# RESEARCH

# Journal of Statistical Distributions and Applications

# **Open Access**

# Meta analysis of binary data with excessive zeros in two-arm trials



Saman Muthukumarana<sup>1\*</sup> , David Martell<sup>2</sup> and Ram Tiwari<sup>3</sup>

\*Correspondence: Saman.Muthukumarana@umanitoba.ca <sup>1</sup>Department of Statistics, University of Manitoba, Machray Hall, Winnipeg, Canada Full list of author information is available at the end of the article

# Abstract

We present a novel Bayesian approach to random effects meta analysis of binary data with excessive zeros in two-arm trials. We discuss the development of likelihood accounting for excessive zeros, the prior, and the posterior distributions of parameters of interest. Dirichlet process prior is used to account for the heterogeneity among studies. A zero inflated binomial model with excessive zero parameters were used to account for excessive zeros in treatment and control arms. We then define a modified unconditional odds ratio accounting for excessive zeros in two arms. The Bayesian inference is carried out using Markov chain Monte Carlo (MCMC) sampling techniques. We illustrate the approach using data available in published literature on myocardial infarction and death from cardiovascular causes. Bayesian approaches presented here use all the data, including the studies with zero events and capture heterogeneity among study effects, and produce interpretable estimates of overall and study-level odds-ratios, over the commonly used frequentist's approaches. Results from the data analysis and the model selection also indicate that the proposed Bayesian method, while accounting for zero events, adjusts for excessive zeros and provides better fit to the data resulting in the estimates of overall odds-ratio and study-level odds-ratios that are based on the totality of the information.

Keywords: Dirichlet process, Model selection, Markov chain Monte Carlo, Simulation

# 1 Introduction

An arm is a standard term for describing clinical trial and it represents a treatment group or a set of subjects. A two-arm study compares a drug with a placebo or drug A with drug B. Sometimes in these studies, the outcome may be binary. A binary outcome is an outcome whose unit can take on only two possible states "0" and "1". For example, outcomes of clinical trials data such as the morbidity and mortality studies are often binary in nature.

The natural distribution for modeling these types of binary data is the binomial distribution given by

$$f(y;p) = \binom{n}{y} p^{y} (1-p)^{n-y} \text{ for } y = 0, 1, \dots, n, \ p \in (0,1).$$

The mean and variance for the binomial random variable are E(Y) = np and Var(Y) = np(1-p) respectively. In a two-arm trial with binary outcomes, it is typically assumed that  $Y_{T_1}, ..., Y_{T_k}$  and  $Y_{C_1}, ..., Y_{C_k}$  are random samples from  $Y_{T_i} \sim Bin(n_{T_i}, P_{T_i})$  and  $Y_{C_i} \sim Bin(n_{C_i}, P_{C_i})$  respectively, where k is the number of studies. In a random effects meta



© The Author(s). 2019 **Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (http://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

analysis of these types of data, the effect size is assumed to vary from study to study. Random effects meta analysis assumes that study effects are a random sample from an underlying relevant distribution of effects, and the combined effect estimates the mean effect of this distribution.

There are a variety of different approaches to analyze these types of data as indicated by some recent literature. See Albert (1995) for various parametrization of binomial models for discrete data within Bayesian settings. Chang et al. (2001) use a mixed effects model to investigate between and within-study variation using rate difference and logit models. Gamalo et al. (2011) propose a Bayesian procedure for testing noninferiority in two-arm studies with a binary primary endpoint that allows the incorporation of historical data on an active control via the use of informative priors but did not consider excessive zeros. Carlin (1992) consider a Bayesian meta-analysis approach for two way contingency table data while Smith et al. (1995) discuss how a full Bayesian analysis can be used to deal with issues in meta-analysis in a natural way using the BUGS language. In this paper, we consider a Bayesian approach for binary data with excessive zeros in two-arm trials. More specifically, we model the excessive zeros using zero inflated binomial distribution and use the Dirichlet process Ferguson (1974) to handle the heterogeneity among studies. There are various zero inflated methods available in the literature. Hall (2000) introduced the framework for count data with many zeros using Poisson and binomial models and likelihood ratio tests based inference for zero inflated Poisson models are discussed in Huang et al. (2014). A Bayesian inference framework for zero inflated Poisson regression models is discussed in Ghosh et al. (2006). A rich class of nonparametric Bayesian priors for study effects and Bayesian nonparametric Polya tree mixture model are developed in Branscum and Hanson (2008) and Burr and Doss (2005).

