

Bayesian Logistic Regression Model on Risk Factors of Type 2 Diabetes Mellitus

Emenyonu Sandra Chiaka (Corresponding author)

Department of Mathematics, Faculty of Science, Universiti Putra Malaysia, 43400 UPM, Malaysia

Email: winchiaka@yahoo.co.uk

Mohd Bakri Adam,

Institute for Mathematical Research and Department of Mathematics, Faculty of Science,

Universiti Putra Malaysia, 43400 UPM, Malaysia

Email: pmbakri@gmail.com

Isthriayagy Krishnarajah

Institute for Mathematical Research and Department of Mathematics, Faculty of Science, Universiti Putra Malaysia, 43400 UPM, Malaysia

Email: isthri@upm.edu.my

Shamarina Shohaimi

Department of Biology, Faculty of Science, Universiti Putra Malaysia, 43400 UPM, Malaysia

Email: sshohaimi@yahoo.com

Chris B Guure

Institute for Mathematical Research, Faculty of Science, Universiti Putra Malaysia, 43400 UPM, Malaysia

Email: bam.bey1@hotmail.com

Abstract

This research evaluates the risk of diabetes among 581 men and women with factors such as age, ethnicity, gender, physical activity, hypertension, body mass index, family history of diabetes, and waist circumference by applying the logistic regression model to estimate the coefficients of these variables. Significant variables determined by the logistic regression model were then estimated using the Bayesian logistic regression (BLR) model. A flat non-informative prior, together with a non-informative non-flat prior distribution were used. These results were compared with those from the frequentist logistic regression (FLR) based on the significant factors. It was shown that the Bayesian logistic model with the non-informative flat prior distribution and frequentist logistic regression model yielded similar results, while the non-informative non-flat model showed a different result compared to the (FLR) model. Hence, non-informative but not perfectly flat prior yielded better model than the maximum likelihood estimate (MLE) and Bayesian with the flat prior.

Keywords: Bayesian approach, Binary logistic regression, Parameter estimate, Prior, MCMC.

1. Introduction

Type 2 diabetes is a non-communicable disease characterised by high blood sugar and relative lack of insulin. According to (Albert et al., 1998), diabetes mellitus is a group of metabolic disorders characterized by excess sugar in the blood over a long period of time which is caused by inadequate secretion of insulin, insulin action or both. Type 2 diabetes mellitus (T2DM) is the commonest form of diabetes which has taken hold of over 90% of the diabetic community throughout the world and the fast upswing in the number of people with diabetes is prominent in the urban and rural regions (Valliyot, 2013). The rise in the prevalence of diabetes has been of great concern globally. In a study on global prevalence of diabetes, (Wild global, 2004) estimate the total number of people with T2DM in the year 2000 at 171 million and anticipate it to rise to 366 million in 2030. The prevalence of type 2 diabetes is higher in the developing nations compared to the developed nations. Therefore, (Mafauzy, 2011) predicts that by 2025, the prevalence of diabetes will be higher by 170 percent in the developing world, compared to a 42 percent increase in prevalence rate in the developed nations. The South East Asian countries like Malaysia, has seen a rapid rise in the prevalence of diabetes. With (Mafauzy2006) showing that in Malaysia, the prevalence of diabetes with in the period of 1986 to 1996 rose from 6.3 percent to 8.2 percent and further predictions by the world health organization, reveals that there will be a total number of 2.48 million people with diabetes by 2030 as compared to 2000, where a figure of 0.94 million was estimated thus, showing a 164 percent increase in prevalence rate.

