

Bayesian Computations for Probit Models and Applications via Chinese Happyness Data

Travis Gubbe

December 17, 2025

Abstract

This report explores the applications of linear and nonparametric models with the addition of Bayesian estimation via Markov Chain Monte Carlo simulations to determine a binary response variable's status. Two spatial probit models - linear and varying coefficient - are constructed to classify a county in China as whether or not it is a high-ranking county based on socioeconomic and population factors. These models contained four explanatory variables: GDP, population per 10,000 residents, and two economic and financial factors. The models created demonstrate the role of spatial dependence in a variable's parameter estimation and the flexibility of regression splines to account for changes in time.

Introduction

The goal of the happyness data set presented is to analyze the various economic conditions and demographic factors leading to a county being ranked in the top 100 counties in China. The data set contains a variety of societal and economic explanatory variables, including GDP, population, savings balances of residents, and numerous fiscal and financial health ratios. The response variable is a binary variable, where if a county is in the top 100 counties, it is coded as a 1 - otherwise it is coded as a 0. The response variable is an unobservable latent variable, which is where the linear probit spatial model is key. A linear probit spatial model and varying coefficient spatial probit model are excellent options for an estimated model, which will be explained in the next section.

This reports aims to demonstrate the linear probit spatial model and varying coefficient spatial probit model to acquire Bayesian estimations and provide Bayesian quantities for the respective models. Bayesian estimations allow for flexibility in the modeling, particularly when accounting for spatial dependence. These models will demonstrate the power of spatial probit models and their relation to Bayesian statistics.

Probit Models

Linear Probit Spatial Model

Linear probit spatial models are useful when the response variable is binary and there is the presence of spatial dependence among the explanatory variables. Spatial dependence is described as a variable

in one location being influenced by a variable in another location, such as the GDP and population of a county being linked together. The model is constructed as follows:

$$y_{it}^* = \rho \sum_{j=1}^n \omega_{ij} y_{jt}^* + \beta_0 + \beta_1 X_{it} + \varepsilon_{it}, \quad \varepsilon_{it} \sim N(0, \sigma_\varepsilon^2)$$

$$y_{it} = \begin{cases} 0, & y_{it}^* < 0, \\ 1, & y_{it}^* \geq 0. \end{cases}$$

ρ : spatial correlation parameter

ω : $N \times N$ spatial weight matrix

β : covariates used to approximate $m(x_{it})$

The spatial weights help define the strength of the units among one another, creating interdependence among the variables. The Markov Chain Monte Carlo Bayesian sampling will help smooth the latent variables and the complex likelihood function, making this a great potential model selection for the happiness data set.

Varying Coefficient Spatial Probit Model

Varying coefficient spatial probit model is similar to the linear spatial probit model, but captures change over time. Since the data set is from 2010-2020, this form of the model can be very useful in its design. The model is constructed as follows:

$$y_{it}^* = \rho \sum_{j=1}^n \omega_{ij} y_{jt}^* + \beta_0(t) + \beta_1(t) X_{it} + \varepsilon_{it}, \quad \varepsilon_{it} \sim N(0, \sigma_\varepsilon^2).$$

where:

ρ : spatial correlation parameter

ω : $N \times N$ spatial weight matrix

β : spline basis functions

This allows the covariates to vary since the relationship among the predictors and outcomes may not be constant over time. Running Bayesian estimates on this model will smooth the covariates and spatial structure, making it a great option for the happiness data set. Unlike the linear version, the β values are spline functions to estimate the time-varying covariates, smoothing the model after many Markov Chain Monte Carlo simulations.

$$m(x_{it}) = \pi^{(T)} * (x_{it}) * \beta = \sum_{j=0}^p \phi_j x_{j,it} + \sum_{j=1}^m \phi_{j+p} (x_{it} - \xi_j)_+^p$$

This means running B-splines in conjunction with the covariates, with the assumption that the baseline of the propensity to be in the top 100 counties changes over time, but the effect of the covariates remains constant over time.

Case Study (Chinese Happyness Data)

Data Preparation

The Chinese happyness data set works well to demonstrate the two spatial probit models, since the response variable (top100) is a binary outcome with multiple variables that are spatially dependent with various socioeconomic variables. Before creating an appropriate model, the explanatory variables were checked for potential multicollinearity. Strong multicollinearity hints at dependency among the explanatory variables, making it difficult to discern the true effect of a variable on the result. Multicollinearity can be investigated by a correlation plot and a Variance Inflation Factor (VIF) plot. High correlations (> 0.8) hint at multicollinearity, while VIF values above 10 are strong indicators of multicollinearity. After running and re-running the variables through these two models, the chosen explanatory variables for modeling are - GDP, population per 10,000 residents, cy, and jr. These predictors are used in both the linear and varying coefficient spatial probit models.