In Section 2, we describe Bayesian model specification used in the paper. The likelihood function and the priors are described. Study effects have a Dirichlet process prior distribution for capturing heterogeneity among studies. We then obtain posterior summary statistics which describe key features in the model. In particular, posterior expectations are approximated through Markov chain Monte Carlo (MCMC) methods. In Section 3, the model is applied to a large dataset available in the literature Nissen and Wolski (2007). We perform the model selection using the log-pseudo marginal likelihood (LPML) comparing the Binomial and zero-inflated Binomial (ZIB). The results suggest that when the data has a high percentage of observed zeros, ZIB model is a more appropriate model to use. Furthermore, the use of Dirichlet process has advantage over the more commonly used random effects model with normally distributed random effects based on DerSimonian-Laird approach DerSimonian and Laird (1986) or a Bayesian approach using normal priors, in terms of its inherent clustering property resulting in the studies with similar effects to cluster, and thus providing more robust estimates. We also test the approach using simulation studies in Section 4 and study the effect of excessive zeros in the ZIB models. We conclude with a short discussion in Section 5.

# 2 Model development

Consider two-arm trials with binary outcomes and let  $Y_{T_i} \stackrel{\text{ind}}{\sim} Bin(n_{T_i}, P_{T_i})$  and  $Y_{C_i} \stackrel{\text{ind}}{\sim} Bin(n_{C_i}, P_{C_i})$ , i = 1, ..., k, where k is the number of studies. Then the joint likelihood  $L = L(y_{T_1}, ..., y_{T_k}, y_{C_1}, ..., y_{C_k} | \mu, P_T, P_C)$  is

$$L = \prod_{i=1}^{k} \left\{ {}^{n_{T_i}} C_{y_{T_i}} P_{T_i}^{y_{T_i}} (1 - P_{T_i})^{n_{T_i} - y_{T_i}} \right\} \prod_{i=1}^{k} \left\{ {}^{n_{C_i}} C_{y_{C_i}} P_{C_i}^{y_{C_i}} (1 - P_{C_i})^{n_{C_i} - y_{C_i}} \right\}.$$
(1)

In random effects meta-analysis formulation, we assume that  $P_T$  and  $P_C$  follow logistic models, and define

 $P_{T_i} = \frac{exp\{\mu+Tr+\alpha_i+e_i\}}{1+exp\{\mu+Tr+\alpha_i+e_i\}} \text{ and } P_{C_i} = \frac{exp\{\mu+e_i\}}{1+exp\{\mu+e_i\}}.$ That is,  $logit(P_{T_i}) = \mu + Tr + \alpha_i + e_i; logit(P_{C_i}) = \mu + e_i, i = 1, \dots, k.$  This gives,

$$logit(P_{T_i}) - logit(P_{C_i}) = Tr + \alpha_i$$
;  $i = 1, ..., k$ ,

where  $logit(p) = log\left(\frac{p}{1-p}\right)$  is the log-odds ratio of p,  $\mu$  is the intercept, Tr is the treatment effect,  $\alpha_i$  and  $e_i$  are the study effects and error terms. As proposed by Muthukumarana and Tiwari in Muthukumarana and Tiwari (2016), consider a Bayesian approach and assume that  $\{\alpha_i; i = 1, ..., k\}$  is a sample from a Dirichlet process with concentration parameter  $\rho$  and the baseline distribution H. We assume that the baseline distribution H is  $N(0, \sigma_H^2)$ . More specifically, we assume that

$$\alpha_{i} \sim DP(\rho, H)$$

$$H \sim N(0, \sigma_{H}^{2})$$

$$e_{i} \sim N(0, \sigma_{e}^{2})$$

$$f(\mu) \propto \text{constant}$$

$$Tr \sim N(0, \sigma_{Tr}^{2})$$

$$\rho \sim U[0.1, 1000]$$
(2)

where hyper parameters  $\sigma_{H}^{2}$ ,  $\sigma_{Tr}^{2}$  and  $\sigma_{e}^{2}$  are assumed to be known. We now obtain the posterior characterizations of parameters using Neal's algorithm Neal (2000) using Gibbs sampling as follows.

$$f(\alpha_c | y_{T_j} : c_j = c) \propto \prod_{j:c_j = c}^k \left( \frac{1}{1 + exp \left\{ \mu + Tr + \alpha_c + e_j \right\}} \right)^{n_{T_j}}$$

$$exp \left\{ \frac{-1}{2\sigma_H^2} \left( \alpha_c - \sigma_H^2 \sum_{j:c_j = c} y_{T_j} \right)^2 \right\}$$
(3)

$$f\left(e_{i}|\underline{y}\right) \propto \left(\frac{1}{1+\exp\left\{\mu+Tr+\alpha_{i}+e_{i}\right\}}\right)^{n_{T_{i}}} \left(\frac{1}{1+\exp\left\{\mu+e_{i}\right\}}\right)^{n_{C_{i}}}$$

$$exp\left\{\frac{-1}{2\sigma_{e}^{2}}\left(e_{i}-\sigma_{e}^{2}\left(y_{T_{i}}+y_{C_{i}}\right)\right)^{2}\right\}$$
(4)

$$f\left(\mu|\underline{y}\right) \propto \prod_{j=1}^{k} \left(\frac{1}{1 + \exp\left\{\mu + Tr + \alpha_j + e_j\right\}}\right)^{n_{T_j}} \left(\frac{1}{1 + \exp\left\{\mu + e_j\right\}}\right)^{n_{C_j}}$$

