

P8160 - Project 3

A Hierarchical Bayesian model for hurricane trajectories

Renjie Wei, Hao Zheng, Xinran Sun
Wentong Liu, Shengzhi Luo

5/7/2022

1. Introduction

1.1 Background and Objectives

Hurricanes are a serious social and economic concern to the United States. Strong winds, heavy rainfall, and high storm surge kill people and destroy property. There is an increasing desire to predict the performance of hurricane, such as its location, speed and so on. In this project, we are interested in modeling the hurricane trajectories to forecast the wind speed achieved by Hierarchical Bayesian Model. The hurricane data contains individual-level-specific effects of each hurricane. Model integration is achieved through a Markov Chain Monte Carlo algorithm.

Also, we present work that is to describe the seasonal difference based on the previous estimated Bayesian model and try to find if there is any evidence supporting that the hurricane wind speed has been increasing over years. Finally, we use additional data which includes the damages and death caused by hurricanes in the United States to build a prediction model. We wish to find the most important factors that affect hurricanes and draw inferences and conclusions based on the model.

2 Data Cleaning and Explorative Data Analysis

2.1 Dataset

hurrican703.csv collected the track data of 702 hurricanes in the North Atlantic area from 1950 to 2013. For all the storms, their location (longitude & latitude) and maximum wind speed were recorded every 6 hours. The data includes the following variables

1. **ID**: ID of the hurricans
2. **Season**: In which **year** the hurricane occurred
3. **Month**: In which **month** the hurricane occurred
4. **Nature**: Nature of the hurricane
 - ET: Extra Tropical
 - DS: Disturbance
 - NR: Not Rated
 - SS: Sub Tropical
 - TS: Tropical Storm
5. **time**: dates and time of the record
6. **Latitude** and **Longitude**: The location of a hurricane check point
7. **Wind.kt** Maximum wind speed (in Knot) at each check point

From the original dataset, we built a new dataset with contains five more variables, including:

1. **Wind_prev**: wind speed at 6 hours ago
2. **Wind_prev_prev**: wind speed at 12 hours ago
3. **Lat_change**: latitude change compared to 6 hours earlier
4. **Long_change**: longitude change compared to 6 hours earlier
5. **Wind_change**: wind speed change at 6 hours earlier compared to 12 hours earlier

These variables will help us to build the model in the following part.

The *hurricanoutcome2.csv* recorded the damages and death caused by 46 hurricanes in the U.S, and some features extracted from the hurricane records. The variables include:

1. **ID**: ID of the hurricans
2. **Season**: In which **year** the hurricane occurred
3. **Month**: In which **month** the hurricane occurred
4. **Nature**: Nature of the hurricane
 - ET: Extra Tropical
 - DS: Disturbance
 - NR: Not Rated
 - SS: Sub Tropical
 - TS: Tropical Storm
5. **Damage**: Financial loss (in Billion U.S. dollars) caused by hurricanes
6. **Deaths**: Number of death caused by hurricanes
7. **Maxspeed**: Maximum recorded wind speed of the hurricane
8. **Meanspeed**: average wind speed of the hurricane
9. **Maxpressure**: Maximum recorded central pressure of the hurricane

10. **Meanpressure:** average central pressure of the hurricane
11. **Hours:** Duration of the hurricane in hours
12. **Total.Pop:** Total affected population
13. **Percent.Poor:** % affected population that reside in low GDP countres (i.e. GDP per Capita \leq 10,000)
14. **Percent.USA:** % affected population that reside in the United States

2.2 EDA

We use a bar plot to examine the number of hurricanes in each month. From Figure 1, we can see that September is the month with the most hurricanes, while there are no hurricanes in February and March. Hurricanes in September also have the highest average wind speed as we can see in Figure 2.

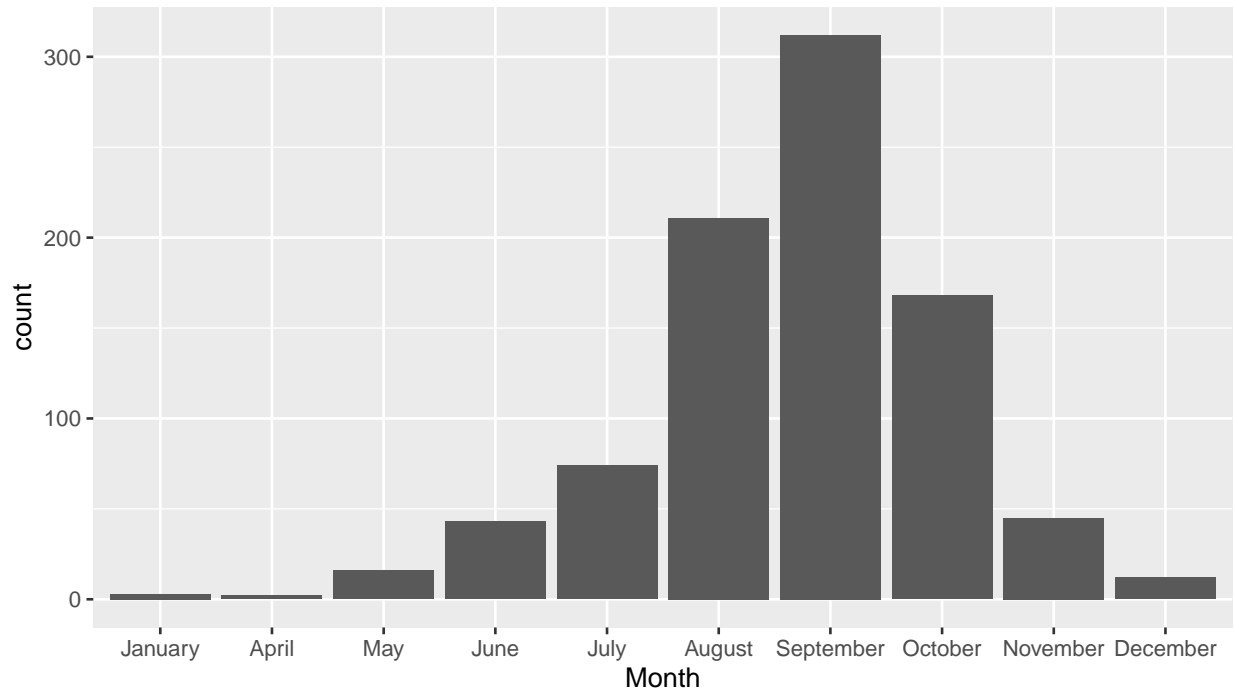


Figure 1. Count of Hurricanes in Each Month

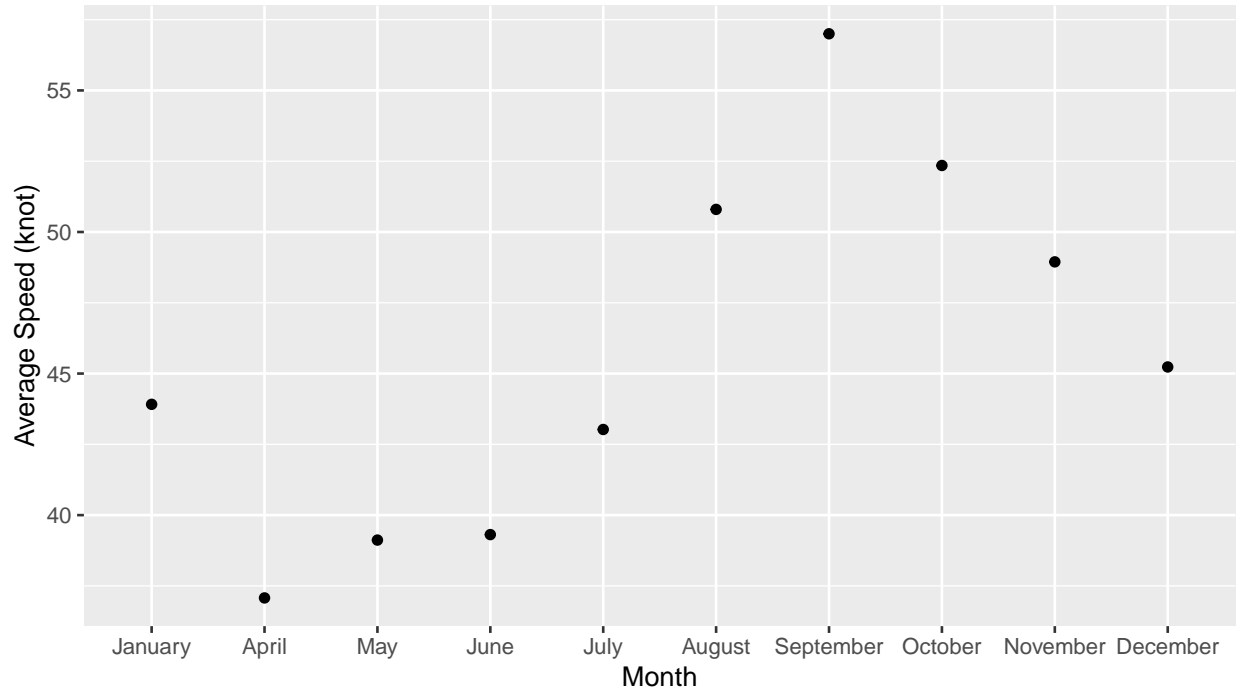


Figure 2. Average Speed (knot) of Hurricanes in Each Month

If we group the hurricanes by years, we can see in general, we have more observations in recently years compared to 50 years ago as shown in Figure 3. However, from Figure 4, the average wind speed seems to have a decreasing trend.