Once the appropriate explanatory variables are selected, the variables are scaled to standardize the values and ensure the variables contribute equally to the model. This is particularly useful since there are high-value variables (GDP, population) and low-value variables (cy, jr), which can cause the low-value variables to be undervalued in the final model interpretation if they are not scaled. In addition, scaling the data can help the model run faster and improve model performance for distance-related models, as is the case with the spatial probit models.

For illustrative purposes of the comparison of the linear probit spatial model and the varying coefficient spatial probit model will be run with the same explanatory variables. This is to demonstrate both model's effectiveness when using the same parameters and their performance to one another for this data set.

Results

First, a linear spatial probit model is created to predict if a county is in the top 100 counties with the four explanatory variables listed (GDP, population, cy, jr):

Variable	Estimate	Std. Dev	95% CI
Intercept	-0.261	0.067	[-0.393, -0.128]
GDP	0.319	0.088	[0.148, 0.491]
Population (per 10K)	0.432	0.087	[0.262, 0.601]
cy	0.357	0.085	[0.189, 0.524]
jr	0.198	0.076	[0.049, 0.346]
ρ	0.619	0.0402	[0.540, 0.698]

Table 1: Summary table for the Linear Spatial Probit Model. All variables are statistically significant since each respective confidence interval does not contain 0.

From the results, GDP, population, cy, and jr all have positive estimates and their credible intervals do not contain zero, which means these factors have a positive impact on whether a county will be ranked in the top 100 and are statistically significant to the model. The estimates also provide the statistical impact of each individual variable - for example, GDP's estimate of 0.319 means, assuming all other variables are fixed, increasing the GDP by one unit increases the probability of the response variable (top 100) by 0.319. The ρ value's estimate value of 0.619 and 95% credible

interval of [0.540, 0.698] indicates there is a positive strong spatial dependence in the model, meaning county's locations to one another can impact their determination of whether they are a top 100 county. This means that a county ranked in the top 100 is more likely to be surrounded by other high-ranking counties. From the model results, these highly-ranked counties may be clustered together geographically, with their ranking being explained by GDP, population, and cy and jr factors.

Since the variables present can change over a period of time, the varying coefficient spatial probit model is utilized to account for the change in time with the use of regression splines. It was assumed the baseline of the propensity to be ranked in the top 100 counties in China would vary over time, while the covariate factors would remain relatively constant for the 10-year time span. This adds only the regression spline values to the model and shortens the risk of overfitting, while keeping the flexibility of the model and makes for easier interpretability. Multiple models were run with varying spline amounts and degrees between 1 and 4. In the end, it was determined to use 1 regression spline in the model with 1 degree to retain statistical significance of the four chosen variables at the 95% significance level. In addition, model performance testing showed very little change in the model performance metrics to confidently change the amount of regression splines and degrees of the model.

From the results table for the varying coefficient model, the covariate estimates differ slightly to the linear model, while the ρ value is equal to the linear model, with both values equal to 0.619. All covariates are statistically significant at the 95% significance level, which can also be confirmed by the credible intervals for each covariate as none contain zero (**Appendix B**). Now that two appropriate spatial probit models have been created, there are a couple of model performance diagnostics that can be run to determine which model performed best for categorizing a top 100 county in China.

Trace plots are also viewed to determine as a diagnostic tool for the spatial probit models. A good trace plot has a fuzzy look, with no clear pattern in the plot. This would indicate good convergence of the MCMC simulations for both the regression and spatial parameters. The trace plots for each spatial model showed random traces and no clear pattern, which indicates good convergence of the MCMC simulations for both models (**Appendix C**).

The density plots for the posterior distributions of the covariates of both models can also be viewed as a diagnostic approach. The goal is to have unimodal density plots and a reasonable spread for each plot (preferably non-zero), indicating each covariate has an association with the response variable and is easy to interpret. The density plots for both the linear spatial probit model and varying coefficient were mostly unimodal in appearance with a good spread in the density as well as what appear to be non-zero covariates, indicating the covariates selected are statistically significant and have an identifiable association with the response variable (**Appendix C**).