$$\exp\left\{\left(\sum_{j=1}^{k} y_{T_j} + y_{C_j}\right)\mu\right\}$$
(5)

$$f\left(Tr|\underline{y}\right) \propto \prod_{j=1}^{k} \left(\frac{1}{1 + exp\left\{\mu + Tr + \alpha_j + e_j\right\}}\right)^{n_{T_j}}$$
$$exp\left\{\frac{-1}{2\sigma_{Tr}^2} \left(Tr - \sigma_{Tr}^2 \sum_{j=1}^k y_{T_j}\right)^2\right\}$$
(6)

$$f\left(\rho|\underline{y}\right) \propto \rho^{r-1}\left(\rho+k\right) B\left(\rho+1,k\right) I_{[0.1,1000]}(\rho) \tag{7}$$

Note that the likelihood in (1) does not account for excessive zeros in the data. For this reason, we now consider a zero inflated binomial model for the data as follows.

$$Y_{T_i} \stackrel{\text{ind}}{\sim} ZIB\left(p_0, n_{T_i}, P_{T_i}\right), \ Y_{C_i} \stackrel{\text{ind}}{\sim} ZIB\left(q_0, n_{C_i}, P_{C_i}\right), i = 1, \dots, k.$$

That is,

$$Y_{T_i} = \begin{cases} 0 & \text{with probability } p_0 \\ Bin\left(n_{T_i}, P_{T_i}\right) & \text{with probability } 1 - p_0. \end{cases}$$

Similarly,

$$Y_{C_i} = \begin{cases} 0 & \text{with probability } q_0 \\ Bin\left(n_{C_i}, P_{c_i}\right) & \text{with probability } 1 - q_0. \end{cases}$$

This modification brings two more extra parameters to the model and we assume that

$$p_0 \sim Beta(a, b)$$

$$q_0 \sim Beta(c, d).$$
(8)

where hyper parameters a, b, c and d are assumed to be known. We obtain the the posterior characterizations of parameters under zero inflated binomial likelihood as follows.

$$f\left(\alpha_{c}|y_{T_{j}}:c_{j}=c\right) \propto \prod_{j:c_{j}=c}^{k} \left[p_{0}+(1-p_{0})\left(\frac{1}{1+exp\left\{\mu+Tr+\alpha_{c}+e_{j}\right\}}\right)^{n_{T_{j}}}\right]^{u_{j}}$$
$$\left[(1-p_{0})\left(\frac{1}{1+exp\left\{\mu+Tr+\alpha_{c}+e_{j}\right\}}\right)^{n_{T_{j}}}exp\left\{\alpha_{c}y_{T_{j}}\right\}\right]^{1-u_{j}} (9)$$
$$exp\left\{\frac{-1}{2\sigma_{H}^{2}}\alpha_{c}^{2}\right\}$$

$$f\left(e_{i}|\underline{y}\right) \propto \left[p_{0} + (1-p_{0})\left(\frac{1}{1+exp\left\{\mu+Tr+\alpha_{i}+e_{i}\right\}}\right)^{n_{T_{i}}}\right]^{u_{i}} \\ \left[(1-p_{0})\left(\frac{1}{1+exp\left\{\mu+Tr+\alpha_{i}+e_{i}\right\}}\right)^{n_{T_{i}}}exp\left\{y_{T_{i}}e_{i}\right\}\right]^{1-u_{i}} \\ \left[q_{0} + (1-q_{0})\left(\frac{1}{1+exp\left\{\mu+e_{i}\right\}}\right)^{n_{C_{i}}}\right]^{w_{i}} \\ \left[(1-q_{0})\left(\frac{1}{1+exp\left\{\mu+e_{i}\right\}}\right)^{n_{C_{i}}}exp\left\{y_{C_{i}}e_{i}\right\}\right]^{1-w_{i}} \\ exp\left\{\frac{-1}{2\sigma_{e}^{2}}e_{i}^{2}\right\}$$
(10)

$$f\left(\mu|\underline{y}\right) \propto \prod_{j=1}^{k} \left[ p_{0} + (1-p_{0}) \left( \frac{1}{1+exp\left\{\mu+Tr+\alpha_{j}+e_{j}\right\}} \right)^{n_{T_{j}}} \right]^{u_{j}} \left[ (1-p_{0}) \left( \frac{1}{1+exp\left\{\mu+Tr+\alpha_{j}+e_{j}\right\}} \right)^{n_{T_{j}}} exp\left\{y_{C_{j}}\mu\right\} \right]^{1-u_{j}} \left[ q_{0} + (1-q_{0}) \left( \frac{1}{1+exp\left\{\mu+e_{j}\right\}} \right)^{n_{C_{j}}} \right]^{w_{j}} \left[ (1-q_{0}) \left( \frac{1}{1+exp\left\{\mu+e_{j}\right\}} \right)^{n_{C_{j}}} exp\left\{y_{C_{j}}\mu\right\} \right]^{1-w_{j}} \right]^{n_{T_{j}}} \left[ f\left(Tr|\underline{y}\right) \propto \prod_{j=1}^{k} \left[ p_{0} + (1-p_{0}) \left( \frac{1}{1+exp\left\{\mu+Tr+\alpha_{j}+e_{j}\right\}} \right)^{n_{T_{j}}} exp\left\{y_{T_{j}}Tr\right\} \right]^{1-u_{j}} \left[ (1-p_{0}) \left( \frac{1}{1+exp\left\{\mu+Tr+\alpha_{j}+e_{j}\right\}} \right)^{n_{T_{j}}} exp\left\{y_{T_{j}}Tr\right\} \right]^{1-u_{j}} \left[ exp\left\{ \frac{-1}{2\sigma_{Tr}^{2}}Tr^{2} \right\}$$