Logistic regression connects a binary dependent variable to a series of independent variables. The dependent variable which is the type 2 diabetes assumes a value of 1 for the probability of occurrence of the disease and 0 for the probability of non-occurrence. Frequentist logistic regression (FLR) makes use of the maximum likelihood estimate (MLE) in order to maximize the probability of obtaining the observed results via the fitted regression parameters. Thus, the FLR brings about point parameter estimates together with standard errors. The uncertainty related to the estimation of parameters is measured by means of confidence interval based on the normality assumption. On the contrary, Bayesian logistic regression (BLR) method makes use of Markov Chain Monte Carlo (MCMC) method in order to obtain the posterior distribution of estimation based on a prior distribution and the likelihood. Thus findings suggest that using the iterative Markov Chain Monte Carlo simulation, BLR provides a rich set of results on parameter estimation. Several studies conclude that BLR performs better in posterior parameter estimation in general and the uncertainty estimation in particular than the ordinary logistic regression. Further reading can be sort from (Lau, 2006) and (Nicodemus, 2001). (Gilks et al., 1996) proposed that in Bayesian, the unknown coefficients β are obtained from posterior distribution, inferences are made based on moment, quantile and the highest density region shown in posterior outcome of the parameter π . In addition, Bayesian approaches in other words can be an alternative to the frequentist approaches. This research aims at applying the BLR model to T2DM to determine the associated risk factors. Uncertainty associated in estimation of the parameters is expressed by means of the posterior distribution. The estimates for the coefficients are obtained by means of FLR, then BLR is also applied on the same variables for coefficient estimation, and the significance of every coefficient estimate is assessed by means of the posterior density generated from the Bayesian analysis. In the present study, factors influencing the occurrence of the disease were determined by applying the Bayesian logistic regression and assuming a non-informative flat prior for every unknown coefficient in the model. On the other hand, several studies also used the BLR method with a non-informative flat prior distribution. However, there have not been many studies on the risk factors of type 2 diabetes mellitus using the Bayesian logistic regression method with a non-informative non-flat prior distribution. Therefore, we decided to assume different prior for the estimation of the parameters which to the best of our knowledge has not been used for the study of T2DM in Malaysia.

2. Materials and Methods

Permission was sought from clinical research centre Kuala Lumpur. The procedure was spearheaded by a family medical specialist who was invited to take part in the study. The main research group organised site feasibility study to recognise clinics that were eligible. Eligibility was based on personal willingness, readiness and agreement to be fully involved and be part of the research group. The research was based on a cluster randomised trial such that the clinics that were selected were done randomly because they met the inclusion criterion. The unit of randomization for the study was the primary health care clinics with males and females ≥ 28 years of age that were diagnosed with T2DM. Individuals with type 1 diabetes and severe hypertension Systolic blood pressure >180 mmHg and Diastolic blood pressure >110 mmHg were excluded. A self-management booklet was shared to all the participants after the training was over and the necessary details were extracted from them. The variables collected during the study were as follows: Demography, social and biological variables and behavioural components.

Logistic regression model will be considered for the occurrence of Type 2 diabetes as a discrete and binary response variable, and factors such as, age, sex, ethnicity, physical activity, family history of diabetes, hypertension, body mass index and waist circumference as explanatory variables. A statistical analysis was carried out to determine the effect of these factors with respect to Type 2 diabetes occurrence. Suppose the Binary logistic regression model is given as:

$$\begin{aligned} \text{Logit}(\pi_i) &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k, \\ \pi &= P(y_i=1 | x_1, \dots, x_k) \end{aligned} \quad (1)$$

Then the estimates of the model can be of the form:

$$\text{Logit}(\hat{\pi}_i) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k \quad (2)$$

Where $\beta = (\beta_0, \beta_1, \dots, \beta_k)$ are estimates of the coefficient β and $\mathbf{x}_i = (x_1, x_2, \dots, x_k)$ are the k independent variables, $\hat{\pi}_i$ is the estimate of the likelihood of type 2 diabetes occurrence.

Given the explanatory variables $x_1, x_2, \dots, x_k, \pi_i$ can be estimated as:

$$\hat{\pi}_i = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)} \quad (3)$$

However, Bayesian framework is the combination of the likelihood function and the prior distribution to yield the posterior distribution. Consequently, the response variable y_i follows a Bernoulli distribution with probability π and is given as:

$$\begin{aligned} y_i &\sim \text{Bernoulli}(\pi_i), \\ \hat{\pi}_i &= \frac{\exp(\mathbf{x}_i \hat{\beta})}{1 + \exp(\mathbf{x}_i \hat{\beta})}. \end{aligned}$$

Where, $\beta = (\beta_0, \beta_1, \dots, \beta_k)$, $\mathbf{y}_i = (y_1, y_2, \dots, y_n)$ and $\mathbf{x}_i = (x_1, x_2, \dots, x_k)$.

The distribution of $(y_i | \mathbf{x}_i \hat{\beta}) = \pi^{y_i} (\pi^{1-y_i})$

For $i = \dots, n$, y_i is the number of successes and $1 - y_i$ is the number of failures.