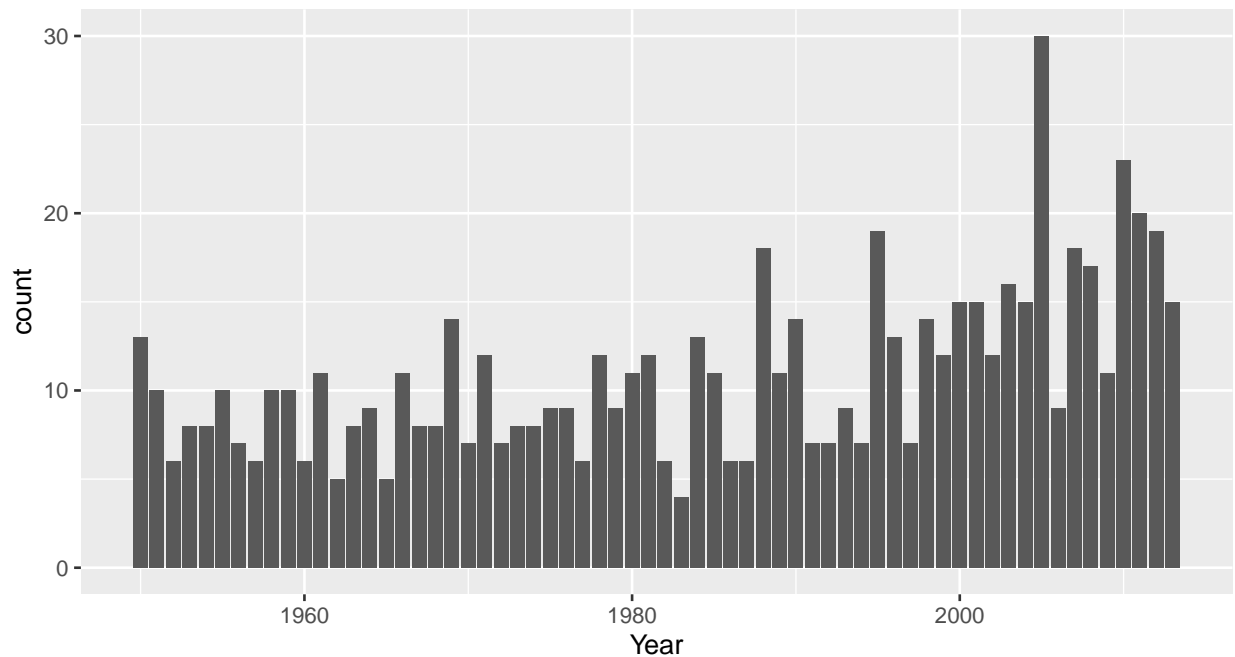


Figure 3. Count of Hurricanes in Each Year

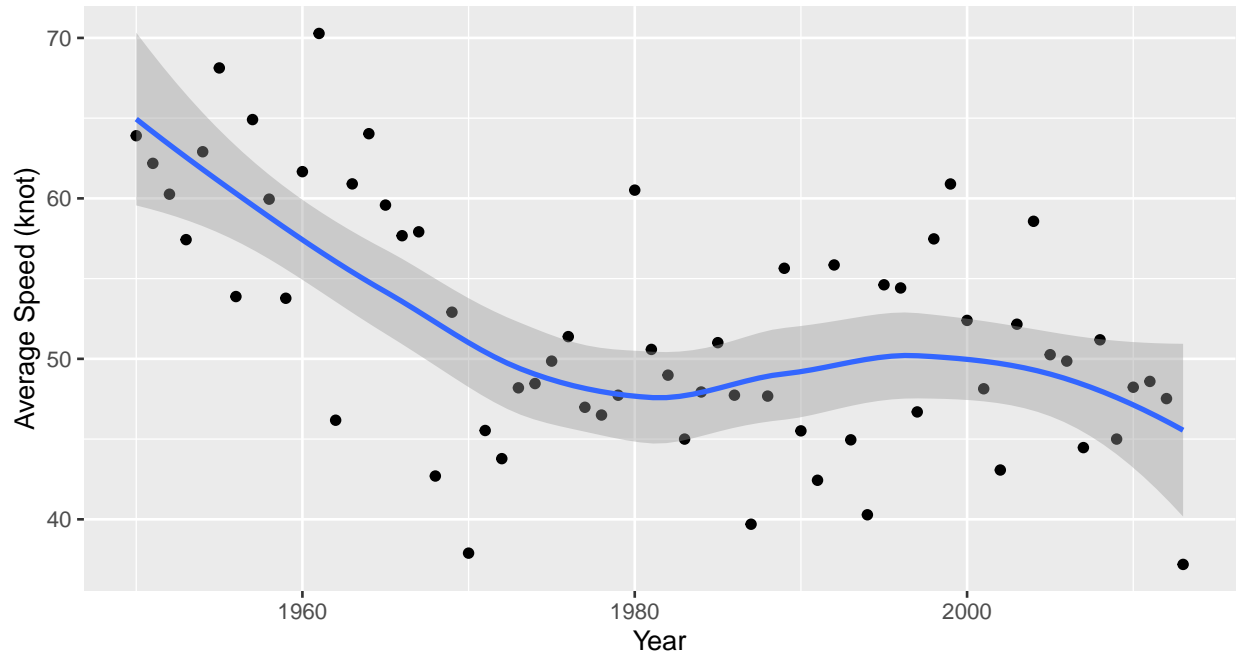


Figure 4. Average Speed (knot) of Hurricanes in Each Year

We also compare the hurricanes with different natures. In our dataset, there are 1214 different nature ratings. This number is larger than the number of hurricanes because some hurricanes are in different natures at different time. From Figure 5, we know that more than half of the natures are in Tropical Storm category. This nature also have the highest average wind speed at about 60 knot, while the disturbance and not rated hurricanes have average wind speed at round 20 knot as Figure 6 illustrates.

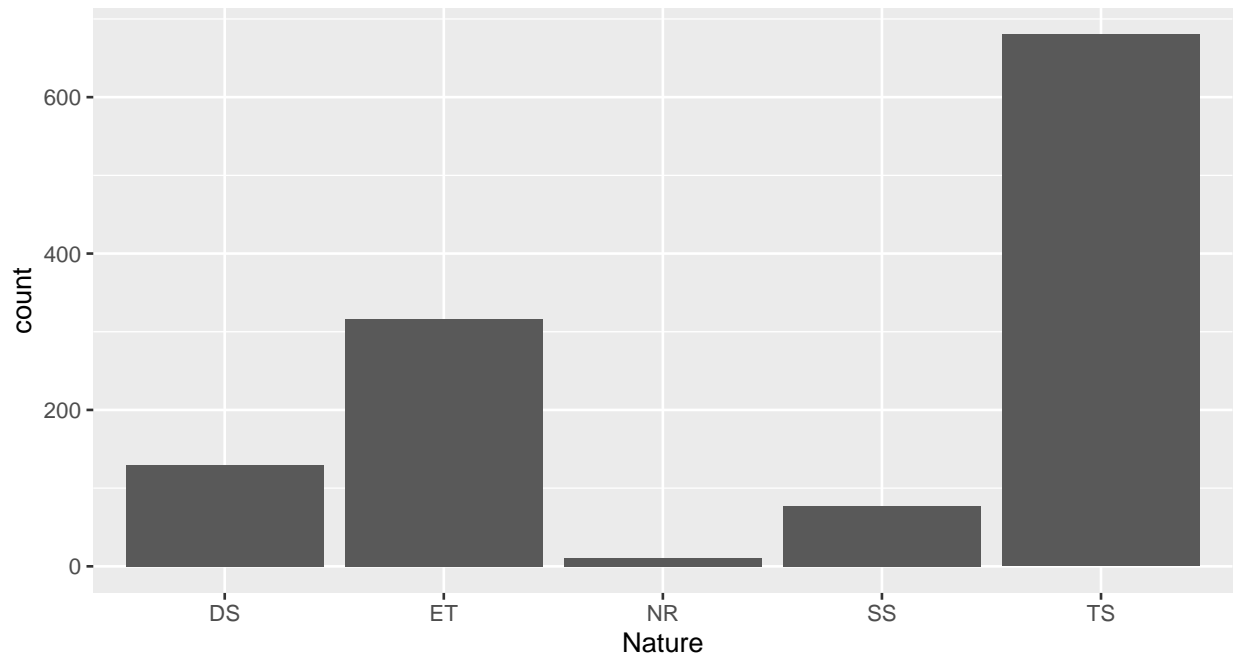


Figure 5. Count of Hurricanes in Each Nature

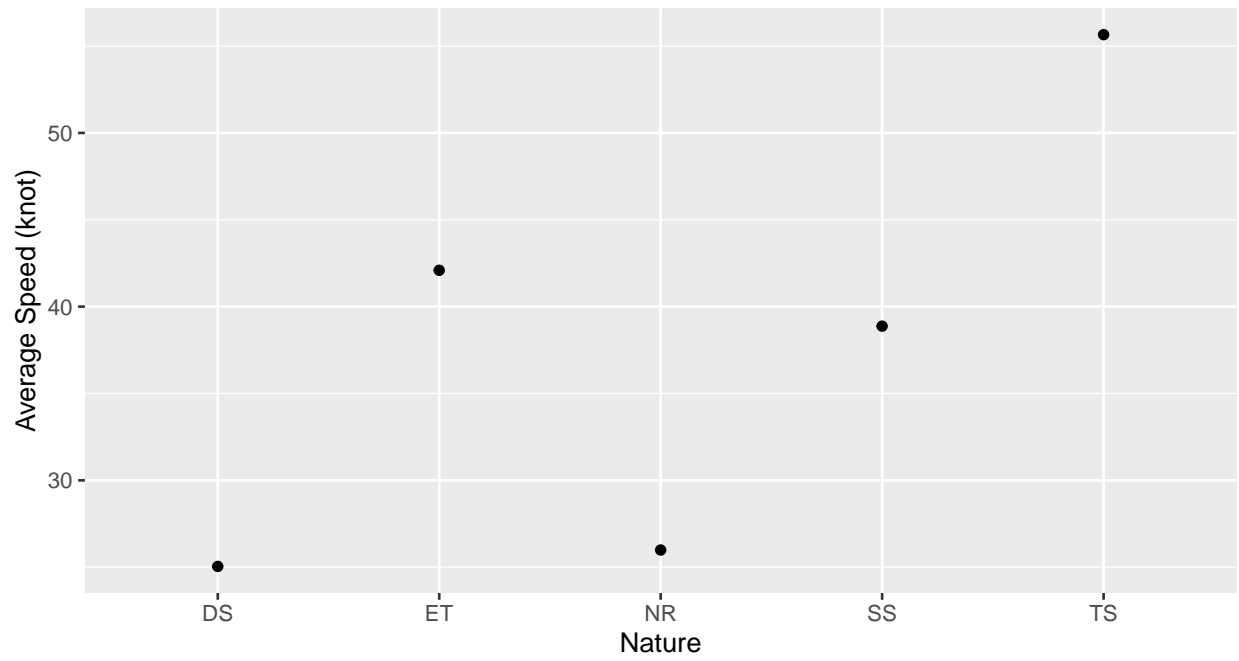


Figure 6. Average Speed (knot) of Hurricanes in Each Nature

3 Bayesian Model for Hurricanes Trajectories

3.1 Markov Chain Monte Carlo (MCMC)

Markov Chain Monte Carlo is combined by two methods, Markov Chain and Monte Carlo Method. Monte Carlo is a random sampling method for approximating a desired quantity, whereas Markov Chain generates a sequence of random variables where the current state only depends on the nearest past in the chain. MCMC algorithm draws samples from Markov Chain successively leading us close to the desired posterior. Two commonly used MCMC algorithm are the Metropolis-Hastings Algorithm and the Gibbs Sampler. Here, we implement the Gibbs Sampler here since we can save much computation cost compared to Metropolis-Hastings Algorithm.

3.2 Gibbs Sampler

Gibbs Sampler is one of Bayesian MCMC approaches with known conditional distributions. By sampling from each random variables given all the others, and changing one random variable at a time, Gibbs Sampler is able to draw parameter samples from the joint distribution. Then given proper starting value, the Markov Chain can reach its stationary distribution.

3.3 Bayesian Hierarchical Modelling

Bayesian hierarchical modelling is a statistical model written in multiple levels (hierarchical form) that estimates the parameters of the posterior distribution using the Bayesian method.

Hierarchical Bayesian Models, which contain both within-group analysis and between-group analysis, are always used to learn about a population from many individual measurements. Therefore, there is natural heterogeneity during the research periods and it could be regarded as subject-specific mean response trajectories for each individual group. To build the model, we split the inference problem into steps, where the full model is made up of a series of sub-models. The Bayesian Hierarchical Model links the sub-models together, correctly propagating uncertainties in each sub-model from one level to the next. MCMC methods work particularly well with hierarchical models, and is the engine that has fueled the development and application of Bayes' theorem.

From the Bayes' theorem:

$$\text{posterior distribution} \propto \text{likelihood} \times \text{prior distribution}$$

$$\pi(\theta|X) \propto \pi(X|\theta) \times \pi(\theta)$$

The Hierarchical Bayes

$$\pi(\theta, \alpha|X) \propto \pi(X|\theta) \times \pi(\theta|\alpha) \times \pi(\alpha)$$

Bayesian Inference is a statistical inference method about parameter. Proper prior distribution of parameter θ is set. After data collection, the belief of parameter θ would be updated by exploring the posterior distribution of θ based on observed data and its pre-assumed likelihood function $L(X; \theta)$. The linear regression model in hierarchical form incorporating with Bayesian inference is implemented with MCMC Integration algorithm for updating the parameter estimation in the final MCMC stationary phase.

3.4 Model Setting

The suggested Bayesian model

$$Y_i(t+6) = \beta_{0,i} + \beta_{1,i}Y_i(t) + \beta_{2,i}\Delta_{i,1}(t) + \beta_{3,i}\Delta_{i,2}(t) + \beta_{4,i}\Delta_{i,3}(t) + \epsilon_i(t)$$

where $Y_i(t)$ the wind speed at time t (i.e. 6 hours earlier), $\Delta_{i,1}(t)$, $\Delta_{i,2}(t)$ and $\Delta_{i,3}(t)$ are the changes of latitude, longitude and wind speed between t and $t-6$, and $\epsilon_{i,t}$ follows a normal distributions with mean zero and variance σ^2 , independent across t .

In the model, $\beta_i = (\beta_{0,i}, \beta_{1,i}, \dots, \beta_{5,i})$ are the random coefficients associated the i th hurricane, we assume that

$$\beta_i \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

follows a multivariate normal distributions with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

We assume the following non-informative or weak prior distributions for σ^2 , $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

$$P(\sigma^2) \propto \frac{1}{\sigma^2}; \quad P(\boldsymbol{\mu}) \propto 1; \quad P(\boldsymbol{\Sigma}^{-1}) \propto |\boldsymbol{\Sigma}|^{-(d+1)} \exp\left(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}\right)$$

d is dimension of β .

3.5 Posterior Distributions

Let $\mathbf{B} = (\beta_1^\top, \dots, \beta_n^\top)^\top$, derive the posterior distribution of the parameters $\Theta = (\mathbf{B}^\top, \boldsymbol{\mu}^\top, \sigma^2, \boldsymbol{\Sigma})$.