$$(12)$$

$$f\left(\rho|\underline{y}\right) \propto \rho^{r-1} \left(\rho+k\right) B\left(\rho+1,k\right) I_{[0.1,1000]}(\rho) \tag{13}$$

$$f\left(p_{0}|\underline{y}\right) \propto \left[p_{0} + (1-p_{0})\left(\frac{1}{1+exp\left\{\mu+Tr+\alpha_{j}+e_{j}\right\}}\right)^{n_{T_{j}}}\right]^{u_{j}} \\ \left[(1-p_{0})\left(\frac{1}{1+exp\left\{\mu+Tr+\alpha_{j}+e_{j}\right\}}\right)^{n_{T_{j}}}exp\left\{y_{T_{j}}\left(\mu+Tr+\alpha_{j}+e_{j}\right)\right\}\right]^{1-u_{j}} \\ p_{0}^{a-1}\left(1-p_{0}\right)^{b-1}$$
(14)

$$f\left(q_{0}|\underline{y}\right) \propto \left[q_{0} + (1 - q_{0})\left(\frac{1}{1 + exp\left\{\mu + e_{j}\right\}}\right)^{n_{C_{j}}}\right]^{w_{j}} \\ \left[\left(1 - q_{0}\right)\left(\frac{1}{1 + exp\left\{\mu + e_{j}\right\}}\right)^{n_{C_{j}}}exp\left\{y_{C_{j}}\left(\mu + e_{j}\right)\right\}\right]^{1 - w_{j}} \\ q_{0}^{c-1}\left(1 - q_{0}\right)^{d-1} \\ \left[1, y_{T_{i}} = 0\right], \qquad \left[1, y_{C_{i}} = 0\right]$$

$$(15)$$

where  $u_j = \begin{cases} 1, y_{T_j} = 0 \\ 0, y_{T_j} = 1 \end{cases}$  and  $w_j = \begin{cases} 1, y_{C_j} = 0 \\ 0, y_{C_j} = 1 \end{cases}$ .

We investigate the suitability of the zero inflated binomial distribution using the log pseudo marginal likelihood (LPML) Gelfand et al. (1992) in Section 4.

# 3 Data analysis

We illustrate the approach discussed in Section 2 using a published data set on counts of the number of people experiencing myocardial infarction from the use of drugs with an active ingredient "rosiglitazone" Nissen and Wolski (2007). The data used in this section provides information on diabetes patients, 42 diabetes trials having zero events in both arms, and possible heart condition or death resulting from the use of rosiglitazone. Rosiglitazone is a treatment used to treat patients with type 2 diabetes. The data

provide information on diabetes patients, 42 diabetes trials, and possible heart condition or death resulting from the use of rosiglitazone. Rosiglitazone is a treatment for diabetes widely used in treating patients with type 2 diabetes. We separately apply the model on myocardial infarction and death from cardiovascular based on these 42 studies. There were 86 myocardial infarctions in the rosiglitazone group and 72 in the control group. There were 39 deaths from cardiovascular causes in the rosiglitazone group and 22 in the control group. Note that the percentages of observed zeros from the 42 studies in the treatment and control arms for myocardial infarction are 23% and 57% respectively. Similar percentages for cardiovascular causes are 50% and 80% respectively. We set the hyper parameters as a = b = c = d = 1,  $\sigma_H^2 = 2$ ,  $\sigma_T^2 = 2$  and  $\sigma_e^2 = 2$ . Note that the choice of these values result sufficiently diffuse priors in the range of logit scale of primary parameters. We implement the models developed in Section 2 using R. The results are based on a MCMC simulation with a burn-in period of 1000 iterations followed by 30,000 iterations using thinning of 5. We use the data from the 42 studies, without stratifying them into small and large studies, as the purpose of the proposed work is an illustration of the method and not in in-depth analysis of the data by using different methods or by slicing and dicing the data. The posterior box plots of study effects under two models on myocardial infarction are given in Figs. 1 and 2. The advantage of using DP prior is the flexibility and also the ability to cluster studies appropriately. The clustering is based on the values assigned to each study effects based on their posterior distributions, which are approximated using MCMC. Those studies that share the same study effects will be considered to belong to the same group. Note that there were 5 clusters in myocardial infarction and 4 clusters in cardiovascular causes based on study effects. To evaluate the performance between Binomial and ZIB models, we use the LPML which is based on Conditional Predictive Ordinates(CPO). A detailed discussion of the CPO statistic and its applications to model selection can be found in Geisser (1993) and Gelfand and Dey (1994). The LPML is computed as  $\sum_{i=1}^{K} logp(y_i|y_{-i})$  where  $y_{-i}$  denotes the observation vector y with the *i*<sup>th</sup>





observation deleted. The model with larger value of LPML is preferred. The estimates of  $\mu$ , *Tr* and the LPML values are given in Table 1. The LPML prefers binomial model over the ZIB model and the two models estimate the parameter *Tr* differently.