2.1 The likelihood function

The likelihood function is the probability density function of the data which is seen as a function of the parameter treating the observed data as fixed quantities.

For a given sample size n, the likelihood function is given as:

$$L(Y|X\beta) = \prod_{i=1}^n F(y_i | \mathbf{x}_i \hat{\beta}).$$

Recall that

$$(y_i | \mathbf{x}_i \hat{\beta}) = \pi^{y_i} (\pi^{1-y_i}).$$

Where

$$\pi_i = \frac{\exp(\mathbf{x}_i \hat{\beta})}{1 + \exp(\mathbf{x}_i \hat{\beta})}.$$

$$1 - \pi_i = \frac{1}{1 + \exp(\mathbf{x}_i \hat{\beta})}.$$

Therefore, the likelihood function is of the form:

$$= \left(\frac{\exp(\mathbf{x}_i \hat{\beta})}{1 + \exp(\mathbf{x}_i \hat{\beta})} \right)^{y_i} \left(\frac{1}{1 + \exp(\mathbf{x}_i \hat{\beta})} \right)^{1-y_i} \quad (4)$$

Hence, the likelihood function can be of the form:

$$= \exp\left(\sum_{i=1}^n y_i x_{i\beta}\right) \prod_{i=1}^n \left(\frac{\exp(x_{i\beta})}{1+\exp(x_{i\beta})}\right) \quad (5)$$

2.2 Prior distribution

After the model for our data has been selected, the specification of our prior distribution for the unknown model parameters is made. We assign a prior distribution to all the unknown parameters. Firstly we assume a non-informative flat prior with mean zero and a large variance to all the parameters. However, we also assume a prior distribution to all the unknown parameters with mean zero and small variance 1, this influences the posterior distribution. In Bayesian analysis, precision is used rather than the variance, a large variance is chosen for it to be considered as non-informative while a small variance makes the prior not to be perfectly flat. Our choice of large variance is 10000 (10^4). We assign a normal distribution as prior to each unknown parameters, and the normal distribution is of the form:

$$P(\beta_j) = \prod_{j=0}^k \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left\{-\frac{1}{2}\left(\frac{\beta_j - \mu_j}{\sigma_j}\right)^2\right\} \quad (6)$$

Each β is assigned with mean zero and precision 0.0001, and it is expressed as

$$\beta_j \sim N(0, 0.0001), j=0, \dots, k.$$

Where β_j includes all the coefficients having normal prior distributions with very large variance. However, to have a prior that is not perfectly flat, using the normal prior distribution we give each unknown parameter a mean of zero and a variance of 1 with a known precision given as:

$$\beta_j \sim N(0, 1), j=0, \dots, k.$$

Where β_j include all the coefficients having normal prior distributions with very small variance.

2.3 Posterior distribution

The posterior distribution of the coefficients β is obtained by multiplying the likelihood function in Equation (5) by the prior distribution in Equation (6). The posterior is given as

$$P(\beta|yx) \propto \prod_{i=1}^n L(y|x_{i\beta}) \times \prod_{j=0}^k P(\beta_j)$$

The above expression can be written as

$$p(\beta|yx) \propto \left\{ \exp\left(\sum_{i=1}^n y_i x_{i\beta}\right) \prod_{i=1}^n \left(\frac{\exp(x_{i\beta})}{1+\exp(x_{i\beta})}\right) \times \prod_{j=0}^k \left[\frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{1}{2}\left(\frac{\beta_j - \mu_j}{\sigma_j}\right)^2\right)\right] \right\} \quad (7)$$

3. Results and discussion

The data used in the analysis consist of eight variables of which result show that five significantly contributed to the occurrence of Type 2 diabetes Mellitus. Factors such as age, family history of diabetes, hypertension, body mass index, and waist circumference were significant, whereas physical activity, ethnicity, and gender showed no significance.

Bayesian logistic regression was applied to type 2 diabetes data in order to draw up inferences about the effects of several risk factors contributing to the disease. Using the non-informative prior (flat), the means of the posterior distribution of every coefficient are similar to the coefficient estimates generated by analysing with the Frequentist Logistic Regression. As a result of the Bayesian analysis making use of non-informative prior basically uses the available information by the sample data. The inclusion of information about parameter values into the analysis through the choice of non-perfectly flat prior had an influence on the model. Owing to the fact that a known variance was used resulting to a known precision. On the other hand, when the standard deviation is small, the sample mean is close to each of the sample point, thereby making the result reliable.

