

# Spatial non-parametric Bayesian clustered regression coefficients

This webpage is created as an online supplementary material for the manuscript **Spatial non-parametric Bayesian clustered regression coefficients**. We present our modeling code using the Nimble package, as well as code to perform posterior inference and clustering.

## Generate Simulated Data

To provide an example, we adopt a spatial layout akin to the arrangement of counties in Georgia, as demonstrated by (Ma et al., 2019), encompassing a total of 159 counties. In our simulated data set, we distribute three observations within each geographical region. This leads us to define the essential parameters, allowing us to subsequently simulate a complete data set. Where the true cluster distribution is give as:

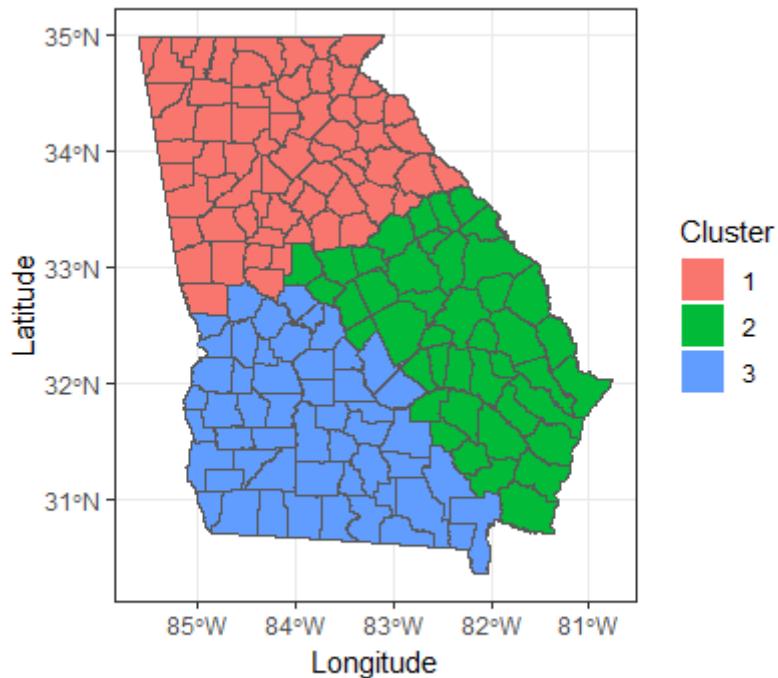


Figure 1: cluster assignment for Georgia counties used for simulation studies.

```
1 library(MASS)
2 library(nimble)
3 library(coda)
4 library(ClusterR)
5 library(mclust)
6 library(ggplot2)
7 library(tidyverse)
```

```

8  library(sf)
9  library(dplyr)
10 library(geosphere)
11 library(ggpubr)
12
13 distMat <- readRDS("./GAcentroidgcd.rds")
14 centroids <- as.data.frame(readRDS("GAcentroids.rds"))
15 dd <- distMat
16
17 # Generate true clustering settings based on centroids
18 # Set a random seed for reproducibility
19 set.seed(123)
20
21 asm <- c()
22 for (i in 1:nrow(centroids)) {
23   if (centroids$x[i] - 2 * centroids$y[i] < -150) {
24     asm[i] <- 1
25   } else if (centroids$x[i] + centroids$y[i] > -51) {
26     asm[i] <- 2
27   } else {
28     asm[i] <- 3
29   }
30 }
31
32 dd <- distMat
33
34 # Create a matrix to store beta values for each cluster
35 betaMat <- t(matrix(nrow = 159, ncol = 6, byrow = TRUE))
36 for (i in 1:159) {
37   ## cluster 1
38   betaMat[,asm == 1] <- c(9, 0, -4, 0, 2, 5)
39   ## cluster 2
40   betaMat[,asm == 2] <- c(1, 7, 3, 6, 0, -1)
41   ## cluster 3
42   betaMat[,asm == 3] <- c(2, 0, 6, 1, 7, 0)
43 }
44
45 # Calculate six different weighted functions using Gaussian kernel
46 # (V1 to V6) based on dd and certain constants
47 V1 <- exp(-dd / (1 * 5))
48 V2 <- exp(-dd / (2 * 0.1))
49 V3 <- exp(-dd / (3 * 0.2))
50 V4 <- exp(-dd / 2 * 4)
51 V5 <- exp(-dd / 30 * 5)
52 V6 <- exp(-dd / 40 * 6)
53
54 # Calculate c1 to c6 based on V1 to V6 and betaMat
55 c1 <- V1[1, ] * (betaMat[1, ])
56 c2 <- V2[1, ] * betaMat[2, ]
57 c3 <- V3[1, ] * betaMat[3, ]
58 c4 <- V4[1, ] * betaMat[4, ]
59 c5 <- V5[1, ] * betaMat[5, ]
60 c6 <- V6[1, ] * betaMat[6, ]