Let

$$\mathbf{X}_i(t)\beta_i^\top = \beta_{0,i} + \beta_{1,i}Y_i(t) + \beta_{2,i}\Delta_{i,1}(t) + \beta_{3,i}\Delta_{i,2}(t) + \beta_{4,i}\Delta_{i,3}(t)$$

where $\mathbf{X}_i(t) = (1, Y_i(t), \Delta_{i,1}(t), \Delta_{i,2}(t), \Delta_{i,3}(t))$, $\beta_i = (\beta_{0,i}, \beta_{1,i}, \beta_{2,i}, \beta_{3,i}, \beta_{4,i})$

then, we can find that

$$Y_i(t+6) \sim N(\mathbf{X}_i(t)\beta_i^\top, \sigma^2)$$

For i^{th} hurricane, there may be m_i times of record (excluding the first and second observation), let

$$\mathbf{Y}_i = \begin{pmatrix} Y_i(t_0+6) \\ Y_i(t_1+6) \\ \vdots \\ Y_i(t_{m_i-1}+6) \end{pmatrix}_{m_i \times 1}$$

denotes the m_i -dimensional result vector for the i^{th} hurricane. Therefore, since $Y_i(t)$'s are independent across t , we can show that the conditional distribution of \mathbf{Y}_i is

$$\mathbf{Y}_i \mid \mathbf{X}_i, \beta_i, \sigma^2 \sim N(\mathbf{X}_i\beta_i^\top, \sigma^2 I)$$

where

$$\mathbf{X}_i = \begin{pmatrix} 1 & Y_i(t_0) & \Delta_{i,1}(t_0) & \Delta_{i,2}(t_0) & \Delta_{i,3}(t_0) \\ 1 & Y_i(t_1) & \Delta_{i,1}(t_1) & \Delta_{i,2}(t_1) & \Delta_{i,3}(t_1) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & Y_i(t_{m_i-1}) & \Delta_{i,1}(t_{m_i-1}) & \Delta_{i,2}(t_{m_i-1}) & \Delta_{i,3}(t_{m_i-1}) \end{pmatrix}_{m_i \times d}$$

and the pdf of \mathbf{Y}_i is

$$\begin{aligned} f(\mathbf{Y}_i \mid \beta_i, \sigma^2) &= \det(2\pi\sigma^2 I_{(m_i \times m_i)})^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\sigma^2 I_{(m_i \times m_i)})^{-1} (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)\right\} \\ &= (2\pi\sigma^2)^{-m_i/2} \exp\left\{-\frac{1}{2}(\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\sigma^2 I_{(m_i \times m_i)})^{-1} (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)\right\} \end{aligned}$$

Since

$$\beta_i \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Therefore the pdf of β_i is

$$\pi(\beta_i | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\beta_i - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\beta_i - \boldsymbol{\mu})^\top\right)$$

Notice that $\mathbf{B} = (\beta_1^\top, \dots, \beta_n^\top)^\top$, i.e.

$$\mathbf{B} = \begin{pmatrix} \beta_{0,1} & \beta_{1,1} & \beta_{2,1} & \beta_{3,1} & \beta_{4,1} \\ \beta_{0,2} & \beta_{1,2} & \beta_{2,2} & \beta_{3,2} & \beta_{4,2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \beta_{0,n} & \beta_{1,n} & \beta_{2,n} & \beta_{3,n} & \beta_{4,n} \end{pmatrix}_{n \times d}$$

So, by using Bayesian rule, we can show the posterior distribution of $\boldsymbol{\Theta}$ is,

$$\begin{aligned} \pi(\boldsymbol{\Theta} | \mathbf{Y}) &= \pi(\mathbf{B}^\top, \boldsymbol{\mu}^\top, \sigma^2, \boldsymbol{\Sigma} | Y) \propto \prod_{i=1}^n f(\mathbf{Y}_i | \beta_i, \sigma^2) \prod_{i=1}^n \pi(\beta_i | \boldsymbol{\mu}, \boldsymbol{\Sigma}) P(\sigma^2) P(\boldsymbol{\mu}) P(\boldsymbol{\Sigma}^{-1}) \\ &\propto \prod_{i=1}^n \left\{ (2\pi\sigma^2)^{-m_i/2} \exp \left\{ -\frac{1}{2}(\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\sigma^2 I)^{-1} (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top) \right\} \right\} \\ &\times \prod_{i=1}^n \left\{ \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2}(\beta_i - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\beta_i - \boldsymbol{\mu})^\top \right\} \right\} \\ &\times \frac{1}{\sigma^2} \times \det(\boldsymbol{\Sigma})^{-(d+1)} \exp \left\{ -\frac{1}{2}\boldsymbol{\Sigma}^{-1} \right\} \end{aligned}$$

3.6 Conditional Posterior Distribution

Due to the high-dimensional problem of the full joint posterior distribution, considering the computational complexity plus we actually well know the form of the joint posterior, we suggest to use Gibbs sampling algorithm instead of Metropolis-Hastings algorithm.

To apply MCMC using Gibbs sampling, we need to find conditional posterior distribution of each parameter, then we can implement Gibbs sampling on these conditional posterior distributions.

Since our suggested model mainly focus on the parameter \mathbf{B} , we decided to derive its conditional posterior distribution.

1. The posterior distribution of \mathbf{B}

Since finding the posterior distribution of \mathbf{B} is the same to find the posterior distribution of β_i , we try to derive the conditional distribution $\pi(\beta_i | \mathbf{Y}, \boldsymbol{\mu}^\top, \sigma^2, \boldsymbol{\Sigma})$

$$\begin{aligned}
\pi(\mathbf{B}|\mathbf{Y}, \boldsymbol{\mu}^\top, \sigma^2, \boldsymbol{\Sigma}) &\propto L_Y(\mathbf{B}^\top, \sigma^2) \times \pi(\mathbf{B} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
&\propto \prod_{i=1}^n f(\mathbf{Y}_i | \beta_i, \sigma^2) \prod_{i=1}^n \pi(\beta_i | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
&\propto \prod_{i=1}^n \left\{ (2\pi\sigma^2)^{-m_i/2} \exp\left(-\frac{1}{2}(\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\sigma^2 I)^{-1} (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)\right) \right\} \\
&\times \prod_{i=1}^n \left\{ \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\beta_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\beta_i - \boldsymbol{\mu})\right) \right\} \\
&\propto \prod_{i=1}^n \exp\left\{-\frac{1}{2}\left((\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\sigma^2 I)^{-1} (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top) + (\beta_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\beta_i - \boldsymbol{\mu})\right)\right\} \\
&= \exp\left\{-\frac{1}{2}\left(\mathbf{Y}_i^\top (\sigma^2 I)^{-1} \mathbf{Y}_i + \beta_i \mathbf{X}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i \beta_i^\top - \mathbf{Y}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i \beta_i^\top \right.\right. \\
&\quad \left.\left. - \beta_i \mathbf{X}_i^\top (\sigma^2 I)^{-1} \mathbf{Y}_i + \beta_i \boldsymbol{\Sigma}^{-1} \beta_i^\top + \boldsymbol{\mu} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}^\top - \boldsymbol{\mu} \boldsymbol{\Sigma}^{-1} \beta_i^\top - \beta_i \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}^\top\right)\right\} \\
&= \exp\left\{-\frac{1}{2}\left(\beta_i (\boldsymbol{\Sigma}^{-1} + \mathbf{X}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i) \beta_i^\top - 2(\mathbf{Y}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i + \boldsymbol{\mu} \boldsymbol{\Sigma}^{-1}) \beta_i^\top + \mathbf{C}\right)\right\}
\end{aligned}$$

where,

$$\mathbf{C} = \mathbf{Y}_i^\top (\sigma^2 I)^{-1} \mathbf{Y}_i + \boldsymbol{\mu} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}^\top$$

By re-writing the conditional posterior distribution, and ignoring some constant terms, we can show that

$$\pi(\mathbf{B}|\mathbf{Y}, \boldsymbol{\mu}^\top, \sigma^2, \boldsymbol{\Sigma}) \propto \prod_{i=1}^n \exp\{(\beta_i^\top - \hat{\beta}_i)^\top (\hat{\boldsymbol{\Sigma}}_{\beta_i})^{-1} (\beta_i^\top - \hat{\beta}_i)\}$$

Hence, each β_i has a conditional posterior multivariate normal distribution

$$\pi(\beta_i|\mathbf{Y}, \boldsymbol{\mu}^\top, \sigma^2, \boldsymbol{\Sigma}) \sim \mathcal{N}(\hat{\beta}_i, \hat{\boldsymbol{\Sigma}}_{\beta_i})$$

where

$$\begin{aligned}
\hat{\beta}_i &= (\boldsymbol{\Sigma}^{-1} + \mathbf{X}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i)^{-1} \mathbf{Y}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i + \boldsymbol{\mu} \boldsymbol{\Sigma}^{-1} \\
\hat{\boldsymbol{\Sigma}}_{\beta_i} &= (\boldsymbol{\Sigma}^{-1} + \mathbf{X}_i^\top (\sigma^2 I)^{-1} \mathbf{X}_i)^{-1}
\end{aligned}$$

2. The posterior distribution of $\pi(\sigma^2|\mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\mu}^\top, \boldsymbol{\Sigma})$

$$\begin{aligned}
\pi(\sigma^2|\mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\mu}^\top, \boldsymbol{\Sigma}) &\propto L_Y(\mathbf{B}^\top, \sigma^2) \times \pi(\sigma^2) \\
&\propto \frac{1}{\sigma^2} \prod_{i=1}^n \left\{ (2\pi\sigma^2)^{-m_i/2} \exp\left(-\frac{1}{2}(\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\sigma^2 I)^{-1} (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)\right) \right\} \\
&\propto \frac{1}{\sigma^2} \left(\frac{\sum_{i=1}^n m_i}{2} + 1\right) \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)\right\}
\end{aligned}$$

which follows the form of pdf of inverse gamma distribution

$$f(x; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{1}{x} \exp\left\{-\frac{\beta}{x}\right\}$$

in this case, x is replaced by σ^2 , α is replaced by $\frac{1}{2} \sum_{i=1}^n m_i$, β is replaced by $\frac{1}{2} \sum_{i=1}^n (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)^\top (\mathbf{Y}_i - \mathbf{X}_i\beta_i^\top)$

i.e.

$$\pi(\sigma^2 | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\mu}^\top, \boldsymbol{\Sigma}) \sim IG(\frac{1}{2} \sum_{i=1}^n m_i, \frac{1}{2} \sum_{i=1}^n (\mathbf{Y}_i - \mathbf{X}_i \boldsymbol{\beta}_i^\top)^\top (\mathbf{Y}_i - \mathbf{X}_i \boldsymbol{\beta}_i^\top))$$

3. The posterior distribution of $\pi(\boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\mu}^\top, \boldsymbol{\sigma}^2)$

$$\begin{aligned} \pi(\boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\mu}^\top, \boldsymbol{\sigma}^2) &\propto \pi(\mathbf{B} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \pi(\boldsymbol{\Sigma}^{-1}) \\ &\propto \prod_{i=1}^n \left\{ \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\boldsymbol{\beta}_i - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\boldsymbol{\beta}_i - \boldsymbol{\mu})^\top\right) \right\} |\boldsymbol{\Sigma}|^{-(d+1)} \exp\left(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}\right) \\ &\propto |\boldsymbol{\Sigma}|^{-(n+d+1+d+1)/2} \exp\left\{-\frac{1}{2}(\boldsymbol{\beta}_i - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\boldsymbol{\beta}_i - \boldsymbol{\mu})^\top - \frac{1}{2}\boldsymbol{\Sigma}^{-1}\right\} \\ &\propto |\boldsymbol{\Sigma}|^{-(n+d+1+d+1)/2} \exp\left\{-\frac{1}{2}\text{tr}(\mathbf{S}\boldsymbol{\Sigma}^{-1})\right\} \end{aligned}$$

where

$$\mathbf{S} = \mathbf{I} + \sum_{i=1}^n (\boldsymbol{\beta}_i - \boldsymbol{\mu})(\boldsymbol{\beta}_i - \boldsymbol{\mu})^\top$$

which is the form of pdf of the inverse wishart distribution Inverse Wishart(\mathbf{V}, \mathbf{S}), where $\mathbf{V} = n + d + 1$, i.e.

$$\pi(\boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\mu}^\top, \boldsymbol{\sigma}^2) \sim IW(n + d + 1, \mathbf{I} + \sum_{i=1}^n (\boldsymbol{\beta}_i - \boldsymbol{\mu})(\boldsymbol{\beta}_i - \boldsymbol{\mu})^\top)$$

4. The posterior distribution of $\pi(\boldsymbol{\mu} | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\sigma}^2, \boldsymbol{\Sigma})$

$$\begin{aligned} \pi(\boldsymbol{\mu} | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\sigma}^2, \boldsymbol{\Sigma}) &\propto \pi(\mathbf{B} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \pi(\boldsymbol{\mu}) \\ &= \prod_{i=1}^n \left\{ \det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\boldsymbol{\beta}_i - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\boldsymbol{\beta}_i - \boldsymbol{\mu})^\top\right) \right\} \\ &\propto \exp\left\{-\frac{1}{2} \sum_{i=1}^n (\boldsymbol{\beta}_i - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\boldsymbol{\beta}_i - \boldsymbol{\mu})^\top\right\} \\ &\propto \exp\left\{-\frac{1}{2} \left(\sum_{i=1}^n \boldsymbol{\beta}_i \boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}_i^\top + n\boldsymbol{\mu}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}^\top - 2 \sum_{i=1}^n \boldsymbol{\beta}_i \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}^\top \right)\right\} \\ &= \exp\left\{-\frac{1}{2} \left(\boldsymbol{\mu}(n\boldsymbol{\Sigma}^{-1})\boldsymbol{\mu}^\top - 2 \left(\sum_{i=1}^n \boldsymbol{\beta}_i \boldsymbol{\Sigma}^{-1} \right) \boldsymbol{\mu}^\top + \mathbf{C}' \right)\right\} \\ &\propto \exp\left\{-\frac{1}{2} \left(\boldsymbol{\mu} - \frac{1}{n} \sum_{i=1}^n \boldsymbol{\beta}_i \right) (n\boldsymbol{\Sigma}^{-1}) \left(\boldsymbol{\mu} - \frac{1}{n} \sum_{i=1}^n \boldsymbol{\beta}_i \right)^\top \right\} \end{aligned}$$

where

$$\mathbf{C}' = \sum_{i=1}^n \boldsymbol{\beta}_i \boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}_i^\top$$

Hence

$$\pi(\boldsymbol{\mu} | \mathbf{Y}, \mathbf{B}^\top, \boldsymbol{\sigma}^2, \boldsymbol{\Sigma}) \sim \mathcal{N}\left(\frac{1}{n} \sum_{i=1}^n \boldsymbol{\beta}_i, \frac{1}{n} \boldsymbol{\Sigma}\right)$$

3.7 Markov Chain Monte Carlo to Generate the Posterior Distributions

Because our hierarchical Bayesian Model exploited non-informative priors for four parameters, the Gibbs Sampling method would be implemented, updating parameters in the following order from their conditional posteriors distributions, \mathbf{B} , σ^2 , Σ and μ .

