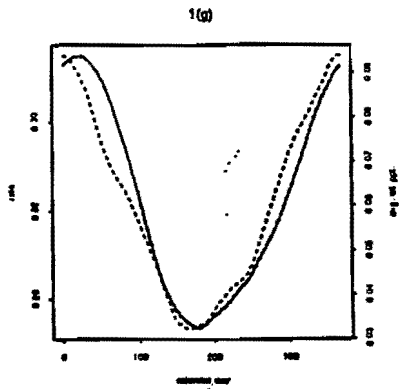
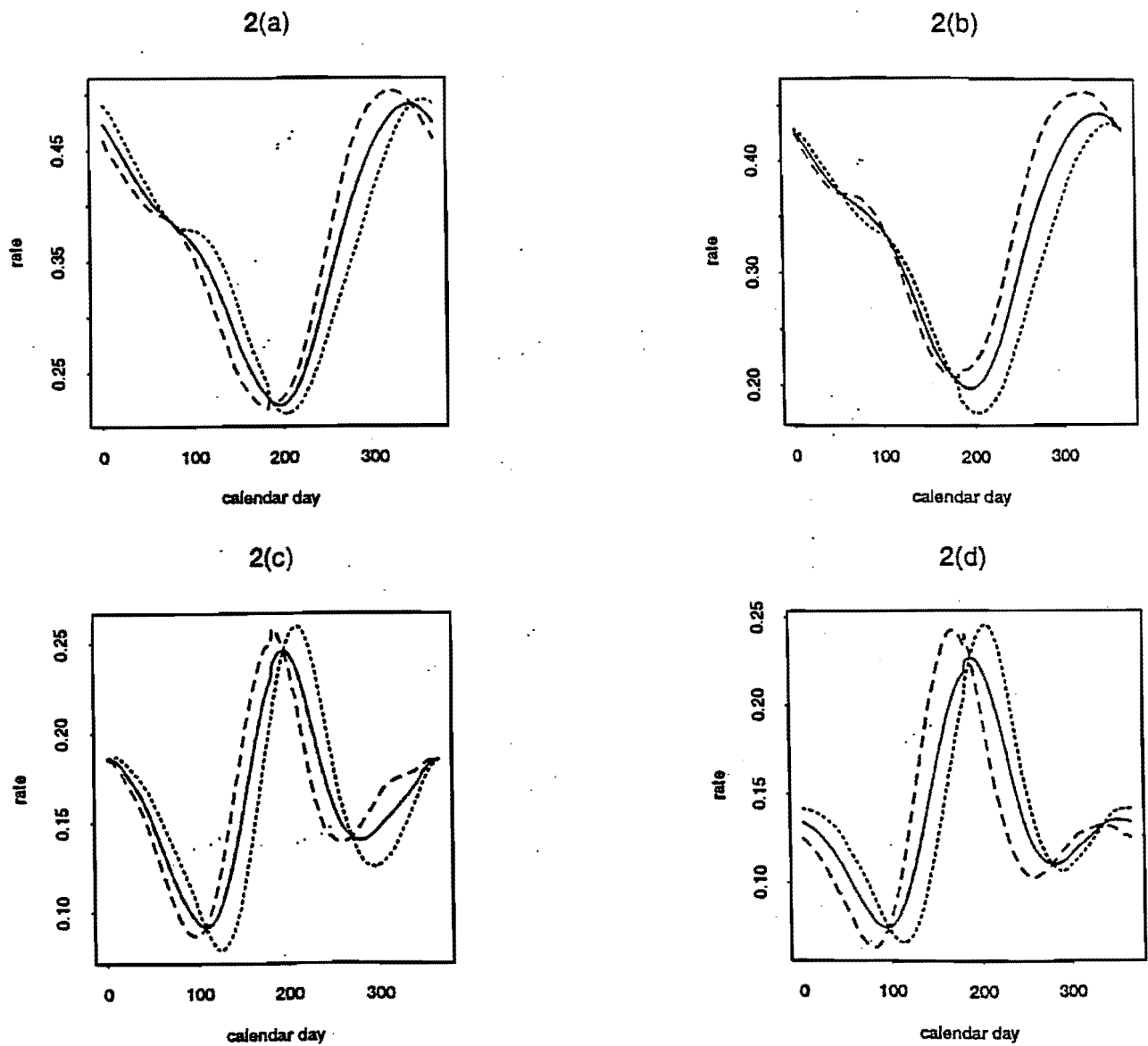


Figure 1. (Cont'd)



**Figure 2.** Average daily rate from the entire historical record (solid line), from the historical record before 1950 (dotted line) and from the historical record after 1950 (dashed line), at (a) Priest River, ID, (b) SandPoint, ID, (c) Miami, AZ and (d) Tucson, AZ.

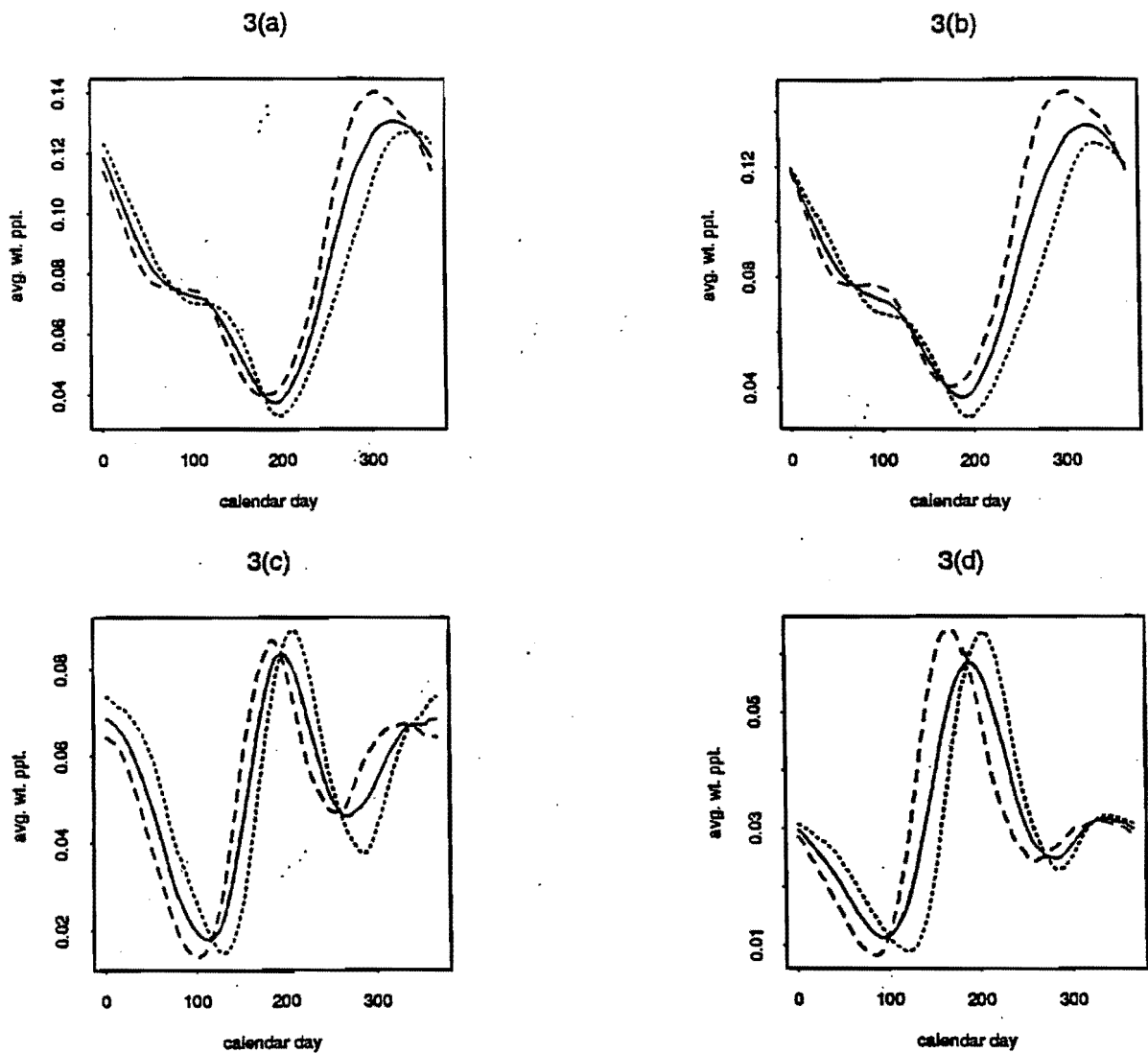


Figure 3. Average weighted precipitation from the entire historical record (solid line), from the historical record before 1950 (dotted line) and from the historical record after 1950 (dashed line), at (a) Priest River, ID, (b) SandPoint, ID, (c) Miami, AZ and (d) Tucson, AZ.