```

```

61
62
63 # Combine c1 to c6 into a matrix Beta.true
64 Beta.true <- cbind(c1, c2, c3, c4, c5, c6)
65
66 # Set the number of samples
67 n <- 159
68
69
70 # Generation of sampling locations Sp
71 Sig.true<- 1
72
73 # Covariates with a given range parameter phi
74 phi <- 0.9
75 dd <- distMat
76 mat <- exp(-dd / phi)
77 x1 <- mvrnorm(1, rep(0, n), mat)
78 x2 <- mvrnorm(1, rep(0, n), mat)
79 x3 <- mvrnorm(1, rep(0, n), mat)
80 x4 <- mvrnorm(1, rep(0, n), mat)
81 x5 <- mvrnorm(1, rep(0, n), mat)
82 x6 <- mvrnorm(1, rep(0, n), mat)
83 X <- data.frame(x1, x2, x3, x4, x5, x6)
84 Mu <- apply(cbind(X) * Beta.true, 1, sum)
85
86 #Create the data
87 Y <- rnorm(n, Mu, Sig.true)
88
89 #Converts the 'x' and 'y' values from the 'centroids' to radians
90 centroids <- centroids
91 cons = centroids$x * pi/180
92 const= centroids$y * pi/180
93
94 new_cen<- data.frame(cons,const)

```

In this chunk, the distance matrix and centroid coordinates obtained from external files. We sets a deterministic seed for reproducibility and assigns data points to one of three clusters based on the conditions defined by the centroid coordinates. The data is then processed to generate beta values specific to each cluster. Using the distance matrix, it creates six different kernel weighted functions. These functions, combined with the beta values, produce six coefficients for each data point. By generating multiple covariates based on a given range parameter, the code ultimately constructs a dataset  $Y$  from GWR model where each data point is drawn from a normal distribution with a mean determined by these coefficients and covariates. The code then transforms the centroid coordinates to radians, preparing them for potential the analysis.

```

1 dnorm_vec2 <- nimbleFunction(
2   run = function(x = double(1), mean = double(1), sd = double(1),
3                 log = integer(0, default = 0)) { returnType(double(0))
4     logProb <- sum(dnorm(x, mean, sd, log = TRUE))
5     if (log) return(logProb)
6     else return(exp(logProb))
7   })
8 registerDistributions('dnorm_vec2')
9
10 # Define the Code for the nimble model