Model Performance

Now that the two models have been created, they must be assessed to determine their performance to one another. The two methods in this paper for model comparison is Deviance Information Criterion (DIC) and a combination of Receiver Operating Characteristic (ROC) with Area Under the Curve (AUC). DIC is a generalization of AIC and is used to compare Bayesian models. The goal is to have a lower DIC value, as this is an indication of model fit and complexity. In contrast, it is best to have a model with higher ROC and AUC scores, as this indicates the model can explain a higher amount of variability and indicates better model performance.

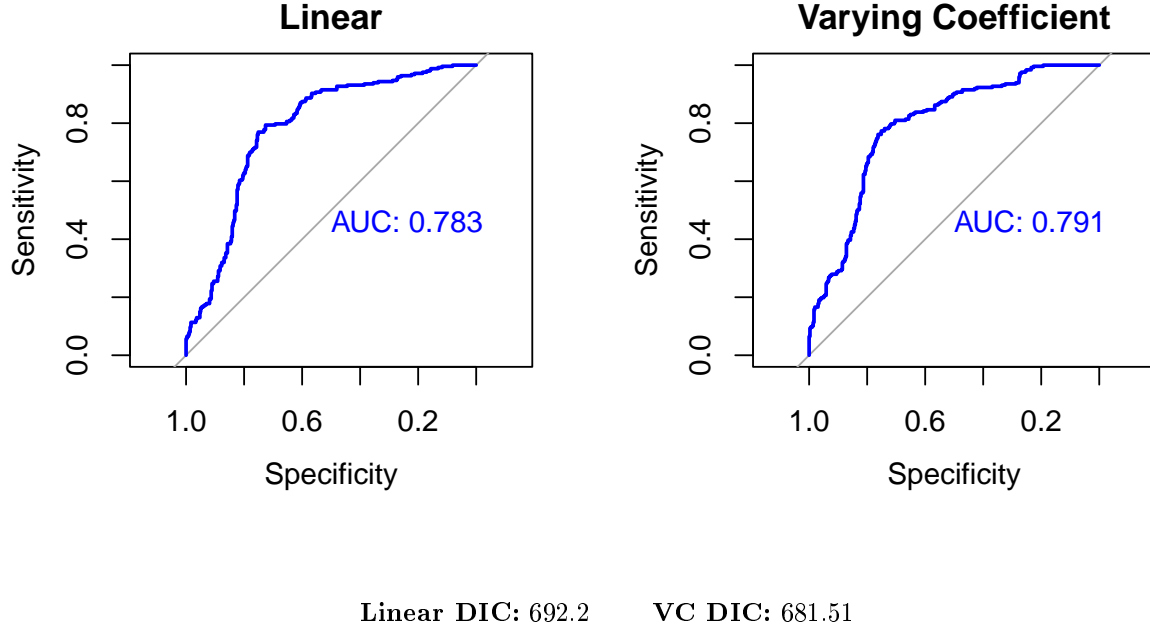


Figure 1: ROC visualizations for the two models with their respective AUC values. The Varying Coefficient model performed slightly better, with an AUC value of 0.791 compared to the Linear Spatial model with an AUC value of 0.783. The DIC values show the Varying Coefficient model performed better than the Linear Spatial model, with a lower DIC value of 680.78 compared to 691.08

From the ROC/AUC results, the Varying Coefficient model had a slightly higher AUC value of 0.791 compared to the Linear spatial model's AUC value of 0.783. An AUC value of close to 0.8 indicates very good model performance, which means both models have performed very well in the prediction of whether a county in China is in the top 100. In addition, with a DIC value of 680.78 compared to a DIC value of 691.08, the Varying Coefficient model performed better than the Linear spatial model in regards to DIC, as the lower DIC value (especially since the DIC difference is greater than 10) indicates better model fit and performance. It is important to note the standard threshold value of 0.5 was chosen for the ROC. Optimization of the ROC threshold can lead to a better AUC as well as a better performance model as it pertains to Type I and Type II error. Doing so can also change the outlook of which model type fits better for the problem at hand, as optimizing the ROC threshold impacts both models' performance. For the purpose of this paper, ROC optimization was not explored since the objective is to provide working probit spatial models.

Conclusion

The results of the linear probit spatial model and varying coefficient spatial probit in conjunction with Bayesian estimation highlight its role in determining spatial dependencies for regional outcomes. The estimated spatial correlation parameter (ρ) poses the observed outcomes in county ranking positively influences its neighboring regions, implying strong spatial interaction. This is an important distinction made by the spatial probit models, as it indicates a county's socio-economic factors can influence neighboring regions. From the two models created, the varying coefficient spatial probit model displayed stronger model performance than its counterpart, as the β -splines

obtain non-parametric relationships between the covariates and the outcome, which is particularly essential for relationships that are not linear.