3.7.1 Algorithm Implementation and Estimation Since we have derived the conditional posterior of the four parameters, we implemented Gibbs sampling algorithm, updating parameters by randomly generating samples from their conditional posterior distributions.

The update of parameters is component wise, at $(t + 1)^{\text{th}}$ step, updating parameters in the following the order:

- Sample $\mathbf{B}^{(t+1)}$, i.e., sample each $\beta_i^{(t+1)}$ from $\mathcal{N}(\hat{\beta}_i^{(t)}, \hat{\Sigma}_{\beta_i}^{(t)})$
- Then, sample σ^2 from $IG(\frac{1}{2} \sum_{i=1}^n m_i, \frac{1}{2} \sum_{i=1}^n (\mathbf{Y}_i - \mathbf{X}_i \beta_i^{(t+1)})^\top (\mathbf{Y}_i - \mathbf{X}_i \beta_i^{(t+1)})$
- Next, sample $\Sigma^{(t+1)}$ from $IW(n + d + 1, \mathbf{I} + \sum_{i=1}^n (\beta_i^{(t+1)} - \mu^{(t)})(\beta_i^{(t+1)} - \mu^{(t)})^\top)$
- Finally, sample $\mu^{(t+1)}$ from $\mathcal{N}(\frac{1}{n} \sum_{i=1}^n \beta_i^{(t+1)}, \frac{1}{n} \Sigma^{(t+1)})$

3.7.2 Train Test Splits For model training as well as the performance evaluation of our Bayesian model, we split the dataset into train and test set. We first drop the hurricanes with less than 3 observations, and removed observations without transformed predictors values. Then we got a dataset with 697 different hurricanes. Our train-test split is within each hurricane's observations, that is, for each hurricane, we randomly set 80% of observations to the training set and left 20% to the test set. We then trained Bayesian model based on training data and evaluate the model performance on test dataset.

3.7.3 Inital Values For a good performance and to speed up the convergence of our algorithm, also keep some uncertainty in the MCMC process, we delicately designed the initial values.

For initial value of \mathbf{B} , we run multivariate linear regressions for each hurricane and use the regression coefficients β_i^{MLR} as the initial value for β_i . Then, the initial value of \mathbf{B} can be represented as $\mathbf{B}_{init} = (\beta_1^{MLR^\top}, \dots, \beta_n^{MLR^\top})^\top$.

For initial value of μ , we take the average of β_i^{MLR} , that is $\mu_{init} = \frac{1}{n} \sum_{i=1}^n \beta_n^{MLR}$

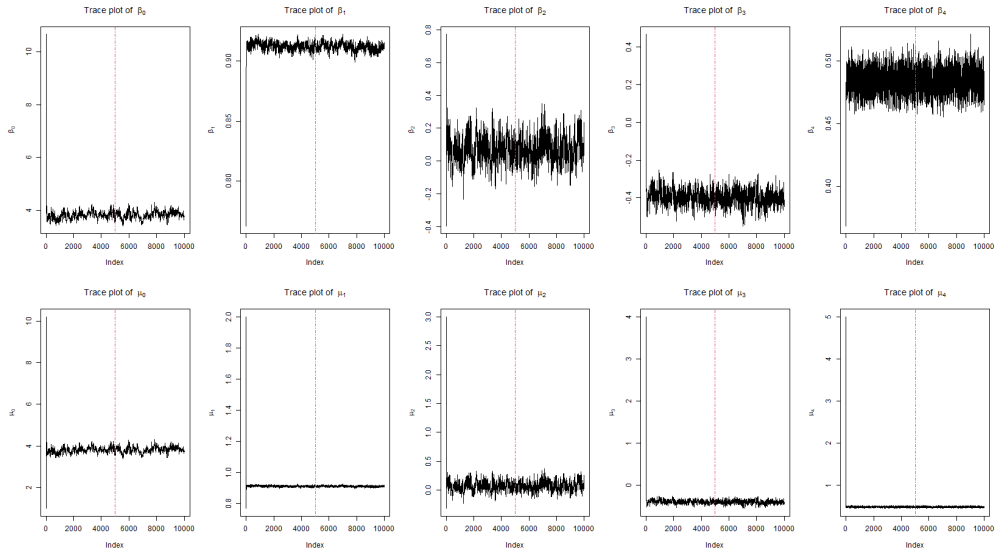
For initial value of σ^2 , we take the average of the MSE for i hurricanes.

For initial value of Σ , we just set it to a simple diagonal matrix, i.e. $\Sigma_{init} = \text{diag}(1, 2, 3, 4, 5)$

3.8 Model Results

3.8.1 Model Convergence The Figures 7 and 8 below show the estimates of each parameters over 10000 iterations. From the trace plots below, we can see that each parameters converge quickly. We take the first 5000 iterations as our burn-in period and use the last 5000 iterations (iterations 5001 to 10000) to do the posterior parameter estimates and inferences. Figures 9 and 10 display the histogram of parameters based on the last 5000 MCMC samples. From the histograms, we found that for β_i and μ , the distributions are relatively normal. However, we also found some skewness in the distributions of Σ , e.g. Σ_{11} , Σ_{12} , Σ_{33} and Σ_{44}

From the corresponding autocorrelation function (ACF) plots, we can get a sense of the degree of serial correlation of the MCMC draws. From these plots we see how the sample autocorrelation between the parameters of the chain decreases as a function of their lag. The ACF plots of μ and β_i shows that autocorrelation is large at short lags, but then goes to around zero pretty quickly. For the ACF of Σ and σ^2 , it seems that for some parameters, the autocorrelation doesn't shrink to zero even after 500 lag (e.g. Σ_{13}).



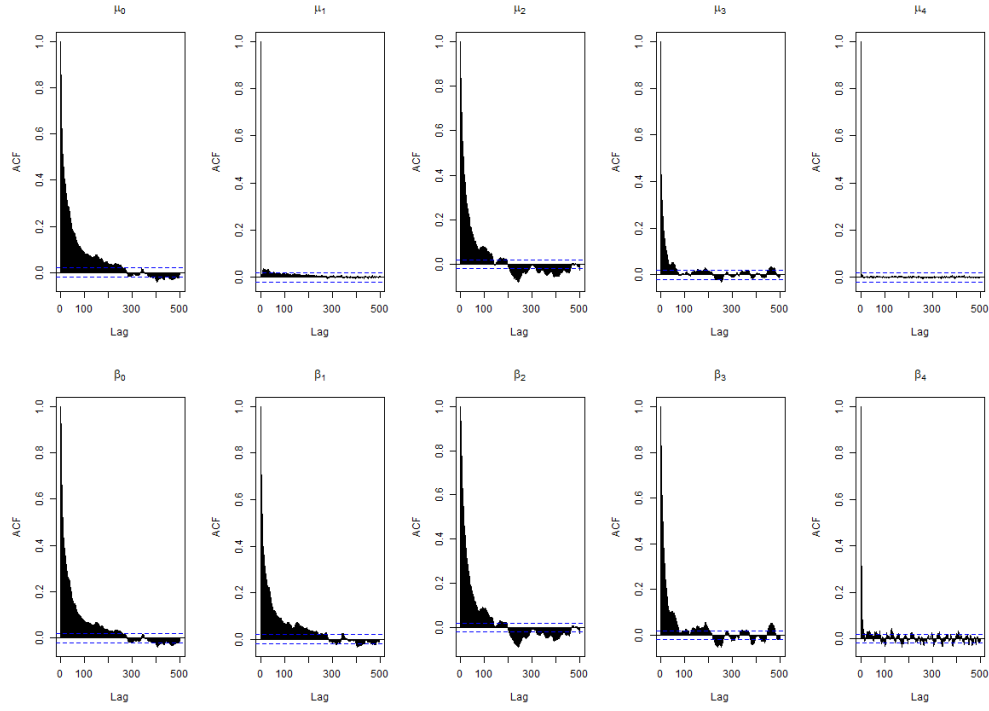
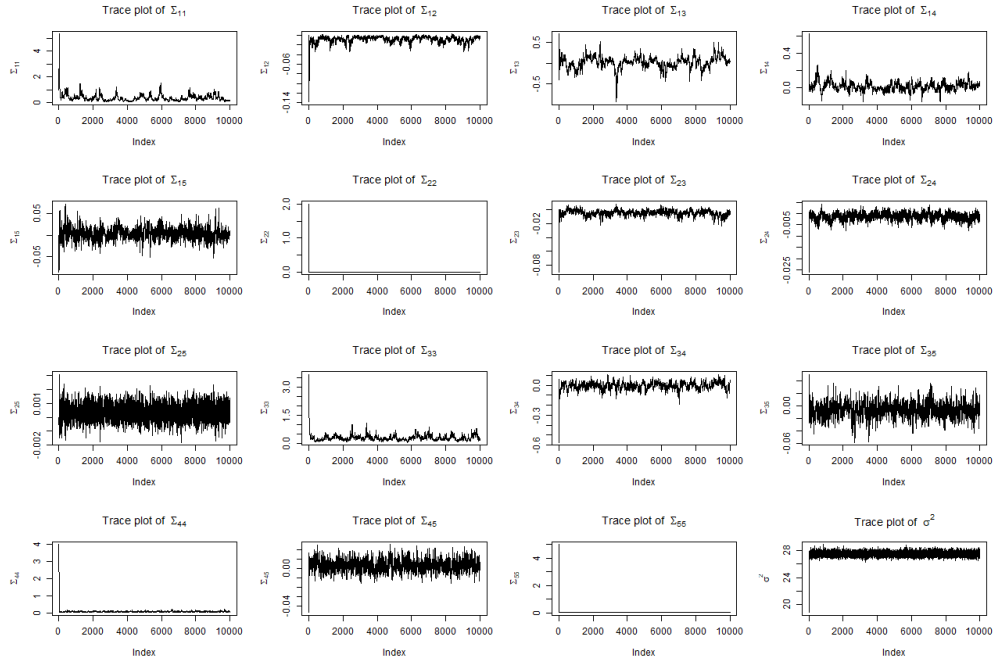


Figure 7. Trace plots and ACF plots of model parameters, based on 10000 MCMC sample