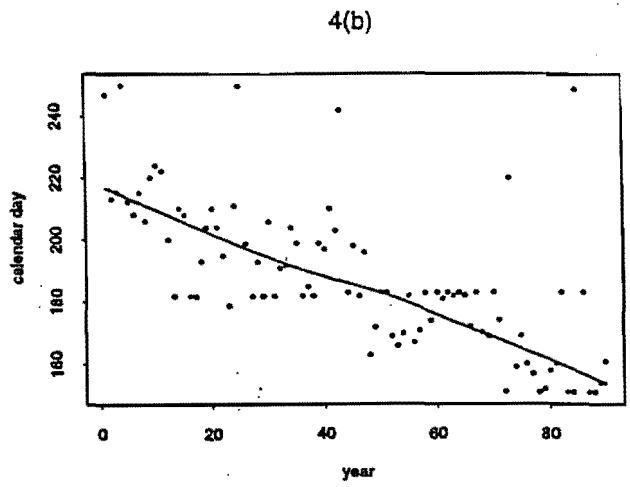
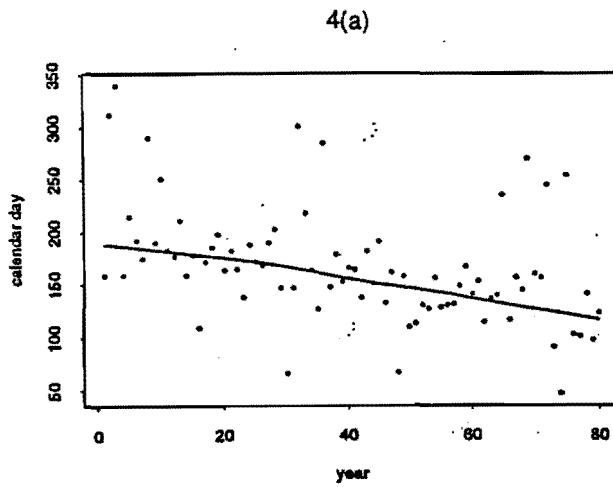


Figure 4. Calendar date of maximum estimated average daily rate in each year (dots), along with a LOWESS smooth (thick line), at (a) Priest River, ID and (b) Tucson, AZ.

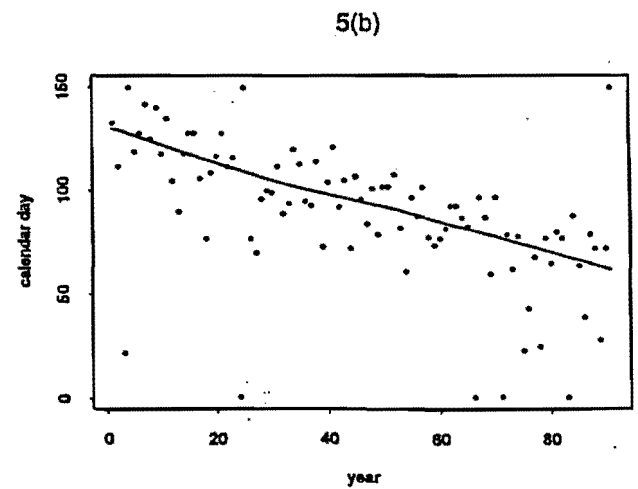
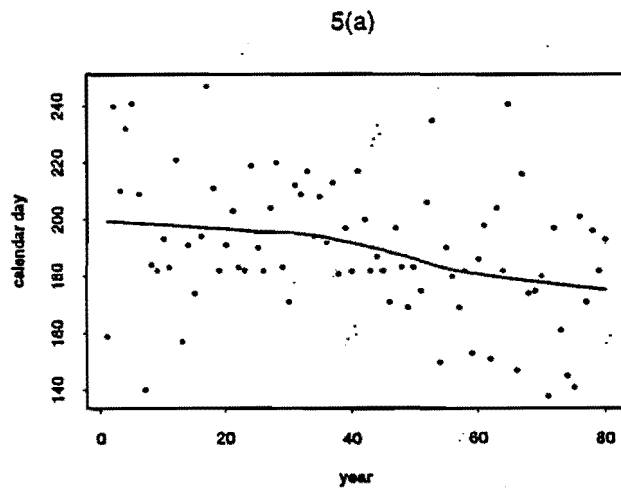


Figure 5. Calendar date of minimum estimated average daily rate in each year (dots), along with a LOWESS smooth (thick line), at (a) Priest River, ID and (b) Tucson, AZ.

# APPENDIX 4D

## A Nonhomogeneous Markov Model for Daily Precipitation Simulation

Balaji Rajagopalan, Upmanu Lall and David G. Tarboton

Utah Water Research Laboratory  
Utah State University  
Logan, Utah, UT-84322-8200

### ABSTRACT

We present a one step nonhomogeneous Markov model for describing daily precipitation at a site. Daily transitions between wet and dry states are considered. The one step, 2x2 transition probability matrix is presumed to vary smoothly day by day over the year. The daily transition probability matrices are estimated nonparametrically. A kernel estimator is used to estimate the transition probabilities through a weighted average of transition counts over a symmetric time interval centered at the day of interest. The precipitation amounts on each wet day are simulated from the kernel probability density estimated from all wet days that fall within a time interval centered on the calendar day of interest over all the years of available historical observations. The model is completely data driven. An application to data from Utah is presented. Wet and dry spell attributes (specifically the historical and simulated probability mass functions (p.m.fs) of wet and dry spell length) appear to be reproduced in our Monte Carlo simulations. Precipitation amount statistics are also well reproduced.

## 1. INTRODUCTION

Markov chains (Gabriel and Neumann, 1962; Todorovic and Woolhiser, 1975; Smith and Schreiber, 1973) have been a popular method for modeling daily precipitation occurrence. Typically a two state (wet or dry), one step model is used, and the state transition probabilities (e.g., transition from wet a day to a wet day, wet day to a dry day) are estimated from the data. One problem with such a description is that the transition probabilities may vary over the year, i.e., the process of precipitation occurrence is nonstationary.

Two approaches are commonly used to address this problem. In the first approach, the year is divided into periods (or seasons) and the transition probabilities are estimated separately for each period. There is an implicit assumption that the occurrence process is stationary over the period. This assumption may not be tenable. The second approach is to consider essentially a nonhomogeneous Markov process by allowing the transition probabilities to vary systematically over the year, and to model such a variation through a Fourier series expansion (Feyerherm and Bark (1965), Woolhiser et al. (1973) and Woolhiser and Pegram (1979)). This can be an effective approach where adequate data is available, and the seasonality in the precipitation process can be captured by a few Fourier series terms. Our nonparametric analyses (Rajagopalan and Lall (1995)) of the seasonality of precipitation for stations along a meridional transect in the Western United States, suggests that sometimes the number of Fourier series terms needed may be large relative to the amount of data available.

In this paper, a nonhomogeneous Markov (NM) model is presented that uses kernel methods to estimate a nonhomogeneous transition probability matrix, and to estimate a corresponding nonstationary probability density function (p.d.f) of daily precipitation amount. Kernel methods are local, weighted averages of the target function (relative frequency of occurrence in this case). Since they are capable of approximating a wide variety of target functions

with asymptotically vanishing error, and use only data from a "small" neighborhood of the point of estimate, they are considered nonparametric. Fourier series methods are shown to be a subset of kernel methods by Eubank (1988, secs. 3.4 and 4.1). A review of hydrologic applications of nonparametric function estimation methods is provided by Lall (1995).

A brief description of the Markov chain and its terminology is first presented as a background to motivate our formulation. The general structure of the NM model proposed is next outlined with the nonparametric estimators for the transition probabilities. The simulation procedure is then outlined. Results from an application of the model to a precipitation data from Utah follow. Musings on the results and discussion on limitations of the approach conclude the paper.