We now investigate the study effects on death from cardiovascular causes. The posterior box plots of study effects under two models on death from cardiovascular causes are given in Figs. 3 and 4. The plots indicate that ZIB model is capable in capturing the heterogeneity of study effects. The estimates of  $\mu$ , *Tr* and the LPML values are given in Table 1. In this case, the LPML strongly prefers ZIB model over the binomial model. This is in agreement with the fact that there are large amount of excessive zeros on death from cardiovascular causes relative to myocardial infarction.

A summary of estimates of odds ratios under Binomial, zero inflated Binomial and DerSimonian- Laird random effects models are given in Figs. 5 and 6. For myocardial infarction, DerSimonian- Laird random effects model gives an overall odds ratio of 1.29 with a 95% confidence interval of (0.9, 1.85). On the other hand, Binomial and zero inflated Binomial models provide an overall summary of odds ratio of 1.04 (0.98, 1.1) and 1.07 (0.97, 1.17) respectively. These estimates and 95% credible intervals for cardiovas-cular causes are 1.2 (0.64, 2.24), 1.03 (0.97, 1.09) and 1.13 (0.93, 1.33) respectively. It is clear that our approach provides overall odds ratios estimates that are slightly lower than that from DerSimonian- Laird overall estimate. Note that DerSimonian-Laird approach is based on the non-zero studies. Also note that Binomial and zero inflated Binomial models identify more heterogeneous study effects than DerSimonian- Laird random effects

Table 1 Paramete	r estimates with	each model	along w	ith LPML
------------------	------------------	------------	---------	----------

	Myocardial infarction		Cardiovascular causes	
Parameter	Binomial model	ZIB model	Binomial model	ZIB model
$\mu$	0.0394 (0.0277)	0.0709 (0.0503)	0.0709 (0.0503)	0.1235 (0.0876)
Tr	-1.1989 (0.3945)	-3.3612 (0.4870)	-3.3612 (0.4870)	-4.5339 (0.5707)
LPML	-173.5474	-179.7584	-156.3964	-125.0447



model. According to Figs. 5 and 6, we notice that zero inflated Binomial model identifies more heterogeneous effects than Binomial model while Binomial model identifies more heterogeneous effects than DerSimonian- Laird approach. DerSimonian- Laird estimated random effects variances are zero for both scenarios and this suggests that our approach is superier than DerSimonian- Laird random effects model when there is heterogeneity among studies and LPML model selection criteria will choose the best model in terms of prediction ability.

We now examine the effects of zero inflated parameters  $p_0$  and  $q_0$  on the analysis. The graphical posterior summaries of  $p_0$  and  $q_0$  on myocardial infarction and cardiovascular causes are given in Figs. 7, 8, 9 and 10. In addition, the numerical posterior summaries of  $p_0$  and  $q_0$  are given in Table 2. It is clear that the posterior distributions of  $p_0$  and  $q_0$ 





and their numerical summaries for myocardial infarction and cardiovascular causes make sense with respect to the percentages of zeros in the data. We also consider a Beta(0.5, 0.5)prior on  $p_0$  and  $q_0$  in order to investigate the prior sensitivity. The numerical posterior summaries of  $p_0$  and  $q_0$  under Beta(0.5, 0.5) prior are given in Table 3. We notice a magnitude change in estimates of  $p_0$  and  $q_0$  in this case but the estimates of primary parameters  $\mu$  and Tr are very close indicating that odds ratios are not sensitive to the choice of prior settings. This indicates that inference on  $p_0$  and  $q_0$  will be sensitive to the choice of priors so one should select these priors carefully based on application specific apriori knowledge on zero inflated parameters.

It is important to look at some convergence assessment plots related to the MCMC simulation as this is a high dimensional problem. The trace plot, histogram and autocorrelation plot of  $\mu$  under binomial model on Myocardial Infarction are given in Fig. 11. The trace plot appears to stabilize immediately and hence provides no indication of lack of convergence in the Markov chain. The autocorrelation plot also appears to dampen quickly. Trace plots of study effects on Myocardial Infarction are given in Fig. 12. The trace plots of study effects on death from cardiovascular causes indicate similar behavior. Similar plots were obtained for all of the parameters under each model and provide the evidence of the convergence of the Markov chains.