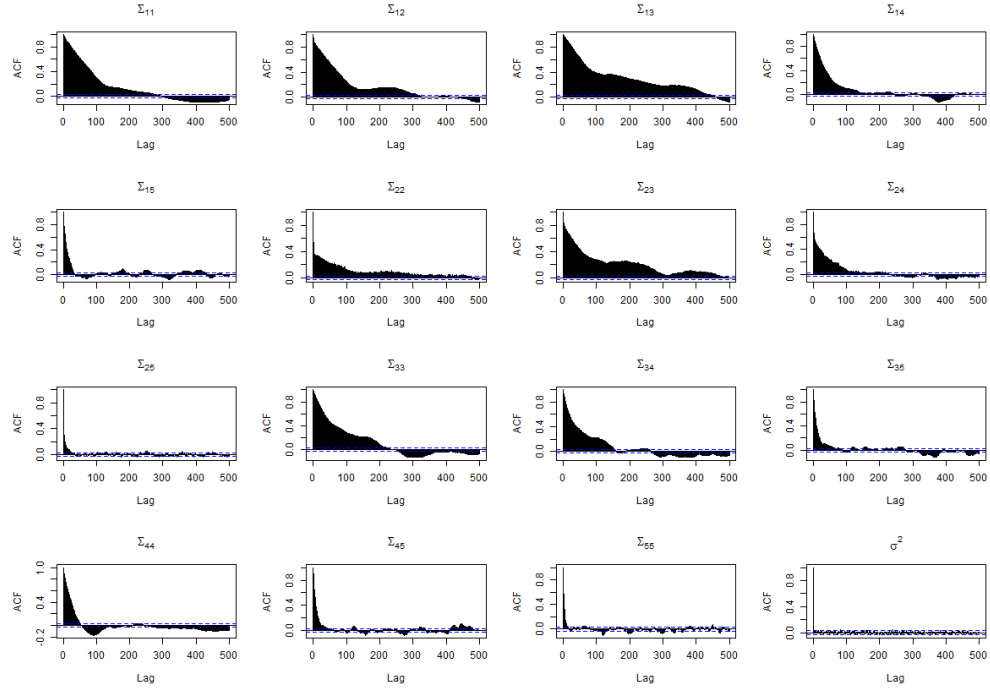


Figure 8. Trace plots and ACF plots of variance parameters, based on 10000 MCMC sample

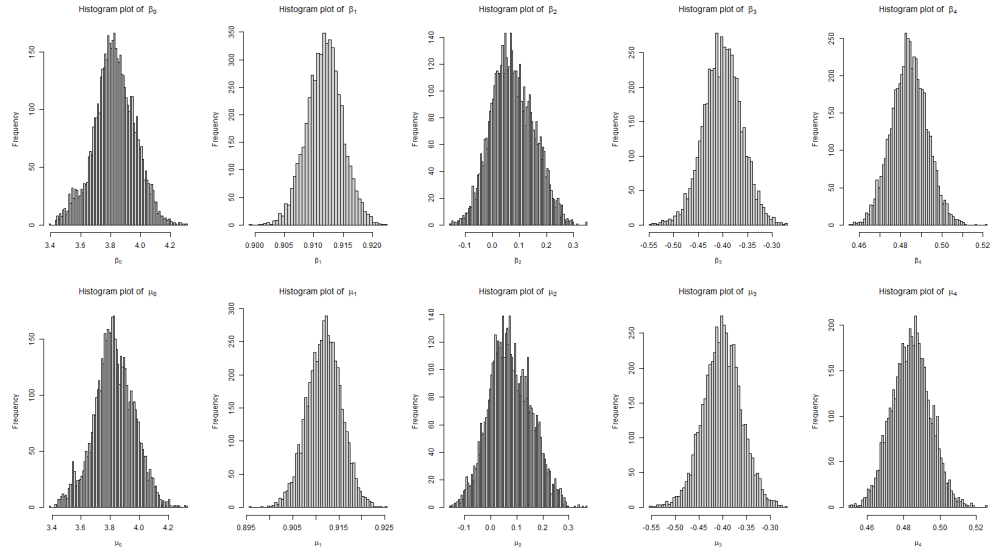


Figure 9. Histograms of model parameters, based on last 5000 MCMC sample

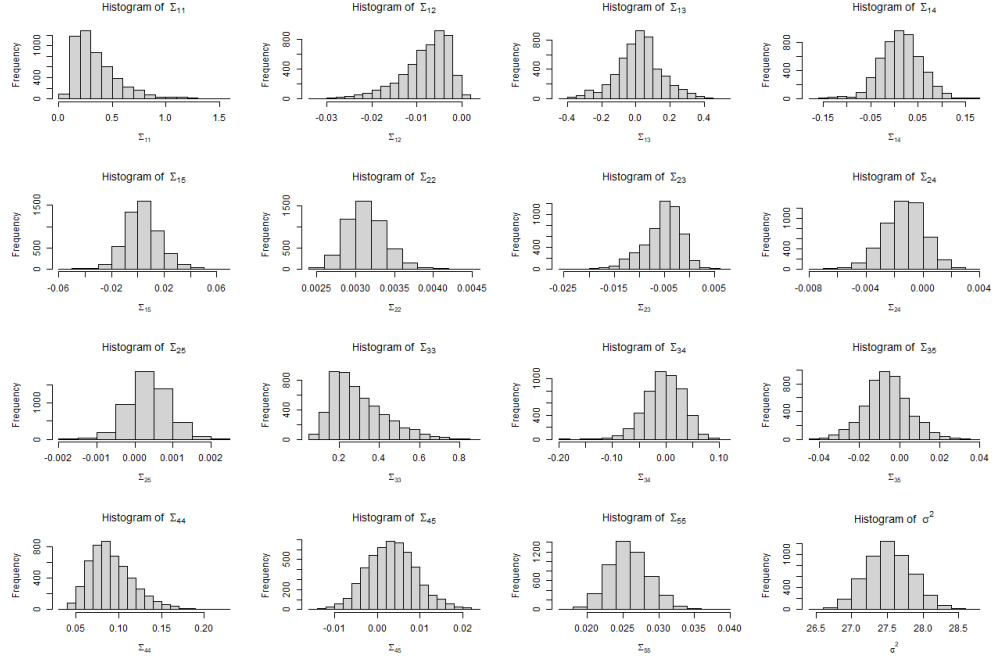


Figure 10. Histograms of variance parameters, based on last 5000 MCMC sample

3.8.2 Parameter Estimates As we stated above, we take the first 5000 iterations as our burn-in period and use the last 5000 iterations (iterations 5001 to 10000) to do the posterior parameter estimates and inferences.

We use the posterior mean as the posterior estimates for each parameter, and the 95% confidence interval is derived based on empirical rule, since the distributions of parameters are almost normal. Notice that the estimates of \mathbf{B} including all β_i 's for all hurricanes, to do the estimation and inference, we take the average value of all β_i 's in each last 5000 MCMC iterations and then take the average of these results as the estimate of \mathbf{B} . So the estimate of \mathbf{B} reflects the sample mean response on the change of the hurricane speed across the predictors. The following tables shows the posterior estimates of the parameters

Variables	$\bar{\beta}_i$	Var($\bar{\beta}_i$)	95% CI of $\bar{\beta}_i$	$\bar{\mu}$	Var($\bar{\mu}$)	95% CI of $\bar{\mu}$
intercept	3.8252	0.0185	(3.789,3.8614)	3.8166	0.0190	(3.7795,3.8538)
Wind_prev	0.9118	0.0000	(0.9118,0.9118)	0.9121	0.0000	(0.9121,0.9122)
Lat_change	0.0744	0.0060	(0.0626,0.0862)	0.0720	0.0065	(0.0594,0.0847)
Long_change	-0.4014	0.0015	(-0.4044,-0.3985)	-0.3968	0.0016	(-0.4,-0.3936)
Wind_change	0.4841	0.0001	(0.484,0.4843)	0.4847	0.0001	(0.4844,0.4849)

Table 1. Posterior Estimates of Model parameters μ and β_i

Parameters	Estimates	Variance	95% CI
Σ_{11}	0.3493	0.0435	(-0.0595,0.7581)
Σ_{12}	-0.0081	0.0000	(-0.0189,0.0027)
Σ_{13}	0.0201	0.0176	(-0.2399,0.2801)
Σ_{14}	0.0131	0.0019	(-0.0725,0.0987)
Σ_{15}	0.0035	0.0002	(-0.0215,0.0285)
Σ_{22}	0.0031	0.0000	(0.0026,0.0036)
Σ_{23}	-0.0053	0.0000	(-0.0125,0.0019)
Σ_{24}	-0.0013	0.0000	(-0.0041,0.0014)
Σ_{25}	0.0004	0.0000	(-7e-04,0.0014)
Σ_{33}	0.2960	0.0176	(0.0362,0.5558)
Σ_{34}	-0.0031	0.0012	(-0.0716,0.0653)
Σ_{35}	-0.0060	0.0001	(-0.0276,0.0156)
Σ_{44}	0.0918	0.0007	(0.0412,0.1424)
Σ_{45}	0.0034	0.0000	(-0.008,0.0148)
Σ_{55}	0.0258	0.0000	(0.0203,0.0313)
σ^2	27.5247	0.1030	(26.8957,28.1538)

Table 2. Posterior Estimates of Model parameters Σ and σ^2

The posterior estimates of μ reflects the population average of our model parameters. For μ_0 , which is the parameter of the intercept, it tells us the average wind speed of all hurricanes. μ_1 is the coefficient associated with the previous wind speed, the estimates of it is 0.912 which is positive, means that an increase in the previous wind speed will causes a higher speed for that in the next time. μ_2 and μ_3 are the coefficients associated with the change in latitude and longitude of the hurricane, it seems that an increase in the change of latitude is associated with the increments in wind speed, on the other hand, an increase in the change of longitude is associated with the decrements in wind speed. Finally, the estimated μ_4 is $0.485 > 0$, indicating that an increase in the change of wind speed will cause an increase in future wind speed.

The posterior estimate of Σ and the corresponding ρ is

$$\Sigma = \begin{pmatrix} 0.349 & -0.008 & 0.020 & 0.013 & 0.004 \\ -0.008 & 0.003 & -0.005 & -0.001 & 0.0004 \\ 0.020 & -0.005 & 0.296 & -0.003 & -0.006 \\ 0.013 & -0.001 & -0.003 & 0.092 & 0.003 \\ 0.004 & 0.0004 & -0.006 & 0.003 & 0.026 \end{pmatrix}, \rho = \begin{pmatrix} 1 & -0.245 & 0.063 & 0.073 & 0.037 \\ -0.245 & 1 & -0.174 & -0.078 & 0.041 \\ 0.063 & -0.174 & 1 & -0.019 & -0.069 \\ 0.073 & -0.078 & -0.019 & 1 & 0.070 \\ 0.037 & 0.041 & -0.069 & 0.070 & 1 \end{pmatrix}$$

From the correlation matrix, we can found that the correlation between $\beta_{j,i}$ $\beta_{k,i}$ ($j \neq k$) is not that strong.

3.9 Model Performance

We evaluate the performance of predictive ability by calculating the RMSE and the R^2 values for each hurricane. The residuals of Bayesian estimates that converged after iterations from MCMC will be used to predict the wind speed of test dataset. The overall R^2 is 0.822 and overall RMSE is 4.51. The valid R^2 is filtered with values between 0 and 1 and we get 77.5% hurricanes (540) indicating that 22.5% of the estimated Bayesian models do not track hurricanes well and have negative R^2 . One of the reason may be the limited number of observations of the hurricanes. Table 3 shows the 10 hurricanes with the least 10 RMSE. R^2 are also large enough to indicates that the estimated model track most hurricanes well and the smallest RMSE is GUSTAV.1996 with R^2 being 0.952.

ID	r_square	rmse
GUSTAV.1996	0.952	0.537

ID	r_square	rmse
LORENZO.2001	0.914	0.733
ERIN.2013	0.878	0.823
JOSE.2011	0.970	0.872
GRETA.1970	0.980	0.876
DELTA.1972	0.825	0.904
EDITH.1967	0.826	0.983
FABIAN.1997	0.955	1.002
DEBBY.2006	0.984	1.045
CRISTOBAL.2002	0.956	1.053

Table 3. R-square and RMSE for prediction result on test data

Figure 11 shows the actual wind speed and the estimated wind speed of randomly selected four hurricanes . We can see that most parts of the two curves overlapped indicating that the predicted values are close to the actual values. In DEBBY.2006, we can see that this is a very good model prediction with small deviation.