Following the application of the generalized linear model (GLM) in the logit link function to the logistic regression, the model shows the estimated parameters, standard errors, and the significance level for all variables

with the intercept estimate in Table.1. The extent of contribution exhibited by the variables in the model is due to their interaction and significance level.

For convergence of each coefficient, the trace plot of 150000 iterations represents the two chains being run for every coefficient that is 75000 iterations for each of the chains. The overlapping of the chains in Figure 2 showed convergence. While the posterior distribution of the model parameters generated from the sampled values reflected kernel densities in Figure 1. Non-informative prior parameter estimates (posterior means) were similar to the estimates obtained by means of the MLE method in Table 2. On the contrary, Convergence monitoring is essential because estimates of the coefficients can only be generated from the iterations, such as the mean, posterior standard deviation, median and quantiles.

Coefficient interpretation for the logistic model were made using (odds ratio) exponential values, For example, estimated coefficient of the variable, family history of diabetes which is 1.217 with an exponential value of 3.378, indicates that the odds of type 2 diabetes are about 3 times higher among persons with family history of diabetes than persons without the history of diabetes in the family, the end points of 95% confidence interval is between 2.320 and 4.974 showing that the change in the odds ratio of type 2 diabetes for the variable family history fall between 2.320 and 4.974 with confidence of 95 percent and thus contributed positively to the development of the disease, while the BLR coefficient estimate for family history of diabetes is 1.233 with exponential value of 3.499 and with the end point of 95 percent posterior interval estimation which falls between 2.349 and 5.067.

For the FLR, the slope variable for waist circumference has an estimate for the coefficient as 0.055 with an exponential value of 1.056 indicating that a one unit increase in waist circumference increases the odds of having type 2 diabetes by a factor of 1.056, with end points of 95 percent confidence interval between 1.033 and 1.082. The estimation indicates the change in the odds ratio of waist circumference for this slope parameter is between 1.033 and 1.082 with 95 percent confidence and the parameter estimate is contributing positively to the disease. While for the BLR coefficient estimate is 0.055 with the exponential value of 1.057 and has a 95 percent posterior interval that falls between 1.032 and 1.082. For body mass index (BMI), the FLR has a coefficient estimate of -0.5097 with exponential value of 0.601 indicating that a unit increase in BMI reduces the odds of developing T2DM by a multiple of 0.601, with end points of 95 percent confidence interval fall between 0.434 and 0.826, while for the BLR the coefficient estimate is -0.515 with exponential value of 0.606 and 95% posterior interval that lies between 0.433 and 0.824. The FLR for the variable age has a coefficient estimate of 0.231 with exponential value of 1.259 also indicating that for a unit increase in age, the odds of developing T2DM increases by a positively correlated multiplicative factor of 1.259, the 95 percent confidence interval is between 1.022 and 1.555, while the BLR for coefficient estimate for age is 0.234 with exponential value of 1.272 and 95 percent posterior interval between 1.022 and 1.563. FLR for the variable hypertension with coefficient estimate of -0.3381 and exponential value of 0.713 with 95 percent confidence interval lying between 0.557 and 0.912 indicating a negative contribution, while the BLR coefficient estimate for hypertension is -0.3429 with exponential value of 0.7154 indicating that a unit increase in hypertension, the odds of developing T2DM reduces by a multiplicative factor of 0.7154, with 95% posterior interval between 0.552 and 0.909.

For the non-informative prior (not perfectly flat). The posterior distribution summaries of the parameters using the non-informative not perfectly flat prior are shown in Table 3. The use of this prior distribution influenced the posterior distribution of the intercept and the regression coefficients. On the other hand, considering the standard deviation and credible interval (which are the Bayesian equivalent of the confidence interval and standard error) for every coefficient assuming a non-informative not perfectly flat prior, shows that each standard deviation is smaller to that of the Bayesian analysis with the non-informative flat prior and the Frequentist analysis implying that the smaller the standard deviation the better the model. In addition, based on the confidence interval estimation for every coefficient in Table 3, since all the coefficients have shorter interval compared to the MLE and the Bayesian model with the non-informative flat prior, this is due to the fact that the approach that gives a shorter interval is considered to be reliable. In other words, the shorter the length of the confidence interval for each coefficient, the better the model. In comparing the non-informative -flat prior with the MLE method on T2DM data. The prior (flat) which is considered will overlap each other because a very large variance for the normal distribution brings about a very small precision which on the other hand yields results that are similar to those of the MLE. So making a comparison between the two methods, one can hardly say with confidence that the model is better than the other. However, with the use of another method, (that is a known variance that results in an informed or known precision) which will still be non- informative but not perfectly flat yielded a better model than the MLE and Bayesian with the flat prior.