Both models had a ρ value equal to 0.619, which indicates a strong positive correlation between neighboring counties in its determination of being ranked in the top 100 counties, demonstrating strong spatial dependency. In other words, counties in the top 100 tend to be regionally close together rather than far apart. Both models also have positive estimate parameters for all selected covariates, meaning a 1 unit increase in each covariate has a positive impact on the outcome. In addition, all covariates' credible intervals do not contain zero, which indicates the covariates are statistically significant for both models. Where the varying coefficient model separated itself from the linear model is its AUC and DIC values. The varying coefficient model had a higher AUC value (0.791 vs. 0.783) and a lower DIC value (680.78 vs. 691.08) which indicates the varying coefficient model captured more variability and better balance when fitting the data.

In the end, both models demonstrate the ability to provide practical insights into spatial dependency for binary regional outcomes. The linear spatial probit model is easier to construct and interpret, while the non-parametric spatial probit model incorporates regression splines to estimate time-varying coefficients. For an alternative model, interaction effects can potentially create a more accurate model if it is believed the impact of one covariate affects the value of another covariate. Overall, probit models are a powerful tool in spatial modeling, and can provide applicable insights in regional outcomes.

Appendix

Appendix A

Variable	VIF
GDP	1.48
Population (per 10K)	1.10
cy	1.29
jr	1.16

Table 2

Appendix B

Variable	Estimate	Std. Dev	95% CI
Intercept	0.149	0.168	[-0.180, 0.478]
GDP	0.472	0.124	[0.230, 0.714]
Population (per 10K)	0.461	0.091	[0.283, 0.638]
cy	0.201	0.099	[0.006, 0.396]
jr	0.280	0.090	[0.103, 0.457]
B1	-0.794	0.325	[-1.430, -0.158]
ρ	0.619	0.041	[0.540, 0.698]

Table 3: Summary table for the Varying Coefficient Spatial Probit Model. All covariates are statistically significant at the 95% confidence level since they do not contain 0.

Appendix C

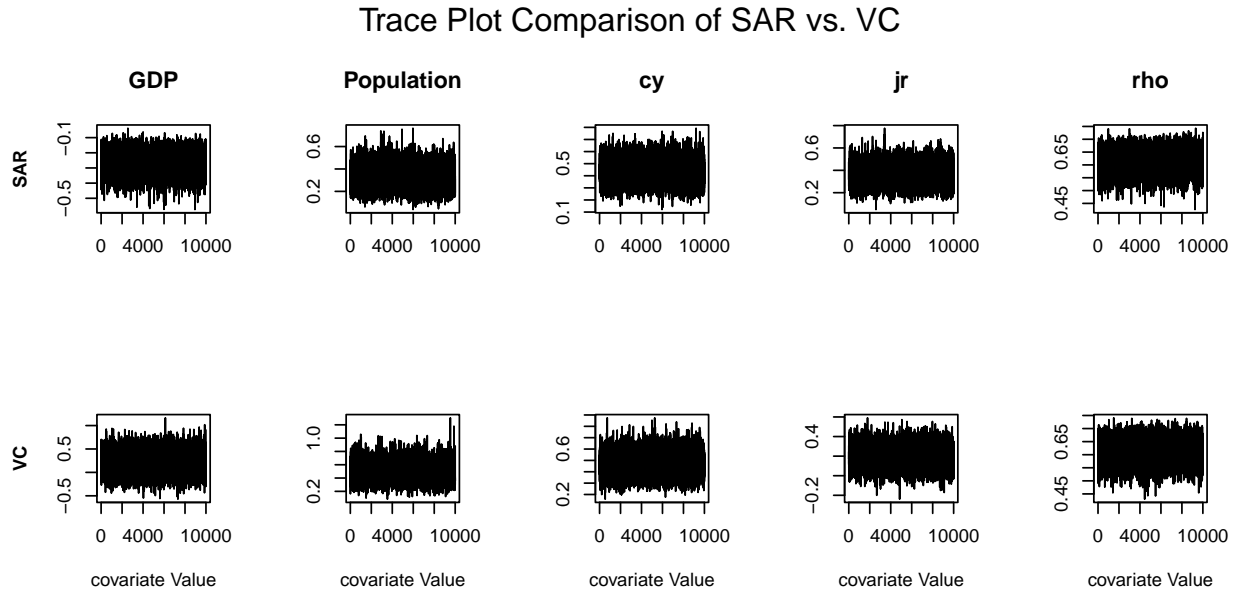


Figure 2: Trace plots of the covariates for the linear spatial probit model and varying coefficient probit spatial model, illustrating good convergence of MCMC simulation for both models.