Note that one can also assign a simpler parametric normal prior on study effects  $\alpha_i$  in place of the DP prior. We now re-analyze the data assuming that study effects are arising from a  $N(0, \sigma_H^2)$  prior distribution. We remark that this is the baseline distribution of the







DP prior in (2). In this case, forest plots of odds ratios for each model are given in Fig. 13. The estimates of primary parameters of interest and LPML values are given in Table 4. The LPML model selection criteria clearly indicates that the DP prior in (2) is superior than the conventional parametric prior.

Note that the overall decision to assess the safety should be based on  $p_0, q_0$  and the overall odds ratio (OR). For example, the treatment can be declared is to be safer than the control, if  $OR \leq 1$ , and  $p_0 > q_0$ . Also notice that estimates of  $(p_0, q_0)$  are independent of the odds ratio because the counts cannot be in "true" zero arms and "Binomial" arms. We combine the two metrics, conditional OR and  $(p_0, q_0)$ , to come up with an overall





unconditional odds ratio. We define it to be modified odds ratio =  $OR \times (1 - p_0)/(1 - q_0)$ . Note that when  $p_0 = q_0$ , modified odds ratio is same as OR. If  $p_0 > q_0$ , this adjusts OR, by multiplying by a factor less than 1, and if  $p_0 < q_0$ , it adjust OR by multiplying by a factor,  $h(p_0, q_0) = (1 - p_0)/(1 - q_0)$  is the ratio of probabilities of observing Bernoulli counts in the two arms, and can be considered as odds for observing Bernoulli counts in the two arms. In frequentist setup,  $h(\hat{p}_0, \hat{q}_0)$  is independent of  $\hat{\mu}$ , and hence independent of conditional odds ratio. In fact,  $\hat{p}_0$  and  $\hat{q}_0$  converge to  $p_0$  and  $q_0$  with probability 1, and hence  $h(\hat{p}_0, \hat{q}_0)$  also converge to  $h(p_0, q_0)$  with probability 1, as h is a continuous function (from Slutsky's theorem). So, the estimated modified odds ratio is a consistent estimator for unconditional odds ratio for various models in Table 5. As estimate of  $p_0$  is less than  $q_0$  for both examples (Myocardial Infarction and cardiovascular causes), the modified OR values are higher than the corresponding OR values.

#### 4 Results from simulation studies

To understand the role of  $p_0$  and  $q_0$  in the model, different simulation studies were carried out. For this purpose, we generate random ZIB values with empirical binomial parameters. We first generate 42 pairs of independent binary, 0 and 1, variables from Bernoulli  $(p_0)$  and Bernoulli  $(q_0)$  where  $p_0$  and  $q_0$  are from the set of values {(0.1, 0.1), ..., (0.9, 0.9)}. We then assign the true-zeros at the places with 1s, and generate binomial outcomes from  $B(\bar{n}_T, \hat{P}_{T_i})$  and from  $B(\bar{n}_C, \hat{P}_{C_i})$ , where  $\bar{n}_T, \bar{n}_C, \hat{P}_{T_i}$  and  $\hat{P}_{C_i}$  are empirical estimates. Then, MCMC sampling scheme described in Section 2 was carried out using *R* to obtain the

Parameter	Myocardial infarction	Cardiovascular causes		
<i>p</i> <sub>0</sub>	0.0495 (0.044)	0.231 (0.117)		
90	0.27 (0.118)	0.566 (0.137)		

**Table 2** Posterior mean and standard deviation (in parentheses) of  $p_0$  and  $q_0$ 

Parameter	Myocardial infarction	Cardiovascular causes
<i>p</i> <sub>0</sub>	0.0254 (0.036)	0.179 (0.125)
90	0.247 (0.125)	0.538 (0.161)
$\mu$	0.073 (0.051)	0.118 (0.084)
Tr	-3.357 (0.444)	-4.707 (0.632)

**Table 3** Posterior mean and standard deviation (in parentheses) of  $p_0$  and  $q_0$  under *Beta*(0.5, 0.5) prior distribution

posterior estimate of  $p_0$  and  $q_0$ . This was done 1000 times for each pair to obtain the mean and standard error of each estimate. For various scenarios of excessive zeros, the results are given in Table 6. The results indicate that when true values of  $p_0$  is small and the observed values of zeros in the simulated data in treatment arm (control arm) is also small (large), the estimated values of  $p_0$  and  $q_0$  are also small (large), whereas when the values of  $p_0$  and  $q_0$  are large the simulated data has large proportion of zeros in both the arms, this results in large estimated values of  $p_0$  and  $q_0$ . In both the situations, the estimated values of  $p_0$  and  $q_0$  are in conformity with the observed percentages of zeros in the simulated data. The estimates of  $p_0$  and  $q_0$  remain high in spite of their true choices from the parameter values. Note that our primary interest is on alphas and on treatment arm not on the control arm, so we may not need to investigate q0 very well as it can be trated as nuisance parameter. In practice, one should have a very good apriori knowledge of q0 which will allow to assign an informative prior as it is reflecting the zeros in the control arm. This indicates that the use of ZIB is more appropriate when there are excessive zeros in the data.