Table 1: Analysis of Maximum Likelihood Estimate for all the variables in the full model.

					95% Confidence interval for odds ratio
Coefficient	Estimate	Standard error	P-value	Odds ratio	(2.5%, 97.5%)
Intercept	-4.294	1.090	<0.001	0.014	0.002, 0.112
Gender	-0.061	0.199	0.762	0.941	0.636, 1.393
Age	0.237	0.107	0.027	1.268	1.028, 1.567
Ethnicity	0.125	0.186	0.499	1.134	0.788, 1.637
Physical activity	-0.077	0.182	0.670	0.925	0.648, 1.323
Hypertension	-0.337	0.127	0.008	0.714	0.555, 0.916
Waist circumference	0.057	0.012	<0.001	1.058	1.034, 1.085
Family history	1.226	0.195	<0.001	3.408	2.336, 5.030
Body mass index	-0.548	0.171	0.001	0.578	0.411, 0.806

Table 2: Point estimate of frequentist logistic regression analysis and posterior distribution summaries of parameter estimates of Bayesian logistic model (non-informative flat prior) for type 2 diabetes occurrence with reference to the significant factors.

Coefficient	Point estimate from FLR	posterior mean	Posterior Standard deviation	Quantiles of posterior distribution		
				2.5%	Median	97.5%
Intercept	-4.241	-4.283	1.061	-6.373	-4.288	-2.248
Age	0.231	0.234	0.108	0.022	0.235	0.447
Hypertension	-0.338	-0.343	0.127	-0.594	-0.342	-0.096
Family history	1.277	1.233	0.196	0.854	1.232	1.623
Waist Circumference	0.055	0.055	0.012	0.032	0.055	0.079
Body mass index	-0.510	-0.515	0.167	-0.838	-0.516	-0.193

Table 3: Posterior distribution summaries of the parameters for type 2 diabetes mellitus occurrence with reference to the significant factors with non-informative (not perfectly flat) prior distribution.

Coefficient	Point estimate from FLR	Posterior Standard deviation	Quantiles of posterior distribution		
			2.5%	Median	97.5%
Intercept	-2.090	0.715	-3.466	-2.097	-0.690
Age	0.107	0.098	0.088	0.108	0.297
Hypertension	-0.372	0.123	-0.614	-0.371	-0.132
Family history	1.109	0.187	0.747	1.109	1.478
Waist Circumference	0.036	0.009	0.018	0.036	0.056
Body mass index	-0.425	0.159	-0.739	-0.424	-0.120