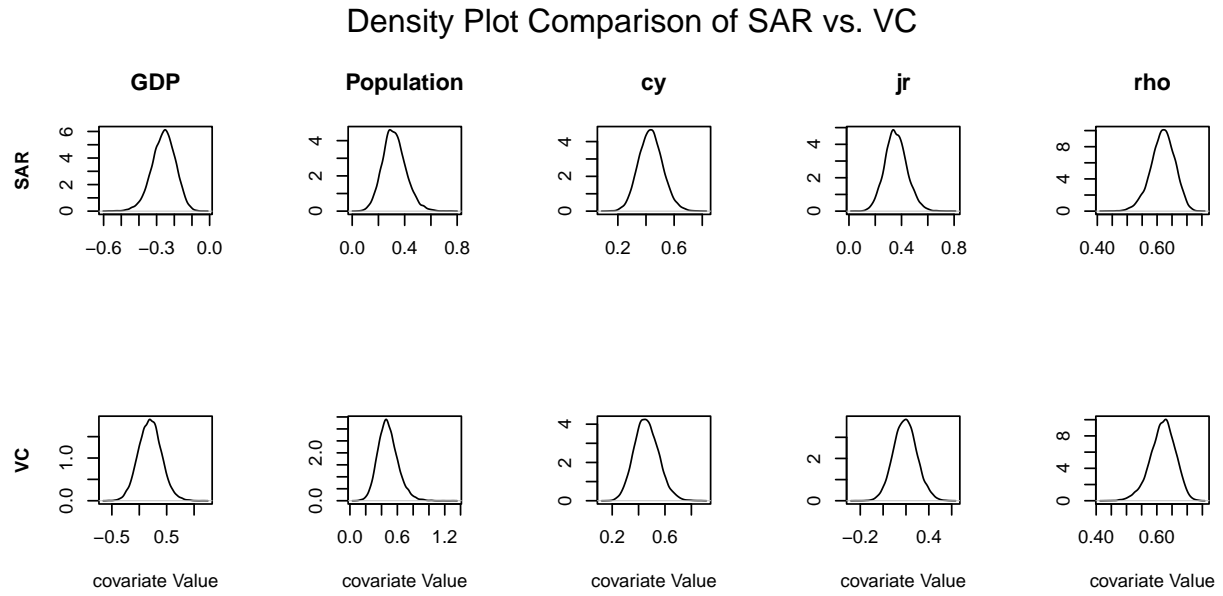


Figure 3: Density plots of the covariates for the linear spatial probit model and varying coefficient probit spatial model, illustrating good convergence of MCMC simulation for both models.