## 2. BACKGROUND

The basic assumption in a two state Markov Chain model is that the present state (wet or dry) depends only on the immediate past. The transition probabilities for transitions (i.e., WW, WD, DW, DD) between the two states (W or D) are estimated directly from the data through a counting process. Two elements of the transition probability matrix are the probability of a dry day following a wet day,  $P_{WD} = a_1$ , and the probability of a wet day following a dry day,  $P_{DW} = a_2$ . The other probabilities, probability of a wet day following a wet day,  $P_{WW}$  and the probability of a dry day following a dry day,  $P_{DD}$  are  $(1 - a_1)$  and  $(1 - a_2)$  respectively.

Seasonal variations in the transition probabilities can be accounted for by expressing the changing transition probabilities through a Fourier series (Woolhiser and Pegram, 1979; Roldan and Woolhiser, 1982). As an illustration, the transition probability  $P(WD)$  can be expressed as:

$$P_{WD}(t) = \bar{P}_{WD} + \sum_{k=1}^m c_k \sin(2\pi tk/365 + \theta_k); \quad t = 1, 2, \dots, 365 \quad (1)$$

where  $m$  = the maximum number of harmonics required to describe the seasonal variability of the transition probability,  $\bar{P}_{WD}$  is the annual mean value of the parameter,  $c_k$  is the amplitude, and  $\theta_k$  is the phase angle in radians for the  $k$ th harmonic.

The means, amplitudes, and phase angles are estimated by numerical optimization of the log likelihood function, as described by Woolhiser and Pegram (1979) and Roldan and Woolhiser (1982). Fourier series representations of parameters of a first-order Markov chain for precipitation have been used (among others) by Feyerherm and Bark (1965) who used least squares techniques for parameter estimation and by Stern and Coe (1984) who formulated the estimation problem as a generalized linear model to obtain maximum likelihood estimators.

The degree of dependence in time is limited by the order (i.e., the number of past days the present state is presumed to depend on) of the Markov chain. Feyerherm and Bark (1967) and Chin (1977) suggest that the order may need to be seasonally variable as well. Lack of parsimony is a drawback of MC models as the order is increased. A number of researchers (Hopkins and Robillard (1964), Haan et al (1976), Srikanthan and McMahon (1983), Guzman and Torrez (1985)) have also stressed the need for multistate MC models that consider the dependence between transition probabilities and rainfall amount. In this paper, we shall consider only a two state, first order Markov Chain. Extensions to other situations follow in the same spirit.

### 3. MODEL FORMULATION

The NM model that we present allows the one step transition probability matrix to change over each day thus capturing the day to day variation in the occurrence process in a natural manner. The daily transition probability matrices are estimated using a discrete kernel estimator, which we describe in the following section. Daily precipitation occurrence sequences are then simulated using the transition probability matrices. To complete the model, precipitation amounts on each wet day



are simulated from the nonparametric probability density estimated from all wet days that fall within a time interval or bandwidth centered on the calendar day of interest over all the years of available historical record. The model is completely data driven.

### 3.1 Transition Probabilities and their Estimation

The precipitation occurrence process is shown in figure 1. From the daily precipitation record we can obtain four types of data, (for illustration refer to figure 1) which are, (1) the day indices  $t_{w1}, t_{w2}, \dots, t_{wnw}$  of  $nw$  wet days; (2) the day indices  $t_{d1}, t_{d2}, \dots, t_{dnd}$  of  $nd$  dry days; (3) the day indices  $t_{wd1}, t_{wd2}, \dots, t_{wdnwd}$  of the  $nwd$  days on which a transition occurs from wet to dry, meaning days  $t_{wd1}, t_{wd2}, \dots$  are wet and days  $t_{wd1}+1, t_{wd2}+1 \dots$  are dry; (4) the day indices  $t_{dw1}, t_{dw2}, \dots, t_{dwnw}$  of the  $ndw$  days on which a transition occurs from dry to wet, meaning days  $t_{dw1}, t_{dw2}, \dots$  are dry and days  $t_{dw1}+1, t_{dw2}+1 \dots$  are wet. A day index refers to a number between 1 to 366, representing the calendar day of the observation. From these we estimate the transition probabilities  $P_{wd}(t)$  (probability of transition from a wet day on calendar day  $t$  to a dry day on calendar day  $t+1$ ),  $P_{dw}(t)$  (probability of transition from a dry day on calendar day  $t$  to a wet day on calendar day  $t+1$ ). The other two transition probabilities (namely  $P_{ww}(t)$  and  $P_{dd}(t)$ ) can be estimated directly from the relations  $P_{wd}(t) + P_{ww}(t) = 1$  and  $P_{dw}(t) + P_{dd}(t) = 1$ . The transition probabilities for calendar day  $t$  are estimated from the data using discrete nonparametric kernel estimators.

For a traditional Markov chain the transition probabilities are estimated simply as the ratio of the number of transitions in the historical record to the number of wet or dry days in the historical record, as appropriate. Here, we try to localize such estimates about the calendar day of interest using kernel estimators. The general idea is that the events (i.e., a wet or dry day, or a state transition) occurring near the calendar day of interest should be given more weightage while the

ones further away should be given a lower weightage. The resulting kernel estimators for the transition probabilities  $P_{wd}(t)$  and  $P_{dw}(t)$  are given as:

$$\hat{P}_{wd}(t) = \frac{\sum_{i=1}^{n_{wd}} K\left(\frac{t - t_{wdi}}{h_{wd}}\right)}{\sum_{i=1}^{n_w} K\left(\frac{t - t_{wi}}{h_{wd}}\right)} \quad (2)$$

$$\hat{P}_{dw}(t) = \frac{\sum_{i=1}^{n_{dw}} K\left(\frac{t - t_{dwi}}{h_{dw}}\right)}{\sum_{i=1}^{n_d} K\left(\frac{t - t_{di}}{h_{dw}}\right)} \quad (3)$$

where  $n_{wd}$  is the number of transitions in the historical record from wet day to dry day,  $n_{dw}$  is the number of transitions in the historical record from dry day to wet day,  $n_d$  is the number of dry days in the historical record,  $n_w$  is the number of wet days in the historical record,  $K(\cdot)$  is the kernel function (or weight function) and  $h(\cdot)$  is a kernel bandwidth,  $t$  is the calendar day of interest and the  $t_{(\cdot)}$ 's have the definitions described earlier. Note that the estimates on any calendar day  $t$  are obtained by using the information from days in the range  $[t - h(\cdot), t + h(\cdot)]$ . Note that the definition of calendar dates is periodic, i.e. day 365 and day 1 are recognized as 1 day apart for a non-leap year. The contribution to the estimate of an event that lies within this range is determined by the kernel or weight function  $K(\cdot)$ , that is described below.

Since we have a discrete situation (i.e. each day being discrete) we use the discrete kernel developed by Rajagopalan and Lall (1995) as:

$$K(x) = \frac{3h}{(1-4h^2)}(1 - x^2) \quad \text{for } |x| \leq 1 \quad (4)$$

where  $x = (t - t_{(.)})/h_{(.)}$ , measures how far an event  $t_{(.)}$ , that lies within a bandwidth  $h_{(.)}$  of the day  $t$ , is from  $t$ ; and  $h_{(.)}$  is an integer.

The kernel in (3) was derived from the consideration that the sum of all weights ascribed to events that lie within a bandwidth  $h_{(.)}$  of  $t$  sum to 1, i.e.,  $\sum_{x=-1}^1 K(x) = 1$ ; that the weights be symmetric on either side of  $t$ , i.e.  $\sum_{x=-1}^1 xK(x) = 0$ ; that each weight be positive; and the resulting estimate of probability have minimum mean square error.

The estimators in equations 2 and 3 are fully defined once the respective bandwidths are specified. We choose the bandwidth using the Least Squared Cross Validation (LSCV) procedure (Scott, 1992, p. 225), where the bandwidth is chosen that minimizes a LSCV function which is given as

$$\text{LSCV}(h) = \frac{1}{n} \sum_{i=1}^n (1 - \hat{P}_{-t_i}(t_i))^2 \quad (5)$$

where  $\hat{P}_{-t_i}(t_i)$  is the estimate of the transition probability ( $\hat{P}_{wd}$  or  $\hat{P}_{dw}$ ) on day  $t_i$  dropping the information on day  $t_i$ ,  $n$  is the number of observations ( $n_{dw}$  or  $n_{wd}$ ). Here we assume a prior probability of transition to be 1 on the days on which transitions have occurred hence the 1 in the equation 5. The bandwidth is searched from 1 to 182 (length of half year). Once the transition probabilities are estimated for each day in the historical record the simulation of the precipitation occurrence for each day using the transition probability matrix of the previous day is possible.