```

```

11 GWRCode <- nimbleCode({
12   for (i in 1:S) {
13     y[1:N, i] ~ dnorm_vec2(b[i, 1] * x1[1:N] + b[i, 2] * x2[1:N] + b[i, 3] *
14                               x3[1:N] + b[i, 4] * x4[1:N] + b[i, 5] * x5[1:N] + b[i, 6] * x6[1:N],
15                               1 / (psi_y[i] * exp(-Dist[1:N, i] / lambda)))
16 
17   psi_y[i] ~ dgamma(100,100)
18 
19 
20   b[i, 1:6] <- bm[latent[i], 1:6]
21   bm[i, 1:6] ~ dmmnorm(mu[ 1:6,latent[i]], Tau[1:6, 1:6,latent[i]])
22   latent[i] ~ dcat(pi[i,1:M])
23 #
24   for (m in 1:M) {
25     mu[ 1:6,m] ~ dmmnorm(mu0[1:6], Tau[1:6,1:6,m])
26     Tau[1:6,1:6,m] ~ dwish(D1[1:6,1:6 ], c)
27     Sigma[ 1:6,1:6,m] <- inverse(Tau[1:6,1:6,m])
28 
29   for (j in 1:6) {
30     mu0[j] ~ dnorm(0, 1)
31   }
32 }
33 
34 
35   for (i in 1:S) {
36     pi[i,1] <- vs[i, 1]
37     for (j in 2 : M){ pi[i, j ] <- vs[i, j ]*prod(vsout[i, 1 : ( j - 1)])
38   }
39 
40 
41   for(j in 1 : (M - 1)){
42     V[j] ~ dbeta(1, 1)
43 
44     for (i in 1:S) {
45       #
46       weight[i,j]<- exp(-(new_cen[i, 1] - knot[j, 1])^2)/r[j,1]^2) * exp(-(new_cen[i, 2]- knot[j, 2])
47 
48       vs[i,j]<- weight[i,j]*V[j]
49       vsout[i,j]<- 1-vs[i,j]
50 
51     }
52     knot[j, 1]~dunif(0,10)
53     knot[j, 2]~dunif(0,10)
54     r[j, 1]<- h^2/2
55     r[j, 2]<- h^2/2
56   }
57 
58   for (k in 1:S){vs[k,M]<- 1}
59 
60   lambda~dunif(0,D)
61   h~dunif(0,5)
62 })

```

## Data List for Model

The next step involves defining the data list for the aforementioned model code. This data list includes the response variable, the covariates, and the distance matrix . It is important to note that the entries in this matrix have been normalized to ensure a maximum value of 10.

```
1 p<- 6
2
3
4 ##Replicate Y to get a square matrix
5 Y <- matrix(Y, nrow = length(Y), ncol = length(Y), byrow = FALSE)
6 # Prepare the data, constants, and initial values for the nimble model
7 GWRdata <- list(y = Y, x1 = X[, 1], x2 = X[, 2], x3 = X[, 3], x4 = X[, 4],
8                  x5 = X[, 5], x6 = X[, 6], Dist = distMat,D1=diag(rep(1, p)),c = p + 1,
9                  new_cen=new_cen)
10
11 GWRConsts <- list(S = 159, M = 10, N = 159, D = 50)
12
13 dim_tau <- c(6, 6, GWRConsts$M)
14
15 # Create the array with all elements set to 1
16 Tau <- array(1, dim = dim_tau)
17
18 # Set the diagonal elements to 1
19 for (i in 1:dim_tau[3]) {
20   Tau[, , i] <- diag(dim_tau[1]) * 1
21 }
22
23 GWRInits <- list( psi_y = rep(100, GWRConsts$S), lambda = 10,
24                     mu0 = rep(0,6),latent = rep(1, GWRConsts $S),
25                     Tau=Tau,knot = matrix(runif((GWRConsts$M-1)*2,0, 1), GWRConsts$M-1, 2),
26                     ,h=1,V = rbeta(GWRConsts$M - 1, 1, 1))
```

## The proposed model

The dnorm\_vec2, using the Nimble framework. The function calculates the sum of log probabilities for a given vector using a normal distribution. After defining the function, it's registered for use in subsequent Nimble models.

The core of the code establishes a Bayesian Geographically Weighted Regression (GWR) model, termed GWRCODE. This model attempts to capture spatial heterogeneity by allowing parameters to vary over space. It uses several probability distributions, including normal (dnorm), gamma (dgamma), multinomial (dmnorm), categorical (dcat), and beta (dbeta). A key feature of the GWR model is that it allows for local relationships between predictors and response variables, using weights based on spatial proximity.

The model incorporates latent variables to capture unobserved heterogeneities and knots for spatial locations that influence local parameter estimates. It calculates the weight for each observation based on its distance to these knots. Additionally, certain parameters, like mu0, are given weakly informative prior distributions.

Lastly, the code processes data (distance matrix, covariates, and response) and prepares it alongside certain constants to be used in the GWR model. The preparation includes the creation of a 3-dimensional array Tau for the covariance structure, where the diagonal elements are set to 1.

