

# Rain

Matthew Wiens, Kelly Kung

December 20, 2017

## Abstract

Predicting precipitation is a key problem for residents of the Pacific Northwest, from its impact on family vacations to landslides. A noteworthy feature of the area is distinct changes of precipitation between the summer and winter. In this report, we propose a model for differentiating the two seasons and understanding the chance of precipitation in each, across different climatic areas of the state. Such a model can be used to predict optimal times to travel from a rainy region to a dry region. We adopt a Bayesian methodology and discuss predictive results within the Bayesian framework.

## 1 Introduction

Residents and visitors to the Pacific Northwest (Oregon and Washington) notice a distinct pattern year over year: most days during the winter are rainy, and then they feel that there is a distinct point during the spring or early summer when the weather becomes dry most days. Similarly, during early fall the nice weather switches back to rain. However, the timing of the switch varies throughout the area, with the prominent Cascade Mountains being a key factor. There are many resorts on the east side of the Cascade Range that cater to residents of the west side looking to escape the rain.

Widmann and Bretherton<sup>1</sup> developed a methodology to control for local variation in topography and precipitation data and generated a dataset of estimated precipitation data over 46 years in a 50km by 50km grid. In comparison, general models of atmospheric weather operate on the order of hundreds of kilometers, so there is significant room to improve the spatial resolution of weather models. In particular, local topographic features are not captured by general models, which is a limitation in regions like the Pacific Northwest. The temporal and spatial correlation is also highly dependent on the topography and difficult to model, and therefore there has been a focus on parametrizing models and reducing the dimensionality.

In this report, we propose a Bayesian approach for two reasons: first, to capture prior beliefs from living in the area of one author, and second to be able to discuss distributions around complicated model parameters and probabilistic beliefs about the state of weather on a specific day of the year. Bayesian approaches for daily precipitation data are not new, for example see [Olson and Kleiber].... With a Bayesian approach, we hope to focus on two main aspects (1) the beginning of dry and wet seasons and (2) the probability of precipitation. These aspects will then allow us to answer further questions we may have regarding precipitation.

This report is organized as follows. In Section 2, we refine our goals and questions we hope to answer with this analysis. This then motivates the creation of several variables in our dataset, which is also outlined in Section 2. In Section 3, we discuss the two models that we used in order to analyze the data. We then begin our analysis of the dataset in Section 4 and we discuss our findings in Section 5. In Section 6, we conclude our report with summaries of our analysis as well as potential future directions.

---

<sup>1</sup>M. Widmann and C. S. Bretherton. Validation of Mesoscale Precipitation in the NCEP Reanalysis Using a New Gridcell Dataset for the Northwestern United States. In *Journal of Climate*, 13(11), 1936-1950. 2000.

## 2 Background

### 2.1 Goals-change name later

Upon doing their analysis, Widmann and Bretherton found that there was some variability in precipitation levels in the different topographical areas in the Northwestern United States. In this report, we further investigate the variability in the different areas in order to build upon their results. In particular, we are interested in finding (1) when the dry and wet seasons start and (2) the probability of precipitation of a day. With this data, we are able to answer further questions to better understand the precipitation patterns in the different areas. Possible questions include how do the lengths of dry/wet periods vary from place to place and where does the dry/wet season first start. In addition, because we are doing a Bayesian model, we can take advantage of this opportunity to obtain probabilities of posterior estimates. One interesting question is what is the probability that one area has precipitation while another area does not.

### 2.2 Dataset

To analyze the precipitation in the Pacific Northwest, we use the same dataset as Widmann and Bretherton. This data contains precipitation logs (measured in millimeters) from 1949 - 1994 that was collected from the National Weather Service, which was then corrected by the National Climatic Data Center<sup>23</sup>. The dataset is a 3-dimensional array with a 2-dimensional grid-cell that consists of 16 longitude subsections and 17 latitude subsections, where each cross-section is approximately  $0.48^\circ$  (latitude) by  $0.62^\circ$  (longitude). The third dimension contains temporal data (16,801 days), which means we have approximately 4 million data points. However, note that there are approximately 32,000 missing values which correspond to areas with few functional weather stations, such as the ocean and some areas in the southeast. Using this dataset, we then transformed and created several variables in order to obtain a model to analyze the precipitation.