### 3.2 Precipitation amount generation

Precipitation amounts for the wet days are generated from a kernel probability density estimated from all wet days that fall within a time interval or bandwidth centered on the calendar

day of interest over all the years of historical record. This amounts to two steps (1) choosing the time interval or bandwidth and (2) generating from the kernel estimated p.d.f.

An appropriate bandwidth for localizing the estimate of the probability density of precipitation amount may be obtained by determining the bandwidth appropriate for estimating the probability that a day is wet. If the probability of daily precipitation is low, the precipitation data will be sparse, and the bandwidth needed for stabilizing the variance of the estimated probability distribution of precipitation will be large. Conversely, as the probability of daily precipitation is high, a large number of days with precipitation will occur and the bandwidth needed to localize the estimate can be smaller.

Consequently, we first consider the smoothing of the proportion of wet days ( $p_t = n_t/NT$ ,  $n_t$  is the number of times calendar day  $t$  was wet;  $NT$  is the total number of calendar day  $t$  in the historical record) on each calendar day  $t = 1, 2, \dots, 366$ . These raw proportions are smoothed using the discrete kernel (DK) estimator of Rajagopalan and Lall (1995) which in this case is:

$$\hat{p}_t = \sum_{j=1}^{366} K\left(\frac{t-j}{h_p}\right) p_j \quad (6)$$

where  $K(\cdot)$  is the discrete kernel as defined by equation 3, and  $h_p$  is the bandwidth that we are interested in. The bandwidth  $h_p$  can be obtained using the LSCV procedure similar to equation (5) as given by Rajagopalan and Lall (1995) as:

$$\text{LSCV}(h_p) = \sum_{t=1}^{366} (\hat{p}_t)^2 - 2 \sum_{t=1}^{366} \hat{p}_{-t} p_t \quad (7)$$

where,  $\hat{p}_{-t}$  is the estimate of the calendar day  $t$ , by dropping the information on that day.

Once we estimate the time interval  $h_p$  the next step is to pick the precipitation amounts on all the wet days that fall within the time interval  $h_p$  from the day of interest in all the years of the historical record. Let us say that the precipitation amounts so picked from the historical records are  $y_1, y_2, \dots, y_{np}$  and  $t_1, t_2, \dots, t_{np}$  are the corresponding calendar day index. The task now is to generate precipitation amount for the calendar day  $t$ , which is a wet day. This can be accomplished by fitting a conditional p.d.f  $f(y|t)$  (see equation 10) and then simulating from it. This step is carried out for each wet day that is simulated. Before describing the simulation procedure we introduce a kernel density estimator for continuous variables which is given as:

$$\hat{f}(y) = \frac{1}{h_y np} \sum_{i=1}^{np} K_c\left(\frac{y - y_i}{h_y}\right) \quad (8)$$

where  $K_c(\cdot)$  is a univariate, continuous kernel, and  $h_y$  is the bandwidth. Here we use the Epanechnikov kernel given by :

$$\begin{aligned} K_c(x) &= 0.75(1 - x^2) \quad \text{for } |x| \leq 1 \\ &= 0. \quad \text{otherwise} \end{aligned} \quad (9)$$

where  $x = \frac{y - y_i}{h_y}$ . For a detailed exposition of kernel density estimation for continuous variables and issues relating to bandwidth selection we refer the reader to Silverman, (1986), Scott (1992), and for kernel density estimation methods with specific application to precipitation modeling we refer to Lall et al. (1995) and Rajagopalan et al. (1995).

A logarithmic transform of the precipitation data prior to density estimation is often considered. Such a transformation is also attractive in the kernel density estimation (k.d.e) context. Since it can provide an automatic degree of adaptability of the bandwidth (in real space).

This alleviates the need to choose variable bandwidths with heavily skewed data, and also alleviates problems that the k.d.e. has with p.d.f. estimates near the boundary (e.g.,, the origin) of the sample space. The resulting k.d.e. can be written as:

$$\hat{f}(y) = \frac{1}{np} \sum_{i=1}^{np} \frac{1}{h_{LY}} K_c\left(\frac{\log(y) - \log(y_i)}{h_{LY}}\right) \quad (10)$$

where  $h_{LY}$  is the bandwidth of the log transformed data. This is chosen using a recursive approach due to Sheather and Jones (1991) (SJ) to minimize the Mean Integrated Square Error (MISE) and recommended by Rajagopalan et al. (1995) typically for precipitation data.

The two step procedure discussed above can be more formally considered through the conditional p.d.f.  $\hat{f}(y|t)$ , defined using a product kernel representation as:

$$\hat{f}(y|t) = \frac{1}{y h_{LY}} \sum_{i=1}^{np} K_c\left(\frac{\log(y) - \log(y_i)}{h_{LY}}\right) K\left(\frac{t - t_i}{h_p}\right) / \sum_{i=1}^{np} K\left(\frac{t - t_i}{h_p}\right) \quad (11)$$

Equation 11, states that the conditional probability density of a rainfall amount  $y$  on calendar day  $t$  is obtained by considering a window of width  $h_p$  centered at  $t$ , weighting the precipitation amounts on wet days that fall within this window using the kernel  $K(\cdot)$ , and then forming a density estimate by further weighting these amounts with the kernel  $K_c(\cdot)$ . Strictly speaking, the bandwidths  $h_p$  and  $h_{LY}$  should be chosen by optimizing a criteria relevant to the conditional density. The description of our procedure given earlier shows that we are essentially choosing these bandwidths independently. McLachlan (1992, p. 306-308), discusses the simultaneous selection of bandwidths in each coordinate, versus the use of the optimal univariate bandwidths in each direction. It is not clear that the additional effort of simultaneous selection of the two bandwidths is justified. Consequently, we choose the bandwidths  $h_{LY}$  and  $h_p$  by the

methods described for the univariate case. Rajagopalan et al (1995) show that bandwidths selected in this way are often satisfactory. For simulation from the kernel estimated p.d.f. (such as equation (11)) it is not necessary to explicitly estimate the density  $\hat{f}(y|t)$ . The estimation of the bandwidths  $h_{LY}$  and  $h_p$  and subsequent perturbation of the historical data is sufficient.

### 3.3 Simulation Procedure

The simulation procedure from the NM model can be described in the following steps.

1. From the historical precipitation sequence evaluate the transition probabilities ( $P_{wd}(t)$ ,  $P_{ww}(t)$ ,  $P_{dw}(t)$  and  $P_{dd}(t)$ ) for each calendar day  $t$  using the estimators described in section 3.1. Similarly evaluate the probability density function for precipitation amount on day  $t$  using the procedure described in section 3.2.
2. Start the simulation with a wet or dry day (deciding by generating a uniform random number  $U$  in  $[0,1]$ , if  $U \leq 0.5$  then wet else dry).
3. The precipitation state for the next day is simulated from the transition probability matrix for the current day (as estimated in step 1).
4. Precipitation amounts on wet days are generated following the process illustrated in figure 2, that is described below:
  - (i) Pick all the wet day precipitation amounts (e.g.,  $y_1, y_2, \dots, y_{np}$ ) from all the years in the historical record that fall within the window  $h_p$  centered on the corresponding calendar day of interest and also the corresponding calendar day indices  $t_1, t_2, \dots, t_{np}$ .
  - (ii) For the calendar day of interest, pick a historical wet day to perturb using the bandwidth  $h_p$  and the kernel  $K(x)$  to specify the resampling metric. Recall that the kernel function describes the weight given to each calendar day that lies within  $h_p$  of calendar day  $t$ , that depend on the "distance" between the two dates relative to the bandwidth  $h_p$ , and the kernel function given in equation (4). Let the weights associated with each of  $np$  wet days that are thus identified be

$wt_1, wt_2, \dots, wt_{np}$ . Now generate a random integer  $j$  between 1 and  $np$  from a probability metric given by these weights.

(iv) The simulated precipitation amount is  $y^* = \exp(\log(y_j) + U h_L Y)$  where  $y_j$  is the precipitation on the historical day point picked to be perturbed. The random variate  $U$  is generated from the probability density corresponding to the kernel function  $K_C(\cdot)$ . As mentioned earlier, we have used the Epanechnikov kernel in this study and simulation from this kernel is easily accomplished using the two step procedure described in Silverman (1986, p. 143)

5. The process (steps 3 and 4) is repeated day by day until the desired length of record is generated.

#### 4. MODEL APPLICATION

The model described was applied to daily rainfall data from Salt Lake City in Utah. Thirty years of daily weather data was available from the period 1961-1991. Salt Lake City is at  $40^{\circ}46'$  N latitude,  $111^{\circ}58'$  W longitude and at an elevation of 1288 m. Most of the precipitation comes in the form of winter snow. Rainfall occurs mainly in Spring, with some in Fall.