In essence, the code captures spatially varying relationships in data using a Bayesian GWR model and sets the stage for model fitting and inference using the Nimble framework.

```

1 # Perform MCMC sampling using nimble
2 mcmc.out_Spatial <- nimbleMCMC(code = GWRCode, data = GWRdata, constants = GWRConsts,
3                                 inits = GWRInits, monitors = c("latent","pi"),
4                                 niter = 5000, nburnin = 2000,
5                                 nchains = 1, setSeed = TRUE)

## |-----|-----|-----|-----|
## |-----|-----|
```

## Run the Model

In Nimble, we can directly run the MCMC engine using the following code. This function is designed to run the Markov Chain Monte Carlo (MCMC) process and typically requires the model code, data, and initial values as input. It offers various options for controlling aspects such as multiple chains, the number of iterations, thinning intervals, and more.

The following code example illustrates the execution of a single MCMC chain with a total of 5000 iterations, where the initial 2000 iterations are used as burn-in. Consequently, the output will provide 3000 posterior samples for the parameters: “pi”, and “latent.”

## Posterior Convergence Diagnostics and Estimation

The coda package (Plummer et al., 2006) offers convenient tools for performing posterior convergence diagnostics. It also provides useful functions for computing various percentiles of the posterior distribution, which are often of great interest in posterior inference. In the next chunk the code perform clustering analysis on the posterior samples from the Bayesian model and visualize the resulting cluster assignments on a map of Georgia. The Dahl’s and mode methods are used for clustering, and the final clusters are plotted on the map using ggplot.