### 2.3 Grouping of Areas

The Pacific Northwest has a varied climate, from temperate rainforests on the Pacific coast to desert in Southeastern Oregon, which has a major impact on the seasons. Therefore, each grid cell is classified into one of nine climatic zones as follows:

1. Washington Coastline (Temperate Rainforest)
2. Oregon Coastline
3. Western Washington / Seattle Metropolitan Area
4. Portland Metropolitan Area / Willamette Valley / Southwestern Oregon
5. Cascade Range
6. Columbia Plateau
7. Selkirk Mountains
8. Blue Mountains
9. Eastern Oregon / High Desert / Great Basin

---

<sup>2</sup>The data was corrected during the 'Validated Historical Daily Data' project, which aimed to verify the data collected.

<sup>3</sup>T. Reek, S. Doty, and T. Owen. A Deterministic Approach to the Validation of Historical Daily Temperature and Precipitation Data from the Cooperative Network. In Bulletin of the American Meteorological Society 73, no. 6 (1992): 753-62.

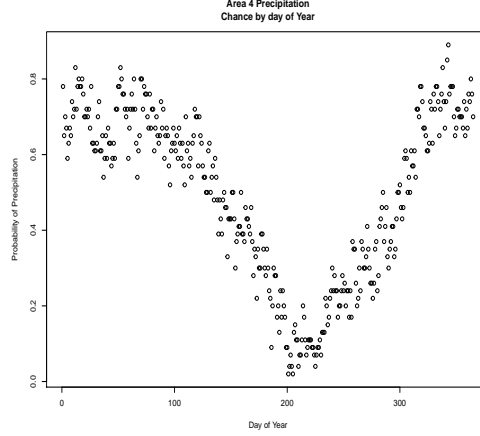


Figure 1: Western Oregon precipitation chance over the year

These zones were decided by considering average precipitation over the time period of the data and the topographic features of the region. The rainfall for the region was aggregated by considering if there was measurable rainfall in each grid by day, and then for each day if the majority of grid cells reported rain, then a value of *rain* was assigned to the area, else it was *dry*.

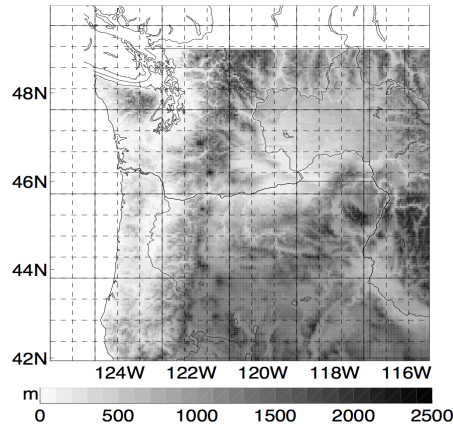


Figure 2: Topography of the region

## 2.4 Indicators for Dry and Wet Season

Using the grouping of areas as mentioned in Section 2.3, we then proceeded to create indicators of when the dry season and wet season starts. To do this, we first created indicators of precipitation for each day in each area, where a 1 indicates that more than 0.1 mm of precipitation occurred. We then aggregated the indicators in each of the 9 areas for each day and determined that there was precipitation in the area for that day if a majority of the cells in the area had more than 0.1 mm of precipitation. The areas with precipitation were deemed to be *wet* and the areas without precipitation were deemed to be *dry*. Afterwards, we created an algorithm that chooses indicators of when the dry season starts ( $X_{min}$ ) and when the wet season starts ( $X_{max}$ ) such that it maximizes the number of *dry* days in the dry season and maximizes the number of *wet* days in the wet season. For example, we expected the  $X_{min}$  days to be approximately 170 and the  $X_{max}$  days to be approximately 260 because these dates correspond to the beginning of summer and the beginning of autumn. Using this algorithm, we found the start dates for the dry and wet seasons

for the 46 years.<sup>4</sup>

### 3 Model

#### 3.1 Model 1

Our first model is motivated by the goal of determining when the dry and wet seasons start. To do so, we use a model with a semi-conjugate multivariate normal prior and multivariate normal likelihood. We chose this model because it has a closed form and we are able to sample from the posterior distribution using a Gibbs sampler on the full conditionals. Furthermore, upon initial exploration of the data, we see that the distribution of the start days of dry and wet seasons are approximately normal as seen in Figure 3. Thus, we proceed with the multivariate normal model with confidence that it represents the data we have.