We shall first list some measures of performance that were used to compare the historical record and the model simulated record, and then outline the experimental design. The aim here is to capture the frequency structure of the events (i.e. the underlying p.d.f) which then amounts to the reproduction of all the statistics. By events we mean the wet spell lengths, dry spell lengths and the wet day precipitation. The wet and dry spell lengths are defined as the successive wet or dry days. Clearly the wet spell lengths and dry spell lengths are defined through the set of integers greater than 1. We look at the model performance both at the seasonal scale and the annual scale. For the seasonal scale comparison we have the year divided into four seasons : Winter or Season 1 (Jan - Mar), Spring or Season 2 (Apr - Jun), Summer or Season 3 (Jul - Sep), and Fall or Season 4 (Oct - Dec).



#### 4.1 *Performance measures*

1. Probability mass function of wet spell length, dry spell length and probability density function of wet day precipitation in each season and annual.
2. Mean of wet spell length, dry spell length and wet day precipitation in each season and annual.
3. Standard deviation of wet spell length, dry spell length and wet day precipitation in each season and annual.
4. Length of longest wet spell and dry spell in each season and annual.
5. Maximum wet day precipitation in each season and annual.
6. Percentage of yearly precipitation in each season and annual.
7. Fraction of wet and dry days in each season annual.

#### 4.2 *Experiment design*

Our purpose here is to test the utility of the NM model. The main steps involved in this are:

1. Thirty sets of synthetic records of thirty years each (i.e. the historical record length) are simulated using the NM model.
2. The statistics of interest are computed for each simulated record, for each season, and are compared to statistics of the historical record using boxplots. The p.m.f.'s of wet and dry spell lengths are estimated using the Discrete Kernel estimator of Rajagopalan and Lall (1995) (same as the estimator in equation (6)) and the p.d.fs of the wet day precipitation is estimated using the estimator in equation (10). The statistics listed in section 4.1 are computed for the simulated record and compared with those of the historical record.

## 5. RESULTS

In this section we present comparative results (using the performance measures listed in section 4.1) of the NM model for the Salt Lake City data. The p.m.f.s/p.d.f.s of the simulated records are compared with those for the historical record using boxplots while other statistics are summarized in Tables 1,2 and 3. A box in the boxplots (e.g., Figure 3) indicates the interquartile range of the statistic computed from thirty simulations, the line in the middle of the box indicates the median simulated value. The solid lines correspond to the statistic of the historical record. The boxplots show the range of variation in the statistics from the simulations and also show the capability of the simulations to reproduce historical statistics. The plots of the p.d.f.'s are truncated to show a common range across seasons and to highlight differences near the origin (mode).

Figure 3 shows the boxplots of kernel estimated p.d.f.'s of simulated data of wet day precipitation and the historical data. It can be seen that the historical p.d.f.'s are very well reproduced by the simulations in all the four seasons. The other statistics are also seen to be well reproduced by the model for all the seasons and also annual, as can be noticed from table 1.

Boxplots of kernel estimated p.m.f.'s of simulated data of wet spell length are found to enclose the p.m.f of the historical data of wet spell length for all the four seasons in figure 4 and for the annual in figure 5. The other statistics are also preserved quite well by the simulations, as seen from table 2. Good performance of the model in reproducing the dry spell statistics can be seen from figures 6 and 7 and also from table 3. The coefficient of skew, and the coefficient of variation, the 25% quantile and the 75% quantile were also preserved for all the three variables, but are not shown here.

The extreme statistics (e.g., longest spell length or maximum wet day precipitation) exhibit a high degree of variability in the simulations (refer tables 1,2 and 3) and an asymmetric sampling distribution, as one would expect.

Note that none of the statistics that we have listed in section 4.1 are explicitly or implicitly considered in the model. Hence the good reproduction of p.d.f.'s/p.m.f.'s of the three variables, is quite heartening.

## 6. SUMMARY AND CONCLUSIONS

A nonhomogeneous Markov model for simulating daily precipitation is presented in this paper. The traditional Markov chain model is extended to consider the a smooth variation in the transition probabilities from day to day, thus attempting to capture the nonstationarity in the precipitation occurrence process. The 2x2 daily transition probability matrix is estimated nonparametrically. The primary intended use of the model is as a simulator that is faithful to the historical data sequence, obviating the need to divide the year into seasons and subsequently fitting the Markov chain parameters separately for each season. Simulations from the model are shown to preserve the frequency structure ( p.d.f/p.m.f) of the wet spell length, dry spell length and wet day precipitation at both the seasonal and annual time scales.

In many cases, the Fourier series approach to addressing seasonal variation in Markov Chain parameters may be just as effective. Recall that the Fourier series approach can be shown to be a subset of the kernel approach with a specific kernel choice. The kernel approach presented here is attractive because it is relatively parsimonious, locally adaptive, and extends quite naturally to localizing the probability density estimation for precipitation amount as well. Extensions to higher order chains or those with more states follow directly. One needs to define the appropriate events as was done here and go through the solution of the corresponding smoothing problem.

## ACKNOWLEDGMENTS

Partial support of this work by the U.S. Forest Service under contract notes, INT-915550-RJVA and INT-92660-RJVA, Amend #1 is acknowledged. The principal investigator of the project is D.S. Bowles. Thanks are due to Alaa Ibrahim Ali for useful discussions.

## REFERENCES

- Chin, E.H., Modeling daily precipitation occurrence process with Markov Chain, *Water Resour. Res.*, 13, 949-956, 1977.
- Eubank, R.L., *Spline smoothing and nonparametric regression*, Marcel Dekker, Inc., New York.
- Feyerherm, A.M. and L.D. Bark, Statistical methods for persistent precipitation patterns, *Journal of Appl. Meteorology*, 4, 320-328, 1965.
- Feyerherm, A.M. and L.D. Bark, Goodness of fit of a markov chain model for sequences of wet and dry days, *Journal of Appl. Meteorology*, 6, 770-773, 1967.
- Gabriel, K.R. and J. Neumann, A Markov chain model for daily rainfall occurrence at Tel Aviv, *Quart. J. Roy. Meteor. Soc.*, 88,90-95, 1962.
- Guzman, A.G. and C.W. Torrez, Daily rainfall probabilities: conditional upon prior occurrence and amount of rain, *Journal of Climate and Applied. Meteorology*, 24(10), 1009-1014 , 1985.
- Haan, C.T., D.M. Allen and J.O. Street, A markov chain model of daily rainfall. *Water Resour. Res.*, 12(3), 443-449, 1976.
- Hopkins, J.W. and P. Robillard, Some statistics of daily rainfall occurrence for the canadian prairie provinces, *Journal of Applied. Meteorology*, 3, 600-602, 1964.
- Lall,U., Nonparametric function estimation: Recent hydrologic applications, *US National report, 1991-1994, International Union of Geodesy and Geophysics*, 1994.
- Lall, U., B. Rajagopalan and D.G. Tarboton, A Nonparametric Wet/Dry Spell model for resampling daily precipitation, *Submitted to Water Resour. Res.*, 1995.
- McLachlan, G.J., *Discriminant analysis and statistical pattern recognition*, John Wiley and Sons, New York, 1992.
- Rajagopalan, B. U. Lall and D.G. Tarboton, Simulation of Daily Precipitation from A Nonparametric Renewal Model, *Submitted to Water Resour. Res.*, 1995.