```

1 mcmc.out <- as.mcmc(mcmc.out_Spatial)
2 latentZMat <- mcmc.out
3
4 # Assuming latentZMat is a data frame or matrix
5 # and you want to extract columns with names containing "latent"
6 latentZMat <- latentZMat[, grep("latent", colnames(latentZMat))]
7
8 membershipList <- purrr::map(1:nrow(latentZMat), .f = function(x) {
9   outer(latentZMat[x, ], latentZMat[x, ], "==")
10 })
11
12 ## the empirical probability matrix
13 bBar <- Reduce("+", membershipList) / length(membershipList)
14
15 ## sum of squared differences
16 lsDist <- purrr::map_dbl(membershipList, ~sum((.x - bBar) ^ 2))
17
18 ## find the optimal iteration, and take as the final inferenced result
19 ## if there are multiple optimal iterations, take the first one
20 mcluster <- which.min(lsDist)
21 finalCluster <- as.numeric(latentZMat[mcluster[1],])
22 finalCluster <- ifelse(latentZMat[mcluster[1],] == 9, 1, ifelse(latentZMat[mcluster[1],] == 10, 2, ifel
23
24 table(finalCluster)

## finalCluster
```

```

## 1 2 3
## 65 42 52
1 # Get the total number of iterations in the chain
2 fossil::rand.index(as.numeric(finalCluster), asm)

## [1] 0.8785925

1 latentPE <- as.numeric(unlist(apply(latentZMat, 2, FUN = function(x) {
2   return(DescTools::Mode(x)[1])
3 })))
4 # Check number of clusters, and number of regions in each cluster
5 latentPE<- ifelse(latentPE == 9, 1, ifelse(latentPE == 10, 2, ifelse(latentPE == 2, 3, latentPE)))
6 fossil::rand.index(asm, latentPE)

## [1] 0.8752488
1 table(latentPE)

## latentPE
## 1 2 3
## 66 39 54

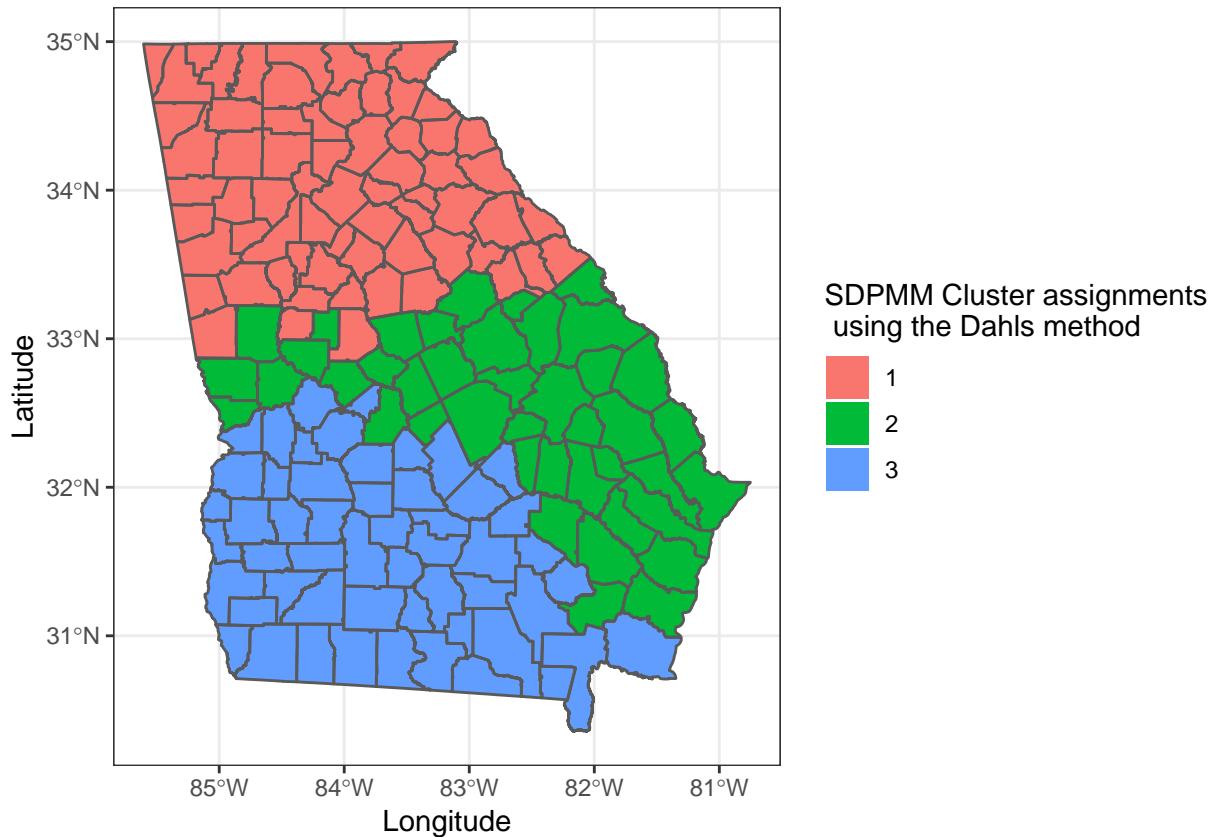