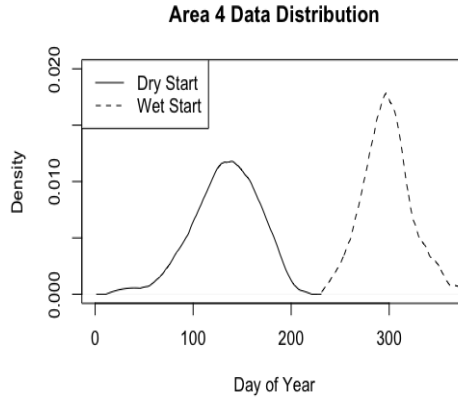


Figure 3: Kernel Density Estimate of the start and end of the dry season for Portland

We define the semi-conjugate prior distributions as follows.

$$\begin{aligned}\theta &\sim MVN(\mu_0, \Lambda_0) \\ \Sigma &\sim Inverse - Wishart_{\nu_0}(\Sigma_0^{-1})\end{aligned}$$

To determine the prior parameters, we set  $\mu_0 = \begin{bmatrix} X_{min} \\ X_{max} \end{bmatrix} = \begin{bmatrix} 170 \\ 260 \end{bmatrix}$  and  $\Lambda_0 = \begin{bmatrix} 30^2 & 0 \\ 0 & 30^2 \end{bmatrix}$  for all 9 area groups. The prior means are chosen according to when the summer and autumn season starts and the prior standard deviations of the  $\theta$ 's were chosen to represent a spread of one month. We choose the  $\Sigma_0$  values by mimicking the methods used during our homework, where we let  $\Sigma_0$  equal the textitvariance - covariance matrices for each area group from the data. The parameters for  $\Sigma_0$  are shown in Table 1.

In order to conduct Gibbs sampling, we need full conditionals of  $\theta$  and  $\Sigma$ , which are shown below. By running Gibbs sampling on these full conditionals, we are effectively sampling from the posterior distribution and so we obtain posterior estimates of  $\theta$  and  $\Sigma$ .

$$\begin{aligned}\theta|\Sigma, y &\sim MVN(\mu_n(\Sigma), \Lambda_n(\Sigma)) \\ \Sigma|\theta, y &\sim Inverse - Wishart_{\nu_n}(\Sigma_n^{-1}(\theta))\end{aligned}$$

---

<sup>4</sup>Note that when creating the indicators for years, we did not take into account of leap years, which resulted in data points that spilled over to the next year. We just take the years in which we have a year's worth of data for.

1		2		3		4		5	
1033.00	-254.79	652.62	32.22	1103.09	-35.55	661.74	17.54	943.12	67.83
-254.79	562.52	32.22	307.61	-35.55	860.19	17.54	197.88	67.83	524.46

6		7		8		9	
162.92	-15.51	778.02	-17.29	1097.19	58.49	288.11	-15.03
-15.51	108.72	-17.29	99.43	58.49	76.20	-15.03	12.34

Table 1: Table of prior  $\Sigma_0$  values for each area group

## 3.2 Model 2

We further developed our model by considering the probability of rain on each day for the summer (dry) and winter (rainy) seasons. This is motivated by the idea of trying to plan outdoor activity in each area of the Pacific Northwest and deciding if there is value in traveling to a different area to escape rain. Therefore, the following model is proposed:

- The data is binary data representing if it rained or not for each day of the year,
- $\pi_w$  and  $\pi_d$  are the probability of rain during the wet season and dry season, respectively.
- $\theta_1$  and  $\theta_2$  are the cutoff points for the wet and dry seasons. The wet season runs from January 1st until the day before  $\theta_1$  and from  $\theta_2$  to December 31st, and the dry season runs from  $\theta_1$  until the day before  $\theta_2$ , and the natural restriction is imposed that  $\theta_1 < \theta_2$
- Let  $y_j$  be the data on day  $j$ . Then  $p(y_j) = \pi_w 1_{j < \theta_1} + \pi_d 1_{\theta_1 \leq j < \theta_2} + \pi_w 1_{j \geq \theta_2}$

## 4 Analysis

Using the semi-conjugate multivariate normal model described in section 3.1, we used a Gibbs sampler to generate samples from the posterior distributions. We report summary statistics of the posterior predictive distribution in Table 2. In particular we note the increased variability of the start of the dry season between regions and the variance within each region. On the other hand, the end of the dry season is much more consistent within and across regions. This suggests larger weather patterns triggering the end of the dry season which are consistent between years. Also of interest is the beginning of dry, summer weather across the Pacific Northwest. In table 3 we report the 95th percentile of the beginning of summer. For regions 3 and 4, which include Seattle and Portland, this date is around the 4th of July weekend, which agrees with common beliefs in the area that one cannot count on summer to begin until then. In addition, a very likely beginning of the dry season is later in the generally drier, Eastern part of the region. This may be due to the drier nature of that region, so distinctly different weather patterns may not develop until later in the summer.

In figure 4 this model is used to explore how one could escape the rainy season in the major metropolitan areas of Seattle, Washington (region 3) and Portland, Oregon (region 4). For example, resort operators in the drier climates directly east of these cities are most desirable when they have dry weather and the cities have rainy weather. We find that there is longer period in the springtime when going to the east side of the Cascades is desirable, and a narrow window during the fall. However, by noting the maximum probability is substantially less than 1 in each peak, we note that the range of this time period is variable year to year, and so making exact future plans would be difficult.

Region	Start Date Mean	Start Date Variance	End Date Mean	End Date Variance
1	133	26	291	14
2	136	15	299	10
3	137	22	293	10
4	134	22	297	13
5	144	22	297	10
6	119	65	300	17
7	126	79	301	20
8	147	49	305	12
9	134	44	309	16

Table 2: Posterior Mean and Variance of mean parameter

Region	Day of Year
1	190
2	179
3	190
4	188
5	197
6	210
7	226
8	226
9	208

Table 3: Date by area such that there is a 95% probability that the dry season has begun

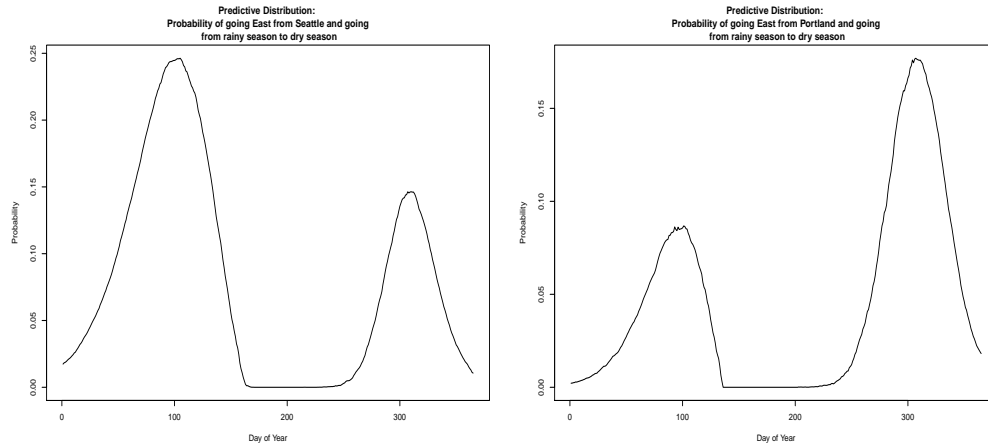


Figure 4: On each day of the year, the probability of escaping the rain by going east for the two major metropolitan areas

## 5 Discussion

## 6 Conclusion