- Rajagopalan, B. and U. Lall, A kernel estimator for discrete distributions, *Journal of Nonparametric Statistics*, (in press), 1995.
- Rajagopalan, B. and U. Lall, Seasonality of Precipitation along a meridian in the western U.S., submitted to *Geophys. Res. Lett.*, 1995.
- Roldan J. and D.A. Woolhiser, Stochastic daily precipitation models 1. A comparison of occurrence processes, *Water Resour. Res.*, 18(5), 1451-1459, 1982.
- Scott, D.W., *Multivariate density estimation: Theory, Practice and Visualization*, Wiley series in probability and mathematical statistics, John Wiley and Sons, New York, 1992.
- Sheather, S.J. and M.C. Jones, A reliable data-based bandwidth selection method for kernel density estimation, *Journal of the Royal Statistical Society, B.* 53, 683-690, 1991.
- Silverman, B.W., *Density estimation for statistics and data analysis*, Chapman and Hall, New York, 1986.
- Smith, J.A. and H.A. Schreiber H.A., Point processes of seasonal thunderstorm rainfall. 1. Distribution of rainfall events, *Water Resour. Res.* 10(3), 418-423, 19743
- Srikanthan, R. and T.A. McMahon, Stochastic simulation of daily rainfall for australian stations. *Transactions of the ASAE*, 754-766, 1983.
- Stern, R.D. and R. Coe, R., A model fitting analysis of rainfall data (with discussion), *Journal of Royal Statistical Society, Series, A.*, 147, 1-34, 1984.
- Todorovic, P and D.A. Woolhiser, Stochastic model of  $n$ -day precipitation, *J. Appl. Meteorology*, 14(1), 17-24, 1975.
- Woolhiser, D.A., E.W. Rovey and P. Todorovic, Temporal and spatial variation of parameters for the distribution of  $n$ -day precipitation, in *Floods and Droughts, Proceedings of the Second International Symposium in Hydrology*, edited by E.F. Schulz, V.A. Koelzer, and K. Mahmood, pp. 605-614, Water Resources Publications, Fort Collins, Colorado, 1973.

Woolhiser, D.A. and G.G.S. Pegram, Maximum likelihood estimation of Fourier coefficients to describe seasonal variations of parameters in stochastic daily precipitation models, *J. Appl. Meteorology*, 18, 34-42, 1979.

Table 1

Statistics of wet day precipitation for Salt Lake City, UT, 1961-1991  
from historical precipitation record and averaged over 30 simulated precipitation records

	mean wet day ppt. (inches)	std. dev. wet day ppt. (inches)	fraction of yearly ppt.	maximum wet day ppt. (inches)
<b>Season 1</b>				
25% quantile	0.16	0.19	0.23	1.26
Median	0.16	0.20	0.23	1.36
75% quantile	0.17	0.21	0.24	1.59
historical	0.15	0.17	0.21	0.92
<b>Season 2</b>				
25% quantile	0.19	0.24	0.26	1.74
Median	0.19	0.25	0.27	1.86
75% quantile	0.20	0.26	0.28	2.18
historical	0.20	0.24	0.28	1.62
<b>Season 3</b>				
25% quantile	0.18	0.27	0.24	1.94
Median	0.18	0.28	0.26	2.3
75% quantile	0.19	0.30	0.26	2.87
historical	0.18	0.29	0.26	2.28
<b>Season 4</b>				
25% quantile	0.16	0.19	0.24	1.37
Median	0.17	0.21	0.24	1.7
75% quantile	0.18	0.23	0.25	2.16
historical	0.17	0.19	0.25	1.23
<b>Annual</b>				
25% quantile	0.18	0.24		2.35
Median	0.18	0.25		2.55
75% quantile	0.19	0.25		3.45
historical	0.17	0.22		2.30



Table 2

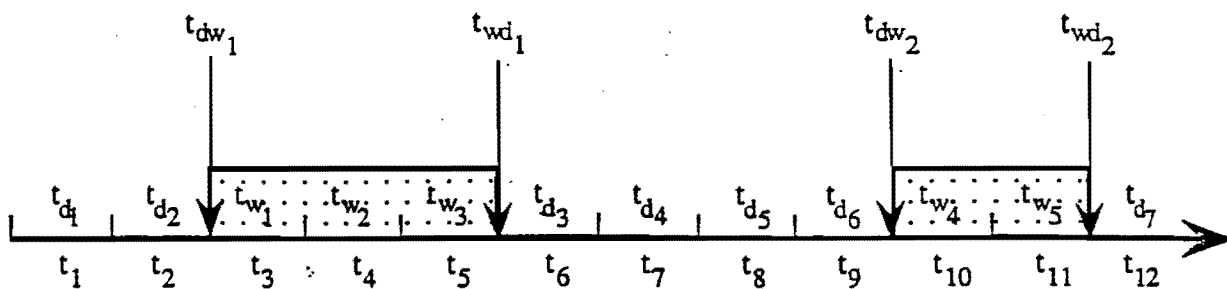
Statistics of wet spell length for Salt Lake City, UT, 1961-1991  
from historical precipitation record and averaged over 30 simulated precipitation records

	mean wet spell length (days)	std. dev. wet spell (days)	fraction of wet days	longest wet spell length (days)
<b>Season 1</b>				
25% quantile	1.89	1.29	0.31	9
Median	1.92	1.37	0.32	10
75% quantile	1.99	1.43	0.33	11.8
historical	1.86	1.29	0.32	10
<b>Season 2</b>				
25% quantile	1.87	1.27	0.25	8
Median	1.91	1.34	0.25	9
75% quantile	1.95	1.41	0.26	10
historical	2.12	1.47	0.27	12
<b>Season 3</b>				
25% quantile	1.79	1.23	0.19	8
Median	1.86	1.29	0.20	9
75% quantile	1.91	1.37	0.20	10
historical	1.60	0.9	0.18	7
<b>Season 4</b>				
25% quantile	1.85	1.27	0.25	8
Median	1.87	1.32	0.26	9
75% quantile	1.92	1.38	0.27	10
historical	1.97	1.36	0.26	9
<b>Annual</b>				
25% quantile	1.88	1.32	0.26	10
Median	1.91	1.36	0.26	11
75% quantile	1.94	1.39	0.26	13
historical	1.91	1.31	0.26	12

Table 3

Statistics of dry spell length for Salt Lake City, UT, 1961-1991  
from historical precipitation record and averaged over 30 simulated precipitation records

	mean dry spell length (days)	std. dev. dry spell (days)	fraction of dry days	longest dry spell length (days)
<b>Season 1</b>				
25% quantile	3.8	3.5	0.67	23
Median	3.92	3.63	0.68	25
75% quantile	4.0	3.75	0.68	27
historical	3.91	3.64	0.68	30
<b>Season 2</b>				
25% quantile	5.21	5.64	0.74	39
Median	5.48	5.91	0.75	46
75% quantile	5.59	6.25	0.76	50
historical	5.5	5.41	0.73	28
<b>Season 3</b>				
25% quantile	6.82	7.12	0.79	44
Median	7.05	7.53	0.80	52
75% quantile	7.26	7.943	0.81	72
historical	6.87	6.92	0.82	55
<b>Season 4</b>				
25% quantile	4.91	5.47	0.73	38
Median	5.09	5.71	0.74	43
75% quantile	5.28	5.91	0.75	51
historical	5.21	5.38	0.74	31
<b>Annual</b>				
25% quantile	5.29	6.13	0.74	58
Median	5.41	6.32	0.74	70
75% quantile	5.54	6.67	0.74	86
historical	5.45	5.99	0.74	61



$t_1, t_2, \dots$  are the day indices

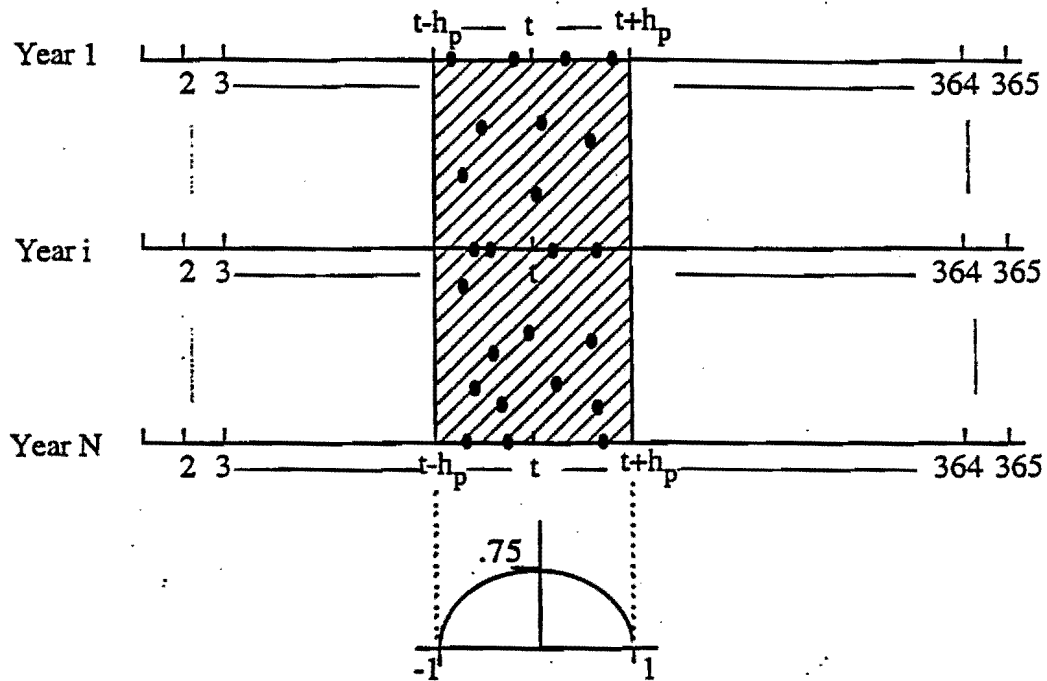
$t_{w1}, t_{w2}, \dots$  are wet day indices