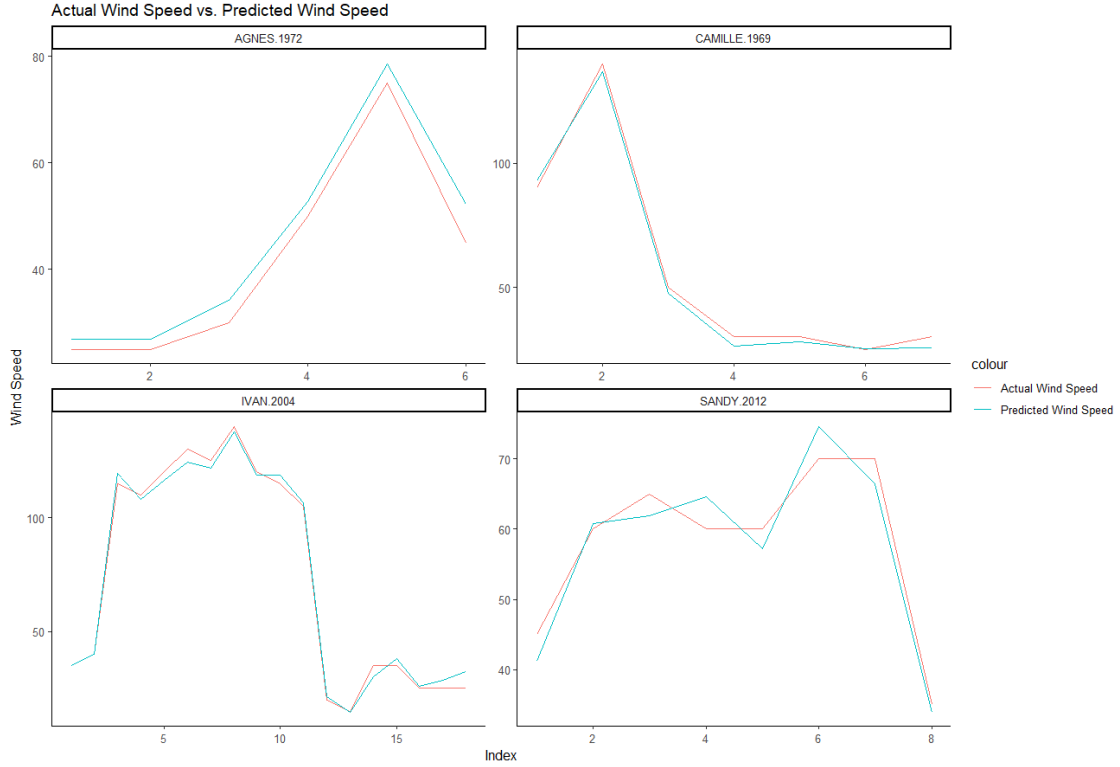


Figure 11. Actual Wind Speed vs. Predicted Wind Speed

4 Explore the seasonal differences and wind speed change

Now based on the estimated Bayesian model from previous questions, we need to explore the seasonal difference. We can fit 5 models using 5 estimated beta values against the three predictors: $X_{i,1}$: the month of the year the i th hurricane started, $X_{i,2}$: the year of the i th hurricane and $X_{i,3}$: the nature of the i th hurricane. The beta values obtained from previous Gibbs Sampler MCMC method contains the mean value of $\beta_{0,i}$, $\beta_{1,i}$, $\beta_{2,i}$, $\beta_{3,i}$ and $\beta_{4,i}$ for each of the 697 unique hurricanes, which is of the size $697 * 5$.

According to the summary, the R squared value for all the five fitted linear models are quite small, which may indicate bad fit. In addition, most coefficients for the model are not significant with a p-value larger than 0.05. However, for those significant coefficients, we could infer a potential relationship between the certain predictors and the beta coefficients respectively. We should consult the previous Bayesian model:

$$Y_i(t+6) = \beta_{0,i} + \beta_{1,i}Y_i(t) + \beta_{2,i}\Delta_{i,1}(t) + \beta_{3,i}\Delta_{i,2}(t) + \beta_{4,i}\Delta_{i,3}(t) + \epsilon_i(t)$$

to interpret the change of the influence on $Y_{i,t+6}$ as the value of the predictor changes.

For the fitted coefficients of β_0 to β_4 , the intercept cannot show information about seasonal difference since they indicate when holding all the predictors zero, the value for the corresponding β . We can only observe that the year is quite significant in the model for β_0 , β_1 with both negative estimates close to zero. Therefore, as the year increase, the coefficient of the intercept and $Y_{i,t}$ may decrease a little, which means for the Bayesian model, the wind speed when holding all the variables zero and the effect the previous wind speed has will decrease over years. Apart from seasonal difference, some other predictors are quite significant, such as **natureET** for β_2 , **natureTS** for β_3 .

	Estimate	Pr(> t)	Estimate	Pr(> t)	Estimate	Pr(> t)	Estimate	Pr(> t)	Estimate	Pr(> t)
(Intercept)	4.481	0.000	1.343	0.000	0.041	0.951	-0.834	0.019	0.289	0.448
monthApril	0.023	0.835	0.015	0.670	0.017	0.931	0.042	0.680	0.036	0.739
monthMay	0.026	0.783	0.000	0.997	0.071	0.660	0.063	0.458	-0.016	0.859
monthJune	0.028	0.765	0.005	0.851	-0.007	0.964	0.056	0.505	0.024	0.792
monthJuly	0.013	0.891	0.015	0.590	-0.009	0.954	0.036	0.664	0.013	0.884
monthAugust	-0.020	0.828	0.023	0.412	-0.052	0.738	0.012	0.881	0.031	0.726
monthSeptember	-0.007	0.938	0.026	0.359	-0.036	0.817	0.021	0.797	0.044	0.618
monthOctober	0.009	0.919	0.021	0.459	-0.029	0.855	0.034	0.680	0.035	0.694
monthNovember	0.015	0.875	0.025	0.393	0.024	0.879	0.026	0.753	0.021	0.817
monthDecember	0.006	0.953	0.009	0.772	-0.054	0.745	0.042	0.633	0.011	0.905
year	0.000	0.072	0.000	0.000	0.000	0.910	0.000	0.203	0.000	0.625
natureET	0.001	0.977	0.004	0.688	-0.070	0.169	-0.026	0.329	-0.021	0.473
natureNR	0.001	0.987	-0.015	0.333	0.006	0.943	0.003	0.944	-0.022	0.646
natureSS	0.014	0.490	-0.003	0.602	-0.001	0.969	0.013	0.496	-0.024	0.234
natureTS	0.012	0.479	-0.006	0.249	-0.015	0.588	-0.023	0.126	-0.017	0.283

Table 4. Coefficients of the fitted β model against three predictors

We also try to represent the months as four seasons and fit a model for β against them. Each model has three dummy variables corresponding to the three seasons except Spring. The latter three rows of estimate shows how the value of β differentiate between Spring and the other three seasons respectively. If with a rather small p-value, we can conclude the existence of seasonal difference. Therefore, by constructing model in this way, we find that β_1 and β_4 will increase a little as season changes from Spring to Summer, then to Autumn, which means a season difference of the effect $Y_{i,t}$ and $\Delta_{i,3}(t)$ has on the wind speed. For β_2 , β_3 , Summer and Autumn may lead to a slightly smaller effect of $\Delta_{i,1}(t)$, $\Delta_{i,2}(t)$ have on the wind speed compared to Spring.

response	coefficient	Estimate	Pr(> t)
β_0	Intercept	3.837	0.000
β_0	seasonSummer	-0.031	0.205
β_0	seasonAutumn	-0.024	0.325
β_0	seasonWinter	-0.019	0.654
β_1	Intercept	0.894	0.000
β_1	seasonSummer	0.015	0.044
β_1	seasonAutumn	0.021	0.005
β_1	seasonWinter	0.003	0.794
β_2	Intercept	0.161	0.000
β_2	seasonSummer	-0.098	0.017
β_2	seasonAutumn	-0.091	0.025
β_2	seasonWinter	-0.098	0.164
β_3	Intercept	-0.350	0.000
β_3	seasonSummer	-0.047	0.034
β_3	seasonAutumn	-0.043	0.046
β_3	seasonWinter	-0.009	0.802
β_4	Intercept	0.442	0.000
β_4	seasonSummer	0.036	0.120
β_4	seasonAutumn	0.049	0.035
β_4	seasonWinter	0.015	0.711

Table 5. Coefficients of the fitted β model against season

Now fit linear models for β against the season variables (corresponding to the year) to seek for potential evidence of the statement :“the wind speed has been increasing over years”. In order to analyze this question, need to inspect on model which corresponds to the wind speed and the year. For β_2 model, the estimate of year is significant, although it's really close to zero. Therefore, we can infer that as the year increases, the impact past wind speed has on the current wind speed may decrease a little, which cannot provide support for the statement. However, it's quite match with the results shown in the figures in the initial EDA session, which indicates the mean wind speed tends to decrease over years.

response	coefficient	Estimate	Pr(> t)
β_0	Intercept	4.514	0.000
β_0	year	0.000	0.050
β_1	Intercept	1.345	0.000
β_1	year	0.000	0.000
β_2	Intercept	-0.106	0.863
β_2	year	0.000	0.776
β_3	Intercept	-1.027	0.002
β_3	year	0.000	0.053
β_4	Intercept	0.305	0.382
β_4	year	0.000	0.607

Table 6. Coefficients of the fitted β model against year

In conclusion, for different months, there is no significant differences observed. Over years, the effect the wind speed 6 hours ago has on the current wind speed may decrease a little. And there is no evidence to support the statement in task 5.

5 Predict the hurricane-induced damage and deaths

Firstly, we plot deaths and financial loss separately. Figure 12. shows the distributions of deaths and damage. We could easily find a few points which are far away from most of the points indicate serious damage of society. In predictions of disasters, these extreme points are important because they enable the model to predict the worst outcome. Therefore, we keep these points in model building.

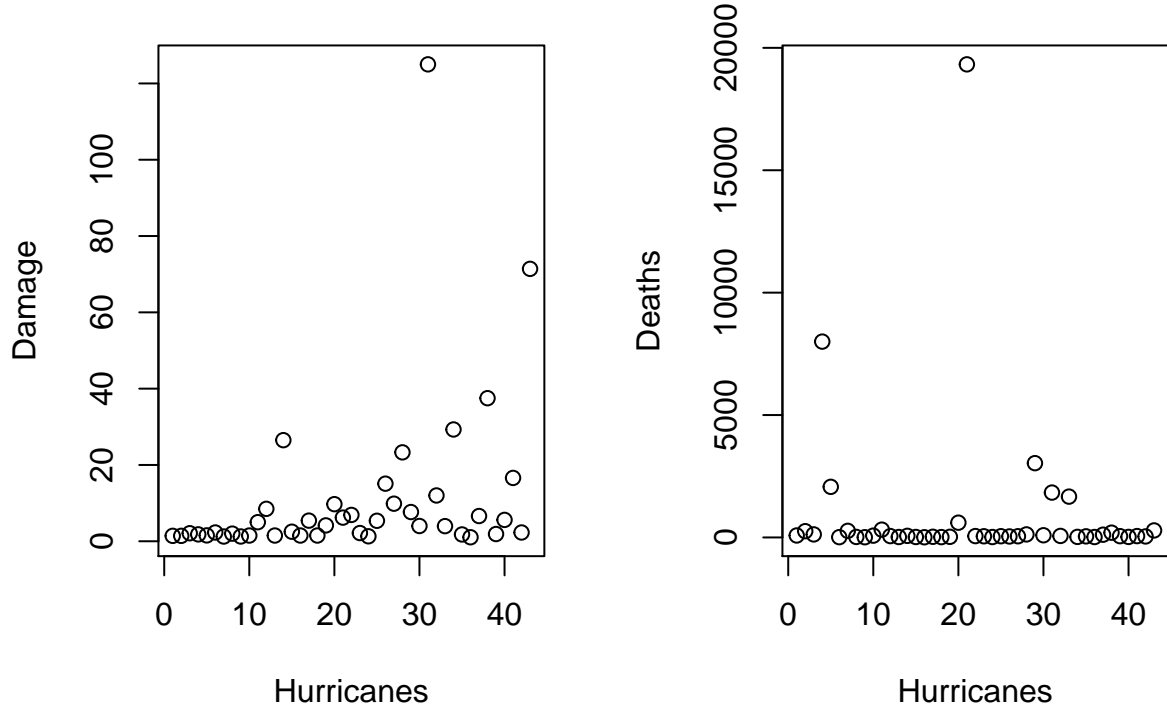


Figure 12. Distributions of Damage and Deaths

In order to build a model that combines information in original data and the estimated coefficients from the Bayesian model, we extract the coefficients from the previous results. By taking the average of β_i at different time points, we obtain $\beta_0 \sim \beta_4$ of each hurricane. Part of the results is shown in Table 7.

id	intercept	beta1	beta2	beta3	beta4
agnes.1972	3.950974	0.9224097	0.0059532	-0.3103372	0.5453543
alex.2010	3.798737	0.9370333	0.0698849	-0.3937358	0.5400187
alicia.1983	3.897408	0.9036878	-0.0748341	-0.3994486	0.5477718
allen.1980	3.687070	0.9655304	0.1306393	-0.5460144	0.5466129
andrew.1992	3.676279	0.9375384	-0.2843257	-0.5782973	0.5370158
betsy.1965	3.808396	0.9513766	-0.4500720	-0.3890718	0.4244575
bob.1991	3.629466	0.9232143	0.0279527	-0.5751636	0.4382048
camille.1969	3.994355	0.9355674	0.0729188	-0.5734830	0.6703910
charley.2004	3.638829	0.9482764	-0.1797332	-0.6955016	0.1818395
david.1979	3.789678	0.9579657	-0.0461134	-0.3823658	0.6853938

Table 7. Coefficients of Each Hurricane

Fortunately, 43 hurricanes recorded in *hurricanoutcome2.csv* are also in *hurrican703.csv*. Thus, we merge two data frame by hurricane id to predict the deaths and damage caused by hurricanes.