R Code

```
1 library(spatialprobit)
2 library(Matrix)
3 library(dplyr)
4 library(splines)
5
6 happy <- read.csv("happy_english.csv", stringsAsFactors = FALSE)
7
8 # Y variable is numeric
9 happy$top100 <- as.numeric(happy$top100)
10
11 # Create IDs for counties and years, then sort (year, county)
12 counties <- sort(unique(happy$county_name))
13 years <- sort(unique(happy$year))
14
15 happy <- happy %>%
16   mutate(
17     county_id = match(county_name, counties),
18     year_id = match(year, years)
19   ) %>%
20   arrange(year_id, county_id)
21
22 N <- length(counties) # number of counties
23 TT <- length(years) # number of years
24 stopifnot(N * TT == nrow(happy))
25
26 # Cross-sectional W_N (N x N) sparse
27 W_N <- Matrix(0, nrow = N, ncol = N, sparse = TRUE)
28
29 for (i in 1:N) {
30   city_i <- happy$city[happy$county_id == i][1]
31   for (j in 1:N) {
32     if (i != j) {
33       city_j <- happy$city[happy$county_id == j][1]
34       if (!is.na(city_i) && !is.na(city_j) && city_i == city_j) {
35         W_N[i, j] <- 1
36       }
37     }
38   }
39 }
40
41 # Zero diagonal
42 diag(W_N) <- 0
43
44 # Row-standardize: rows sum to 1 when neighbors exist
45 rs <- rowSums(W_N)
46 rs[rs == 0] <- 1
47 W_N <- Diagonal(x = 1 / rs) %*% W_N
48
49 # Panel W = I_T x W_N (block-diagonal, sparse)
50 I_T <- Diagonal(TT)
51 W_panel <- kronecker(I_T, W_N)
52
53 selected_vars <- c("GDP", "registered_pop_10k", "cy", "jr")
54
```

```

55 Xvars <- selected_vars
56 # 1) Extract predictors as numeric matrix
57 X_no_int <- as.matrix(happy[, Xvars])
58
59 # Make sure everything is numeric
60 X_no_int <- apply(X_no_int, 2, as.numeric)
61
62 # 2) SCALE the predictors: center = TRUE, scale = TRUE
63 X_scaled <- scale(X_no_int, center = TRUE, scale = TRUE)
64
65 # 3) Add intercept AFTER scaling
66 X <- cbind(Intercept = 1, X_scaled)
67 colnames(X)[1] <- "Intercept"
68
69 # Response
70 y <- happy$top100
71
72 set.seed(123)
73
74 # Linear Spatial Probit Model
75 #####
76 fit_sar <- sar_probit_mcmc(
77   y      = y,
78   X      = X,
79   W      = W_panel,
80   ndraw  = 10000,
81   burn.in = 2000,
82   thinning = 5,
83   prior  = NULL,
84   showProgress = TRUE
85 )
86
87 # 'fit_sar' is an object of class 'sarprobit'
88 summary(fit_sar)
89 # diagnostic plots: chains, densities, ACF, etc.
90 plot(fit_sar)
91
92 mean(fit_sar$rho) #p rho = 0.619
93
94 rho_mean <- 6.187e-01
95 rho_sd   <- 4.024e-02
96
97 rho_ci_norm <- rho_mean + qnorm(c(0.025, 0.975)) * rho_sd
98 rho_ci_norm #[0.540, 0.698]
99
100 intercept_ci_norm <- -0.26051 + qnorm(c(0.025, 0.975)) * 0.06785
101 gdp_ci_norm <- 0.31924 + qnorm(c(0.025, 0.975)) * 0.08762
102 pop_ci_norm <- 0.43169 + qnorm(c(0.025, 0.975)) * 0.08662
103 cy_ci_norm <- 0.35673 + qnorm(c(0.025, 0.975)) * 0.08540
104 jr_ci_norm <- 0.19793 + qnorm(c(0.025, 0.975)) * 0.07576
105
106 intercept_ci_norm #(-0.393, -0.128)
107 gdp_ci_norm #(0.148, 0.491)
108 pop_ci_norm #(0.262, 0.601)
109 cy_ci_norm #(0.189, 0.524)
110 jr_ci_norm #(0.049, 0.346)

```

```

111
112
113 #Varying Coefficient Spatial Probit Model
114 #####
115 # Response is numeric
116 happy$year <- as.numeric(happy$year)
117
118 # Scale time to [0, 1]
119 t_min <- min(happy$year)
120 t_max <- max(happy$year)
121
122 happy$t_scaled <- (happy$year - t_min) / (t_max - t_min)
123
124 # Number of splines (B1)
125 K <- 1
126 #1 Degree
127 B <- bs(happy$t_scaled, df = K, degree = 1)
128 B <- as.matrix(B)
129 colnames(B) <- paste0("B", 1:K)
130
131 # Attach to data frame
132 for (k in 1:K) {
133   happy[[paste0("B", k)]] <- B[, k]
134 }
135
136 spline_terms <- paste0("B", 1:K)
137
138 vars_static <- c("GDP", "registered_pop_10k", "cy", "jr")
139
140 # Extract and coerce to numeric
141 X_static_raw <- happy[, vars_static, drop = FALSE]
142 X_static_raw[] <- lapply(X_static_raw, as.numeric)
143
144 # Scale for numerical stability
145 X_static_scaled <- scale(as.matrix(X_static_raw))
146 colnames(X_static_scaled) <- vars_static
147
148 # GDP as numeric
149 happy$GDP <- as.numeric(happy$GDP)
150
151 # Scale GDP
152 X_vc <- cbind(
153   Intercept = 1,
154   X_static_scaled,           # constant effects
155   as.matrix(happy[, spline_terms]) # varying intercept basis
156 )
157
158 X_vc <- as.matrix(X_vc)
159 storage.mode(X_vc) <- "double"
160
161 fit_vc <- sar_probit_mcmc(
162   y      = y,
163   X      = X_vc,
164   W      = W_panel,
165   ndraw  = 10000,
166   burn.in = 2000,