$t_{d1}, t_{d2}, \dots$  are dry day indices

$t_{dw1}, t_{dw2}, \dots$  are day indices of transition from a dry day to wet day

$t_{wd1}, t_{wd2}, \dots$  are the day indices of transition from a wet day to dry day

Figure 1. Precipitation occurrence process.



$t$  is the calendar day on which precipitation is required  
 $h_p$  is the time interval around the calendar day  $t$   
 $1, \dots, N$  are the years in the historical record  
 Thick dots are the rainy days in the historical record  
 The kernel function shown at the bottom is used to weight the rainfall amounts on each of the rainy day.

Figure 2. Precipitation generation process.

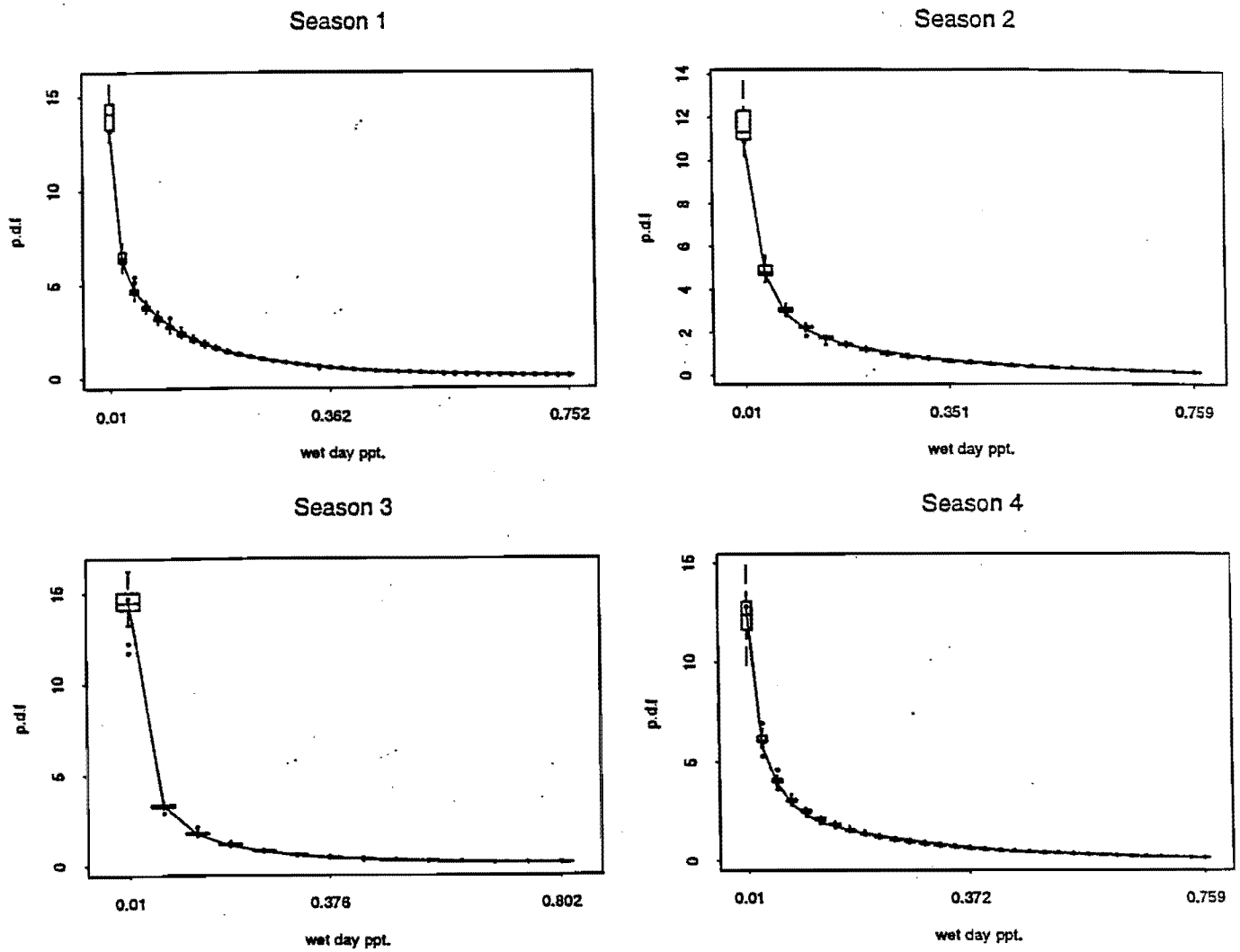


Figure 3. Boxplots of p.d.f. of wet day precipitation in each season, for model simulated records along with the historical values.

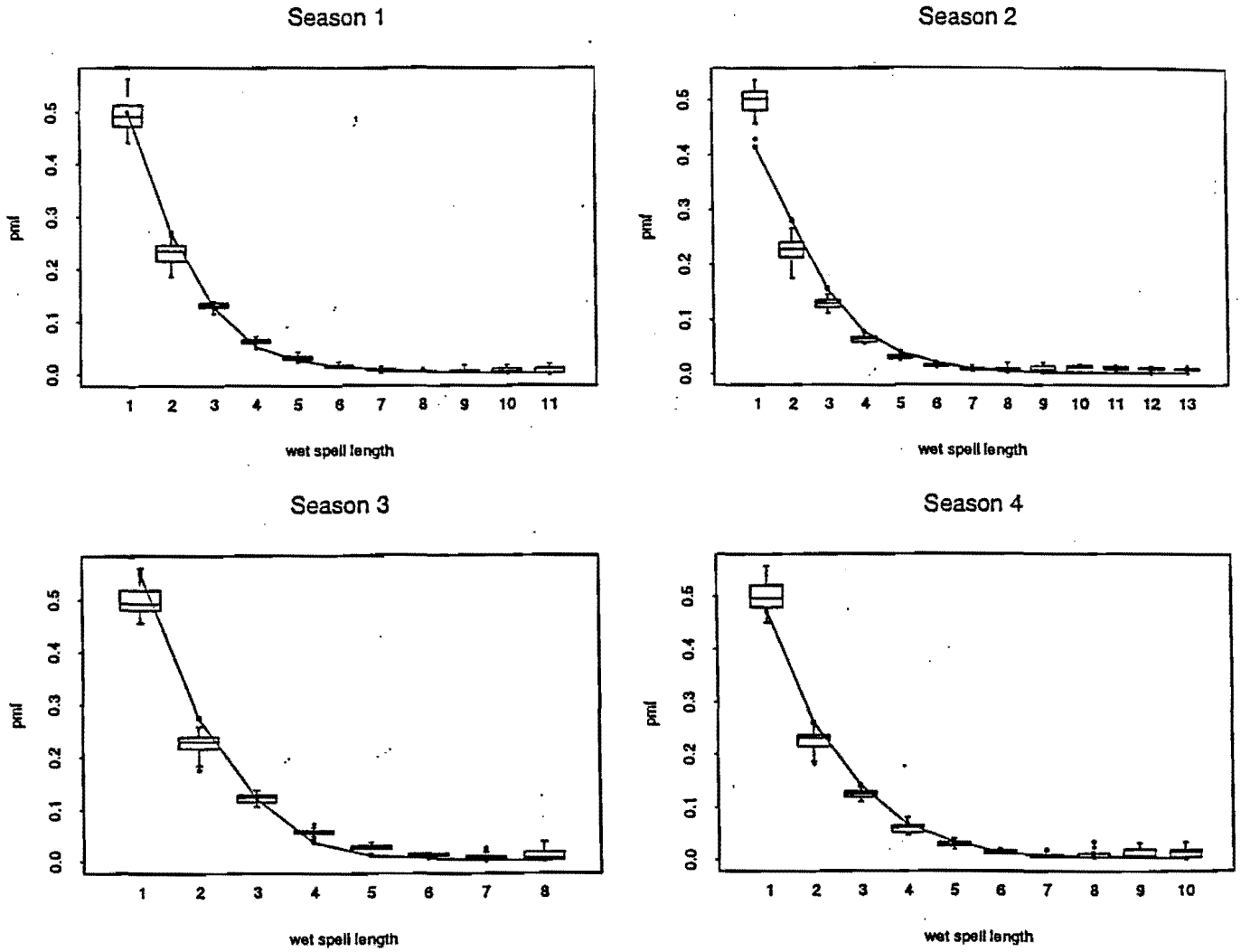
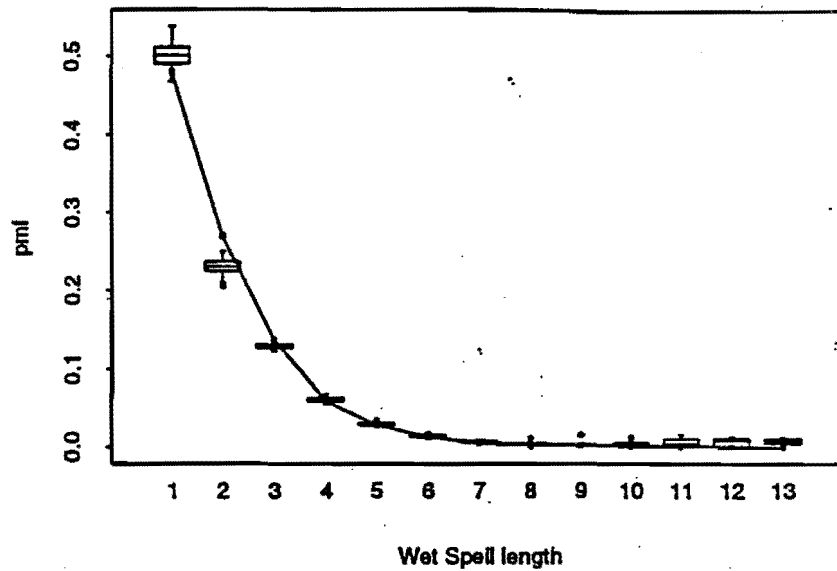