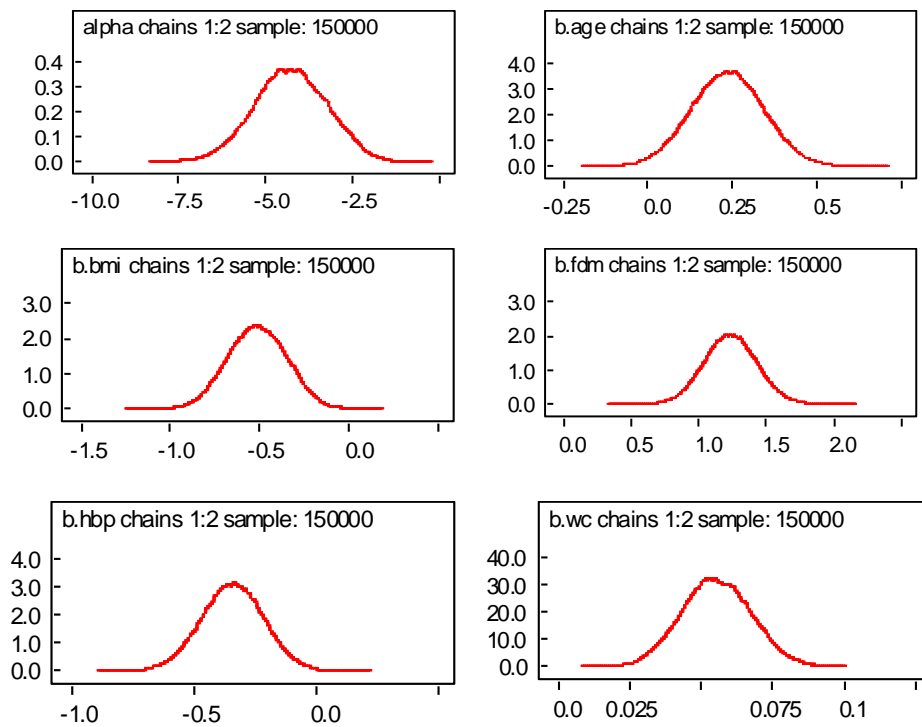


Figure. 1: Density distribution of the corresponding posterior estimates of the intercept, age, and body mass index, family history of diabetes and waist circumference respectively .

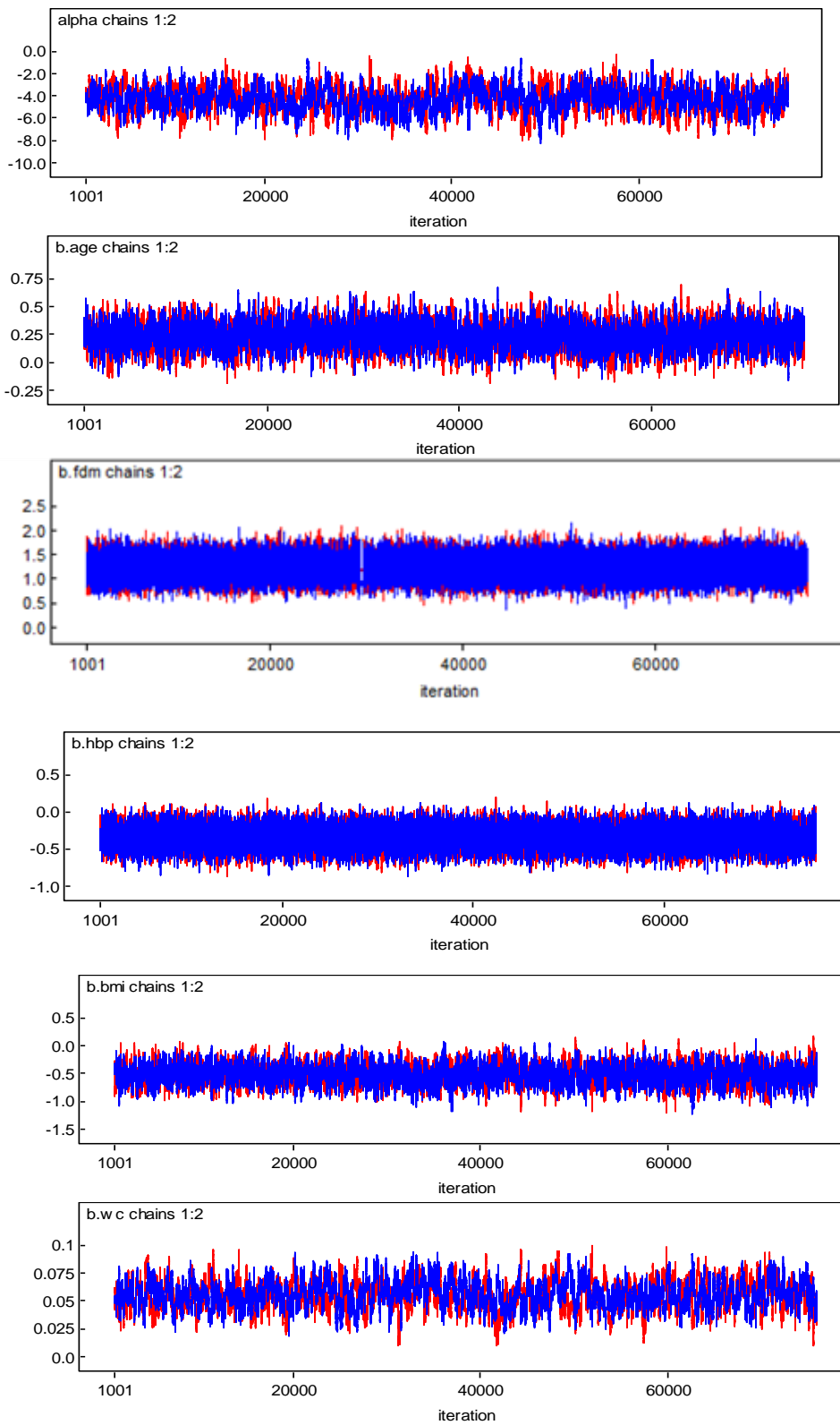


Figure 2: History of trace plots of the given posterior coefficients estimates for the variables of interest. History of trace plots indicating the coefficient values of iterations for the two chains being run.