```

```

167   thinning = 5,
168   prior = NULL,      # or your prior list
169   showProgress = TRUE
170 )
171
172 summary(fit_vc)
173
174 plot(fit_vc)
175
176 fit_vc$bdraw
177 #p_rho = 0.620
178 rho_mean_vc <- 6.19e-01
179 rho_sd_vc <- 4.0128e-02
180
181 rho_ci_norm_vc <- rho_mean_vc + qnorm(c(0.025, 0.975)) * rho_sd_vc
182 rho_ci_norm_vc #[0.540, 0.698]
183
184 int_ci_norm_vc <- 0.14883 + qnorm(c(0.025, 0.975)) * 0.16784
185 int_ci_norm_vc #[-0.180, 0.478]
186
187 gdp_ci_norm_vc <- 0.47206 + qnorm(c(0.025, 0.975)) * 0.12366
188 gdp_ci_norm_vc #[0.230, 0.714]
189
190 pop_ci_norm_vc <- 0.46081 + qnorm(c(0.025, 0.975)) * 0.09051
191 pop_ci_norm_vc #[0.283, 0.638]
192
193 cy_ci_norm_vc <- 0.20094 + qnorm(c(0.025, 0.975)) * 0.09928
194 cy_ci_norm_vc #[0.006, 0.396]
195
196 jr_ci_norm_vc <- 0.28026 + qnorm(c(0.025, 0.975)) * 0.09023
197 jr_ci_norm_vc #[0.103, 0.457]
198
199 b1_ci_norm_vc <- -0.79404 + qnorm(c(0.025, 0.975)) * 0.32471
200 b1_ci_norm_vc #[-1.430, -0.158]
201
202 #####
203 library(pROC)
204
205 predict_prob_probit_approx <- function(fit, X) {
206   beta_hat <- colMeans(fit$bdraw)
207   eta_hat <- as.vector(X %*% beta_hat)
208   p_hat <- pnorm(eta_hat)
209   pmin(pmax(p_hat, 1e-12), 1 - 1e-12)
210 }
211
212 p_lin <- predict_prob_probit_approx(fit_sar, X)
213 roc_lin <- roc(response = y, predictor = p_lin)
214 auc_lin <- auc(roc_lin) #0.7833
215 plot(roc_lin, plot = TRUE, col = "blue", print.auc = TRUE, main = "Linear")
216
217 p_vc <- predict_prob_probit_approx(fit_vc, X_vc)
218 roc_vc <- roc(response = y, predictor = p_vc)
219 auc_vc <- auc(roc_vc) #0.791
220
221 plot(roc_vc, plot = TRUE, col = "blue", print.auc = TRUE, main = "Varying
    Coefficient")

```

```

222 #Follow-up
223 #####
224 compute_dic_probit_approx <- function(fit, X, y, ndraw_eval = 1000) {
225   # X: design matrix used in fitting (with intercept)
226   # y: 0/1 outcome vector
227
228   nd <- min(ndraw_eval, nrow(fit$bdraw))
229   idx <- seq(1, nrow(fit$bdraw), length.out = nd)
230
231   dev_vec <- numeric(nd)
232
233   for (k in seq_along(idx)) {
234     j <- idx[k]
235     beta_j <- fit$bdraw[j, ]
236
237     eta_j <- as.vector(X %*% beta_j)
238     p_j <- pnorm(eta_j)
239     p_j <- pmin(pmax(p_j, 1e-12), 1 - 1e-12)
240
241     loglik_j <- sum(y * log(p_j) + (1 - y) * log(1 - p_j))
242     dev_vec[k] <- -2 * loglik_j
243   }
244
245   # Posterior mean deviance
246   D_bar <- mean(dev_vec)
247
248   # Deviance at posterior mean beta
249   beta_hat <- colMeans(fit$bdraw)
250   eta_hat <- as.vector(X %*% beta_hat)
251   p_hat <- pnorm(eta_hat)
252   p_hat <- pmin(pmax(p_hat, 1e-12), 1 - 1e-12)
253
254   loglik_hat <- sum(y * log(p_hat) + (1 - y) * log(1 - p_hat))
255   D_hat <- -2 * loglik_hat
256
257   # DIC
258   DIC <- 2 * D_bar - D_hat
259
260   list(D_bar = D_bar, D_hat = D_hat, DIC = DIC)
261 }
262
263 DIC_lin <- compute_dic_probit_approx(fit_sar, X, y)
264 DIC_lin$DIC #691.08
265
266 DIC_vc <- compute_dic_probit_approx(fit_vc, X_vc, y)
267 DIC_vc$DIC #680.78
268
269 #Plots
270 #####
271 #Trace plots for SAR and VC
272 par(mfrow=c(2,5), oma=c(0,0,2,0))
273 plot(fit_sar$bdraw[,1], type = "l",
274      main = "GDP",
275      xlab = "",
276      ylab = expression(bold("SAR")))
277 plot(fit_sar$bdraw[,2], type = "l",