The death variable is a count variable, so we decided to use Poisson regression to analysis relationship between death and other variables excluding damage. We use **Total.Pop** and **Hours** as the offset, since the outcome of deaths is proportional and the results would be different in some dimension (different populations, different duration). The Poisson regression is:

$$\log(E(Deaths)) = \beta_i X_i + offset$$

Where X_i indicates all predictors included in the model. We use **glm** function to achieve the Poisson model. The coefficients result is in Table 8.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	272.6548414	11.9772996	22.764300	0.0000000
intercept	12.8670811	0.2681866	47.978086	0.0000000
beta1	139.3302799	2.2527075	61.850143	0.0000000
beta2	6.2673041	0.1222299	51.274709	0.0000000
beta3	11.6385648	0.3090069	37.664416	0.0000000
beta4	-15.5433361	0.3093214	-50.249788	0.0000000
nobs	-0.0259659	0.0011052	-23.494615	0.0000000
Season	-0.0101636	0.0022226	-4.572917	0.0000048
MonthJuly	-2.6409629	0.1476034	-17.892292	0.0000000
MonthJune	-0.3043892	0.0911374	-3.339894	0.0008381
MonthNovember	-3.0007313	0.1534637	-19.553361	0.0000000
MonthOctober	-2.1031926	0.0625508	-33.623762	0.0000000
MonthSeptember	-0.6480044	0.0473666	-13.680606	0.0000000
NatureNR	2.1051263	0.1265128	16.639630	0.0000000
NatureTS	4.2677929	0.1266627	33.694166	0.0000000
Maxspeed	0.0099355	0.0010283	9.661642	0.0000000
Meanspeed	-0.0634120	0.0033070	-19.175107	0.0000000
Maxpressure	-0.4345774	0.0087466	-49.685178	0.0000000
Meanpressure	0.0092073	0.0002035	45.253527	0.0000000
Percent.Poor	0.0520424	0.0008541	60.929497	0.0000000
Percent.USA	-0.0206136	0.0005851	-35.231400	0.0000000

Table 8. Coefficients of Deaths Prediction

From the results, $\beta_0 \sim \beta_4$ indicate the relatively strong association. Especially, β_1 , which represents the earlier wind speed has the biggest coefficient. We could conclude that high wind speed of hurricane more easily leads to serious casualties. Also, months seem to be an important factor in prediction. Comparing to June and September, July, November and October have lower proportion of death given all other variables constant.

In order to obtain the integer data, we transform the units of **Damage** from billion to million. Thus, **Damage** could be regarded as a count variable which could also be fitted by Poisson regression. In order to adjust the exposure, we use **Hours** as the offset.

$$\log(E(Damage * 1000)) = \beta_i X_i + offset$$

Where X_i presents all predictors included in the model. We use **glm** function to achieve the Poisson model. The coefficients results is in Table 9.

term	estimate	std.error	statistic	p.value
(Intercept)	-206.5342295	2.0185627	-102.317470	0.00e+00
intercept	4.7792873	0.0279968	170.708107	0.00e+00
beta1	60.3768193	0.4513469	133.770326	0.00e+00
beta2	-1.0909958	0.0132414	-82.392801	0.00e+00
beta3	3.6396116	0.0259580	140.211430	0.00e+00
beta4	-1.6090883	0.0335258	-47.995549	0.00e+00
nobs	0.0308929	0.0002582	119.662746	0.00e+00
Season	0.0757660	0.0003973	190.691693	0.00e+00
MonthJuly	0.4806325	0.0187755	25.598871	0.00e+00
MonthJune	-3.2641224	0.0241511	-135.154056	0.00e+00
MonthNovember	-1.8378419	0.0249958	-73.526098	0.00e+00
MonthOctober	-1.3044789	0.0094414	-138.165832	0.00e+00
MonthSeptember	-1.7755988	0.0077805	-228.211660	0.00e+00
NatureNR	-4.2815939	0.0356423	-120.126833	0.00e+00
NatureTS	-1.9548375	0.0143805	-135.936287	0.00e+00
Maxspeed	0.0507898	0.0002142	237.112410	0.00e+00
Meanspeed	-0.0644452	0.0004911	-131.235042	0.00e+00
Maxpressure	-0.0140045	0.0012224	-11.456463	0.00e+00
Meanpressure	-0.0001769	0.0000407	-4.351453	1.35e-05
Total.Pop	0.0000003	0.0000000	62.112535	0.00e+00
Percent.Poor	-0.0384685	0.0001886	-203.998906	0.00e+00
Percent.USA	-0.0049715	0.0000729	-68.209623	0.00e+00

Table 9. Coefficients of Damage Prediction

The results of coefficients in predicting damage also show the importance of β_1 . From the model, we can see that serious casualties are also accompanied by serious financial losses which are strongly influenced by earlier wind speed and are slightly affected by months, latitude change, longitude change and wind speed change. $\beta_0 \sim \beta_4$ are generally powerful in damage and deaths prediction.

6 Discussions

6.1 Conclusions

Our MCMC algorithm successfully estimates the high-dimensional parameters. Firstly, All the parameters converges quickly under a good initial values setting. Secondly, The overall R^2 is relatively large, our model fits the data well.

Based on posterior estimates of μ , an increase in current wind speed and the change in wind speed is associated with increase in the wind speed in the upcoming future. And the posterior estimate of Σ shows that the correlations between β_i 's are relatively small.

In exploring how the year, month and the nature of hurricanes affect the wind speed, we found that for different months, there is no significant differences observed. Over years, the effect the wind speed 6 hours ago has on the current wind speed may decrease a little. The nature of the hurricanes does play a important role in wind speed.

Finally, when focusing on predicting the death and loss caused by hurricanes, the β_i coefficients estimated from the Bayesian model have great importance. Especially β_1 , which is the estimated coefficient associated with the previous wind speed, is significant and relatively strong in both damage and loss model, which may indicates that the wind speed of a hurricane strongly affect the damage and financial loss caused by it.

6.2 Limitations

Due to the computational cost of MCMC algorithm, we do not have the chance to test our algorithm stability by setting different starting values. We did some exploration about the starting value (see Appendix), even if the starting value is ridiculous, our MCMC algorithm successfully converged and the posterior estimates are pretty close to the ones with good initial value. More effort should be done to test the model reliability in the future.

In model training and testing, we dropped several hurricanes' data with few observations (<3). However, the lack of observations may be associated with the nature of the hurricanes and hence affect the damage and financial loss caused. Carefully studies should be done on these short-term hurricanes data.

From the ACF plots of components in the variance covariance matrix, we can see that for some parameters, not only autocorrelation is large at short lags, but it also dies out very slowly. We should future examine the serial correlation between successive draws.

6.3 Group Contribution

Xinran contributed to the data cleaning and EDA of the hurricanes data and the introduction of the project. Renjie contributed to the math derivation of the posterior distributions and designed the MCMC algorithm, and assesing the results from MCMC. Shengzhi contributed to the model evaluation and interpretation of the model results. Hao focusing on exploring the seasonal and annual difference in the change of the wind speed. Wentong helped to analysis the important predictors in predicting the damage and financial loss caused by hurricanes.

Appendix

The code of data cleaning and EDA is in the .Rmd file of the report. For more information see https://github.com/MefiMefi/P8160_Project3_MCMC/

Posterior Distribution Sampling

```
# posterior of B
# each beta_i 1*5
# B n*5

post.B <- function(dt, muvec, sigma2, Sigma){
  n = length(dt)
  B = NULL
  RSS = NULL
  m = NULL
  for (i in 1:n){
    # stuffs to define the distribution
    X = as.matrix(dt[[i]][,-1])
    y = as.vector(dt[[i]][,1])
    V = sigma2^(-1) * t(X) %*% X + solve(Sigma)
    M = sigma2^(-1) * t(y) %*% X + muvec %*% solve(Sigma)
    mean_bi = solve(V) %*% t(M) # we need a 5*1
    vcov_bi = solve(V)
    bi = mvrnorm(1, mu = mean_bi, Sigma = vcov_bi)
    B = rbind(B, bi)
    # calculate RSS, m_i for sigma2
    RSS = rbind(RSS, sum((y - X %*% bi)^2))
    m = rbind(m, nrow(X))
  }
  return(list(B = B, RSS = RSS, m = m))
}
```

```

}

# test passed
#testB = post.B(dt = dt_mtx, muvec = rep(0,5), sigma2 = 2, Sigma = diag(2,5,5))

# sampling sigma2 from a inverse gamma, need sum(m_i) and SSR
post.sigma2 <- function(m, RSS){
  alpha_ = sum(m)/2
  beta_ = sum(RSS)/2
  sigma2 = rinvgamma(1, alpha = alpha_, beta = beta_)
  return(sigma2)
}

# test passed
#testsigma2 = post.sigma2(testB$m, testB$RSS)

# sampling Sigma from inverse wishart distribution
post.Sigma <- function(B, muvec){
  n = nrow(B)
  S.matrix = diag(1, 5, 5)
  for (i in 1:n){
    beta_i = B[i,]
    S.matrix = S.matrix + (beta_i-muvec) %*% t(beta_i-muvec)
  }
  v = n + 5 + 1
  Sigma = rinvwishart(1, nu = v, Omega = S.matrix, checkSymmetry = F)
  return(Sigma[,1])
}

# test passed
#testSigma = post.Sigma(testB$B, rep(1,5))

# sampling muvec from multivariate normal distribution
post.mu <- function(B, Sigma){
  n = nrow(B)
  mean_mu = colMeans(B)
  vcov_mu = 1/n * Sigma
  muvec = mvrnorm(1, mean_mu, vcov_mu)
  return(muvec)
}

# test passed
#testmuvec = post.mu(testB$B, testSigma)

```

MCMC Algorithm

```

MCMC.Gibbs <- function(dt, init.B, init.muvec, init.sigma2, init.Sigma, max.iter = 1e4){
  B.res = list()
  sigma2.res = list()
  Sigma.res = list()
  muvec.res = list()

```

```

# initialize B
init.Bstuff = post.B(dt, init.muvec, init.sigma2, init.Sigma)

cur.stuff = NULL
cur = list(B = init.B, sigma2 = init.sigma2, muvec = init.muvec, Sigma = init.Sigma, RSS = init.Bstuff$RSS)

B.res[[1]] = init.B
sigma2.res[[1]] = cur$sigma2
Sigma.res[[1]] = cur$Sigma
muvec.res[[1]] = cur$muvec

prev.stuff = cur.stuff
prev = cur

# start iteration
for (i in 2:max.iter) {
  cur.stuff = post.B(dt, prev$muvec, prev$sigma2, prev$Sigma)
  sigma2 = post.sigma2(cur.stuff$m, cur.stuff$RSS)
  Sigma = post.Sigma(cur.stuff$B, prev$muvec)
  muvec = post.mu(cur.stuff$B, Sigma)

  # update current parameters
  cur = list(B = cur.stuff$B, sigma2 = sigma2, muvec = muvec, Sigma = Sigma, RSS = cur.stuff$RSS)
  prev.stuff = cur.stuff
  prev = cur

  # append to final result
  B.res[[i]] = cur$B
  sigma2.res[[i]] = cur$sigma2
  Sigma.res[[i]] = cur$Sigma
  muvec.res[[i]] = cur$muvec
}
return(list(
  B = B.res,
  sigma2 = sigma2.res,
  Sigma = Sigma.res,
  muvec = muvec.res
))
}

```

MCMC implementation

```

mu0 = colMeans(ols.B, na.rm = T)
sigmasq0 = mean(ols.residual.res)

# Good initial value
MCMC.res.2 = MCMC.Gibbs(dt_train, init.B = ols.B, init.muvec = mu0, init.sigma2 = sigmasq0, init.Sigma = diag(5, 697))
save(MCMC.res.2, file = "MCMC.res.2.RData")

# Bad starting value
bad.B = matrix(seq(1:(5*697)), 697)
bad.mu = colMeans(bad.B)

```

```
MCMC.res.bad = MCMC.Gibbs(dt_train, init.B = bad.B, init.muvec = bad.mu, init.sigma2 = sigmasq0*1000,
save(MCMC.res.bad,file = "MCMC.res.bad.RData")
```

Model Performance

```
beta.res.postmean = beta.res.postmean %>% rename(beta_0 = intercept,
                                                beta_1 = Wind_prev,
                                                beta_2 = Lat_change,
                                                beta_3 = Long_change,
                                                beta_4 = Wind_change)

dt_res = merge(dt_test_id, beta.res.postmean, by = "ID")

dt_res = dt_res %>%
  mutate(Wind_kt_pred = beta_0*intercept+beta_1*Wind_prev
          +beta_2*Lat_change+beta_3*Long_change+beta_4*Wind_change) %>%
  group_by(ID) %>%
  mutate(r_square = 1-(sum((Wind_kt_pred-Wind.kt)^2))/(sum((Wind.kt-mean(Wind.kt))^2)),
         rmse = rmse(Wind.kt,Wind_kt_pred))
```

```
dt_rmse=
dt_res %>%
dplyr::select(ID, r_square, rmse) %>%
distinct() %>%
mutate(r_square = round(r_square, 3),
       rmse = round(rmse,3)) %>%
filter(r_square > 0 && r_square < 1) %>%
arrange(rmse)
mean(dt_rmse$r_square)
mean(dt_rmse$rmse)
```