4. Conclusion

In this study, the type 2 diabetes and its associated risk factors were addressed by the use of Bayesian logistic regression (BLR) model. The Bayesian method incorporated a non-informative flat prior distribution and also another prior, which is still non-informative but not perfectly flat. These models allowed us to analyse the uncertainty associated with the parameter estimation. Comparison between the frequentist logistic regression and Bayesian logistic regression models revealed a similarity in the model results owing to the use of non-informative flat prior distribution. Therefore, our study shows that the use of non-informative but not perfectly flat yielded better model than the MLE and Bayesian with the flat prior.

Acknowledgement

We wish to thank the staff from Kuala Lumpur clinical research centre for providing the data used in this research.

References

- Albert, K., Davidson, M.B., Defronzo, R.A., Drash, A., Genuth, S., Harris, M.I., Kahn, R., Keen, H., Knowler, W.C., Lebovitz, H., et al, (1998). Report of the expert committee on the diagnosis and classification of diabetes mellitus. *Diabetes Care* 21: 5-20
- Gilks, W.R., S. Richardson, and D.J. Spiegelhalter, (1996). *Introducing Markov Chain Monte Carlo*. Pages 1-19. Springer.
- Lau, E., Leung, P., Kwok, T., Woo, J., Lynn, H., Orwell, E., Cummings, S., and Cauley, J., (2006). The determinants of bone mineral density in Chinese men-results from mr.os (hong kong), the first cohort study on osteoporosis in Asian men. *Osteoporosis International*, 17(2): 297-303.
- Mafauzy, M., (2006). Diabetes mellitus in Malaysia. *Medical Journal of Malaysia*, 61(4): 397-398.
- Mafauzy, M., Z. Hussein and S. Chan, (2011). The status of diabetes control in Malaysia: results of diabcare 2008. *Med. J. Malaysia*, 66(3): 175-181.
- Nicodemus, K.K. and A.R. Folsom, (2001). Type 1 and type 2 diabetes and incident hip fractures in postmenopausal women. *Diabetes Care*, 24(7): 1192-1197.
- Valliyot, B., Sreedharan, J., Muttappallymyalil, J., Valliyot, S.B., et al, (2013). Risk factor of type 2 diabetes in the rural population of north Kerala, India: a case control study. *Diabetologia Croatica*, 42(1): 33-40.
- Wild, S. H., Roglic, G., Green, A., Sicree, R., and H. King, (2004). Global prevalence of diabetes: estimate for the year 2000 and projections for 2030 response to rathman and giani. *Diabetes care*, 27(10): 2569-2569.

The IISTE is a pioneer in the Open-Access hosting service and academic event management. The aim of the firm is Accelerating Global Knowledge Sharing.

More information about the firm can be found on the homepage:

<http://www.iiste.org>

CALL FOR JOURNAL PAPERS

There are more than 30 peer-reviewed academic journals hosted under the hosting platform.

Prospective authors of journals can find the submission instruction on the following page: <http://www.iiste.org/journals/> All the journals articles are available online to the readers all over the world without financial, legal, or technical barriers other than those inseparable from gaining access to the internet itself. Paper version of the journals is also available upon request of readers and authors.

MORE RESOURCES

Book publication information: <http://www.iiste.org/book/>

Academic conference: <http://www.iiste.org/conference/upcoming-conferences-call-for-paper/>

IISTE Knowledge Sharing Partners

EBSCO, Index Copernicus, Ulrich's Periodicals Directory, JournalTOCS, PKP Open Archives Harvester, Bielefeld Academic Search Engine, Elektronische Zeitschriftenbibliothek EZB, Open J-Gate, OCLC WorldCat, Universe Digital Library, NewJour, Google Scholar