### 5 Discussion

Binary data naturally arise in clinical trials in health sciences. In some cases, they arise with excessive zeros. In this paper, we have provided a random effects meta analysis





approach for binary data with excessive zeros in two-arm trials. The suitability of the binomial and zero inflated binomial model was assessed in the presence of Dirichlet process as the prior for the study effects. The approach can be used as a template for meta analysis of binary data and a user may choose the proper model using log pseudo marginal likelihood. We have shown that our approach is superior than DerSimonian- Laird random



			< 117	
Model	μ	Tr	LPML	Overall odds ratio with 95% C.I.
Myocardial - Bin	0.039 (0.028)	-1.148 (0.369)	-179.5883	1.04 (0.979, 1.101)
Myocardial - ZIB	0.073 (0.052)	-3.311 (0.435)	-182.5765	1.07 (0.961, 1.194)
Cardiovascular - Bin	0.036 (0.025)	-1.892 (0.427)	-161.3744	1.04 (0.983, 1.094)
Cardiovascular - ZIB	0.123 (0.089)	-4.585 (0.568)	-125.5402	1.13 (0.917, 1.356)

**Table 4** Parameter estimates for various models under  $N(0, \sigma_{\mu}^2)$  prior on study effects

effects model when there is heterogeneity among studies and LPML model selection criteria can be used to selection the best model among the Bayesian models (not including DerSimonian-Laird model) for a given data set.

The Bayesian approaches discussed in this paper allowed to incorporate the zerostudies in the likelihood, and we found that the point estimates of the overall odds-ratio from these methods, were lower than the estimates reported in the literature Nissen and Wolski (2007). The use of ZIB model was to identify the percentage of excessive zeros, that is, the studies where the events could not occur, from the (Binomially) modeled zeros where the zero events occurred. Note that under ZIB, some zeros are observed with probability  $p_0$  and some from Binomial model, making the probability of zero-event to be  $p_0 + (1 - p_0)(1 - P_T)^{n_T}$  in the treatment arm. With the use of ZIB model, the Bayes estimates of the odds-ratio went slightly up than with the use of Binomial model, but still they were lower than the results from DerSimonian-Laird random effects model and the resulting estimates in Nissen and Wolski (2007). Note also that DP model being discrete with probability 1, has a clustering property, where the study effects, that are alike, fall in the same cluster. We also investigated the suitability of the DP prior over the conventional parametric normal prior on study effects. The LPML model selection indicated that DP prior is superior than the conventional parametric normal prior. Finally, as the results from ZIB model on the parameters  $p_0$ ,  $q_0$  and OR need to be interpreted together, a modified OR was introduced.

As a future direction of research, we would like to extend the approach discussed in this article for ordinal category data. For example, in some applications, the clinical trial end point could be a response variable in an ordinal scale with multiple categories such as Good/Moderate/Critical etc. This type of ordinal response data can be viewed as multivariate responses arising from continuous latent variables with cut-points. We assume that there is a continuous latent outcome behind these ordinal outcomes such that  $X_i = (X_{i1}, \ldots, X_{im})' \sim \text{Normal}(\mu, \Sigma)$  where *X*'s are the latent outcomes and *m* is the number of ordinal categories. Then the latent variables  $X_{ij}$ 's can be converted to the observed  $Y_{ij}$  using a cut-point vector  $\lambda$ . However the choice of cut-points and their priors need to be carefully selected as there are two arms and the counts on categories could be

**Table 5** Modified odds ratios, standard deviations (in parentheses) and credible intervals under DP and normal priors

Model	Modified odds ratio	95% credible interval		
Myocardial - DP Prior	1.448 (0.277)	(1.05, 2.11)		
Myocardial - Normal Prior	1.446 (0.275)	(1.05, 2.11)		
Cardiovascular - DP Prior	2.209 (0.830)	(1.15, 4.27)		
Cardiovascular - Normal Prior	2.224 (0.837)	(1.16, 4.33)		

Initial pair ( $p_0, q_0$ )	Mean of the pos- terior means of <i>p</i> <sub>0</sub> estimates	Standard error of $p_0$ estimates	Mean of the pos- terior means of q <sub>0</sub> estimates	Standard error of $q_0$ estimates
(0.1,0.1)	0.19888206	0.0907763	0.52069774	0.08299897
(0.2,0.2)	0.29206103	0.10530334	0.59059405	0.05131143
(0.3,0.3)	0.33836514	0.09115085	0.62706227	0.06062431
(0.4,0.4)	0.3587838	0.10839263	0.63652	0.06102938
(0.5,0.5)	0.4977337	0.1267577	0.7200964	0.04900527
(0.6,0.6)	0.6180761	0.1255134	0.7705962	0.06610487
(0.7,0.7)	0.6525938	0.1885614	0.8096615	0.07434429
(0.8,0.8)	0.796206	0.09240875	0.8696273	0.06536744
(0.9,0.9)	0.8958989	0.0501579	0.935196	0.03310256

Table 6 Simi	ulation	studies for	the m	vocardial	infarction	data
--------------	---------	-------------	-------	-----------	------------	------

sparse. In this case, one can consider an objective Bayes approach following the development in Bayarri et al. (2008). Yet another extension of the proposed model is where there are multinomial data with some particular cell(s) being observed excessively. This kind of data may arise from trials with patient reported outcomes.