```

```

278     main = "Population",
279     xlab = "",
280     ylab = "")
281 plot(fit_sar$bdraw[,3], type = "l",
282      main = "cy",
283      xlab = "",
284      ylab = "")
285 plot(fit_sar$bdraw[,4], type = "l",
286      main = "jr",
287      xlab = "",
288      ylab = "")
289 plot(fit_sar$pdraw, type = "l",
290      main = "rho",
291      xlab = "",
292      ylab = "")
293 plot(fit_vc$bdraw[,1], type = "l",
294      main = "",
295      xlab = "covariate Value",
296      ylab = expression(bold("VC")))
297 plot(fit_vc$bdraw[,2], type = "l",
298      main = "",
299      xlab = "covariate Value",
300      ylab = "")
301 plot(fit_vc$bdraw[,3], type = "l",
302      main = "",
303      xlab = "covariate Value",
304      ylab = "")
305 plot(fit_vc$bdraw[,4], type = "l",
306      main = "",
307      xlab = "covariate Value",
308      ylab = "")
309 plot(fit_vc$pdraw, type = "l",
310      main = "",
311      xlab = "covariate Value",
312      ylab = "")
313 mtext("Trace Plot Comparison of SAR vs. VC",
314      outer = TRUE, cex = 1.2)
315 #Density plots for SAR and VC
316 par(mfrow=c(2,5), oma=c(0,0,2,0))
317 plot(density(fit_sar$bdraw[,1]),
318      main = "GDP",
319      xlab = "",
320      ylab = expression(bold("SAR")))
321 plot(density(fit_sar$bdraw[,2]),
322      main = "Population",
323      xlab = "",
324      ylab = "")
325 plot(density(fit_sar$bdraw[,3]),
326      main = "cy",
327      xlab = "",
328      ylab = "")
329 plot(density(fit_sar$bdraw[,4]),
330      main = "jr",
331      xlab = "",
332      ylab = "")
333 plot(density(fit_sar$pdraw),

```

```

334     main = "rho",
335     xlab = "",
336     ylab = "")
337 plot(density(fit_vc$bdraw[,1]),
338      main = "",
339      xlab = "covariate Value",
340      ylab = expression(bold("VC")))
341 plot(density(fit_vc$bdraw[,2]),
342      main = "",
343      xlab = "covariate Value",
344      ylab = "")
345 plot(density(fit_vc$bdraw[,3]),
346      main = "",
347      xlab = "covariate Value",
348      ylab = "")
349 plot(density(fit_vc$bdraw[,4]),
350      main = "",
351      xlab = "covariate Value",
352      ylab = "")
353 plot(density(fit_vc$pdraw),
354      main = "",
355      xlab = "covariate Value",
356      ylab = "")
357 mtext("Density Plot Comparison of SAR vs. VC",
358      outer = TRUE, cex = 1.2)
359 #Diagnostics
360 #####
361 X <- happy[, Xvars]
362 X[] <- lapply(X, as.numeric)
363
364 drop_vars <- c("top100", "county_name", "city", "province", "county_id", "year_id")
365 drop_vars <- intersect(drop_vars, names(happy))
366
367 Xvars <- setdiff(names(happy), drop_vars)
368
369 Xvars # check predictors
370
371 round(cor(X_scaled), 2)
372 # correlation matrix
373 cor_X <- cor(X, use = "pairwise.complete.obs")
374 round(cor_X, 2)
375
376 library(car)
377
378 fm2 <- lm(
379   as.numeric(top100) ~ .,
380   data = cbind(top100 = happy$top100, happy[, selected_vars])
381 )
382
383 vif_values2 <- vif(fm2)
384
385 sort(vif_values2, decreasing = TRUE)
386
387 #Scaling the selected variables and checking for multicollinearity
388 sel_vars <- happy[, selected_vars] # your chosen predictors
389 sel_vars[] <- lapply(sel_vars, as.numeric) # make sure numeric

```



```

390 sel_vars_scaled <- scale(sel_vars)
391
392 sel_vars_scaled_df <- as.data.frame(sel_vars_scaled)
393
394 lm_scaled_2 <- lm(
395   as.numeric(top100) ~ .,
396   data = cbind(top100 = happy$top100, sel_vars_scaled_df)
397 )
398
399 vif_vals2 <- vif(lm_scaled_2)
400 vif_vals2

```

Listing 1: R Source Code