Seasonal Difference

```
load("./dt_long.RData")
load("./ID_in.RData")
load("./beta.res.postmean.RData")

dt_season <-
  dt_long %>%
  drop_na() %>%
  filter(ID %in% ID_in) %>%
  distinct(ID, .keep_all = TRUE) %>%
  select(ID, Season, Month, Nature) %>%
  mutate(Month = factor(Month, levels = month.name))
```

```
season_diff <-
  merge(dt_season, beta.res.postmean, by = c("ID")) %>%
  janitor::clean_names()
colnames(season_diff)[2] <- "year"

# Beta0
intercept.fit <- lm(intercept ~ month + year + nature, data = season_diff)
# Beta1
```

```

wind_prev.fit <- lm(wind_prev ~ month + year + nature, data = season_diff)
# Beta2
lat_change.fit <- lm(lat_change ~ month + year + nature, data = season_diff)
# Beta3
long_change.fit <- lm(long_change ~ month + year + nature, data = season_diff)
#Beta4
wind_change.fit <- lm(wind_change ~ month + year + nature, data = season_diff)

summary(intercept.fit)
summary(wind_prev.fit)
summary(lat_change.fit)
summary(long_change.fit)
summary(wind_change.fit)

sum0 <- summary(intercept.fit)$coefficients[,c(1,4)]
sum1 <- summary(wind_prev.fit)$coefficients[,c(1,4)]
sum2 <- summary(lat_change.fit)$coefficients[,c(1,4)]
sum3 <- summary(long_change.fit)$coefficients[,c(1,4)]
sum4 <- summary(wind_change.fit)$coefficients[,c(1,4)]

kable(cbind(sum0, sum1, sum2, sum3, sum4)) %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed")) %>%
  add_header_above(c(" " = 1, "Beta 0" = 2, "Beta 1" = 2, "Beta 2" = 2, "Beta 3" = 2, "Beta 4" = 2))

# Try to fit the beta model only with the four seasons
season_diff <- as.data.frame(season_diff) %>%
  mutate(month = recode(month, April = "Spring"),
         month = recode(month, May = "Spring"),
         month = recode(month, June = "Summer"),
         month = recode(month, July = "Summer"),
         month = recode(month, August = "Summer"),
         month = recode(month, September = "Autumn"),
         month = recode(month, October = "Autumn"),
         month = recode(month, November = "Autumn"),
         month = recode(month, December = "Winter"),
         month = recode(month, January = "Winter"),
         month = factor(month, levels = c("Spring", "Summer", "Autumn", "Winter")))
colnames(season_diff)[3] <- "season"

# Beta0
intercept.fit.2 <- lm(intercept ~ season, data = season_diff)
# Beta1
wind_prev.fit.2 <- lm(wind_prev ~ season, data = season_diff)
# Beta2
lat_change.fit.2 <- lm(lat_change ~ season, data = season_diff)
# Beta3
long_change.fit.2 <- lm(long_change ~ season, data = season_diff)
# Beta4
wind_change.fit.2 <- lm(wind_change ~ season, data = season_diff)

sum0_2 <- summary(intercept.fit.2)$coefficients[,c(1,4)]
sum1_2 <- summary(wind_prev.fit.2)$coefficients[,c(1,4)]
sum2_2 <- summary(lat_change.fit.2)$coefficients[,c(1,4)]
sum3_2 <- summary(long_change.fit.2)$coefficients[,c(1,4)]

```

```

sum4_2 <- summary(wind_change.fit.2)$coefficients[,c(1,4)]

kable(cbind(sum0_2, sum1_2, sum2_2, sum3_2, sum4_2)) %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed")) %>%
  add_header_above(c(" " = 1, "Beta 0" = 2, "Beta 1" = 2, "Beta 2" = 2, "Beta 3" = 2, "Beta 4" = 2))

# Try to fit the beta model only with the year
# Beta0
intercept.fit.new <- lm(intercept ~ year, data = season_diff)
# Beta1
wind_prev.fit.new <- lm(wind_prev ~ year, data = season_diff)
# Beta2
lat_change.fit.new <- lm(lat_change ~ year, data = season_diff)
# Beta3
long_change.fit.new <- lm(long_change ~ year, data = season_diff)
#Beta4
wind_change.fit.new <- lm(wind_change ~ year, data = season_diff)

summary(intercept.fit.new)
summary(wind_prev.fit.new)
summary(lat_change.fit.new)
summary(long_change.fit.new)
summary(wind_change.fit.new)

sum0.new <- summary(intercept.fit.new)$coefficients[,c(1,4)]
sum1.new <- summary(wind_prev.fit.new)$coefficients[,c(1,4)]
sum2.new <- summary(lat_change.fit.new)$coefficients[,c(1,4)]
sum3.new <- summary(long_change.fit.new)$coefficients[,c(1,4)]
sum4.new <- summary(wind_change.fit.new)$coefficients[,c(1,4)]

kable(cbind(sum0.new, sum1.new, sum2.new, sum3.new, sum4.new)) %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed")) %>%
  add_header_above(c(" " = 1, "Beta 0" = 2, "Beta 1" = 2, "Beta 2" = 2, "Beta 3" = 2, "Beta 4" = 2))

```

Damage and Death

```

data_deaths = data_res %>%
  select(-c(id, Damage))
offset_pop = data_deaths %>% pull(Total.Pop)
offset_dur = data_deaths %>% pull(Hours)
deaths.fit = glm(Deaths ~ . + offset(log(offset_pop)) + offset(log(offset_dur)), data = data_deaths %>%
deaths.tidy = summary(deaths.fit) %>% na.omit()
knitr::kable(deaths.tidy$coefficients, digits = 3)

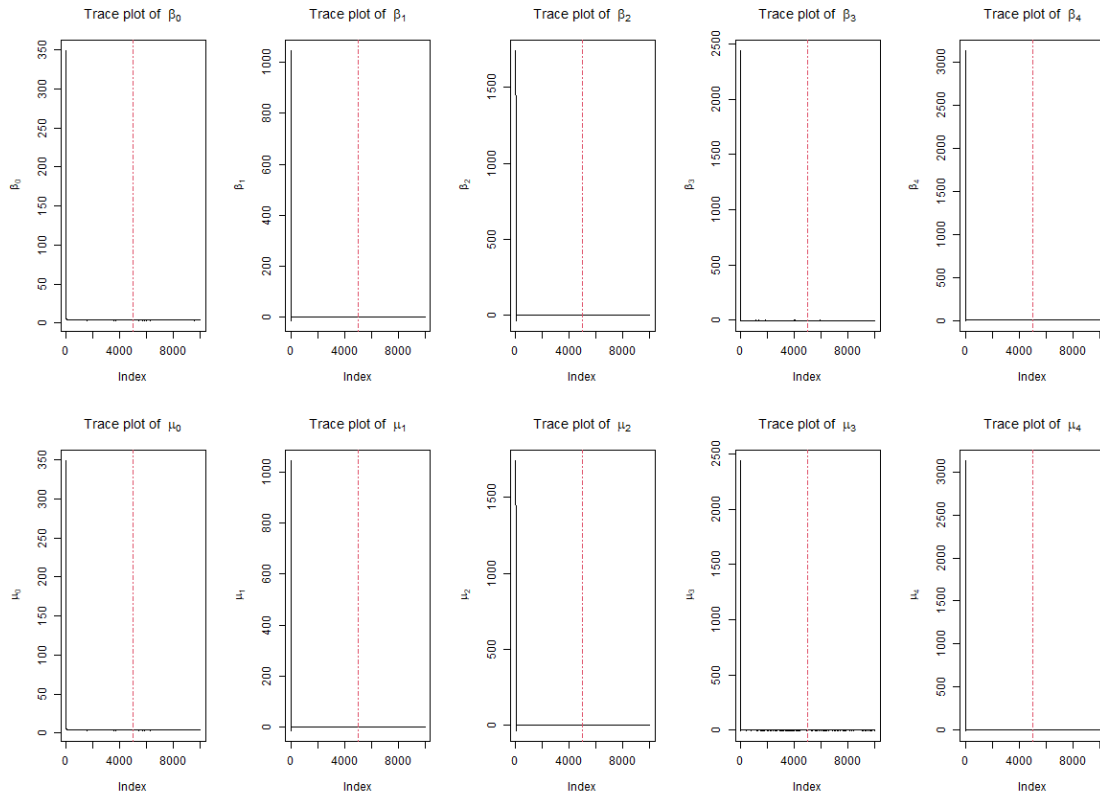
data_damage = data_res %>%
  select(-c(id, Deaths)) %>%
  mutate(Damage = 1000*Damage)

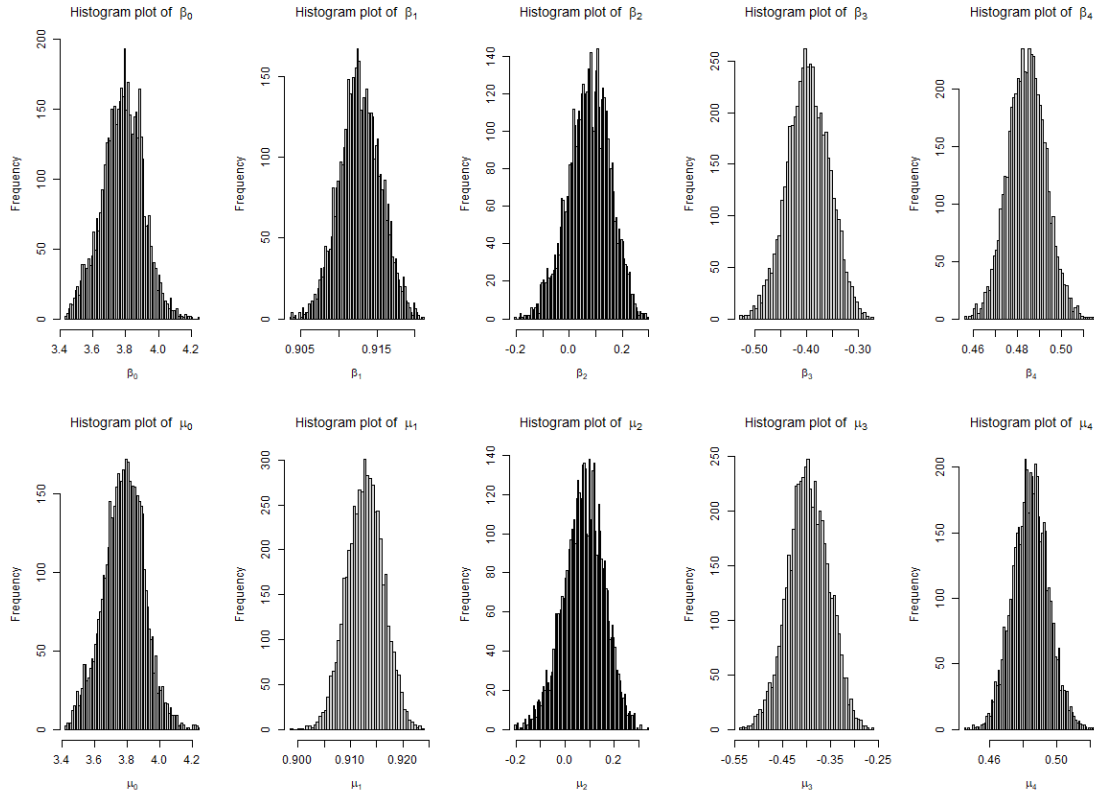
offset_dur_damage = data_deaths %>% pull(Hours)
require(broom)
damage.fit = glm(Damage ~. + offset(log(offset_dur_damage)), data = data_damage %>% select(-Hours), family = "poisson")
damage.tidy = tidy(damage.fit)

```

```
table6 = knitr::kable(damage.tidy, digits = 3)
table6
```

Bad initial value





Citations

1. Taboga, Marco (2021). "Markov Chain Monte Carlo (MCMC) diagnostics", Lectures on probability theory and mathematical statistics. Kindle Direct Publishing. Online appendix.
2. Elsner, J. B., Niu, X., & Jagger, T. H. (2004). Detecting Shifts in Hurricane Rates Using a Markov Chain Monte Carlo Approach, *Journal of Climate*, 17(13), 2652-2666. Retrieved May 9, 2022, from https://journals.ametsoc.org/view/journals/clim/17/13/1520-0442_2004_017_2652_dsihru_2.0.co_2.xml