#### Acknowledgments

The authors thank Editor-in-Chief and three anonymous reviewers whose comments helped to improve the manuscript. This article reflects the views of the authors and should not be attributed to FDA's views or policies.

#### Authors' contributions

All authors have contributed equally to the work and approved the final version of the paper.

#### Funding

Muthukumarana's research has been partially supported by a Discovery grant from the Natural Sciences and Engineering Research Council of Canada. Martell's research internship was funded by Mitacs Globalink program.

#### Availability of data and materials

Data and code can be requested by contacting the authors.

#### **Competing interests**

The authors declare that they have no competing interests.

#### Author details

<sup>1</sup> Department of Statistics, University of Manitoba, Machray Hall, Winnipeg, Canada. <sup>2</sup>ITAM, Mexico City, Mexico. <sup>3</sup>Office of Biostatistics, Center for Drug Evaluation and Research, Food and Drug Administration, 10903 New Hampshire Ave,Silver Spring, USA.

# Received: 13 September 2018 Accepted: 2 July 2019

Published online: 24 July 2019

#### References

- Albert, J.: Teaching Inference about Proportions Using Bayes and Discrete Models. J. Stat. Educ. 3 (1995). https://doi.org/ 10.1080/10691898.1995.11910494
- Bayarri, M. J., Berger, J. O., Datta, G. S.: Objective Bayes testing of Poisson versus inflated poisson models. Inst. Math. Stat. 3, 105–121 (2008)

Branscum, A. J., Hanson, T. E.: Bayesian nonparametric meta-analysis using Polya tree mixture models. Biometrics. 64, 825–833 (2008)

Burr, D., Doss, H.: A Bayesian semiparametric model for random-effects meta-analysis. J. Am. Stat. Assoc. **100**, 242–251 (2005)

Carlin, J. B.: Meta-analysis for 2×2 tables: A bayesian approach. Stat. Med. 11, 141–158 (1992)

Chang, B. H., Waternaux, C., Lipsitz, S.: Meta-analysis of binary data: which within study variance estimate to use? Stat. Med. 20, 1947–1956 (2001)

DerSimonian, R., Laird, N.: Meta-analysis in clinical trials. Control. Clin. Trials. 7, 177–188 (1986)

Ferguson, T. S.: Prior distributions on spaces of probability measures. Ann. Stat. 2, 615-629 (1974)

Gamalo, M., Wu, R., Tiwari, R.: Bayesian approach to noninferiority trials for proportions. J. Biopharm. Stat. **21**, 902–919 (2011)

Geisser, S.: Predictive Inference: An Introduction. Chapman and Hall, London (1993)

Gelfand, A. E., Dey, D. K.: Bayesian Model Choice: Asymptotics and Exact Calculations. J. R. Stat. Soc. Ser. B. 56, 501–514 (1994)

Gelfand, A. E., Dey, D. K., Chang, H.: Model determination using predictive distributions with implementation via sampling-based methods (with discussion). *Bayesian Statistics 4*(Bernardo, J. M., Berger, J. O., Dawid, A. P., Smith, A. F. M., eds.) Oxford University Press (1992)

Ghosh, S. K., Mukhopadhyay, P., Lu, J. C.: Bayesian analysis of zero-inflated regression models. J. Stat. Plan. Infer. **136**(4), 1360–1375 (2006)

Hall, D. B.: Zero-Inflated Poisson and Binomial Regression with Random Effects: A Case Study. Biometrics. 56, 1030–1039 (2000)

Huang, L., Zheng, D., Zalkikar, J., Tiwari, R.: Zero-inflated Poisson model based likelihood ratio test for drug safety signal detection. Stat. Methods Med. Res. (2014). https://doi.org/10.1177/0962280214549590

Muthukumarana, S., Tiwari, R.: Meta-analysis using dirichlet process. Stat. Methods Med. Res. 25(1), 352–365 (2016)

- Neal, RM: Markov Chain Sampling Methods for Dirichlet Process Mixture Models. J. Comput. Graph. Stat. **9**(2), 249–265 (2000)
- Nissen, S. E., Wolski, K.: Effect of rosiglitazone on the risk of myocardial infarction and death from cardiovascular causes. New Eng. J. Med. 356, 2457–2471 (2007)

Smith, T. C., Spiegelhalter, D. J., Thomas, A.: Bayesian approaches to random-effects meta-analysis: a comparative study. Stat. Med. 14, 2685–2699 (1995)

#### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

# Submit your manuscript to a SpringerOpen<sup>®</sup> journal and benefit from:

- Convenient online submission
- ► Rigorous peer review
- ► Open access: articles freely available online
- ► High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at > springeropen.com