1 library(plyr)
2 library(dplyr)
3 library(sf)
4 Georgia <- read_sf("Georgia_dat.shp") %>% filter(!st_is_empty(.))

# Create a new dataset with finalCluster assignments
mydata_and_myMap <- cbind(Georgia, finalCluster,latentPE,asm)

# Create a ggplot map visualization
A<- ggplot() +
  xlab("Longitude") +
  ylab("Latitude") +
  theme_bw() +
  theme(legend.position = "right") +
  labs(fill = 'SDPMM Cluster assignments \n using the Dahls method ') +
  geom_sf(data = mydata_and_myMap, aes(fill = factor(mydata_and_myMap$finalCluster)), color = NA) +
  geom_sf(data = mydata_and_myMap, fill = NA)

A

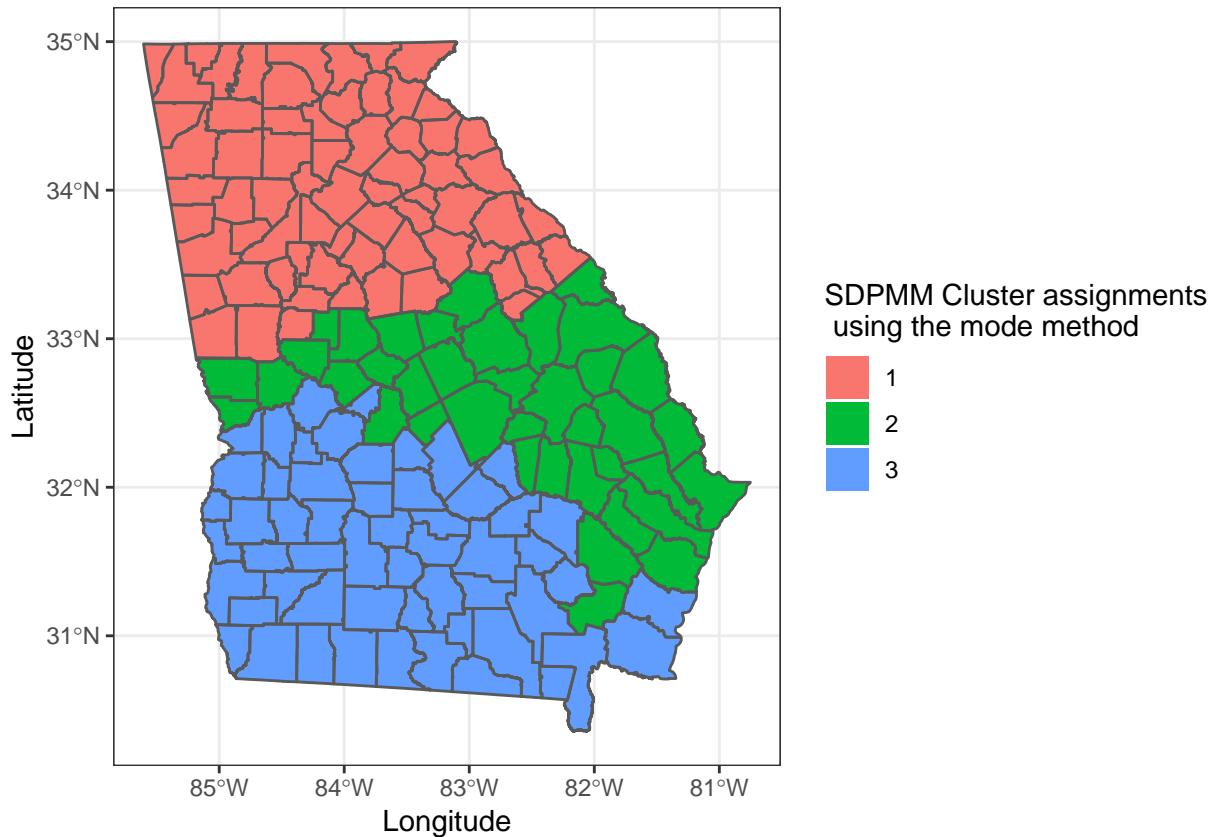
```



```

1 B<- ggplot() +
2   xlab("Longitude") +
3   ylab("Latitude") +
4   theme_bw() +
5   theme(legend.position = "right") +
6   labs(fill = 'SDPMM Cluster assignments \n using the mode method ') +
7   geom_sf(data = mydata_and_myMap, aes(fill = factor(mydata_and_myMap$latentPE)), color = NA) +
8   geom_sf(data = mydata_and_myMap, fill = NA)
9
10
11 B

```



## References

- Plummer, M., Best, N., Cowles, K., & Vines, K. (2006). CODA: Convergence Diagnosis and Output Analysis for MCMC. *R News*, 6(1), 7-11. (<https://journal.r-project.org/archive/>)
- Sugasawa, S., & Murakami, D. (2022). Adaptively Robust Geographically Weighted Regression. *Spatial Statistics*, 48, 100623. (<https://www.sciencedirect.com/science/article/pii/S2211675322000185>)
- Valpine, P. D., Turek, D., Paciorek, C. J., Anderson-Bergman, C., Temple Lang, D., & Bodik, R. (2017). Programming with Models: Writing Statistical Algorithms for General Model Structures with NIMBLE. *Journal of Computational and Graphical Statistics*, 26(2), 403-413.