Figure 4. Boxplots of p.m.f. of wet spell length in each season, for model simulated records along with the historical values.



**Figure 5.** Boxplots of p.m.f. of wet spell length over the whole year, for model simulated records along with the historical values.

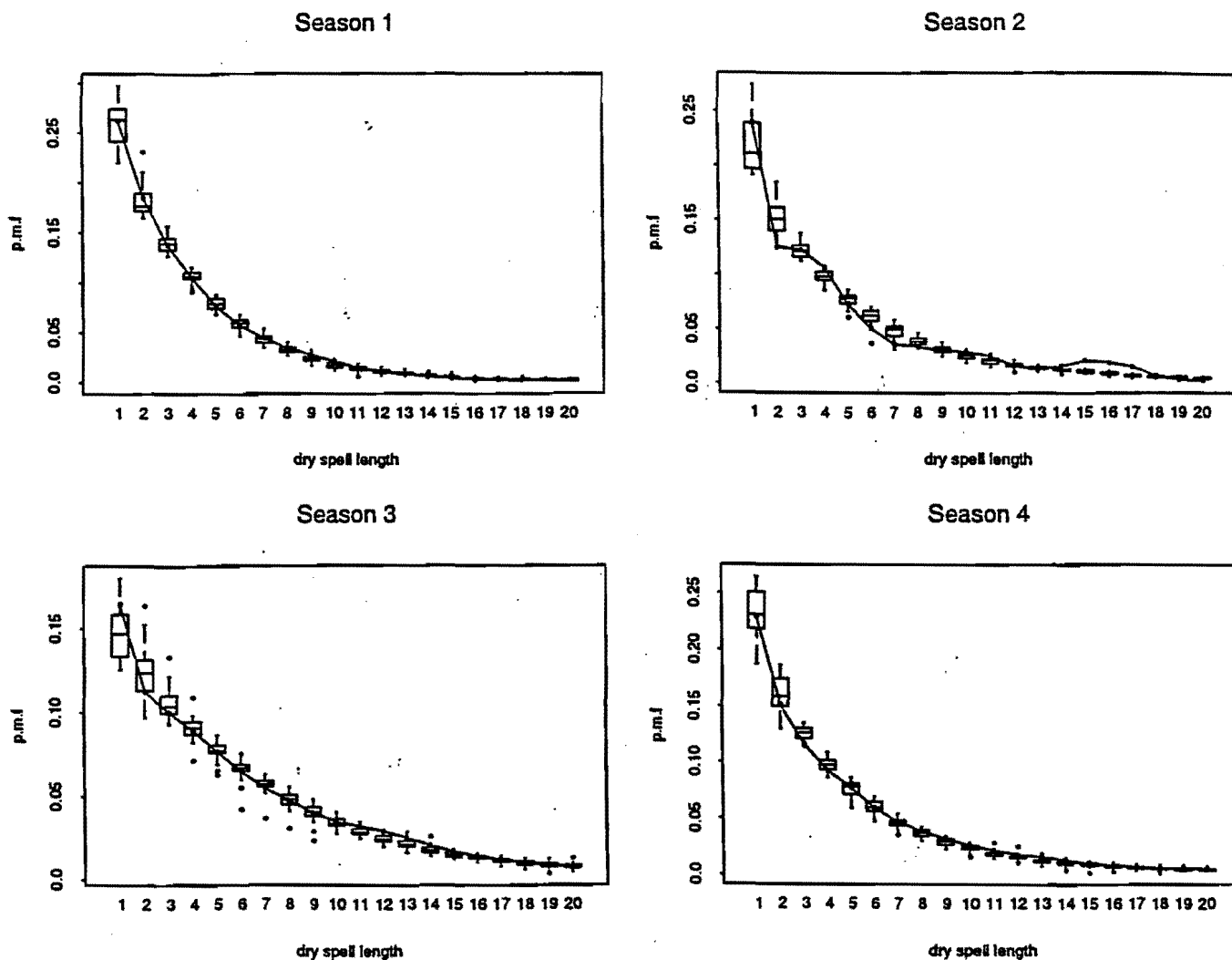
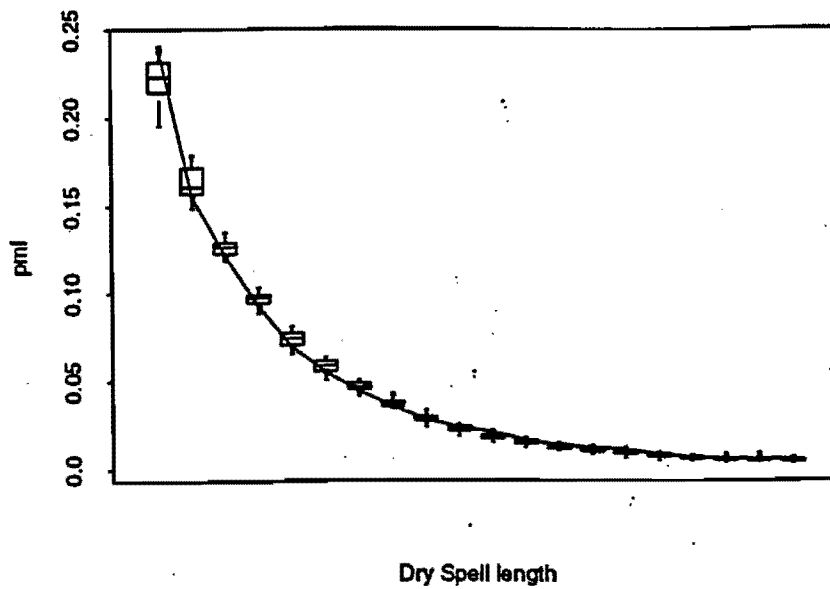


Figure 6. Boxplots of p.m.f. of dry spell length in each season, for model simulated records along with the historical values.





**Figure 7.** Boxplots of p.m.f. of dry spell length over the whole year, for model simulated records along with the historical values.